

STRUCTURAREA DATELOR DE PE WEB FOLOSIND XML

As. Ing. Drd. Gabriel Dacian CUREA

As. Ing. Drd. Daniel MORARIU

**UNIVERSITATEA “LUCIAN BLAGA” SIBIU,
FACULTATEA DE INGINERIE “HERMANN OBERTH”,
CATEDRA CALCULATOARE ȘI AUTOMATIZĂRI**

Abstract

Standardele deschise și concentrarea spre comunicare și colaborare între oameni și aplicații au creat mediul necesar trecerii web-ului de la o orientare pe document la o orientare pe date. Structurarea informației cu ajutorul marcajelor generale puse la dispoziție de XML au condus la dezvoltarea serviciilor web, care au devenit platforma pentru integrarea aplicațiilor, și a web-ului semantic, care încearcă organizarea datelor web, atribuindu-le un înțeles, făcând astfel posibilă procesarea acestora.

Introducere

Nevoia schimbului de informație și a accesului la informație a dus la apariția Internetului ca rezultat al unui proiect de cercetare al Departamentului Apărării Statelor Unite, realizându-se astfel suportul fizic de transmitere al informației între sistemele de calcul. Organizarea și structurarea informației folosind sistemele de calcul au dus la apariția tehnologiilor

hypermedia. Cercetătorii de la Organizația Europeană pentru Cercetări Nucleare (Couseil Européan pour la Recherche Nucléaire - CERN) au fost cei care au combinat tehnologiile hypermedia și Internetul, având ca rezultat World Wide Web (WWW).

Conceptul de marcaj generalizat, apărut din nevoia structurării informației, implică utilizarea etichetelor pentru identificarea porțiunilor de informație. Marcajul generalizat cere ca informația să fie cuprinsă între etichete de început și etichete de sfârșit, acest lucru permițând imbricarea informației și structurarea sa într-o manieră ierarhică.

Următorul exemplu utilizează marcaje pentru a indica faptul că o carte are un titlu, câțiva autori și o editură:

```
<carte>
  <titlu>Structuri de date și algoritmi în Java</titlu>
  <autori>
    <autor>Mitchell WAITE</autor>
    <autor>Robert LAFORE</autor>
  </autori>
  <editură>Teora</editură>
</carte>
```

Utilizarea marcajului pentru reprezentarea informației, face ca aceasta să poată fi ușor de procesat software, deoarece etichetele de început și de sfârșit delimitează clar locul unde bucățile de informație încep și unde se termină. Reprezentarea informației prin marcaje face ca aceasta să fie extensibilă.

SGML (Standard Generalized Markup Language – Limbajul de Marcare Generalizat Standard) este prima încercare de definire a marcajului generalizat; acesta însă a produs o specificație foarte complexă. SGML este un metalimbaj deoarece el definește modul în care orice limbaj de marcare dat poate fi specificat formal.

HyperText Markup Language (HTML)

HTML (Limbaaj de marcare hypertext) este o aplicație SGML, fiind totodată limbajul de marcare care domină Web-ul. HTML definește structura și modul de așezare al informației într-o pagină web (document HTML), prin intermediul marcajelor și al atributelor acestora. Specificația HTML este deținută de World Wide Web Consortium (W3C). [1]

Hypertext-ul este un tip special de bază de date, în care obiectele (texte, imagini, sunete etc.) pot fi legate unele de altele. Hypertext-ul permite deplasarea de la un obiect la altul, chiar dacă aceste obiecte sunt total diferite (deplasare de la un text la o imagine etc.).

Structura unui document HTML este următoarea:

```
<html>
  <head>
    <title> Document HTML</title>
  </head>
  <body>
    Pagina web.
  </body>
</html>
```

Marcajele folosite de HTML sunt grupate în câteva categorii diferite: etichete de structurare text, etichete de formatare, etichete de legare și includere și etichete de introducere date.

Principalele dezavantaje ale HTML-ului sunt: setul de marcaje nu poate fi extins cu marcaje ce ar putea fi folositoare programatorilor; lipsa de semantică a informației.

Extensible Markup Language (XML)

XML (Limbaaj de marcare extensibil) este creația celor de la W3C și a fost creat din nevoia de a simplifica SGML și nevoia controlării evoluției HTML.

XML este similar cu SGML prin aceea că păstrează noțiunea de marcaj generalizat, totuși specificația acestuia este mult mai simplă. Avantajul major este acela că utilizatorul își poate defini propriile marcaje.

Tehnologiile XML au două arii largi de aplicație: aplicațiile centrate pe document și aplicațiile centrate pe date.

În aplicațiile centrate pe document XML-ul este folosit ca un mecanism pentru reprezentarea documentelor semistructurate (manuale tehnice, cataloage de produse etc.). În continuare este prezentat un exemplu de document semistructurat folosind marcaje XML:

```
<H1>Cerințe pentru participarea la cursul de microelectronica</H1>
<P>Pentru participarea la cursul de microelectronică trebuie:</P>
<LISTA>
<ITEM>să cunoașteți limba engleză</ITEM>
<ITEM>să fiți liberi sâmbăta și duminica</ITEM>
<ITEM>să fi absolvit cursul de Bazele electronicii</ITEM>
</LISTA>
<P>Dacă îndepliniți condițiile de mai sus puteți continua cu
<LINK HREF="Program.xml">"Programul cursului"</LINK>
</P>
```

În cazul aplicațiilor centrate pe date XML-ul este utilizat pentru a marca informație "înalt structurată" (reprezentarea textuală a datelor relaționale din bazele de date, structurile de date ale limbajelor de programare etc.). În aceste cazuri documentele XML sunt generate de către aplicații și sunt destinate prelucrării de către aplicații. În continuare prezentăm un exemplu (un ordin de cumpărare) de document XML centrat pe date:

```
<ord_cump id="2398" transmis="2003/10/07">
  <adresa_livrare>
    <societate>S.C. PARTICULARUL S.R.L.</societate>
    <strada>Eroilor</strada>
```

```

<nr>7</nr>
<localitate>Mediaș</localitate>
<județ>Sibiu</județ>
<cod_poștal>56890</cod_poștal>
</adresa_livrare>
<ordin>
  <element cod="P4190" cantitate="2">
    <descriere>Ștampilă rotundă tip P4190</descriere>
  <element cod="R1133" cantitate="2">
    <descriere>Ștampilă dreptunghiulară tip R1133</descriere>
  </ordin>
</ord_cump>

```

În acest exemplu se poate vedea că marcajul este utilizat pentru descrierea unei unități de informație, în loc de modul în care trebuie prezentată aceasta unei persoane. Exemplu anterior de utilizare centrată pe date a XML ar putea fi reprezentat, într-un limbaj obiectual, ca o structură de date astfel:

```

Class ord_cump
{
  int id;
  Date transmis;
  Adresa adresa_livrare;
  Element ordin[];
}

```

La început Web-ul s-a dezvoltat foarte rapid ca un mediu al documentelor utile pentru oameni, în detrimentul datelor și informației care pot fi procesate automat. XML-ul stă la baza a două tehnologii folosite pentru prelucrarea automată a informației de pe Web: Web-ul semantic (Semantic Web) și servicii web (Web Services). [2]

Web-ul semantic încearcă să organizeze datele și informațiile web într-un mod natural pentru oameni și în același timp facil pentru procesarea automată. Pentru dezvoltarea web-ului semantic sunt folosite două elemente: XML-ul și RDF-ul (Resource Description Framework – cadru de descriere de resurse).

XML permite utilizatorilor să adauge o structură arbitrară documentelor , dar nu spune nimic despre semnificația structurii. RDF codifică semnificații în seturi de triplete, fiecare triplet conținând un subiect, un verb și un obiect al unei secvențe elementare. RDF este un cadru pentru descrierea și interschimbarea metadatelor (informații despre informații). [1]

RDF este construit pe următoarele reguli: o resursă poate fi orice structură care poate avea un URI (Universal Resource Identifier – Identificator de resurse universal); o proprietate este o resursă care are un nume și o declarație care este o combinație între o resursă, o proprietate și o valoare.

Două baze de date pot folosi identificatori diferiți pentru ceva ce reprezintă de fapt același concept. Un program trebuie să aibă o cale să descopere astfel de înțelesuri comune, indiferent de bazele de date întâlnite. O ontologie este un document sau un fișier care în mod formal definește relațiile între termeni. Cel mai tipic fel de ontologie pentru web are o taxonomie și un set de reguli logice de inferență. Taxonomia definește clasele de obiecte și relațiile dintre ele.

Web-ul semantic este foarte flexibil, două programe pot ajunge să-și pună în comun înțelesurile prin schimbarea de ontologii, care furnizează vocabularul necesar pentru discuție.

Serviciile web sunt construite pentru realizarea procesării distribuite pe Internet. Serviciile web XML prezintă o funcționalitate folosită de utilizatorii de web prin intermediul unui protocol web standard (în cele mai multe cazuri SOAP – Simple Object Access Protocol – Protocol simplu de acces la obiecte). SOAP este un protocol simplu folosit pentru schimbul de informații într-un mediu distribuit. Are la bază XML-ul și conține trei părți: un pachet care

definește un cadru pentru descrierea a ceea ce se află în mesaj, și cum se procesează; un set de reguli de codare pentru exprimarea instanțelor de tipuri de date definite de aplicație; o convenție pentru reprezentarea răspunsurilor și apelurilor de procedură la distanță.

Serviciile web XML pune la dispoziție o cale de descriere a interfețelor sale în detalii suficiente pentru a permite utilizatorului să creeze o aplicație client care să interacționeze cu el (în mod reactiv, chiar proactiv). Această descriere este pusă la dispoziție în mod uzual într-un document WSDL (Web Services Description Language – Limbaj de descriere a serviciilor web).

Serviciile web sunt înregistrate astfel încât utilizatorii potențiali să le poată găsi ușor. Această înregistrare este făcută cu ajutorul UDDI (Universal Discovery Description and Integration – Descoperire, descriere și integrare universale).

CONCLUZII

Structurarea informației de pe Web, folosind marcasele generale puse la dispoziție de XML, deschide drumul procesării ei pentru realizarea sistemelor distribuite și pentru extragerea de cunoștințe.

Întreaga putere a Web-ului semantic va fi pusă în valoare atunci când oamenii vor crea programe care vor colecta conținutul web-ului din diverse surse, vor procesa informația și vor face schimburi de rezultate cu alte programe.

Aplicațiile vor putea fi construite folosind servicii web XML multiple din surse variate, care vor lucra împreună indiferent de locul unde se află și cum sunt implementate.

În continuare autorii își propun studierea căutării informației pe web și extragerea informației utile din paginile web.

MULȚUMIRI

Acest articol a fost elaborat în cadrul programului de doctorat urmat de autori la Universitatea “Lucian Blaga” din Sibiu (conducător științific prof. univ. dr. ing. VINȚAN Lucian).

Autorii își exprimă întreaga grațitudine companiei SIEMENS AG, CT IC, München, Germania în special domnului vicepreședinte Dr. H.C. mat. Hartmut Raffler pentru sponsorizarea programului de doctorat precum și pentru sugestiile profesionale extrem de utile.

BIBLIOGRAFIE

- [1] **Thuraisingham B. – “XML Databases and the Semantic Web”, CRC Press, 2002**
- [2] **Graham S., Simeonov S., Boubez T., Davis D., Daniels G., Nakamura Y., Neyama R. – “Building Web Services with Java™ – Making sense of XML, SOAP, WSDL and UDDI”, SAMS Publishing, 2002**
- [3] **Fensel D., Hendler J., Liberman H., Wahlster W. – “Spinning the Semantic Web – Bringing the World Wide Web to its Full Potential”, MIT Press, 2003**
- [4] **Cagle K., Dix C., Hunter D., Kovack R., Pinnock J., Rafter J. – “Beginning XML 2nd Edition”, Wrox Press Ltd., 2001**