

XML, l'eXtensible Markup Language

Ingénierie Documentaire
<http://doc.crzt.fr>



Stéphane Crozat

Contributions : Bruno Bachimont, Erik Gebers

Table des matières



I - Introduction à XML	5
1. Préambule : XML selon le W3C	5
2. Définition du XML	6
3. XML : un langage à balise	7
4. XML : un méta-langage	7
5. Exemple : Un document XML	8
6. Historique : de SGML à XML	8
7. Discussion : HTML et XML	9
II - Langages XML	11
1. Langages standard et langages locaux	11
2. Langages XML orientés données	11
3. Documents structurés et documents formatés	12
4. Langages XML orientés documents	15
5. Caractéristiques recherchées pour un format XML	16
6. Outillage XML	16
III - Syntaxe de base XML	19
1. Document bien formé	19
2. Balise	19
3. Élément	20
4. Attribut	21
5. Structure générale d'un fichier XML	22
6. Exemple de fichier XML - Le courriel	22
7. La syntaxe XML en résumé	23
IV - Syntaxe avancée XML	25

1. Commentaires	25
2. Namespace	25
3. Syntaxe XML et espaces	26
4. Encodage des caractères	27
V - Introduction aux schémas XML	28
1. Notion de document valide	28
2. Document Type Definition	29
3. W3C XML Schema	29
4. Regular Language for XML Next Generation	30
VI - Quelques bonnes questions sur XML	32
VII - Bibliographie	35
VIII - Compléments	36
1. Balises et poignées de calcul	36
2. La structuration logique	36
3. Exemple de structuration logique	37
4. Définition des chaînes éditoriales XML	38
5. Définition de Scenari	39
6. Un langage pour publier les documents XML	40
7. Définition de XSL-XSLT	40
IX - Exercices	41
1. Exercices : Culture XML	42
1.1. Exercice	42
1.2. Exercice : Galaxie XML	43
1.3. Exercice : Historique XML	43
1.4. Exercice :	43
2. Exercices : Syntaxe XML	44
2.1. Exercice	44
2.2. Exercice	45
2.3. Exercice	45
Questions de synthèse	46
Solutions des exercices	48
Glossaire	52

Abréviations	53
Bibliographie	54
Webographie	55
Index	56

Introduction à XML

I

1. Préambule : XML selon le W3C

XML * est un standard du W3C * : Extensible Markup Language (XML) 1.0 (Fifth Edition) * .

Les définitions suivantes sont directement copiées du site du W3C * (novembre 2009).

Définition : What is XML?

The Extensible Markup Language (XML) is a *simple text-based format* for representing *structured information*: documents, data, configuration, books, transactions, invoices, and much more. It was derived from an older standard format called *SGML* (ISO 8879), in order to be more suitable for Web use.

Définition : What is XML Used For?

XML is one of the most widely-used formats *for sharing structured information* today: between programs, between people, between computers and people, both locally and across networks.

Exemple : Extrait d'un lexique

```

1 <lexique>
2   ...
3   <entree id="1976">
4     <mot>Essuie-glace</mot>
5     <description>L'essuie-glace retire automatiquement la pluie du <voirAussi ref="1965">pare-brise
</voirAussi> lorsqu'elle tombe dessus. Il possède une <voirAussi ref="1977">lame caoutchoutée</voirAussi>
qui peut être changée séparément en cas de besoin.</description>
6   </entree>
7   <entree id="1977">
8     <mot>Lame caoutchoutée</mot>
9     ...
10  </entree>
11  ...
12 </lexique>

```

2. Définition du XML

Définition : XML

L'eXtensible Markup Language XML * est un *méta-langage* permettant de définir des langages à balises.

Il standardisé comme une recommandation par le W3C * depuis 1998.

Exemple

- LaTeX et HTML sont des langages à balises (mais ce ne sont pas des langages XML)
- XHTML est un langage XML

Fondamental

En tant que méta-langage XML sert à définir des *formats* informatiques, c'est à dire des façons de représenter de l'information.

Les bonnes caractéristiques d'XML - *non-ambiguïté, lisibilité, passivité* - en font un très bon candidat pour de très nombreux usages, il est donc utilisé dans des secteurs très variés de l'informatique.

On distingue généralement deux grandes catégories d'utilisation : *les langages orientés données et les langages orientés documents*.

Complément : Versions de XML

La version historique de XML est la version 1.0, qui reste aujourd'hui la plus largement utilisée.

Il existe également une version 1.1 de XML, qui propose des différences mineures et nécessaires uniquement dans des contextes particuliers. Cette nouvelle version a été nécessaire pour des raisons de compatibilité des outils existants. Les évolutions principales de XML 1.1 sont de nouveaux caractères permis pour les noms d'éléments (pour suivre l'évolution d'Unicode depuis 1998), de nouvelles conventions pour les caractères de fin de ligne (pour la compatibilité avec des ordinateurs *main frame*) et de nouveaux caractères de contrôle dans le contenu.

Complément : Voir aussi

XML : Un langage à balise (cf. p.7)

XML : un méta-langage (cf. p.7)

Langages XML orientés données (cf. p.11)

Langages XML orientés documents (cf. p.15)

Caractéristiques recherchées pour un format XML (cf. p.16)

3. XML : un langage à balise

Définition : Langage à balises

Une balise est un descripteur ajouté au contenu en vue d'un traitement informatique.

Un langage à balises est un langage permettant d'associer à un contenu (généralement du texte) des balises explicites (par exemple pour rendre compte de la structure du texte).

Exemple : Balises Lambda

```

1 <lambda>
2   <body>
3     <p>This is an example</p>
4   </body>
5 </lambda>
```

Ici le texte est "This is an example", il est balisé afin d'ajouter les informations suivantes : c'est est un paragraphe (balise `p`) du corps (balise `body`) d'un document Lambda (balise `lambda`).

Complément : Voir aussi

Balises et poignées de calcul (cf. p.36)

4. XML : un méta-langage

Définition : Méta-langage

Un méta-langage est un langage permettant de définir d'autres langages, les langages ainsi définis permettent à leur tour la description d'informations respectant ces langages.

Exemple : Langage Lambda

Si l'on désire définir informatiquement un langage Lambda, en disant qu'il peut contenir un élément de type `body` qui lui même contient des éléments de type `p`, il faut spécifier les contraintes de ce langage grâce à un méta-langage permettant donc de définir :

- Le vocabulaire : `lambda`, `body`, `p`
- La grammaire : `lambda` contient exactement un `body` qui contient un ou plusieurs `p`

XML en tant que méta-langage ne contient pas les mots du langage Lambda, en revanche il dispose des mécanismes permettant de les définir.

Définition : Notion de schéma

La définition d'un langage XML particulier se fait grâce à un schéma qui permet de lister les mots du langage (vocabulaire) et de définir leur utilisation (grammaire).

On parle également de *format*.



Exemple : Schéma Lambda

```

1 <!ELEMENT lambda (body) >
2 <!ELEMENT body (p+) >
3 <!ELEMENT p (#PCDATA) >

```

Ce code exprime formellement (avec la syntaxe DTD) le langage, ou format, Lambda.

5. Exemple : Un document XML



Exemple : Un mail

```

1 <?xml version='1.0' encoding='iso-8859-1'?>
2 <mail>
3   <de>stephane.crozat@utc.fr</de>
4   <a>fabrice.issac@utc.fr</a>
5   <objet>Demande d'information</objet>
6   <date>01-03-2001</date>
7   <corps>
8     <paragraphe>Bonjour Fabrice,</paragraphe>
9     <paragraphe>Peux tu m'expliquer ce qu'est le langage XML ?</paragraphe>
10    <paragraphe>Il me faudrait en effet ton avis éclairé sur le sujet.</paragraphe>
11    <paragraphe>J'attends ta réponse, à bientôt</paragraphe>
12  </corps>
13 </mail>

```

6. Historique : de SGML à XML

SGML

XML hérite historiquement (et fonctionnellement) du méta-langage SGML, développé au début des années 1980 et largement utilisé dans le cadre des systèmes documentaires, notamment dans les industries de hautes technologies pour lesquelles la fiabilité de la documentation est cruciale. SGML restait un formalisme insuffisamment explicite et trop complexe, ce qui rendait difficile la programmation d'applications SGML dans des cadres plus larges. SGML laissait notamment des ambiguïtés en terme d'interprétation du balisage. XML est un langage respectant la syntaxe SGML, mais en y ajoutant des contraintes supplémentaires, afin d'en lever les ambiguïtés.

- *A l'origine SGML* :
 - norme ISO (1986) créé pour la représentation de documents structurés
 - bien adapté aux systèmes documentaires massifs
 - exemple de langage SGML : HTML
- *Un formalisme pas assez explicite* :
 - rend difficile la programmation d'applications SGML (navigateurs, etc.)
 - XML (98) : hérite de SGML en contraignant la syntaxe

☞ *Exemple : Ambiguïté en SGML*

En SGML une balise ouvrante n'est pas obligatoirement fermée par une balise fermante (sa fermeture pouvant rester implicite), tandis qu'en XML toute balise ouvrante est obligatoirement fermée.

Le code SGML `<A>` est ambigu car on ne sait pas s'il faut l'interpréter comme :

- `<A>` ou
- `<A>`.

En XML seule l'une des deux formulations non ambiguës est autorisée.

☞ *Exemple*

- HTML est un langage SGML
- DocBook est à l'origine un langage SGML, porté depuis en XML

🌸 *Fondamental*

XML hérite de SGML en contraignant la syntaxe, tout document XML est donc un document SGML.

📦 *Complément : Voir aussi*

Documents structurés et documents formatés (cf. p.12)

7. Discussion : HTML et XML

🌸 *Fondamental : XML et HTML ne sont pas au même niveau*

XML est un méta-langage et HTML est un langage, XML et HTML ne sont donc pas directement comparables. En particulier XHTML est un langage XML, c'est donc du XML !

XML et le Web


Une question pertinente est en revanche de se demander pourquoi HTML ou XHTML ne suffisent pas pour les usages Web, et donc pourquoi XML est utile en dehors de son instanciation via XHTML.

La réponse est que si HTML, et *a fortiori* XHTML, fournissent une solution de publication efficace, ils n'offrent pas le même potentiel en terme d'*ingénierie documentaire*. XML possède donc un potentiel de *manipulation* plus important (par exemple pour le multi-supports, la rééditorialisation, ...).

☞ *Exemple : Exemple de limites du HTML*

Le HTML *, langage de référence du web qui a su s'imposer par sa simplicité d'écriture, possède également des inconvénients qui ont motivé les recherches du W3C * sur le XML * :

- *Jeu de balises restreint* : HTML est un langage, les balises sont définies *a priori* , il est impossible de définir ses propres balises et de créer un vocabulaire propre à un domaine spécifique. Ce manque d'*extensibilité* le rend mal approprié pour l'indexation de documents variés, le multimédia, ...
- *Mélange des descripteurs physiques et logiques* : Les balises HTML fournissent des informations relatives à la typographie ou la mise en page d'un texte (CENTER, B, ...) mais également relatives à son rôle dans le document (H1, H2, ...). Cette non séparation stricte entre contenu et présentation nuit à la publication multi-supports et la rééditorialisation par exemple.
- *Liens* : En HTML, les liens sont généralement des pointeurs directs vers une autre ressource (A). Cette mise en œuvre limite fortement les possibilités de réutilisation, notamment par le risque de lien cassé.

 *Complément : Voir aussi*

Documents structurés et documents formatés (cf. p.12)

Langages XML

II

1. Langages standard et langages locaux

Définition : Format standard

On appelle langage (ou format) standard un langage dont le schéma est défini par une organisation reconnue et multipartite (W3C, ISO, OASIS, ...) pour les besoins génériques d'une large communauté, en général internationale.

Remarque : Standard de fait

On appelle standard de fait les langages imposés par une organisation unique, mais à une telle échelle que leur usage est incontournable. C'est typiquement l'exemple des formats de la suite bureautique Microsoft Office. (.doc, .xls, .ppt, ...).

Définition : Format local

On appelle langage (ou format) local - par opposition au format standard - un langage défini par une organisation seule ou une communauté restreinte pour ses besoins propres (un service, une entreprise, voire un réseau d'organisations dans un pays, ...).

Définition : Format ouvert ou fermé

Un format est dit ouvert lorsqu'il est public et que n'importe qui peut l'utiliser sans contrainte. Il est au contraire fermé si son accès est protégé, payant ou soumis à d'autres contraintes.

Remarque : Format propriétaire

On parle de format propriétaire pour désigner un format à la fois local et fermé.

2. Langages XML orientés données

Définition

Ils permettent d'enregistrer et de transporter des données informatiques structurées (comme par exemple des données gérées par des bases de données) selon des formats ouverts (c'est à dire dont on connaît la syntaxe) et faciles à manipuler (les structures arborescentes XML étant plus riches que des fichiers à plat par exemple).

Les langages XML orientés données servent surtout à l'échange ou la sérialisation * des données des programmes informatiques.

Remarque

L'utilisation d'XML est en fait ici assez accessoire, d'autres formalismes s'y substituent sans difficulté, souvent moins explicites pour la manipulation humaine et souvent plus performant pour la manipulation informatique : CSV, JSON, ... (voir par exemple : <http://www.xul.fr/ajax-format-json.html> pour une comparaison JSON / XML).

Remarque : Format verbeux

XML est en général plus verbeux que ses alternatives (il contient plus de caractères), c'est pourquoi il est plus explicite pour un humain (ce qui n'est pas toujours utile), et moins performant informatiquement (ce qui n'est pas toujours un problème).

Exemple : Formats XML orientés données standards

- MathML, ChemML, ...
- ATOM, RSS
- Dublin Core
- RDF, OWL
- SVG
- ...

Exemple : Formats XML orientés données locaux

XML est utilisé pour de très nombreux langages orientés données locaux, en fait chaque fois qu'un format XML est défini pour les besoins spécifique d'un programme.

```

1 <myVector>
2   <x>5</x>
3   <y>19</y>
4 </myVector>
```

Complément : Langages XML de programmation

Le formalisme XML est également utilisé pour définir des langages de programmation, à l'instar du C ou du Java. Les langages de programmation écrits en XML sont généralement à vocation déclarative et adressent le plus souvent des traitements eux mêmes liés à XML.

XSL-XSLT est un exemple de langage de programmation écrit en XML. On peut également citer par exemple le langage de script ANT (<http://ant.apache.org/>) ou XUL pour la réalisation d'IHM (<http://www.xul.fr/>)

3. Documents structurés et documents formatés

Définition : Document formaté

On appelle document formaté un document dont le fichier informatique source décrit la façon de le mettre en forme. C'est la façon la plus courante de traiter avec les documents informatiques, telle qu'elle est mise en œuvre dans les traitements de texte ou sur le Web avec HTML.

Un document formaté nécessite un logiciel capable d'en interpréter le format pour permettre de le lire.

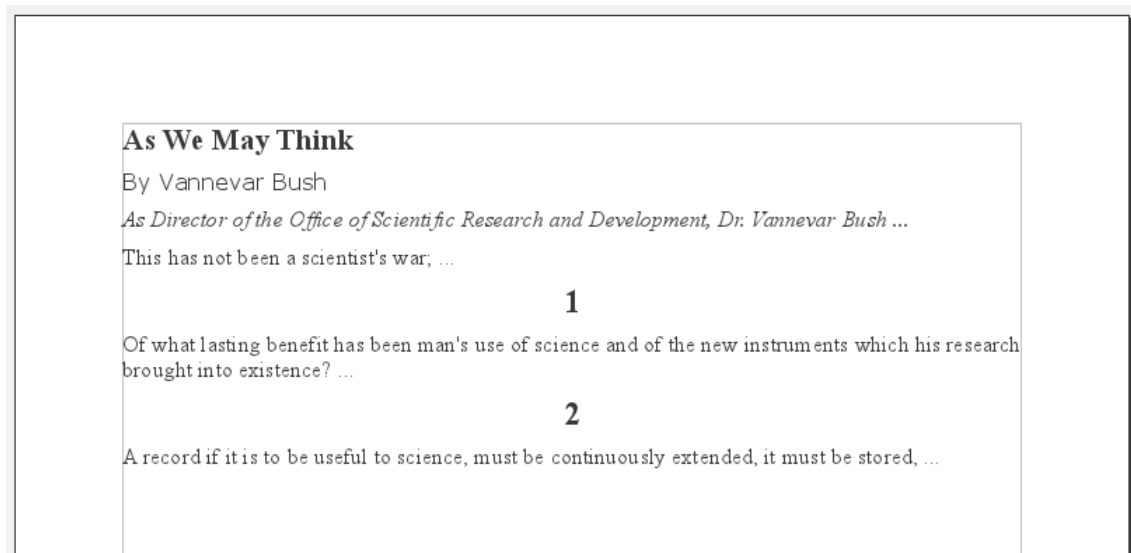
XML est un excellent candidat à la sérialisation * de documents formatés.

☞ *Exemple : Document formaté (format OpenDocument : extrait simplifié)*

```

1 <document-content>
2 <automatic-styles>
3   <style name="P1" family="paragraph"><text-properties font-name="Times New Roman" fo:font-size=
4     "12pt"/></style>
5   <style name="P2" family="paragraph"><text-properties font-name="Times New Roman" fo:font-size=
6     "16pt" fo:font-weight="bold"/></style>
7   <style name="P3" family="paragraph"><paragraph-properties fo:text-align="center"/><text-properties
8     font-name="Times New Roman" fo:font-size="18pt" fo:font-weight="bold" /></style>
9   <style name="P4" family="paragraph"><text-properties font-name="Times New Roman" fo:font-size=
10    "12pt" fo:font-style="italic"/></style><style name="T1" family="text"><text-properties font-name=
11    "Verdana"/></style>
12 </automatic-styles>
13 <body>
14   <text>
15     <p style-name="P2">As We May Think</p>
16     <p style-name="P1"><span style-name="T1">By Vannevar Bush</span></p>
17     <p style-name="P4">As Director of the Office of Scientific Research and Development, Dr. Vannevar
18     Bush ...</p>
19     <p style-name="P1">This has not been a scientist's war ...</p>
20     <p style-name="P3">1</p>
21     <p style-name="P1">Of what lasting benefit has been man's use of science ...</p>
22     <p style-name="P3">2</p>
23     <p style-name="P1">A record if it is to be useful to science ...</p>
24   </text>
25 </body>
26 </document-content>

```



Visualisation dans OpenOffice.org Writer de l'extrait de l'article "As We May Think"

☞ *Définition : Notion de documents structurés*

On appelle document structuré un document dont la structure logique est décrite plutôt que la mise en forme physique (Structured documents *, p.7).

Après SGML qui avait été inventé pour cela, XML est aujourd'hui le candidat quasi-unique pour la réalisation de documents structurés.

Exemple : Document structuré (format DocBook, légèrement simplifiée)

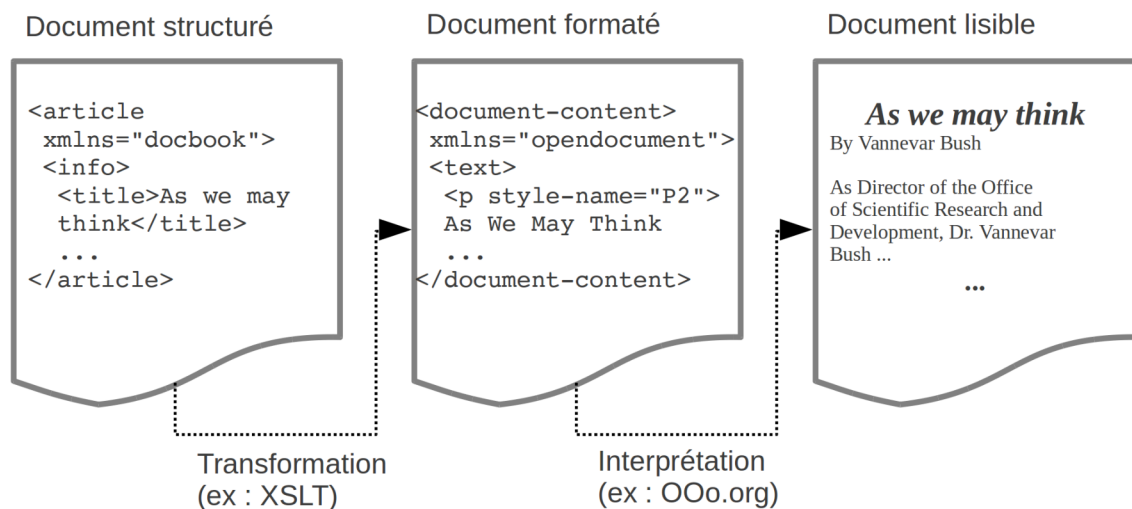
```

1 <article xmlns="http://docbook.org/ns/docbook">
2   <info>
3     <title>As we may think</title>
4     <author><personname>Vannevar Bush</personname></author>
5   </info>
6   <epigraph><para>As Director of the Office of Scientific Research and Development, Dr. Vannevar Bush
...</para></epigraph>
7   <para>This has not been a scientist's war ...</para>
8   <sect1>
9     <para>Of what lasting benefit has been man's use of science ...</para>
10  </sect1>
11  <sect1>
12    <para>A record if it is to be useful to science ...</para>
13  </sect1>
14 </article>

```

Définition : Transformation des documents structurés

Un document structuré n'est pas destiné à être directement utilisé pour la lecture humaine, il doit être d'abord transformé dans un document formaté pour être utilisé (visualisé, imprimé, ...).



Processus de transformation d'un document structuré en document formaté

Remarque : Chaîne XML complète

La technologie XML peut en fait être utilisée tout au long de la chaîne :

- Langage XML orienté document structuré en amont (DocBook, DITA, TEI, langage local, ...)
- Langage XML de programmation pour la transformation (XSL-XSLT)
- Langage XML orienté document formaté pour l'aval (OpenDocument, XHTML, ...)

 *Complément : Voir aussi*

La structuration logique (cf. p.36)

Exemple de structuration logique (cf. p.37)

Architecture des chaînes éditoriales XML (cf. p.)


Un langage pour publier les documents XML (cf. p.40)

Définition de XSL-XSLT (cf. p.40)

 *Complément : Bibliographie*

Technologie de la GED * : Structure logique et structure physique (pp58-61)

4. Langages XML orientés documents

 *Définition*

Ils permettent de représenter informatiquement des documents numériques. Les formats de documents faisant appel à des représentations fortement arborescentes, XML est un candidat idéal pour cet usage. En fait SGML, l'ancêtre de XML, a été inventé pour l'informatique documentaire. Aujourd'hui la très grande majorité des formats de représentation de document sont en XML.

 *Définition : Langages XML orienté documents formatés*


Ils définissent des formats de mise en forme de document (structure physique), en général pour un support donné.

 *Exemple : Langages XML orientés documents formatés*

- XHTML
- XSL-FO
- SMIL
- OpenDocument
- OOXML
- ...

 *Définition : Langages XML orientés documents structurés*

Ils définissent des formats de structuration de document (structure logique), en général pour un métier donné, plus ou moins précis selon la généralité du langage.

 *Exemple : Langages XML orientés documents structurés*

- DocBook
- TEI
- DITA
- ...

Exemple : Format local orienté document structuré

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <document type="Lorem ipsum">
3   <paragraphe>Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque sit amet
4 libero ac mauris egestas venenatis nec vitae sapien.</paragraphe>
5   <paragraphe>Donec a lectus sed augue pellentesque accumsan eu ac justo. Etiam est urna, sagittis
6 ac cursus nec, placerat quis velit.</paragraphe>
7 </document>
```

Complément : Voir aussi

Documents structurés et documents formatés (cf. p.12)

5. Caractéristiques recherchées pour un format XML

Fondamental : Non ambiguïté

Les règles syntaxiques strictes d'XML rendent son interprétation informatique univoque.

Fondamental : Lisibilité

Un format XML a pour objectif d'être lisible par un être humain, afin de plus facilement en appréhender le langage.

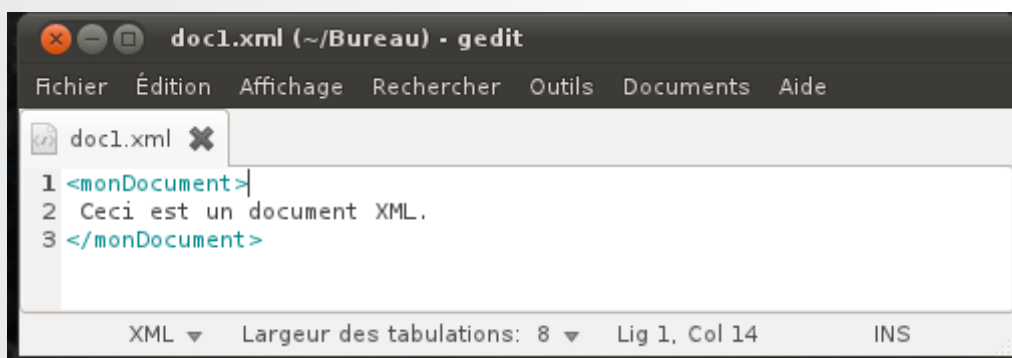
Fondamental : Passivité

Un format XML est passif, il dépend de programmes informatiques qui le transforment.

6. Outillage XML

Fondamental

Un éditeur de texte suffit pour faire du XML.



Document XML créé avec un simple éditeur de texte (gedit)

Parseurs

Un *parser* (ou parseur) XML est un programme capable d'analyser un document XML afin de :

- vérifier qu'il est bien formé
- éventuellement vérifier sa validité par rapport à un schéma (DTD, W3C Schema, RelaxNG)
- éventuellement en fournir une représentation mémoire permettant son traitement (SAX, DOM)

Par exemple Xerces est un parser Java (<http://xerces.apache.org/xerces2-j/>). Tous les langages de programmation proposent aujourd'hui des parsers XML.

Éditeur validant

Un éditeur validant est un éditeur spécialisé dans l'écriture du XML (on parle simplement d'éditeur XML), permettant de :

- vérifier dynamiquement que le fichier est bien formé,
- de charger des schémas pour vérifier dynamiquement la validité,
- de proposer des fonctions d'auto-complétion,
- ...

Éditeurs XML orientés développeurs

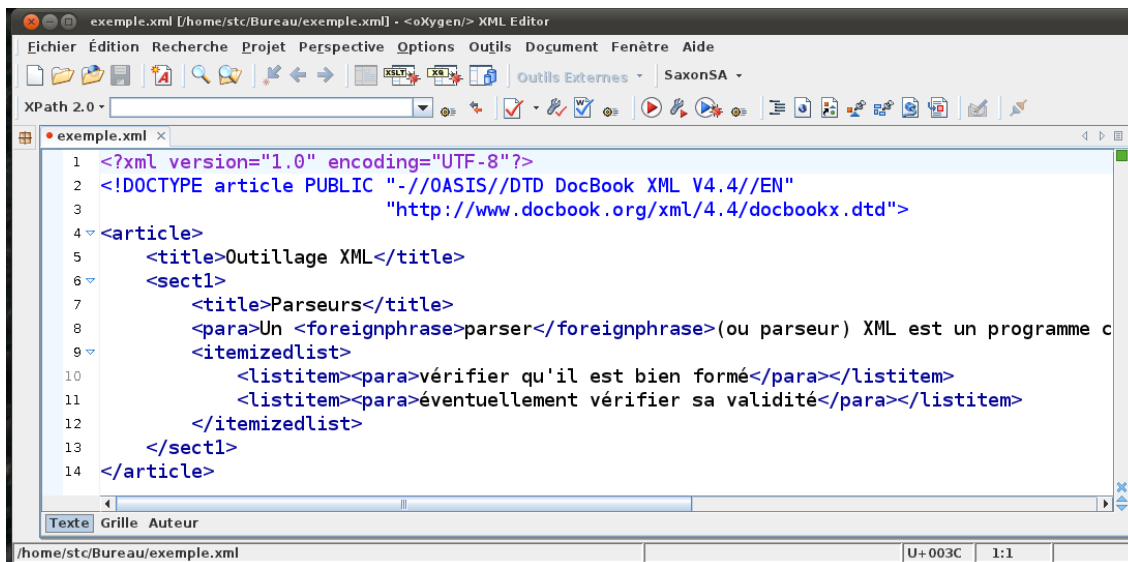
Les éditeurs spécialisés orientés développeurs sont des IDE * permettant l'écriture :

- des schémas,
- des transformations,
- des fichiers XML de test,
- ...

Ils sont en général peu adaptés à la production des documents XML par les utilisateurs finaux (rédacteurs).

On peut citer par exemple :

- oXygen (<http://www.oxygenxml.com/>), un produit commercial complet accessible à moins de 100€ pour les usages non commerciaux, multi-plateformes (Linux, Mac, Windows) ;
- XML Spy, disponible uniquement sous Windows
- ...



Exemple de fichier XML dans l'éditeur oXygen

Éditeurs XML orientés rédacteurs

Les éditeurs XML orientés rédacteurs prennent en général en entrée un schéma, des fichiers permettant de configurer l'éditeur, des programmes de transformation, et offre un environnement dédié à l'écriture et la publication de document XML.

On peut citer :

- Jaxe, libre et multi-plateformes (<http://jaxe.sourceforge.net/fr/>)
- Adobe Framemaker, sous Windows (<http://www.adobe.com/products/framemaker.html>)
- Arbortext, sous Windows (<http://www.ptc.com/products/arbortext/>)
- XMetal, sous Windows (<http://na.justsystems.com/content-xmetal>)
- Scenari, chaîne éditoriale complète libre et multi-plateformes (<http://scenari-platform.org/>)

Complément : Voir aussi

Définition des chaînes éditoriales XML (cf. p.38)

Définition de Scenari (cf. p.39)

Syntaxe de base XML



1. Document bien formé

Définition

Un document est dit bien formé lorsqu'il respecte les règles syntaxiques du XML *.

Fondamental

XML * ne pose pas de sémantique à priori mais uniquement des règles syntaxiques.

Syntaxe

Les règles syntaxiques à respecter sont :

- Il n'existe qu'un seul élément père par document, cet élément contenant tous les autres. Il est également appelé *racine*.
- Tout élément fils est inclus complètement dans son père (pas d'éléments croisés).

Attention

XML ne définit pas le comportement ni la manière dont doit être traité un document, mais uniquement un *format*.

Complément

Balise (cf. p.19)

2. Balise

Définition : Balise

Les balises sont les composants fondamentaux permettant l'écriture de documents XML *.

Elles se distinguent du contenu en utilisant les caractères <, > et /. Elles n'ont pas de présentation ou de signification définie par le langage mais elles auront un sens pour les applications et/ou le lecteur.

Syntaxe

Il existe trois types de balises :

- balise d'ouverture : `<nom_balise>`
- balise de fermeture : `</nom_balise>`
- balise vide : `<nom_balise/>`

Exemple : Exemple de balise XML

```
1 <adresse>12, rue de Paris</adresse>
```

3. Élément

Définition : Élément

Un élément XML * est un extrait d'un fichier XML comprenant :

- une balise ouvrante
- une balise fermante
- le contenu compris entre les deux balises, qui peut être du *texte* et/ou d'autres *éléments*, ou *rien* (élément vide)

Le nom d'un élément peut être composé de tout caractère alphanumérique plus `_`, `-` et `.`. De plus, ils doivent commencer par un caractère alphabétique ou `_` et ceux commençant par la chaîne `xml` sont réservés (que ce soit en minuscules, majuscules ou un mélange des deux).

Attention : Case-sensitive

XML * différencie les majuscules des minuscules, il faut donc respecter la casse.

Attention

Deux éléments *ne peuvent pas* avoir un contenu croisé : `<nom1> ... <nom2> ... </nom1> ... </nom2>` est interdit.

Définition : Élément vide

On appelle élément vide un élément qui ne comporte rien entre sa balise ouvrante et sa balise fermante.

Syntaxe : Élément vide

Les syntaxes `<element></element>` et `<element/>` sont strictement équivalentes.

4. Attribut

Définition

Un attribut est une information supplémentaire attachée à un élément, on parle de métadonnée.

Syntaxe

Les attributs d'un élément sont formés d'une suite d'affectations séparées par des espaces :
`attribut1='valeur' attribut2='valeur' ...`

Ils sont ajoutés à la balise ouvrante ou à une balise vide (jamais à une balise fermante) :

- `<nom_element [attributs]>`
- `<nom_element [attributs]/>`

La valeur est indiquée entre apostrophes ou guillemets (au choix, mais pas de mélange des deux) :

- `attribut1='valeur' ou`
- `attribut1="valeur"`

Méthode

Utilisez des apostrophes si la valeur de l'attribut inclut des guillemets et vice et versa.

Attention

Un élément ne peut pas contenir deux attributs ayant le même nom.

Syntaxe

Le nom d'un attribut est soumis aux mêmes contraintes que les noms d'éléments.

La valeur de l'attribut quant à elle peut contenir tout caractère à l'exception de ^, % et &.

Remarque : Équivalence attribut / élément

Un attribut peut toujours être représenté alternativement par un élément fils de l'élément qu'il caractérise, avec une signification du même ordre :

- `<element attribut="x"/>` et
- `<element><attribut>x</attribut><element>` sont similaires.

Il est donc tout à fait possible de faire du XML sans utiliser d'attribut.

Méthode : Usage des attributs

On utilise généralement les attributs :

- Pour différencier le contenu destiné à être affiché dans le document lisible des métadonnées qui ne le seront pas (version, date de création, ...)
- Pour simplifier l'écriture du document

- Pour ajouter des identifiants et des références

5. Structure générale d'un fichier XML

Syntaxe

Un document XML * est constitué de :

- Un *prologue*
Il est facultatif et comprend une déclaration XML *, indiquant la version du langage XML utilisé et le codage des caractères dans le document. Chacune de ces informations est optionnelle mais leur ordre est obligatoire

```
<?xml version="numéro de version" encoding="encodage des caractères"?>
```
- Un *arbre d'éléments* contenant au moins un élément (l'élément racine)

Exemple : Prologue

```
1 <?xml version="1.0" encoding="UTF-8"?>
```

Indique que le document est codé en utilisant un langage XML * de version 1.0, avec des caractères codés selon la norme UTF-8.

Exemple : Arbre d'éléments

```
1 <lettre>
2   <expediteur>moi</expediteur>
3 </lettre>
```

6. Exemple de fichier XML - Le courriel

Exemple : Un courriel imprimé

Date: Mar, 28 Oct 2003 14:01:01 +0100 (CET)

De : Marcel <marcel@ici.fr>

A : Robert <robert@labas.fr>

Sujet: Hirondelle

Salut,

Pourrais-tu m'indiquer quelle est la vitesse de vol d'une hirondelle transportant une noix de coco ?

A très bientôt,

Marcel

Représentation possible de ce message en XML

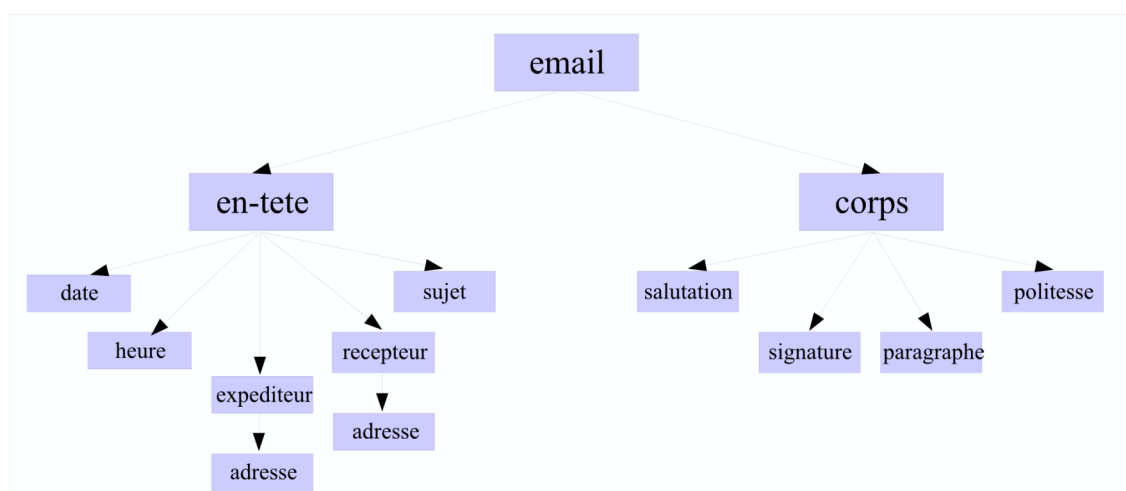
```
1 <?xml version="1.0" encoding="ISO-8859-1"?>
```

```

2 <email>
3 <entete>
4 <date type='JJMMAAAA'>28102003</date>
5 <heure type='24' local='(GMT+01 :00)'+>14:01:01</heure>
6 <expediteur><adresse mail='marcel@ici.fr'>Marcel</adresse></expediteur>
7 <recepteur><adresse mail='robert@labas.fr'>Robert</adresse></recepteur>
8 <objet>Hirondelle</objet>
9 </entete>
10 <corps>
11 <salutation>Salut,</salutation>
12 <paragraphe>Pourrais-tu m'indiquer quelle est la vitesse de vol d'une hirondelle transportant une noix de
coco ?</paragraphe>
13 <politesse>À très bientôt,</politesse>
14 <signature>Marcel</signature>
15 </corps>
16 </email>

```

 Complément : Arbre d'éléments du document email correspondant



Représentation d'un document XML sous forme d'arbre

7. La syntaxe XML en résumé

Les règles de syntaxe d'XML sont très simples : un document XML comporte un en-tête permettant de l'identifier en tant que document XML et un élément racine contenant tous les autres éléments. Un élément est caractérisé par une balise ouvrante et une balise fermante et contient d'autres éléments et/ou du contenu.

Aspect	Syntaxe	Explication
Entête	<code><?xml version='1.0' encoding='UTF-8'?></code>	Un document XML peut contenir une ligne d'entête qui sert à identifier le document en tant que document XML, et généralement à spécifier le système d'encodage de caractères utilisé dans le document.
Élément racine	<code><NomElementRacine> ... </NomElementRacine></code>	Tout document XML contient un et un seul élément, dit racine, contenant tous les autres.
Élément	<code><NomElement> ... </NomElement></code>	Un élément est caractérisé par une balise ouvrante et une balise fermante obligatoire, les éléments sont donc inclus les uns dans les autres et ne peuvent se croiser. Les balises tiennent compte de la casse des caractères.
Attribut	<code><NomElement NomAttribut='valeur'></code>	Un attribut associé à une valeur peut être spécifié dans la balise ouvrante d'un élément.
Élément vide	<code><NomElementVide/></code>	Un élément peut être vide, il est alors noté par une unique balise à la syntaxe particulière. Un élément vide contient en général des attributs.
Contenu	<code><NomElement> contenu ici </NomElement></code>	Le contenu est placé à l'intérieur des éléments, c'est à dire entre les balises.

Aspects principaux de la syntaxe XML

Syntaxe avancée XML

IV

1. Commentaires

Syntaxe

Un fichier XML peut contenir des commentaires, ils peuvent être insérés n'importe où dans le document avec la syntaxe :

```
<!-- Ceci est un commentaire -->
```

Exemple

```
1 <document>
2   Ceci est un contenu.
3   <!-- Ceci est un commentaire -->
4 </document>
```

2. Namespace

Principe

Un *namespace* (ou espace de noms en français) est une mécanique qui permet d'assurer l'unicité des noms des éléments utilisés au sein des fichiers XML, dans l'objectif de pouvoir « mélanger » différents schémas.

- Soit un extrait de schéma S1 qui définit un élément de syntaxe en informatique comme contenant du code : ... *{element syntaxe {code {...}}}**
- Soit un extrait de schéma S2 qui définit un élément de syntaxe en mathématique comme contenant une équation : ... *{element syntaxe {equation {...}}}**

Si je souhaite, dans un schéma S3 réutiliser les éléments syntaxe issus de S1 et S2, je rencontre un conflit de noms : en effet, deux éléments portant le même nom, « syntaxe », définissent en fait des éléments différents.

Le *namespace* va me permettre de différencier ces deux éléments, en associant un nom unique aux schémas S1 et S2, par exemple une adresse web.

Si j'associe *www.utc.fr/S1* au premier schéma et *www.utc.fr/S2* au second, j'obtiens alors deux noms de balises différents :

- *www.utc.fr/S1:syntaxe*

- www.utc.fr/S2:syntaxe

Préfixe

Cette écriture étant quelque peu fastidieuse, il est également possible d'associer un préfixe au namespace. La correspondance entre le namespace et le préfixe est déclarée dans le fichier XML, ce qui permet à un programme informatique de remplacer les préfixes par les namespaces.

Finalement, on obtient :

- Préfixe s1 associé au namespace www.utc.fr/S1 et écriture XML `s1:syntaxe`
- Préfixe s2 associé au namespace www.utc.fr/S2 et écriture XML `s2:syntaxe`

Remarque : Namespace par défaut

Il est possible de définir un namespace par défaut ce qui permet d'avoir un *namespace* pour chaque balise, sans avoir à utiliser de préfixe.

Syntaxe

```
1 <elementRacine
2   xmlns:prefixe1="namespace1"
3   xmlns:prefixe2="namespace2" ...
4   xmlns:prefixeN="namespaceN"
5   xmlns="namespaceDesÉlémentsNonPréfixés"
6 >
```

Définition : Nom développé

Un *nom développé* (*expanded name*) est le couple constitué par un *nom d'espace de nommage* (*namespace name*) et par un *nom local* (*local name*).

Complément

- <http://www.w3.org/TR/2004/REC-xml-names11-20040204/>
- Traduction : <http://www.yoyodesign.org/doc/w3c/xml-names11/>

3. Syntaxe XML et espaces

Soit les trois syntaxes XML suivantes (les espaces sont symbolisés par des ~ et les retours chariot par des §) :

1. `<a>bonjour~aurevoir`
2. `<a>bonjour~~~~aurevoir`
3. `<a>bonjour§
aurevoir~~~~`
4. `<a>bonjour~aurevoir§
`

Ces trois syntaxes sont logiquement équivalentes.

Donc :

- À l'intérieur d'un élément qui ne contient que du texte : 1 espace = N espaces = 1 retour chariot = N retours chariot
- Entre deux éléments : 1 espace = N espaces = 1 retour chariot = N retours chariot = Rien (absence de caractère)

Dans la pratique seuls les cas de *mixed content* ont besoin que les espaces soient préservés. Les deux exemples ci-après ne sont pas équivalents :

- `<p>Ceci est <i>important</i></p>`
- `<p>Ceci est<i>important</i></p>`

Mais sans schéma il est impossible de décider a priori si le contenu d'un élément qui contient des éléments fils et des espaces est *mixed* ou si ces espaces servent juste à la mise en forme du XML.

Espaces conservés

Par défaut, un programme doit considérer que tous les espaces sont signifiants et doivent être conservés lors des traitements.

Remarque : Équivalence entre séparateurs

XML considère comme équivalents les caractères de séparation :

- Espace ' '
- Tabulation \t
- Retour chariot \r
- Fin de ligne \n

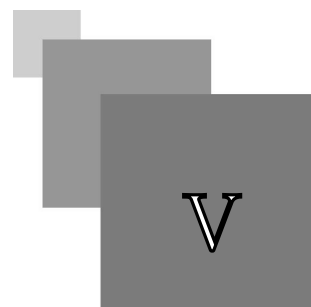
4. Encodage des caractères

L'encodage des caractères définit l'équivalence entre des caractères et leur valeur numérique en machine.

L'ASCII est l'ancêtre et la base de tous les systèmes d'encodage actuels : il définit par exemple que la lettre A à la valeur 65.

L'encodage universel Unicode UTF-8 est supporté par XML et aujourd'hui la plupart des outils informatiques associés, il est à utiliser dans le cas général.

Introduction aux schémas XML



1. Notion de document valide

Définition : Schéma

Un schéma est une description de la structure que doit respecter un document lui faisant référence, c'est à dire qu'il établit la liste des éléments XML autorisés (avec leurs attributs), ainsi que l'agencement possible de ces éléments.

On parle aussi de *grammaire*, au sens où le schéma définit l'enchaînement autorisé des balises et vient en *complément de la syntaxe XML* (qui elle est indépendante d'un schéma particulier).

Définition : Document valide

Un document XML * bien formé est dit valide pour un schéma donné s'il respecte les règles structurelles imposées par ce schéma.

Ce contrôle de la structure permet :

- De s'assurer l'homogénéité structurelle des documents de même type.
- Le traitement automatique d'un ensemble de documents de même type (mise en forme, stockage, extraction d'informations...).
- La création de formats standard et leur respect.

Exemple : Exemples de langages de schéma

Il existe plusieurs langages de définition schéma, mais les trois principaux sont :

- Document Type Définition (W3C) : Un langage hérité de SGML qui fait partie du standard XML
- W3C XML Schema (W3C) : Une alternative aux DTD destiné à moderniser et compléter ce langage historique
- Relax NG (OASIS, ISO) : Une autre alternative, compromis entre W3C XML Schema et DTD

2. Document Type Definition

Le formalisme de définition de schéma DTD est le premier qui a été introduit dès la première version du standard XML. Il est en fait intégré au standard W3C de XML.

Il est directement hérité de la norme SGML.

Les DTDs utilisent un langage spécifique (non XML) pour définir les règles structurelles. Un fichier de DTD peut contenir principalement deux types de déclarations :

- *des déclarations d'éléments*,
indiquent les éléments pouvant être inclus dans un document et l'organisation du contenu de chaque élément (éléments fils ou texte).
- *des déclarations d'attributs*,
définissent les attributs pouvant être associés à un élément ainsi que leur type.

Exemple : Exemple de DTD

```
1 <!ELEMENT document (paragraphe+)>
2 <!ATTLIST document type CDATA #REQUIRED>
3 <!ELEMENT paragraphe (#PCDATA)>
```

Exemple : Exemple de document XML valide

```
1 <?xml version='1.0' encoding='iso-8859-1'?>
2 <!DOCTYPE document SYSTEM "document.dtd">
3 <document type='memo'>
4   <paragraphe>Lorem ipsum dolor sit amet.</paragraphe>
5   <paragraphe>Consectetur adipiscing elit.</paragraphe>
6   <paragraphe>Sed do eiusmod tempor.</paragraphe>
7 </document>
8
```

3. W3C XML Schema

Les XML Schema ont été proposés par le W3C pour permettre de dépasser les limites des DTD.

<http://www.w3.org/XML/Schema>

On notera en particulier :

- une syntaxe XML
- l'extension de l'expression des règles d'organisation structurelle (héritage, réutilisation, etc.)
- l'ajout d'un langage de typage des éléments (particulièrement utile pour les format XML orientés données)

👉 Exemple : Exemple de DTD

```

1 <!ELEMENT document (paragraphe+)>
2 <!ATTLIST document type CDATA #REQUIRED>
3 <!ELEMENT paragraphe (#PCDATA)>

```

👉 Exemple : Exemple de W3C XML Schema correspondant

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
3   <xs:element name="document">
4     <xs:complexType>
5       <xs:sequence>
6         <xs:element maxOccurs="unbounded" ref="paragraphe"/>
7       </xs:sequence>
8       <xs:attribute name="type" use="required"/>
9     </xs:complexType>
10  </xs:element>
11  <xs:element name="paragraphe" type="xs:string"/>
12 </xs:schema>

```

4. Regular Language for XML Next Generation

RelaxNG (REgular LAnguage for XML Next Generation) est un langage de schéma XML.

- RelaxNG est une alternative aux DTD et à W3C XML Schema, qui combine les avantages de ces deux autres langages.
- RelaxNG est un standard OASIS et une norme ISO/CEI.
- Deux syntaxes : une syntaxe XML (alternative à W3C Schema) et une syntaxe compacte (alternative aux DTD).
- RelaxNG ne définit que la structure (comme les DTD) et utilise W3C XML Schema pour le typage des données.

<http://relaxng.org/>

📦 Complément

Le standard est porté par James Clark depuis ses travaux sur Trex (il est issu de la fusion de Trex et Relax de Murata Makoto).

👉 Exemple : Exemple de schémas publics définis en Relax NG

- OpenDocument (format bureautique)
- DocBook (format documentaire)
- Atom (syndication)

👉 Exemple : Exemple de DTD

```

1 <!ELEMENT document (paragraphe+)>
2 <!ATTLIST document type CDATA #REQUIRED>

```

```
3 <!ELEMENT paragraphe (#PCDATA)>
```

☞ *Exemple : Exemple de schéma RelaxNG correspondant (syntaxe XML)*

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <grammar xmlns="http://relaxng.org/ns/structure/1.0">
3   <start>
4     <element name="document">
5       <attribute name="type"/>
6       <oneOrMore>
7         <element name="paragraphe">
8           <text/>
9         </element>
10        </oneOrMore>
11      </element>
12    </start>
13 </grammar>
```

☞ *Exemple : Exemple de schéma RelaxNG correspondant (syntaxe compacte)*

```
1 start = element document {
2   attribute type {text},
3   element paragraphe {text}+
4 }
```

☞ *Exemple : Autre exemple de schéma RelaxNG correspondant (syntaxe XML, patterns nommés)*

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <grammar xmlns="http://relaxng.org/ns/structure/1.0">
3   <start>
4     <ref name="document"/>
5   </start>
6   <define name="document">
7     <element name="document">
8       <attribute name="type"/>
9       <oneOrMore>
10        <ref name="paragraphe"/>
11      </oneOrMore>
12    </element>
13  </define>
14  <define name="paragraphe">
15    <element name="paragraphe">
16      <text/>
17    </element>
18  </define>
19 </grammar>
```

☞ *Exemple : Autre exemple de schéma RelaxNG correspondant (syntaxe compacte, patterns nommés)*

```
1 start = document
2 document = element document {attribute type {text}, paragraphe+ }
3 paragraphe = element paragraphe {text}
```

Quelques bonnes questions sur XML

VI

Réalisé à partir d'une interview de Michaël Tartar dans un *chat* de LeJournalDuNet en 2001.

http://www.journaldunet.com/chat/retrans/011206_tartar.shtml

Quel est l'avantage de ce langage par rapport au HTML ?

HTML est un langage de mise en forme. XML est beaucoup plus générique. Il permet de définir des langages alors que HTML définit un jeu de balises limitées et donc non-extensibles.

Pourquoi le XML n'a-t-il pas encore remplacé le HTML ?

Le XML n'a pas vocation à remplacer le HTML. XML permet aux architectes de sites de se concentrer sur l'information à fournir à l'internaute, et dans un deuxième temps de la mettre en forme pour des besoins spécifiques : PDF pour des documents à imprimer, HTML pour les navigateurs, ...

Quel est l'intérêt de l'utilisation d'XML sur le web en opposition aux fichiers CSS ?

CSS * ne s'applique qu'à la définition de règles graphiques de restitution d'une information, alors que XML permet de structurer les informations. A ce titre, sur le Web, ils sont complémentaires.

Qui est le mieux placé pour monter en puissance sur XML ? Un webmaster ou un développeur ?

Les deux : Le développeur devra mettre à disposition du webmaster une information XML prête à être intégrée au site. Le webmaster sélectionnera les morceaux d'informations qu'il souhaite afficher en fonction du contexte.

Pour le stockage du contenu XML, dans quels cas faut-il s'appuyer sur une base de données et dans quels cas peut-on se satisfaire de conserver les fichier XML à proprement dit ?

C'est le même problème avec n'importe quel type d'information. Qu'elle soit structurée en XML ne change rien : les bases de données apportent des services de confidentialité, d'intégrité, de montée en charge, ... que n'apportent pas les simples systèmes de fichiers.

Comment les outils de recherche vont-ils s'adapter à ce langage ?

XML permet de définir un vocabulaire et donc une sémantique particulière à chaque information. Les outils de recherche devront s'adapter et proposer des méthodes d'accès aux informations complémentaires à la recherche plain-texte (microformats, Web sémantique, ...).

Y a-t-il un lien entre XML et le Web sémantique ?

Le Web sémantique est un Web sur lequel les informations disponibles ont un sens pour les humains et pour les machines. XML a alors une place de choix pour sa mise en œuvre dans la mesure où le principe du balisage mêle justement information pour l'homme et pour la machine. Mais ce n'est pas le seul domaine de développement de XML, dont les usages dépassent largement le domaine du Web.

Les jeux de balises extensibles, c'est bien, mais trop de jeux différents ne risque-t-il pas de nuire à l'universalité que veut véhiculer XML ?

Cette question est centrale concernant XML. Il faut d'abord distinguer les langages XML spécifiques à chaque organisation, ceux que chacun se crée sur mesure pour ses propres besoins, des langages XML standard.

Pour les seconds, il faut bien comprendre que chaque vocabulaire sera spécialisé pour une application particulière. Par exemple SVG pour modéliser les documents vectoriels en 2D, XSL-FO pour modéliser des mises en pages à destination de l'impression papier, etc. Donc, sous réserve que deux langages distincts ne répondent pas aux mêmes besoins, ces langages resteront complémentaires.

Concernant les langages "maison" de chacun, il est évident qu'ils se recoupent les uns les autres et participent d'une certaine « babélisation ». Et c'est une des raisons pour lesquelles les solutions d'EAI * sont amenées à se développer, ne serait-ce que pour assurer la transformation des informations à priori identiques. Sur le Web, des services de traduction assurent déjà la réconciliation sémantique de vocabulaires frères. Les technologies comme XSL-XSLT de transformation automatique permettent de concilier une approche spécifique à chaque contexte et des échanges dans via langages standards ou des traductions de dialecte à dialecte.

Via quels outils XML va-t-il se généraliser dans les échanges entre entreprises ?

Clairement, les outils d'EAI * prendront une part importante, ne serait-ce que pour traiter les transformations des différents langages frères, ou encore assurer la gestion transactionnelle des échanges.

Quelles relations y a-t-il entre XML et EDI ?

L'EDI * propose une infrastructure d'échange d'informations entre entreprises, ce qui comprend le format des données (par exemple EDIFACT) et le réseau à valeur ajoutée, qui assure la sécurité de la communication entre les entreprises. XML ne propose que la structuration. Un support transactionnel reste nécessaire, avec, par exemple, Internet pour le transport et ebXML * pour la modélisation des processus d'échange d'informations.

Est-ce que l'évolution de l'Internet mobile va favoriser l'utilisation du XML dans le développement d'applications ?

Certainement : dès que l'information doit être restituée sur plusieurs supports, XML est un choix d'excellence.

Est-ce que les XML Schema sont d'un usage répandu ? Peut-on se concentrer dessus et oublier DTD, ou l'usage de DTD est-il encore trop répandu ?

Clairement, les DTD * sont très répandues et sont amenées à le rester pour encore un bon moment, ne serait-ce que par leur adaptation à l'aspect documentaire de XML. Concernant l'aspect données, les limites des DTD sont telles que le passage par les schémas est nécessaire. Le fait que le W3C ait passé les schémas à l'état de recommandation permettra d'assurer un support de cette technique de définition des vocabulaires XML dans les différents produits du marché.

Si XML a l'air de résoudre autant de problèmes, pourquoi ne l'a-t-on pas inventé plus tôt ?

Pouvait-on l'inventer plus tôt ? Les premiers problèmes à résoudre étaient surtout liés à l'interopérabilité des systèmes. XML se place un cran au dessus en permettant de structurer l'information de manière standardisée, et surtout acceptée par tous.

XML est-il LA réponse au stockage structuré de n'importe quel document ?

Les documents graphiques binaires, tels que les photos ou les films, resteront dans des formats binaires, bien que certainement décrits par des fichiers XML (MPEG-7 pour la vidéo par exemple). Quant aux informations amenées à être réutilisées, recherchées, partagées, XML est effectivement *la* réponse.

Bibliographie



VII

- Michard A., "XML, langage et applications" (1997) *
- SELFHTML, "Formation XML" (2003) *
- Hoizey N., "Faisons le point sur les langages de schéma XML" (2004). *

Compléments

VIII

1. Balises et poignées de calcul

L'ingénierie documentaire met à profit deux thèses complémentaires :

- le contenu est numérisé dans sa forme signifiante : il est manipulable par la machine mais indépendamment de sa signification qui lui reste inaccessible ;
- le contenu est enrichi par des balises qui sont connues syntaxiquement et sémantiquement par la machine ; elle sait quoi en faire.

Fondamental

Le principe du balisage consiste à enrichir un contenu numérisé (dans sa forme sémiotique), sans l'altérer, pour lui ajouter des poignées qui vont être manipulables par l'ordinateur (logiquement).

Remarque

Le contenu est donc interprétable par l'homme et la machine, chacun via ce qui lui est destiné :

- l'humain interprète le contenu signifiant numérisé ;
- la machine interprète les balises ;

Exemple : XML

XML est une illustration de ce principe, puisque l'on va coupler une information sémiotique (texte, image, etc.) destinée à l'interprétation humaine, avec un ensemble de balises qui permettent de décrire formellement une structure documentaire qui sera alors manipulable par le calcul.

2. La structuration logique




Un document peut être décrit comme une collection d'objets comportant des objets de plus haut niveau composés d'objets plus primitifs. Les relations entre ces objets représentent les relations logiques entre les composants du document. Par exemple [...] un livre est divisé en chapitres, chaque chapitre en sections, sous-sections, paragraphes, etc. Une telle organisation documentaire est appelée représentation de document structuré. (traduit depuis la préface de Structured documents *)



 *Définition : Structuration logique*

On appelle structuration logique d'un contenu une inscription explicitant la structure de ce contenu en fonction de son organisation et des attributs intrinsèques qui le caractérisent et non en fonction de propriétés de présentation sur un support.

 *Définition : Document abstrait*

Un document décrit par sa structure logique est appelé document abstrait, on parle aussi de *document structuré*.

 *Définition : Structuration physique*

On appelle structuration physique ou mise en forme d'un contenu une inscription décrivant la façon dont ce contenu doit être présenté sur un support donné.

 *Définition : Document formaté*

Un document décrit par sa structure physique est appelé document formaté, c'est en général également ce dont on parle quand on parle simplement de document.

 *Remarque*

Il est possible de calculer une ou plusieurs structurations physiques pour une même structuration logique. Il est possible de calculer d'autant plus de structurations physiques que la structuration logique est indépendante de ses supports de présentation.


 *Remarque*

La structuration logique est associée au fond du contenu, tandis que la structuration physique est associée à sa forme sur un support.

 *Complément : Voir aussi*

Langages XML orienté documents (cf. p.15)

3. Exemple de structuration logique

 *Exemple : Un exercice structuré logiquement*

Soit la structuration logique d'un exercice :

- 1 Exercice = {Enonce, Question, Indice, Solution}
- 2 avec
- 3 Enonce = Soit un triangle rectangle disposant d'un angle de 30 degrés.
- 4 Question = Donner la valeur des autres angles du triangle.
- 5 Indice = La somme des angles d'un triangle est égale à 180 degrés.
- 6 Solution = 90 et 60 degrés.

Il est possible à partir de cette représentation de calculer différentes présentations. Pour l'écran on peut générer une présentation HTML, en laissant la solution en hyperlien cliquable. Pour le papier on peut générer une présentation PDF, en affichant la solution sur une page séparée de l'énoncé. Pour un usage multimédia on pourra générer une présentation SMIL, avec affichage de l'énoncé, lecture de la question, et affichage de la solution après un temps de pause.

Notons que si l'on avait choisi une des représentations physiques, plutôt que la représentation logique, il n'aurait pas été possible de générer les autres représentations.

Exemple : Un exercice mis en forme

Soit la mise en forme en HTML du même exercice :

```

1 <HTML>
2 <BODY>
3 Soit un triangle rectangle disposant d'un angle de 30 degrés. </BR>
4 <B> Donner la valeur des autres angles du triangle. </B> </BR>
5 <A HREF="ex001i01.html"> Vous avez besoin d'aide ? </A> </HR>
6 <A HREF="ex001s01.html"> Vérifier votre réponse ! </A>
7 </BODY>
8 </HTML>

```

On voit que dans ce format la structure logique n'apparaît plus explicitement et qu'il n'est plus possible d'identifier l'énoncé, la question et la solution sans comprendre le contenu.

Fondamental

L'exemple montre que l'on peut calculer la mise en forme à partir de la structure logique, *mais non l'inverse*.

4. Définition des chaînes éditoriales XML

Définition : Chaîne éditoriale XML

Une chaîne éditoriale XML est un système informatique permettant la production de documents structurés en XML, grâce à :

- un éditeur *WYSIWYM* ("What You See Is What You Mean" ou "ce que vous voyez est ce que vous voulez dire"),
- la publication *polymorphe* (Web, papier, diaporama, etc.),
- et la *ré-éditorialisation* sans recopie (dés-agrégation / ré-agrégation de documents).

Complément : Implémentations

On peut distinguer trois grandes modalités pour implémenter une chaîne éditoriale :

- L'implémentation « sur mesure » consiste à mobiliser différentes briques logicielles et à programmer un logiciel spécifique pour un besoin particulier. Par exemple en couplant un éditeur du marché avec un outil de stockage et en réalisant des moteurs de publication sous la forme de feuilles XSL-XSLT.

- L'implémentation « modèle dédié » consiste à implémenter en dur un modèle de chaîne éditoriale pour un besoin assez transversal, tel que la publication de livres ou de modules de cours universitaires.
- L'implémentation « générique » consiste à réaliser un système paramétrable, indépendant d'un modèle en particulier, et à le configurer en fonction des besoins. Il s'agit alors de ce que l'on peut appeler un Système de Gestion de Chaînes Éditoriale par analogie aux Systèmes de Gestion de Bases de Donnée

Complément : Historique : de la recherche au développement économique

Le concept de chaîne éditoriale trouve son origine dans l'édition et dans la presse : il consiste à organiser les tâches de production et de publication, en séparant les métiers intervenant dans le processus. Une chaîne éditoriale est donc un processus avant d'être un outillage, qui est né d'un besoin d'industrialisation. Historiquement, les deux technologies ayant permis l'implémentation de chaînes éditoriales numériques sont LaTeX (1982) et SGML (norme ISO depuis 1986). XML (standard W3C depuis 1998) est aujourd'hui la technologie de référence pour la réalisation de chaînes éditoriales. C'est sa maturité qui a permis de sortir ce procédé des domaines auxquels il était confiné jusque là (documentation technique stratégiques dans l'aviation ou publication scientifique typiquement).

En 1999, l'Université de Technologie de Compiègne réalise Scenari, premier environnement intégré de conception de chaînes éditoriales XML proposant un langage déclaratif de modélisation et de publication de haut niveau (Scenari, la chaîne éditoriale libre ^{*}, Ingénierie des connaissances et des contenus ^{*}).

5. Définition de Scenari

Définition : Scenari

Scenari est un logiciel libre permettant de créer des contenus multimédia structurés selon une approche innovante : celle de la chaîne éditoriale XML.

Il permet de répondre aux enjeux actuels de la création et la gestion des documents numériques :

- réduction des coûts de production et de maintenance
- maîtrise de la qualité des publications
- multiplication des canaux de diffusion
- diversification des modalités de communication
- pérennisation de l'information

Complément

<http://scenari-platform.org>

6. Un langage pour publier les documents XML

XML lorsqu'il est utilisé pour définir des formats documentaire métier est un format de *représentation* de l'information, et non un format de *publication* (de présentation) de cette information : donc un tel fichier XML *n'est pas utilisable tel que par un lecteur*.

XML ne peut donc être utilisé pour des langages abstraits que si l'on est capable de transformer les documents sources en documents publiés lisibles grâce à un format de présentation : HTML par exemple dans le cas de publication Web, ou PDF pour l'impression.

7. Définition de XSL-XSLT

Définition : XSL-XSLT

XSL-XSLT est une partie du standard W3C XSL qui a trait à la transformation des documents XML (l'autre partie étant XSL-FO).

XSL-XSLT est un langage de programmation déclaratif écrit en XML (un programme XSL-XSLT est un document XML).

- XSL-XSLT est langage de manipulation de document XML (fondé sur XPath et sur le modèle arborescent de représentation des documents XML)
- XSL-XSLT est utilisé pour transformer un document XML source dans un autre format, typiquement HTML, mais aussi tout autre format codé sur des caractères dont la syntaxe est connue.
- XSL-XSLT est aussi utilisé pour faire des changements de schéma XML (export d'un XML vers un autre XML structuré différemment) par exemple pour échanger des données selon un standard.

Remarque : XSL-XSLT, XSL-FO, XSLT, XSL, FO

On parle souvent (par simplification) de XSL ou de XSLT pour désigner XSL-XSLT et de FO pour désigner XSL-FO.

XSL utilisé seul désigne donc par convention XSL-XSLT (et non XSL-FO).

Exercices



1. Exercices : Culture XML

1.1. Exercice

[solution n°1 p.48]

1.1.1. Exercice

Historique d'XML : Cocher toutes les propositions correctes.

- XML a été créé dans la continuité du SGML
- XML est une norme ISO
- XML est un standard W3C
- XML est plus récent que HTML

1.1.2. Exercice

XML versus HTML : Cocher toutes les propositions correctes.

- HTML est un langage de mise en forme
- XML est un langage de mise en forme
- HTML définit un jeu de balises fixé a priori et donc non-extensibles
- XML est beaucoup plus générique, il permet de définir des langages

1.1.3. Exercice

Qu'est ce qui est nécessaire pour avoir un document XML bien formé ?

- Un entête XML
- Une référence à une DTD
- Une référence à un XML Schema
- Une unique balise racine qui encadre toutes les autres
- Que chaque balise ouverte soit fermée

1.1.4. Exercice

Que permet un schéma XML ?

- D'indiquer la mise en page pour un type de document et un format cible.
- De vérifier automatiquement qu'un document correspond à une structure type.
- De définir une grammaire documentaire.

1.2. Exercice : Galaxie XML

[solution n°2 p.49]

Ranger les termes suivants dans les bonnes catégories.

Open Document Format	SGML	XHTML	SMIL	XML	W3C	HTML
SVG	DocBook	OASIS				

Méta-langage	Format de publication XML	Format de publication SGML	Langage de structuration documentaire	Organisme de standardisation
--------------	---------------------------	----------------------------	---------------------------------------	------------------------------

1.3. Exercice : Historique XML

[solution n°3 p.49]

Indiquer pour chaque format en quelle année il a été standardisé.

1995	2000	1998	1986
------	------	------	------

SGML	XML	HTML 2.0	XHTML 1.0
------	-----	----------	-----------

1.4. Exercice :

Discutez les propositions suivantes en 5 à 10 lignes, illustrez chacune à l'aide d'un exemple.

Question 1

[solution n°4 p.49]

Pourquoi est-il utile d'utiliser des schémas pour valider les documents XML ?

Question 2

[solution n°5 p.49]

Pourquoi le XML n'a-t-il pas encore remplacé le HTML ?

2. Exercices : Syntaxe XML

2.1. Exercice

[solution n°6 p.50]

Soit le fichier ci-après, quelles sont les assertions vraies ?

```

1 <papier type="scientifique">
2   <titre>Réinterroger les structures documentaires</titre>
3   <sousTitre>De la numérisation à l'informatisation</sousTitre>
4   <auteur>Stéphane Crozat</auteur>
5   <auteur>Bruno Bachimont</auteur>
6   <resume>Nous proposons dans cet article d'aborder ...</resume>
7   <abstract>In this paper we define...</abstract>
8   <motsCles>
9     <terme>Ingénierie des connaissances</terme>
10    <terme>XML</terme>
11    <terme>Document</terme>
12  </motsCles>
13  <keywords>
14    <word>Knowledge engineering</word>
15    <word>XML</word>
16    <word>Document</word>
17  </keywords>
18  <publication date="2004-07-05"/>
19  <version maj='1' min='0'/>
20  <ressource uriSrc="http://archivesic.ccsd.enrs.fr/docs/00/06/23/97/PDF/sic_00001015.pdf"/>
21 </papier>

```

Ce fichier n'est pas un fichier XML

Ce fichier est un fichier XML, il possède un élément racine unique qui contient tous les autres.

Ce fichier est un fichier XML, tous ses éléments sont inclus dans les uns dans les autres.

Ce fichier est un fichier XML, toutes ses balises ouvertes sont fermées.

Ce fichier est un fichier SGML.

2.2. Exercice

[solution n°7 p.50]

Compléter les trous avec les balises *fermantes* adéquates.

```
<?xml version="1.0"?>
<document>
  <entete>
    <titre>Document à corriger</titre>
    <auteur>Stéphane Crozat [ ]
    <date>01-03-01 [ ]
    <version numero="1"/>
    [ ]
  <corps>
    <division titre="Première partie">
      <paragraphe>Ce texte doit être <important>corrigé [ ]
    </paragraphe> <paragraphe>Le contenu est sans importance ici</paragraphe>
    </division>
    <division titre="Seconde partie">
      <paragraphe>Ai-je le droit de mettre du texte ici ?</paragraphe>
    [ ]
    [ ]
    [ ]
  </corps>
</document>
```

2.3. Exercice

[solution n°8 p.51]

Fermer les balises du fichier suivant afin qu'il soit un document XML bien formé.

```
<docAbstrait>
  <contenu>
    <paragraphe>Lorem ipsum dolor sit amet, consectetur adipiscing elit.</ [ ] >
    <Texte>Nulla erat tellus, molestie a ultrices sed, gravida eget ipsum.
    <para>Phasellus lacinia, ipsum vitae interdum tincidunt.
    <P>Nullam pulvinar diam et tellus ullamcorper eleifend.</ [ ] ></ [ ] ></
    [ ] ></ [ ] >
    <ligne>Nunc eu lectus in diam tempus adipiscing in rhoncus elit.</ [ ] ></ [ ] >
```

Questions de synthèse



HTML est-il un langage XML ?

.....

.....

.....

.....

.....

.....

.....

SGML est plus libre que XML, par exemple il autorise l'omission des balises fermantes. Quels avantages XML présente-t-il d'un point de vue informatique ?

.....

.....

.....

.....

.....

.....

.....

Quels outils sont nécessaires pour faire du XML ?

.....

.....

.....

.....

.....

.....

.....

Les autres média que le texte sont-ils aussi adaptés à une structuration en XML ?

.....

.....

.....

.....

.....

.....
.....



Solutions des exercices



> Solution n°1

Exercice p. 42

Exercice

Historique d'XML : Cocher toutes les propositions correctes.

- XML a été créé dans la continuité du SGML
- XML est une norme ISO
- XML est un standard W3C
- XML est plus récent que HTML

Exercice

XML versus HTML : Cocher toutes les propositions correctes.

- HTML est un langage de mise en forme
- XML est un langage de mise en forme
- HTML définit un jeu de balises fixé a priori et donc non-extensibles
- XML est beaucoup plus générique, il permet de définir des langages

Exercice

Qu'est ce qui est nécessaire pour avoir un document XML bien formé ?

- Un entête XML
- Une référence à une DTD
- Une référence à un XML Schema
- Une unique balise racine qui encadre toutes les autres
- Que chaque balise ouverte soit fermée

Exercice

Que permet un schéma XML ?

- D'indiquer la mise en page pour un type de document et un format cible.
- De vérifier automatiquement qu'un document correspond à une structure type.
- De définir une grammaire documentaire.

> **Solution n°2**

Exercice p. 43

Ranger les termes suivants dans les bonnes catégories.

Méta-langage	Format de publication XML	Format de publication SGML	Langage de structuration documentaire	Organisme de standardisation
XML	SVG	HTML	DocBook	W3C
SGML	SMIL			OASIS
	XHTML			
	Open Document Format			

> **Solution n°3**

Exercice p. 43

Indiquer pour chaque format en quelle année il a été standardisé.

SGML	XML	HTML 2.0	XHTML 1.0
1986	1998	1995	2000

NB : La première version d'HTML a été spécifiée en 1993, mais le premier standard officiel (W3C) est HTML 2.0 en 1995.

> **Solution n°4**

Exercice p. 43

Si l'on n'utilise pas de schéma pour valider les documents XML, ceux-ci sont hétérogènes et difficiles à exploiter automatiquement.

Par exemple, on ne peut pas utiliser de XSLT pour les transformer automatiquement.

> **Solution n°5**

Exercice p. 43

Le XML n'a pas vocation à remplacer le HTML. XML permet aux architectes de sites de se concentrer sur l'information à fournir, et dans un deuxième temps de la mettre en forme pour des besoins spécifiques : sur un téléphone portable, pour imprimer, pour les navigateurs Web, ...

> **Solution n°6**

Exercice p. 44

Soit le fichier ci-après, quelles sont les assertions vraies ?

```

1 <papier type="scientifique">
2  <titre>Réinterroger les structures documentaires</titre>
3  <sousTitre>De la numérisation à l'informatisation</sousTitre>
4  <auteur>Stéphane Crozat</auteur>
5  <auteur>Bruno Bachimont</auteur>
6  <resume>Nous proposons dans cet article d'aborder ...</resume>
7  <abstract>In this paper we define...</abstract>
8  <motsCles>
9    <terme>Ingénierie des connaissances</terme>
10   <terme>XML</terme>
11   <terme>Document</terme>
12 </motsCles>
13 <keywords>
14   <word>Knowledge engineering</word>
15   <word>XML</word>
16   <word>Document</word>
17 </keywords>
18 <publication date="2004-07-05"/>
19 <version maj='1' min='0'/>
20 <ressource uriSrc="http://archivesic.ccsd.cnrs.fr/docs/00/06/23/97/PDF/sic_00001015.pdf"/>
21 </papier>

```

Ce fichier n'est pas un fichier XML



Ce fichier est un fichier XML, il possède un élément racine unique qui contient tous les autres.

Ce fichier est un fichier XML, tous ses éléments sont inclus dans les uns dans les autres.

Ce fichier est un fichier XML, toutes ses balises ouvertes sont fermées.

Ce fichier est un fichier SGML.

> **Solution n°7**

Exercice p. 45

Compléter les trous avec les balises *fermantes* adéquates.

<?xml version="1.0"?>

<document>

```

<entete>
  <titre>Document à corriger</titre>
  <auteur>Stéphane Crozat</auteur>
  <date>01-03-01</date>
  <version numero="1"/>
</entete>
<corps>
  <division titre="Première partie">
    <paragraphe>Ce texte doit être <important>corrige</important>
  </paragraphe>
  <paragraphe>Le contenu est sans importance ici</paragraphe>
  </division>
  <division titre="Seconde partie">
    <paragraphe>Ai-je le droit de mettre du texte ici ?</paragraphe>
  </division>
</corps>
</document>

```

> Solution n°8

Exercice p. 45

Fermer les balises du fichier suivant afin qu'il soit un document XML bien formé.

```

<docAbstrait>
  <contenu>
    <paragraphe>Lorem ipsum dolor sit amet, consectetur adipiscing elit.</paragraphe>
    <Texte>Nulla erat tellus, molestie a ultrices sed, gravida eget ipsum.
    <para>Phasellus lacinia, ipsum vitae interdum tincidunt.
    <P>Nullam pulvinar diam et tellus ullamcorper eleifend.</P></para></Texte></contenu>
    <ligne>Nunc eu lectus in diam tempus adipiscing in rhoncus elit.</ligne></docAbstrait>

```

Les balises doivent être fermées dans l'ordre dans lequel elles sont ouvertes, pour que chaque élément soit inclus dans son père. La syntaxe XML est sensible à la casse : P est différent de p.

Glossaire



EAI

L'EAI (Enterprise Application Integration ou Intégration des applications d'entreprise) désigne à la fois les solutions et les méthodes destinées à assurer l'intégration des différentes composantes du système d'information. Il s'agit de gagner en souplesse et de baisser les coûts de maintenance des interfaces inter-applicatives. Les chantiers d'EAI sont souvent un préalable nécessaire à des projets e-business (source : Le Journal Du Net).

Sérialisation

Processus consistant à enregistrer des données en mémoire vive (par exemple des objets) sous une forme permettant leur persistance, typiquement sur une mémoire secondaire.

Abréviations

CSS : Cascading Style Sheet

DTD : Document Type Definition

EAI : Enterprise Application Integration

ebXML : electronic business XML

EDI : Electronic Data Interchange ou Echange de Données Informatisé

HTML : HyperText Markup Language

IDE : Integrated Development Environment (Environnement de Développement Intégré)

W3C : World Wide Web Consortium

XML : eXtensible Markup Language

Bibliographie



Stéphane Crozat, *Scenari, la chaîne éditoriale libre*, Accès Libre, Eyrolles, 2007.

Gérard Dupoirier, *Technologie de la GED : Techniques et management des documents électroniques*, Hermes, 1995.

Jacques André, Richard Furuta, Vincent Quint, *Structured documents*, Cambridge University Press, 1989.

Bruno Bachimont, *Ingénierie des connaissances et des contenus : le numérique entre ontologies et documents*, Lavoisier, Hermès, 2007

Michard Alain, *XML, langage et applications*, Eyrolles, 1997.

Webographie



Nicolas Hoizey, *Faisons le point sur les langages de schéma XML*, <http://www.clever-age.com/veille/clever-link/faisons-le-point-sur-les-langages-de-schema-xml.html>, 2004.

SELFHTML, *Formation XML*, <http://fr.selfhtml.org/xml/>, 2003.

W3C, *Extensible Markup Language (XML) 1.0 (Fifth Edition)*, <http://www.w3.org/TR/2008/REC-xml-20081126/>, consulté en 2008.

XML, <http://www.w3.org/XML>.

Index



Attribut

p. 20, 28

DTD

p. 28

Méta-langage

p. 7

Valide

p. 28

Balise

p. 6

Elément

p. 28

SGML

p. 8

W3C

p. 5

Bien formé

p. 19

HTML

p. 9

Syntaxe

p. 8, 22, 22

XML

p. 5, 5, 6, 7, 16