

## UNELE ASPECTE PRIVIND CODIFICAREA ADNOTĂRII

*N.Dulgheru*

*Universitatea Tehnică a Moldovei*

### 1. FORMALISME DE ADNOTARE PENTRU PROCESAREA LIMBAJULUI NATURAL

Adnotarea corpului de date cere alegerea unui format pentru reprezentarea textului și a adnotărilor sale în formă electronică. Formatul dat trebuie să permită folosirea și re folosirea corpului de date adnotat de programele soft disponibile în diferite locuri de cercetare. Corpul de date și documentele de adnotare de asemenea trebuie să fie codate ca să facă accesul datelor cât mai ușor și flexibil. Limbajul extins de marcări (Extended Markup Language, XML) asigură un cadru de codificări standarde pentru adnotări, care satisface aceste necesități. Folosind XML ca bază, Standardul XML de Codificare a Corpului de Date (XML Corpus Encoding Standard, XCES) a fost realizat ca parte a proiectului EAGLES, cu scopul de a oferi un cadru pentru codificarea și organizarea corpului de date și a adnotărilor sale într-un format standard, flexibil și reutilizabil.

Acest raport aduce la cunoștința studentului unele standarde de suport și înrudite cu XML, dezvoltate în cadrul XML, cum sunt XSLT (Limbajul de Transformare XML/ XML Transformation Language) și RDF (Cadrul de Definiții Resurse/ Resource Definition Framework). Aceste implementări oferă împreună atât mijloace de codificare a datelor și a adnotărilor lor, folosind arhitectura documentului XCES, cât și posibilitatea de a manipula și accesa aceste date. În plus, aici, se schițează problemele și preocupările pentru reprezentarea datelor adnotate și se face o privire generală asupra unui model abstract al corpului de date și asupra adnotărilor lor.

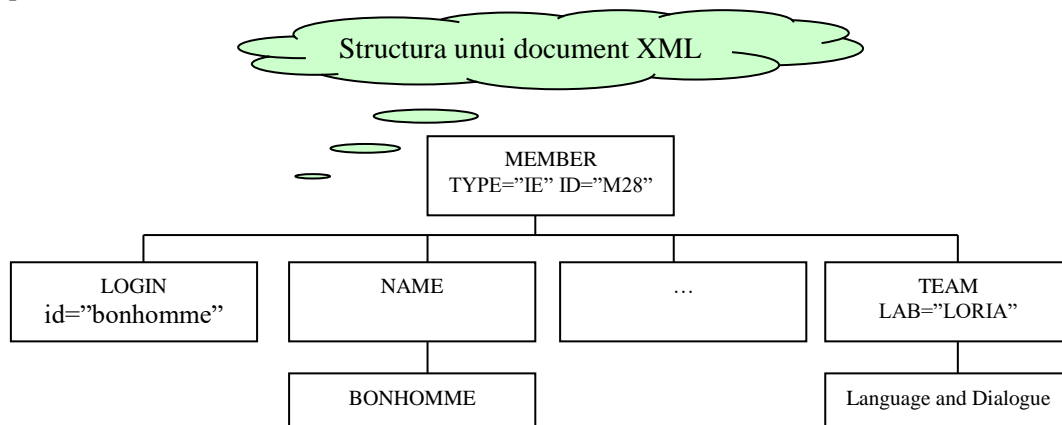
Ce este XML? XML, o abreviere pentru Extended Markup Language (Limbaj Extins de Marcări), este un metalimbaj care permite crearea unui limbaj propriu de marcări. XML prezintă o simplificare a standardului SGML, ultimul fiind destinat reprezentării structurii “logice” a unui document, iar HTML fiind conceput ca o aplicație a lui SGML.

Un document XML este un document SGML cu unele diferențe mici (dar esențiale). XML este tot atât de expresiv, însă nu atât de complex în comparație cu SGML, și face posibilă transmiterea în rețea a documentelor structurate și a bazelor de date.

Cu ajutorul XML putem modela date, publica date structurate în rețea, separa structura logică a unui document de prezentarea sa reală, integra date din surse eterogene. Nu putem evita folosirea XML din cauza simplității sale, care face posibilă integrarea acestui limbaj în orice fel de aplicații. Plus la aceasta, o varietate largă de aplicații este deja implementată în industrie (publicare, baze de date, catalogare, businessul electronic etc.), știință și cercetări (astronomie, matematică etc.), iar la dispoziție avem multe programe soft: editoare, analizatoare sintactice (parsers), punți de la și spre mediul de editare existent.

XML are următoarele proprietăți: accentul trebuie pus pe “semantica” unui document, modelul de bază are structura unui arbore, posibilitatea de a imagina un limbaj-schiță pentru a accesa orice parte a unui document XML, XML suportă codificările de caractere Unicode.

Un document XML reprezintă o structură ierarhică (fig. 1):



**Figura 1**

## 2. CODIFICAREA ADNOTĂRII CU XML

### Ce prezintă un corp de date adnotat?

- Niște date primare, inițiale...
  - un text, un semnal de vorbire etc.
- ...asociate cu o însemnare (notă explicativă) despre proprietățile sale lingvistice;
  - categoria morfo-sintactică pentru fiecare cuvânt, categoriile sintactice și structura, structura discursului, coreferințe etc.

### Unde sunt adnotările?

- Deseori, în același document /fișier cu date primare

#### *Penn Treebank*

```
((S ((NP-SBJ-1 Jones)
      (VP followed)
      (NP him)
      (PP-DIR into)
      (NP the front room))
      ,
      (S-ADV (NP-SBJ* -1)
            (VP closing)
            (NP the door)
            (PP behind)
            (NP him))))))
```

```
<TEXT>
<p>
<s>
  <lex pos=DT>The</lex>
  <lex pos=NNP>Federal</lex>
  <lex pos=NNP>Aviation</lex>
  <lex pos=NNP>Administration</lex>
  <lex pos=VBD>underestimated</lex>
  <lex pos=DT>the</lex>
  <lex pos=NN>number</lex>
  <lex pos=IN>of</lex>...
```

### De ce aceasta nu este o idee prea bună?

- Este greu de avut câteva tipuri diferite de adnotări în același corp de date;
- Este greu de avut adnotări alternative de același tip;
- Este greu de menținut;
- ...și multe alte cauze.

### XCES rezolvă problema!

- XCES = Standardul XML de codare a corpului de date;
- “Însemnări aparte”(stand-off markup);
  - Adnotări în documente XML separate, legate de original (fig. 2).

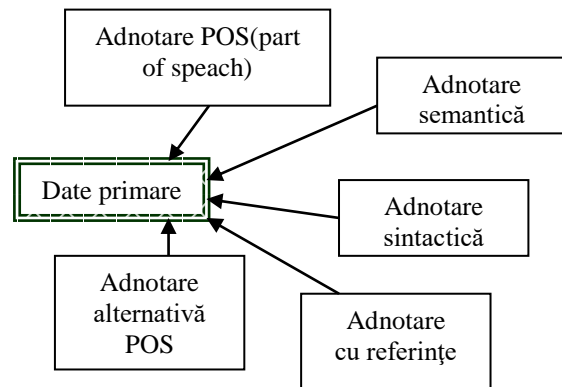


Figura 2

### Pointerii XML ușurează lucrul

- Se poate de indicat un element, caracter, șir de caractere etc. în același document sau în altul (fig. 3).

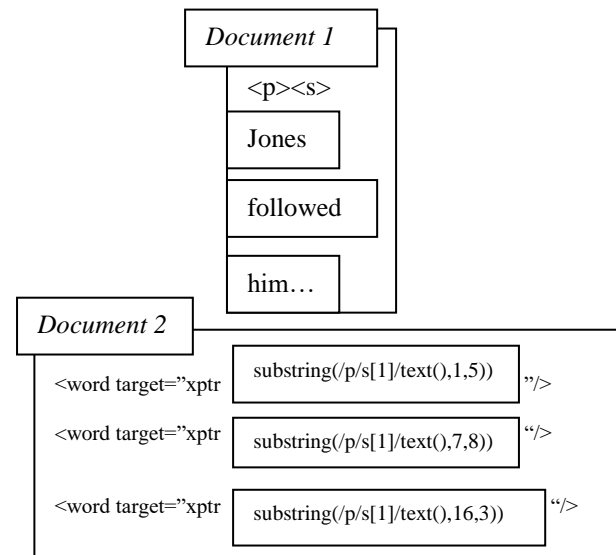


Figura 3

### Ce fel de notări se fac asupra datelor primare?

- XCES: “Trăsături sugerate de tipografie”
  - Structura logică generală (capitol, paragraf, titlu, referință etc.)
  - Propoziții?
  - Cuvinte?

Unde se începe  
“adnotarea?”

### Deci cum să alcătuiesc o schemă de codificare pentru adnotări?

- Răspuns simplu:
  - Oricum vrei, doar să urmezi un model coerent de date.

### Ce este de dorit să conțină o schemă de adnotări?

- Cuprinde caracteristicile de care ai nevoie;

- Ușor de procesat atât pentru tine, cât și pentru alții;
- Compatibilă cu alte scheme, ca să permită fuzionarea, compararea etc.;
- Permite ca adnotările să poată fi *utilizate din nou* și să fie *extensibile*.

### Pasul întâi

- Găsește *categoriile lingvistice și trăsăturile* pe care dorești să le reprezinți, fără să iei în considerare problema *formatului fizic*

- E.g.

❖ Parte de vorbire: cât de detaliată este informația despre categorii? Care sunt lemele specificate?

❖ Sintaxa: care sunt categoriile (NP,VP)? Este funcția sintactică (subiect, complement) sau informația tematică (agent, pacient) inclusă?

### Pasul doi

- Dezvoltă o modalitate de a reprezenta informația;
- Preferabil în XML, dar nu numaidecât;
- Care poate să fie procesată de software-ul tău;
- Nici o informație pierdută;
- Atât timp cât informația este reprezentată, nu este nevoie să te îngrijezi de formatul specific XML.

### Exemplu (Parte de vorbire)

```
<w><orth>dogs</orth>
  <cat>NNP</cat>
  <lemma>dog</lemma>
</w>
```

```
<seg target="xptr(substring(/p/s[1]/text(),1,5))"/>
  lex="MyCats#DX51"/>
```

```
<w cat="NNP" lemma="dog">dogs</w>
```

Figura 4

### De ce nu este important formatul exact (modul de aranjare și prezentare)?

- Pentru că există XSLT!
- Este ușor de a trece de la un format la altul (...Atât timp cât informația există);
- Adnotatorii pot să continue să lucreze după formatele lor interne;
- Trebuie să fie construit după un model de date comun;

### De ce avem nevoie?

- De un model comun, abstract care poate exprima (reprezenta) informația adnotării, ne luând în considerare codarea fizică;
- De o exemplificare concretă generică XML a modelului, spre care și din care pot fi traduse formate concrete folosind XSLT;
- De un set comun de categorii de date, la care adnotatorii să facă referință și de care să se folosească.

### Deci cum lucrează aceasta? (fig. 5)

#### Registrul Categoriilor de Date (RCD)

- Inventarul categoriilor de date;
- Specificate de colectivul de cercetători;
- Poate fi în continuu extinsă;
- Menținut pentru referire și folosire universală;
- Nu este un standard, ci un punct de referință;
- Arată cum pot fi găsite subcategoriile specifice în categoriile din RCD:

#### Instanțierea Registrului Categoriilor de Date

- Este definit folosind Cadrul de Definiții Resurse (Resource Definition Framework, RDF);
- Realizat de comunitatea XML (Consortiul W3);
- Schemele RDF specifică semantica datelor bazându-se pe XML (obiecte și relațiile lor, ierarhii ale tipurilor derivate);
- Descrierile RDF dau exemple de obiecte.

#### Scheletul structural

- Un model general, fundamental care informează practica curentă;
- Un Instrument de trasare generică (Generic Mapping Tool, GMT) furnizează realizarea XML a scheletului structural;
- *Lingua franca* folosită să compare și să unească adnotările, să facă posibil designul instrumentelor generice pentru vizualizare, editare, extragere etc.

#### Specificarea categoriei de date

- Set de categorii de date într-o schemă de adnotare dată;
- Subclasă de categorii din RCD + (opțional) categorii specifice în aplicație definite suplimentar;
- Descrie constrângerile asupra categoriilor de date;
- Restricții asupra valorilor;

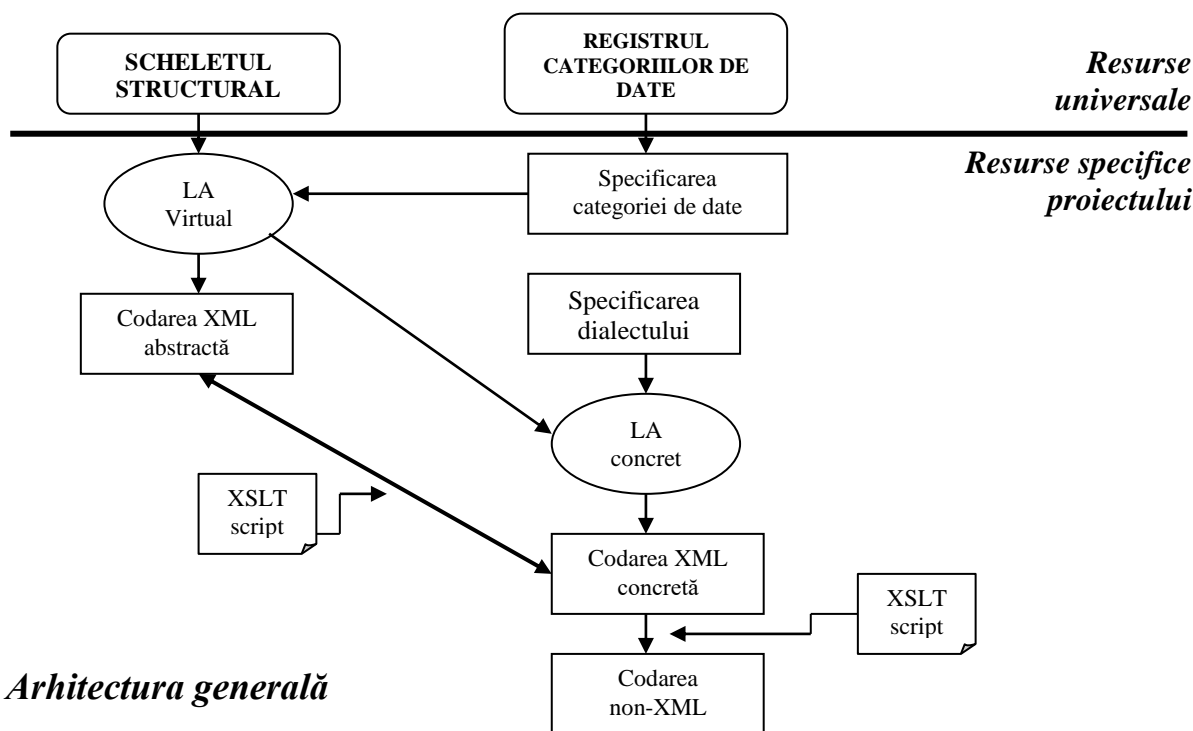


Figura 5

- Restricții asupra locului unde pot apărea categoriile de date (de ex., nivelul într-o ierarhie structurală).

### Specificarea dialectului

- Definește formatul XML specific proiectului;
- Folosește scheme XML, XSLT scripts, foi de stil XSL;

- Include;

- *Stiluri de instanțiere a categoriilor de*

*date*

```
<NounPhrase>
<cat type="NounPhrase">
```

- *Stiluri ale vocabularului de categorii de*

*date*

```
<cat type="NounPhrase">
<cat type="NP">
<cat type="SN">
```

- *Structuri de lărgire (expansiune)*

- ❖ Modifică structura exprimată folosind scheletul structural;

- ❖ Ex. Creează două subnoduri sub un nod dat pentru a grupa diferite tipuri de informație.

### Limbajul de Adnotare (LA) / (Annotation Markup Language, AML)

- Scheletul structural + SCD(Specificarea Categoriilor de Date)=**LA Virtual**;

- LA Virtual + Specificarea Dialectului folosită pentru a genera în mod automat **LA concret**

- Se conformează formatului specific proiectului în Specificarea Dialectului;

- Filtrele XSLT traduc

- ❖ Reprezentările adnotării în LA concret și în Instrumentul de Trasare Generică (ITG) / (GMT, Generic Mapping Tool);

- ❖ Din LA concret în formate non-XML (ex. Penn Treebank).

### Bibliografie

1. A.Mengel, W.Lezius. *An XML –based representation format for syntactically annotated corpora/ Proceeding of the 2<sup>nd</sup> Conference on Language Resources and Evaluation LREC, Athens, Greece, 2000.*

2. XCES.2000. *Corpus Encoding Standard for XML.2000.* Vassar College. LORIA/CNRS (<http://www.cs.vassar.edu/XCES/>).

3. W.3C.2000.XSL Transformation (XSLT). *Version 1.0* (<http://www.w3c.org/TR/xlst>).