



Interacțiunea Om-Calculator

Laborator 1

XML

Cuprins

1	Obiective laborator	1
2	Despre XML – Scurta istorie	1
3	Sintaxa	2
4	Cazuri de utilizare	3
5	Tooluri	5
6	Șabloane de proiectare a fișierelor XML	5
6.1	Șabloane pentru folosirea Metadatelor (date despre date)	5
6.2	Abstractizare	5
6.3	Organizarea datelor	6
6.4	Altele	6
7	Concluzii	6
8	Referinte	7
	Anexa 1 – Exemplu de document XML	8
	Anexa 2 Fisierul document.xml pentru un document Microsoft Office .docx in care apare textul „Hello World”	9

1 Obiective laborator

- Familiarizarea cu standardul XML si utilitatea acestuia

2 Despre XML – Scurta istorie

XML – (eXtensible Markup Language) este un limbaj de adnotare/structurare a datelor. Istoria aparitiei sale porneste de la SGML (Standard Generalized Markup Language) aparut in anii ‘80. Acest limbaj isi propunea la aparitie sa creeze documente ce pot fi analizate de catre masini prin introducerea de marcaje (sau taguri). SGML nu a fost un succes deoarece era foarte complex si astfel crearea de programe care sa-l analizeze.



În anii 90 la CERN a apărut HTML-ul (HyperText Markup Language). La baza acestui limbaj a fost SGML-ul și scopul său a fost marcarea documentelor astfel încât să poată fi transmise prin rețea. Unul din primele documente care au stat la baza WWW-ului și respectiv HTML-ului le puteți găsi la [1]. Una din ideile acestui document este că browserele trebuie să ignore marcasele și atributele pe care nu le înțeleg. Această idee a avut ca efect din câte cred eu apariția unui limbaj mai simplu care să poată fi parsat rapid de către browsere. Și astfel s-a dezvoltat HTML-ul.

Succesul avut de Web a cauzat o dezvoltare rapidă a browserelor care analizau marcasele HTML. Astfel au apărut o gamă largă de marcase și atribute care puteau fi scrise fără prea multe constrângeri. Browserele au trebuit să țină cont de aceste constrângeri și au devenit foarte complexe. Pe de altă parte s-a sesizat utilizarea HTML-ului pentru adnotarea documentelor și o slăbiciune de-a sa – faptul că nu se pot adăuga marcase noi.

În acest context în decembrie 1997 W3C (World Wide Web Consortium [2]) a lansat ca recomandare propusă versiunea 1.0 de XML [3]. În acest document se afirmă că XML este un subset al SGML care păstrează caracteristici cum ar fi posibilitatea creării de marcase, posibilitatea de validare, posibilitatea de a fi citit și înțeles de către oameni. Totodată se afirmă că XML este creat pentru a fi folosit pentru adnotări, schimb de date, publicare documente și prelucreare automată de către clienți (agenți) inteligenți. Pastrand un subset al SGML împreună cu un set de constrângeri documentele XML vor putea fi prelucrate foarte rapid de către parsere neavând limitele date de complexitatea SGML. Vom vedea care sunt cazurile de utilizare actuale pentru XML după ce îi vom înțelege sintaxa – după care îi vom înțelege mai clar și avantajele și slăbiciunile.

3 Sintaxa

Am vorbit în prima parte a acestei prezentări despre marcase și atribute. Acum putem să încercăm să le și definim. Pentru a ilustra sintaxa vom folosi exemplul din Anexa 1. Un marcaj este un text care descrie sensul sau structura unor date. Acest text este inserat între semnele „<” și „>”. Marcasele pot avea început și sfârșit ca de exemplu:

```
<nume>Popescu</nume> (a)
```

sau pot fi vide (fără text în interiorul marcajului)

```
<casatorit/>
```

Marcasele sunt case-sensitive. La întâlnirea unui marcaj de tipul <nume>Popescu</NUME> parserul va da o eroare.

Un element este tot ceea ce se află între marcajul de început și marcajul de sfârșit. (a) este un element. Numele elementului este „nume”. Aceste nume trebuie să îndeplinească câteva condiții [4].

- Pot să conțină litere, cifre și alte caractere
- Nu pot începe cu cifre și cu semne de punctuație
- Nu pot începe cu „xml”
- Nu pot conține spații

Un element poate conține mai multe elemente. Acest lucru determină apariția unor relații de rudenie între elemente și a unor constrângeri suplimentare. Cea mai importantă dintre aceste constrângeri este aceea că într-un fișier XML trebuie să existe un element rădăcină care să le



contina pe toate celelalte. Elementul radacina in Anexa 1 este <persoana>. O a doua constrangere este ca elementele nu pot fi incrucisate. Astfel urmatoarea constructie este incorecta <tranzactie><data>2005-11-05</tranzactie></data>. Cu alte cuvinte marcajele trebuie inchise invers fata de cum sunt deschise.

Relatiile de rudenie pot fi parinte-copil sau frate-frate. Un exemplu de relatie parinte-copil in cadrul unui element este urmatoarea.

```
<tranzactie>
  <data>20-07-2005</data>
  <ora>13:00:00</ora>
</tranzactie>
```

In acest exemplu elementul <tranzactie> este parinte al elementelor <data> si <ora> care sunt “frati”.

Un element poate sa contina unul sau mai multe atribute. Atributele se specifica in cadrul marcajului de inceput al elementului si sunt de tip nume_atribut=”valoare”. Un exemplu in Anexa 1 este in cadrul marcajului urmator: <suma moneda=”RON”> 1000 </suma>. Atributul in acest caz este moneda iar valoarea atributului este RON. Valorile atributelor se pun intotdeauna intre ghilimele.

Ar mai fi de mentionat ca intotdeauna un document XML incepe cu un rand in care se mentioneaza versiunea limbajului si codificarea folosita. In exemplul din Anexa 1 se specifica faptul ca se foloseste versiunea 1.0 si ca se foloseste codificarea ISO-8859-2 care este codificarea folosita pentru limba romana.

Respectand aceste simple reguli de sintaxa vom putea obtine ceea ce se numeste un document XML **bine format**. Acest document se valideaza la randul lui cu ajutorul unui alt fisier in care se definesc tagurile folosite. Acest alt fisier poate fi de tip **DTD** sau **XML Schema** si va fi prezentat in laboratoarele viitoare.

Astfel, un document XML **“bine formatat”** este un document ce foloseste o sintaxa corecta si anume:

- Documentele XML trebuie sa aiba un element radacina (root)
- Elementele XML trebuie sa aiba tag-uri de sfarsit
- Tag-urile XML sunt sensibile litere mari/litere mici
- Elementele XML trebuie sa fie corect imbricate
- Atributele XML trebuie sa fie incluse intre ghilimele

Un document XML **valid** este un document “bine formatat” care se supune si regulilor DTD (Document Type Definition) sau XML Schema.

4 Cazuri de utilizare

Din introducerea si din scurta trecere in revista a sintaxei am putut observa ca XML-ul este un limbaj care permite structurarea si adnotarea datelor intr-un mod lizibil atat de catre oameni cat si de catre calculatoare. Avand aceste calitati putem oarecum deduce unde se poate folosi in mod practic si unde este mai putin recomandat.



Poate cel mai evident caz de utilizare este cel al fișierelor de configurare. În momentul de față foarte multe aplicații își păstrează fișierele de configurare în XML. Motivele sunt cele enunțate mai sus – **parsare rapidă pentru sisteme automate, lizibilitate pentru oameni și structura logică a datelor.**

Un alt caz în care se folosește XML și în care se pare că va deveni un standard este cel al stocării informațiilor din fișiere tip office. Inițiativa a aparținut dezvoltatorilor de soluții office open-source care au standardizat un format numit Open Document[5][6]. Acest format presupune salvarea oricărui document de tip office – document, foaie de calcul (spreadsheet), prezentare – într-un format XML. Informația este stocată în mai multe fișiere care sunt arhivate și astfel utilizatorul poate vedea numai un singur fișier de tip *.sxw sau *.odt, etc. Începând cu Suita Office 2007 și Microsoft a abordat o arhitectură bazată pe XML în reprezentarea internă a fișierelor (în Anexa 2 este prezentat un fișier .docx în care este scris textul „Hello World”).

De ce s-a ales XML pentru un astfel de format? Pentru că oferă avantajele structurării informației și mai ales pentru că s-au dezvoltat deja foarte multe alte limbaje și standarde pe baza lui. Aceste limbaje (XHTML, SVG, XSL, SMIL, XLink, XForms, MathML, Dublin Core) sunt folosite și ele în acest format și deoarece toate au la bază XML interoperabilitatea este garantată. XML are și un dezavantaj care este dat de faptul că introduce informație redundantă (overhead). Informația redundantă este reprezentată de numele marcărilor care se repetă în document. Această problemă a fost rezolvată de cei de la Open Office prin arhivarea fișierelor XML într-o arhivă zip. Astfel s-a rezolvat problema spațiului suplimentar apărut prin introducerea unei mici pierderi în viteză.

Un alt caz de utilizare este în stocarea efectivă a datelor – mai precis baze de date. În momentul de față XML poate fi folosit cu succes pentru baze de date mici. În momentul în care bazele de date cresc ca dimensiune overheadul devine inacceptabil. Din cauza spațiului mare ocupat și a necesității accesului foarte rapid la date arhivarea fișierelor XML cu date nu poate fi o soluție și în acest caz. Oricum, multe interfețe pentru baze de date permit exportarea bazei de date în format XML – atât definiția tabelelor și a relațiilor între ele cât și datele propriu zise. Două exemple de aplicații care au această facilitate sunt Microsoft Access și PHPMyAdmin[7]. Un alt exemplu de caz în care XML-ul nu a avut succes este legat tot de stocarea unor cantități mari de date și de acces rapid la acestea – mai precis păstrarea fișierelor jurnal ale serverelor. În acest domeniu este folosit printre altele un standard numit CLF – Common Logfile Format – care are avantajul că nu introduce informație suplimentară [8].

XML stă la baza multor tehnologii și limbaje noi. Astfel pentru prezentare de conținut pe web există XHTML și WML (pentru dispozitive mobile), pentru descrierea unor fișiere XML există XSchema, pentru transformarea fișierelor XML pentru a fi reprezentate există XSL, pentru realizarea unor prezentări există SMIL, pentru descrierea obiectelor grafice există SVG, pentru reprezentarea semanticii unor domenii există RDF și/sau OWL, pentru schimb de informații între aplicații există SOAP și așa mai departe, lista aceasta fiind departe de a cuprinde toate tehnologiile/limbajele existente bazate pe XML. De ce au apărut atâtea limbaje bazate pe XML? Pentru că are o structură simplă și totuși strictă, pentru că este foarte ușor de analizat și foarte ușor de definit și dezvoltat.



5 Tooluri

Există o largă varietate de instrumente pentru editarea fișierelor XML. Din seria celor comerciale cele mai cunoscute sunt Altova (<http://www.altova.com/xml-editor/>), Oxygen XML (<http://www.oxygenxml.com/>) și StylusStudio (<http://www.stylusstudio.com/xml/editor/>). Ambele oferă versiuni de „trial” pentru a vedea ce facilități oferă.

Datorită simplității sintaxei o gamă largă de editoare au facilități pentru crearea fișierelor XML: Notepad++ (<http://notepad-plus.sourceforge.net>), din gama editoarelor mai simple, și Eclipse sau Netbeans din gama editoarelor mai complexe.

6 Șabloane de proiectare a fișierelor XML

Site-ul XMLPatterns.com prezintă un mare număr de șabloane de construcție a fișierelor XML. Folosirea acestor șabloane permite crearea unor documente mai ușor de citit, folosit și prelucrat, permite gestionarea mai facilă a informațiilor oferite.

În această secțiune vom rezuma șabloanele cele mai importante.

6.1 Șabloane pentru folosirea Metadatelor (date despre date)

Când există metadata într-un document pot fi folosite mai multe soluții pentru integrarea lor:

- **Head + Body** – varianta folosită și în HTML – datele despre date se pun într-o secțiune specială a documentului numită head și datele efective se vor plasa în secțiunea body.
- **Metadatale se pun înaintea datelor în document.** Motivul este destul de simplu și anume în condițiile în care procesarea XML-ului se face secvențial s-ar putea ca un proces, după ce interpretează meta-datele să nu mai dorească să prelucreze și restul documentului și implicit să nu-l mai încarce sau parseze
- În cazul în care există un mare număr de **metadata**, poate fi util ca acestea să fie **plasate într-un fișier separat**. Acest lucru este recomandabil mai ales în cazul în care este vorba de metadata comune mai multor fișiere (vezi ontologii și RDF)

6.2 Abstractizare

Cum se pot alege și respectiv grupa elementele astfel încât structura XML-ului format să fie cât mai inteligibilă.

- **Containere** – în containere se pot pune mai multe elemente diferite dar care caracterizează un „părinte” comun dintr-un anumit punct de vedere. De exemplu pentru un om putem pune într-un container „date de identificare” elementele „nume”, „prenume”, „CNP”
- **Colecții** – când în cadrul unui element pot apărea ca și copii mai multe elemente de același tip acestea este recomandabil să se grupeze ca și copii al unui element separat. Exemplu un om are mai mulți copii ceea ce putem exprima în mai multe feluri:

```
<om>
```

```
  <date_identificare>...</date_identificare>
```

```
  <copil></copil>
```



```
<copil></copil>
...
</om>
Sau
<om>
  <date_identificare>...</date_identificare>
  <copii>
    <copil></copil>
    <copil></copil>...
  </copii>
...
</om>
```

A doua variantă presupune adăugarea elementului „copii” care reprezintă o colecție de elemente de tip „copil”

- **Elementele este bine să aibă o corespondență în cadrul domeniului descris** (v. Descrierea claselor și a proprietăților la programarea orientată obiect)

6.3 Organizarea datelor

- **Declararea elementelor înainte de a fi referite.** Cu ajutorul atributelor de tip IDREF elemente pot fi referite în cadrul altor elemente. Este recomandabil ca întâi să apară în document elementul și apoi referința (tot pentru ușurința procesării)
- **Referirea datelor și nu repetarea lor.** Datele pot fi declarate folosind declarații de tip Entity sau pot fi identificate în mod unic folosind atribute de tip ID și referite prin atribute de tip IDREF
- **Clasificare folosind atribute** – când pot clasifica anumite elemente în mai multe feluri se pot introduce atribute pentru a nu complica ierarhia. În exemplul de mai sus elementului copil i se poate introduce atributul „sex” care practic introduce o clasificare suplimentară fără a adăuga elemente suplimentare.

6.4 Altele

- **Reutilizare atribute** – dacă mai multe elemente au aceleași atribute nu trebuie definite atributele de mai multe ori
- **Nume scurte dar inteligibile** – „loc_de_munca” este preferat elementului “serv”.

Pentru mai multe detalii și pentru mai multe șabloane vizitați <http://xmlpattern.com>

7 Concluzii

XML-ul este un limbaj cu o sintaxa simplă și care permite doar structurarea datelor într-o manieră proprie prin definirea propriilor taguri. Această facilitate de structurare a datelor a permis folosirea sa pentru a dezvolta limbaje noi precum și pentru a fi folosit în noi standarde de stocare a datelor. Nu în ultimul rând XML-ul poate fi folosit pentru a schimba date între aplicațiile care au nevoie de a comunica într-un limbaj comun.



Chiar si in ciuda dezavantajelor pe care le are (overhead) calitatile XML mai sus mentionate ii permit sa fie folosit cu succes in multe domenii si sa fie piatra de temelie pentru si mai multe domenii.

8 Referinte

- [1] „The first version of HTML” - <http://www.w3.org/History/19921103-hypertext/hypertext/WWW/MarkUp/MarkUp.html>
- [2] World Wide Web Consortium – <http://www.w3c.org>
- [3] Press Release XML 1.0 - <http://www.w3.org/Press/XML-PR>
- [4] Tutorial XML – http://www.w3schools.com/xml/xml_elements.asp
- [5] Standardizarea Open Document – siteul web OASIS - http://www.oasis-open.org/news/oasis_news_05_23_05.php
- [6] Oasis OpenDocument Essentials - <http://books.evc-cit.info/odbook/book.html>
- [7] PHPMyAdmin – web site - http://www.phpmyadmin.net/home_page/index.php
- [8] CLF – formatul de inregistrare al jurnalelor - http://www.w3.org/Daemon/User/Config/Logging.html#common_logfile_format
- [9] Șabloane XML – www.XMLPattern.com



Anexa 1 – Exemplu de document XML

```
<?xml version="1.0" encoding="ISO-8859-2"?>
  <persoana>
    <nume>Popescu</nume>
    <prenume>Ion</prenume>
    <informatii_personale>
      <oras>Bucuresti</oras>
      <strada>Calea Mosilor</strada>
      <nr>100</nr>
      <tel_fix>0213003000</tel_fix>
      <tel_mobil></tel_mobil>
      <e-mail>popescu.ion@gmail.com</e-mail>
      <casatorit />
    </informatii_personale>
    <tranzactii>
      <tranzactie>
        <data>20-07-2005</data>
        <ora>13:00:00</ora>
        <suma moneda="RON">1000</suma>
        <directia>venit</directia>
        <descriere>salariu</descriere>
      </tranzactie>
      <tranzactie>
        <data>20-07-2005</data>
        <ora>13:05:00</ora>
        <suma moneda="RON">35</suma>
        <directia>plata</directia>
        <descriere>telefon</descriere>
      </tranzactie>
    </tranzactii>
  </persoana>
```




Anexa 2 Fisierul document.xml pentru un document Microsoft Office .docx in care apare textul „Hello World”

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<w:document
xmlns:wpc=http://schemas.microsoft.com/office/word/2010/wordprocessingCanvas
xmlns:mo=http://schemas.microsoft.com/office/mac/office/2008/main
xmlns:mc=http://schemas.openxmlformats.org/markup-compatibility/2006
xmlns:mv="urn:schemas-microsoft-com:mac:vml"
xmlns:o="urn:schemas-microsoft-com:office:office"
xmlns:r=http://schemas.openxmlformats.org/officeDocument/2006/relationships
xmlns:m=http://schemas.openxmlformats.org/officeDocument/2006/math
xmlns:v="urn:schemas-microsoft-com:vml"
xmlns:wp14=http://schemas.microsoft.com/office/word/2010/wordprocessingDrawing
xmlns:wp=http://schemas.openxmlformats.org/drawingml/2006/wordprocessingDrawing
xmlns:w10="urn:schemas-microsoft-com:office:word"
xmlns:w=http://schemas.openxmlformats.org/wordprocessingml/2006/main
xmlns:w14=http://schemas.microsoft.com/office/word/2010/wordml
xmlns:wpg=http://schemas.microsoft.com/office/word/2010/wordprocessingGroup
xmlns:wpi=http://schemas.microsoft.com/office/word/2010/wordprocessingInk
xmlns:wne=http://schemas.microsoft.com/office/word/2006/wordml
xmlns:wps=http://schemas.microsoft.com/office/word/2010/wordprocessingShape
mc:Ignorable="w14wp14">
  <w:body>
    <w:p w:rsidR="00007A72" w:rsidRDefault="003C523F">
      <w:r>
        <w:t>Hello World</w:t>
      </w:r>
      <w:bookmarkStart w:id="0" w:name="_GoBack"/>
      <w:bookmarkEnd w:id="0"/>
    </w:p>
    <w:sectPr w:rsidR="00007A72" w:rsidSect="006575D2">
      <w:pgSz w:w="11900" w:h="16840"/>
      <w:pgMar w:top="1440" w:right="1800" w:bottom="1440"
w:left="1800" w:header="708" w:footer="708"
w:gutter="0"/>
      <w:cols w:space="708"/>
      <w:docGrid w:linePitch="360"/>
    </w:sectPr>
  </w:body>
</w:document>
```