

Înregistrările de autoritate ca date interconectate: un experiment autohton

DAN MATEI

Institutul Național al Patrimoniului
București

Abstract

Authority Records as Linked Data: A Local Experiment

The phrase „linked data” was set in 2006 by Tim Berners-Lee, the „inventor” of the web. It designates a (relatively) new paradigm, which implies the association of web assertions in the manner in which the pages on/ between the sites are connected. Their reason is what is known as „semantic web,” that is semantic associations between the entities on the web to allow the software agents (e.g. search engines) to make logical inferences. And the idea is that these linked data should (also) be done between assertions coming from different sources.

Keywords: *linked data, web, assertions, index, associations*

1. Rememorare: datele interconectate

Sintagma „date interconectate” [*linked data*] a fost propusă în 2006 de Tim Berners-Lee, „inventatorul” webului.¹ Sintagma designează o paradigmă (relativ) nouă, care presupune asocierea de aserțiuni pe web, în maniera în care sunt conectate paginile pe / între sauturi. Rațiunea lor este ceea ce se cheamă „webul semantic,” adică asocieri semantice între entități pe web,² care să permită agenților soft (e.g. motoarelor de căutare) să facă inferențe

¹ A se vedea celebrele-i principii la: www.w3.org/DesignIssues/LinkedData.html

² Sau, cum zice sloganul Google Knowledge Graph: „things, not strings” [lucruri, nu șiruri (de caractere)]. <http://googleblog.blogspot.com/2012/05/introducing-knowledge-graph-things-not.html>

logice. Și ideea e ca aceste interconectări să se facă (și) între aserțiuni provenite din surse diferite.

În Figura 1 se ilustrează (foarte simplificat) cum se pot interconecta aserțiuni provenind din surse diferite (sugerate prin nuanțe diferite). Adică, cineva aserțiază că **Război și pace** (o lucrare, în terminologia FRBR [Functional Requirements for Bibliographic Records]³) are drept creator pe Tolstoi. Altcineva adaugă aserțiuni despre o expresie a lucrării în limba engleză, iar altcineva aserțiuni despre o expresie în limba română. În fine, din alte surse provin antroponimele lui Tolstoi în engleză, respectiv în rusă.

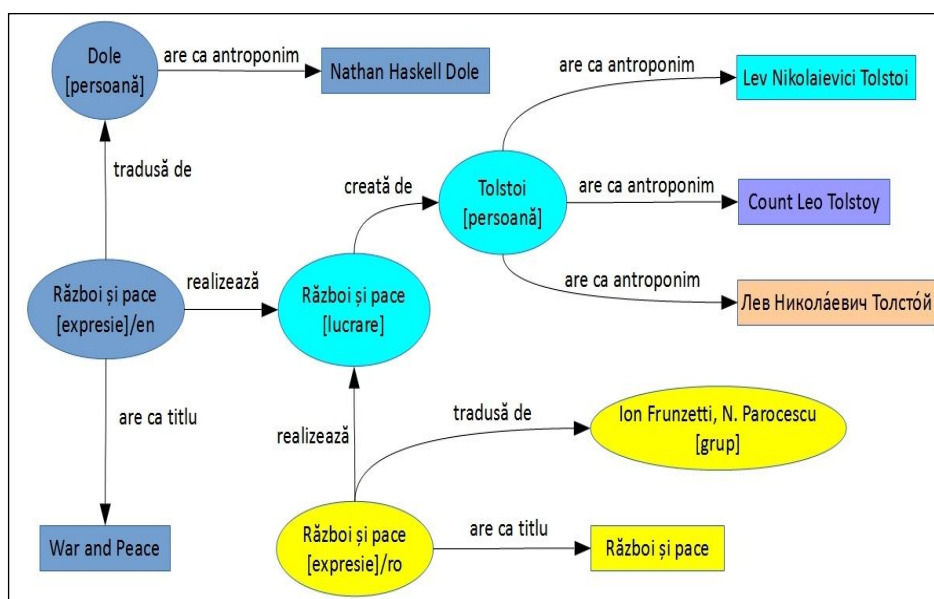


Fig. 1. Exemplu de date bibliografice interconectate

Esențial în această manieră de a trata datele este ca entitățile invocate în aserțiuni să fie referite prin identificatori „rezolvabili” (resolvable URI

³ www.ifla.org/publications/functional-requirements-for-bibliographic-records

[Uniform Resource Identifiers]⁴) pe web, adică care să identifice neambiguu resurse pe web.

Modelul conceptual „clasic” pentru datele interconectate este RDF [Resource Description Framework]⁵ care, în esență, definește tripletele subiect-predicat-obiect. Din punct de vedere practic, bazele de date ce implementează modelul RDF⁶ au și avantajul că au tabele (abstracte și) puține și permit tratarea unitară a claselor și proprietăților. Adică, e posibilă adăugarea succesivă de noi clase și proprietăți (care pot fi rafinări sau abstractizări ale celor deja existente), cu alte cuvinte se pot aduce modificări taxonomiei subiacente, fără a se modifica structura bazei de date. Așadar administratorul bazei de date poate face asta fără a apela la programatori. De pildă, dacă avem clasa „organizație,” oricând se poate adăuga o subclasă a acesteia, „persoană juridică.” Similar, dacă avem proprietatea „are drept contributor pe,” se poate adăuga o subproprietate a acesteia, „are drept traducător pe.” Dezavantajul acestui gen de baze de date pare a fi complexitatea sporită pe care o impune interogărilor.⁷ Deci, este sarcina proiectanților interfețelor-utilizator să „ascundă” această complexitate sub prezentări simple și intuitive.

În exemplul grafic din Figura 1, s-ar putea extrage trei entități (lucrarea, expresia și autorul) și să se exprime în formalismul RDF câte trei proprietăți ale fiecăreia:

⁴ „A URI whose resource has one or more representations available via invoking HTTP GET on the URI”: www.w3.org/TR/2010/WD-sparql11-http-rdf-update-20100126/

⁵ www.w3.org/TR/2004/REC-rdf-primer-20040210

⁶ „Triplestore”: www.answers.com/topic/triplestore

⁷ www.museumsandtheweb.com/mw2012/papers/a_new_framework_for_querying_semantic_networks

subiect	predicat	obiect	limba
id-1	e de tip	expresie	
id-1	are ca titlu	Război și pace	ro
id-1	realizează	id-2	
id-2	e de tip	lucrare	
id-2	are ca titlu	Război și pace	ro
id-2	creată de	id-3	
id-3	e de tip	persoană	
id-3	are ca antroponim	Lev Nikolaievici Tolstoi	ro
id-3	are ca antroponim	Count Leo Tolstoy	en

Deja instituții importante - cum ar fi British Library⁸ (Figura 2: 2,6 milioane de înregistrări, care au generat 84.961.180 triplete) sau British Museum⁹ - își oferă informațiile catalografice sub formă de date interconectate deschise (adică atât gratuite, cât și liber reutilizabile !).

Figura 2. Pagina cu datele interconectate a British Library



⁸ www.bl.uk/bibliographic/datafree.html#lod

⁹ <http://collection.britishmuseum.org>

2. Ce înseamnă „date interconectate” (în context bibliografic)?

- identificarea entităților (și referirea la ele prin identificatori, nu prin etichete) și a proprietăților lor;
- deconstrucția înregistrărilor existente - fie ele MARC sau alt format - în aserțiuni elementare: subiect - predicat - obiect;
- conectarea între entități, chiar dacă provin din surse diferite.

3. Proiectul: indexul antroponimelor

Idea e să tratăm proprietățile asociate unei persoane ca date interconectate și să expunem pe web indexul antroponimelor care apar în bazele de date ale Institutului Național al Patrimoniului întreținute de colectivul CIMEC.¹⁰ Și asta ca un prim pas spre expunerea tuturor entităților. Figura 3 ilustrează grafic entitățile evidențiate (entitatea „persoană” fiind centrală), împreună cu proprietățile lor (relațiile între ele și atributele lor¹¹).

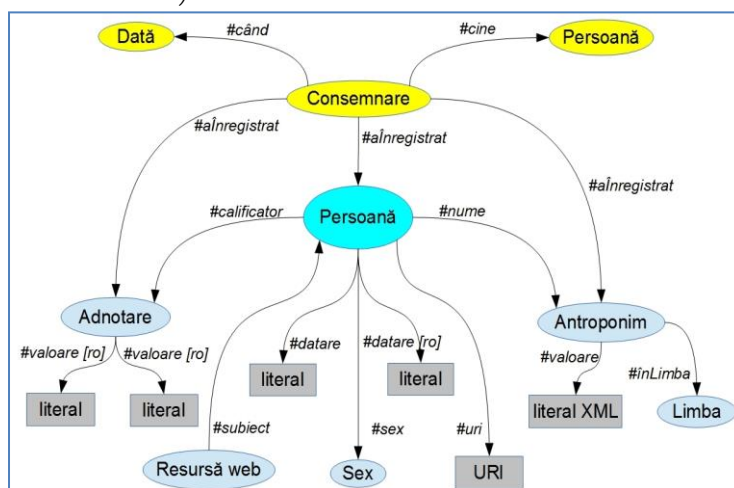


Figura 3.
Asocierile unei persoane (elipsele reprezintă entități, iar dreptunghiurile, atribute;
URI = Uniform Resource Identifier¹²).

¹⁰ <http://cimec.ro/Baze-date-online.html>

¹¹ Eticheta prefixată cu # reprezintă un identificator.

¹² http://en.wikipedia.org/wiki/Uniform_resource_identifier

Entitățile și proprietățile consemnate fac parte din ontologia noastră internă RIDE [Rich Information Describing Entities] (denumirea face aluzie la ontologia LIDO [Lightweight Information Describing Objects],¹³ care a fost dezvoltată în cadrul proiectului European Athena,¹⁴ la care CIMEC a fost partener). Ontologia RIDE este în curs de elaborare, pentru a folosi la proiectarea bibliotecii digitale naționale *culturalia.ro*. Ea este o extensie a ontologiilor EDM [Europeana Data Model]¹⁵ (dezvoltată pentru Biblioteca Digitală Europeană *europa.eu*) și CIDOC-CRM [Conceptual Reference Model]¹⁶ (model conceptual elaborat de CIDOC - Comitetul de Documentare al ICOM - Consiliul Internațional al Muzeelor și oficializat ca standardul ISO 21127:2006¹⁷).

Se observă că entitățile evidențiate sunt:

entitatea	reprezintă	identificatorul RIDE
persoană	obiectul înregistrării de autoritate	crm:E21_Person
antroponim	numele persoanei	ride:Anthroponym
adnotare	calificatorul persoanei	ride:Annotation
concept	sexul persoanei	crm:E55_Type
limba	limba antroponimului	crm:E56_Language
document	pagina web al cărei subiect este persoana	crm:E31_Document
document de referință	lucrarea de referință din care face parte pagina web	crm:E32_Authority_Document
consemnare	evenimentul redactării înregistrării	ride:Recording
persoană	catalogatorul care a redactat înregistrarea	crm:E21_Person
dată	data la care a fost redactată înregistrarea	crm:E2_Temporal_Entity

¹³ <http://network.icom.museum/cidoc/working-groups/data-harvesting-and-interchange/what-is-lido/>

¹⁴ <http://www.athenaeurope.org/>

¹⁵ <http://pro.europeana.eu/edm-documentation>

¹⁶ <http://www.cidoc-crm.org/>

¹⁷ http://www.iso.org/iso/catalogue_detail?csnumber=34424

[antroponim]		[limbă]	e
#AlbertCelMare [persoană]	#uri	http://viaf.org/viaf/...	ride:has_URI
#Bărbat [concept]	"valoare" [ro]	"bărbat"	rdfs:label
#TeologSiFilozofGerman [adnotare]	"valoare" [ro]	"teolog și filozof german"	rdfs:label
#TeologSiFilozofGerman [adnotare]	"valoare" [en]	"German theologian and philosopher"	rdfs:label
http://ro.wikipedia.org... [document]	#subiect	#AlbertCelMare [persoană]	crm:P129_is_about
#WikipediaRo [document de referință]	#componentă	http://ro.wikipedia.org... [document]	crm:P148_has_component
# {o anume consemnare} [consemnare]	#aÎnregistrat	*	ride:recorded
# {o anume consemnare} [consemnare]	#cine	# {un anume indexator} [persoană]	crm:P14-carried_out_by
# {o anume consemnare} [consemnare]	#când	# {o anume dată} [dată]	ride:has_timing

De notat că entitatea de tip „consemnare” (care e un tip specific de „eveniment,” i.e. de crm:E5_Event) se relaționează cu fiecare entitate și aserțiune. Cu alte cuvinte, se consemnează „paternitatea” fiecărei afirmații. Pe de altă parte, la o privire atentă se observă în fig. 4 că aserțiunea:

<#AlbertCelMare> <#uri> http://viaf.org/viaf/...

are ea însăși o proprietate, și anume are #tip = #VIAF. În jargonul RDF, această atribuire de proprietăți unei aserțiuni se numește reificare.

Practic, procesarea antroponimelor se desfășoară astfel: se extrag numele de persoane din bazele de date, iar apoi indexatorii analizează fiecare nume, identifică entitățile și redactează setul de proprietăți ale acestora.

Sursele datelor sunt (deocamdată doar):

- Inventarul bunurilor culturale mobile clasate;¹⁸
- Catalogul Colectiv al Cărții Vechi Românești;
- Catalogul Incunabilelor;¹⁹
- Repertoriul Teatral;²⁰
- Premiere Muzicale și Coregrafice.²¹

Procesul de catalogare/indexare revine la a completa formula (Access) exemplificată în Figura 5.

cheia	Albertus Magnus	cheie-sort-A-text-1	Albert cel Mare	ro
apelatiune	Sfântul ^1Albert cel Mare, zis ^2Albertus Magnus	cheie-sort-A-numar-1		
calificator	teolog și filozof german	cheie-sort-A-text-2		
datare	1200 - 1280	cheie-sort-A-numar-2		
femeie ?	<input type="checkbox"/>	Intrare principală ?	<input checked="" type="checkbox"/>	
URI-Wikidata	http://www.wikidata.org/wiki/Q60059	cheie-sort-B-text-1	Albertus Magnus	la
URI-VIAF	http://viaf.org/viaf/88125532/	cheie-sort-B-numar-1		
URI-Answers	http://www.answers.com/topic/albertus-magnus	cheie-sort-B-text-2		
URI-Wikipedia-Ro	http://ro.wikipedia.org/wiki/Albertus_Magnus	cheie-sort-B-numar-2		
URI-ULAN		Intrare principală ?	<input type="checkbox"/>	
editor	Matei	data	16.07.2013	
surse	http://www.calendarcatholic.ro/Sfinti/tabid/66/articoleType/ArticleView/articoleId/3450/Sf-Albert-cel-Mare-ep-inv-.aspx	note	^3Sfântul (i.e. ar merita o intrare la sfînti)	
istoricul modificărilor				
id	1873			
copie-a		problematică ?	<input type="checkbox"/>	

Figura 5. Formula de catalogare/indexare a unei persoane.

Așadar, plecând de la forma primară a numelui persoanei (așa cum se găsește ea în baza de date din care provine și vizibilă în rubrica „cheie”), catalogatorul/indexatorul:

¹⁸ <http://clasate.cimec.ro>

¹⁹ <http://cimec.ro/Carte/Catalogul-colectiv-incunabilelor-Schatz-Stoica.pdf>

²⁰ www.cimec.ro/Teatre/Star_Home.htm

²¹ www.cimec.ro/scripts/Muzica/Premiere/selPREM.asp

- elaborează designarea persoanei destinată afișării, cu cele trei componente ale ei, adică:
 - antroponimul complet,
 - un calificator care să încadreze persoana profesional și cultural,
 - anii de viață sau de activitate ai persoanei, pentru a o încadra temporal;
- identifică (și consemnează) eventualele URI-uri [Uniform Resource Identifier] ale persoanei (deocamdată) în trei liste de autoritate consacrate: Wikidata²², VIAF [Virtual International Authority File]²³ și ULAN [Union List of Artist Names];²⁴
- identifică (și consemnează) eventualele articole descriptive despre respectiva persoană, în două enciclopedii online consacrate: Wikipedia (în română)²⁵ și Answers.com;²⁶
- specifică două chei de sortare/căutare pentru apelațiune și o asociază pe fiecare cu un segment al acesteia.

O mostră din rezultatul preliminar (vizualizat cu Chrome) se vede în Figura 6.

²² <http://www.wikidata.org>

²³ <http://viaf.org/>

²⁴ <http://www.getty.edu/research/tools/vocabularies/ulan/>

²⁵ http://ro.wikipedia.org/wiki/Pagina_principal%C4%83

²⁶ <http://wiki.answers.com/Q/Special:search>

Sex	Designare	nume	calif.	dat.	Wikidata	VIAF	ULAN	Wikipedia (ro)	Answers
♂	Hans von Aachen [German painter] (1552-1615)								
♂	Petru Pavel Aaron [episcop al Bisericii Române Unite cu Roma din Transilvania, 1752-1764] (1709-1764)								
♀	Elena Ailincăi [țesătoare româncă din Moldova] (fl. 1930)								
♂	Johannes de Aingre (sec XV-XVI)								
♂	Sechel Aipotochiței [grefier român] (fl. 1838)								
♂	Aka Muhammad Khan [Persia] (1779 - 1797)								
♂	Albano [filozof și astrolog italian, profesor de medicină la Padova] (1250-1316)								
♂	Sfântul Albert cel Mare, zis Albertus Magnus [teolog și filozof german] (1200-1280)								
♂	Albert de Saxa-Coburg și Gotha [prinț consort al reginei Victoria a Marii Britanii] (1819-1861)								
♂	Albert Lebourg [pictor francez] (1849-1928)								
♂	Albert Marionnet [French sculptor] (1852-1910)								
♂	Albert Nagy [pictor român] (fl.sec. XX)								
♂	Alberti Cherubino [pictor italian] (1553-1615)								
♂	Albertus de Placentia (15th century)								
♂	Sfântul Albert cel Mare, zis Albertus Magnus [teolog și filozof german] (1200-1280)								
♂	Petrus Albigannus Trecius [editor științific] (1545)								
♂	Albius Tibullus [poet elegiac latin] (55 î.Hr.-19 î.Hr)								
♂	Paolo Antonio Alboni [pictor italian] (1665-1735)								
♂	Albrecht Dürer [pictor, grafician, teoretician al artei german] (1471-1528)								

Figura 6. Expunerea indexului lexicografic (HTML5).

Indexul expune designările în ordine alfabetică, de regulă de două ori: atât la nume, cât și la prenume. Aceasta pentru a facilita utilizatorului poziționarea pe antroponimul căutat,²⁷ cu alte cuvinte, poziționarea pe baza unui prefix al prenumelui să fie la fel de ușoară ca și poziționarea pe baza unui prefix al numelui.

Cheile de sortare (de regulă prenumele, respectiv numele) sunt aliniate²⁸ pentru a facilita înțelegerea logicii ordonării.

Iconurile asociate fiecărei designări sunt linkuri către cele trei fișiere

²⁷ Când aceasta va fi posibilă, adică atunci când se va implementa un mecanism de căutare/poziționare/paginare - sperăm că în viitorul apropiat.

²⁸ Pe baza a două segmente pe care le evidențiază catalogatorul în cadrul antroponimului.

de autoritate, respectiv cele două enciclopedii unde apare persoana respectivă.

De dragul ilustrării, se expun (când e cazul) și iconurile limbilor celor trei componente ale designării: antroponimul, calificatorul și datarea. În mod normal, limbile nu vor fi evidențiate în expunerea online a indexului, dar vor putea fi folosite la filtrări. De pildă, ca să vedem doar antroponimele în limba latină.²⁹

4. Indexul antroponimelor: comentarii

Deocamdată, expunerea lexicografică a indexului este rudimentară (fig. 6): este un fișier HTML5 static, ce conține 1.881 de antroponime (cu 3.527 de intrări distincte)³⁰ și are un volum de 3,3 MB, deci se încarcă greu, chiar și într-o rețea locală.

Patru decizii de proiectare diferențiază acest index de cele tradiționale:

- forma de afișare nemaifiind și cheie de sortare, s-a abandonat tradiția bibliografică de inversare a numelui cu prenumele, cu alte cuvinte, antroponimul este expus în ordinea firească dată de cultura sa de origine;
- se consemnează - de regulă - pe lângă tradiționala datare și un calificator care poziționează persoana într-o cultură și un domeniu;
- se specifică sexul persoanei (ceea ce va permite și filtrarea după sex);
- (unde e cazul) se consemnează limbile numelui, a calificatorului și chiar și a datării (dacă formularea ei este specifică unei culturi, e.g. cu „i.Hr.”, formulare specifică limbii române; astfel se obține un index „independent de limbă”, cu alte cuvinte multilingv).

Consemnarea trimiterii la Wikipedia are ca efect colateral interesant obținerea și a unui index alfabetic al unui subset al Wikipediei!

5. Indexul antroponimelor: planuri

Planul de termen scurt este expunerea online (reală) a acestui index,

²⁹ În lipsa unui icon consacrat pentru limba latină, am folosit drapelul Sfântului Scaun, unde latina este limbă oficială.

³⁰ Valabil la 22.11.2013.

ceea ce presupune rafinarea softului care-l generează. Pe termen mediu, planul ar fi:

- să se programeze o interfață web flexibilă de alimentare/corectare a înregistrărilor, ceea ce ar permite comunităților profesionale să contribuie, dar mai ales ar facilita (deoarece structura internă permite):

- asocierea mai multor antroponime unei persoane (cum ar mai fi „Albert, Count von Bollstädt, known as Albert the Great” asociat persoanei #AlbertCelMare, sau „Nenea Iancu” asociat persoanei #ILCaragiale);

- împărțirea unui antroponim în mai multe segmente, necesară pentru apelățiuni luxuriante, cum ar fi „Pieter Bruegel the Elder, called Peasant Bruegel, also called Peer den Drol”;

- asocierea mai multor „chei-umbră” unui segment, adică chei de căutare care să nu fie expuse, dar care să permită regăsirea cu a) prefixe „nestandard,” cum ar fi „Breughel” și „Breugel” pentru Bruegel sau „Göthe” pentru Goethe sau b) prefixe fără diacritice, cum ar fi „Brancusi” pentru Brâncuși;

- să se mai adauge proprietăți asociate unei persoane, ca de pildă datele de naștere, pentru a permite și o ordonare cronologică;

- să se consemneze și URI-ul persoanei din Freebase,³¹ pentru a se facilita absorbția ei în Google Knowledge Graph;³²

- să se trateze și identitățile bibliografice, cum ar fi „Karol Józef Wojtyła” vs. „Ioan Paul al II-lea” sau „Otilia Valeria Coman” vs. „Ana Blandiana”.

Pe termen lung, planul este mult mai ambițios, și anume dezvoltarea unui index lexicografic cuprinzător al (viitoarei) bibliotecii digitale a României (*Culturalia.ro*) și al viitorului catalog național partajat³³ - și care să joace și rolul de fișier de autoritate „deschis” -, care să includă și:

- locuri;

³¹ www.freebase.com

³² www.google.com/insidesearch/features/search/knowledge.html

³³ Idee obsesivă, la care n-am renunțat încă.

- concepte;
- perioade;
- obiecte;
- lucrări;
- expresii;
- manifestări;
- exemplare.

Pe de altă parte, ne dorim să oferim informațiile din acest index ca date interconectate deschise, în format RDF/XML (poate și JSON-LD), expunându-le (și) ca seturi de date procesabile pe situl guvernului dedicat datelor deschise *data.gov.ro*.³⁴

Anexa A.

Fragmentul din ontologia RIDE folosit la elaborarea indexului

Taxonomia claselor

rdfs:Resource

- crm:E1_CRM_Entity
 - edm:NonInformationResource
 - crm:E39_Actor
 - crm:E21_Person
 - crm:E18_Physical_Thing
 - crm:E19_Physical_Object
 - crm:E20_Biological_Object
 - crm:E21_Person
 - crm:E77_Persistent_Item
 - crm:E80_Thing
 - crm:E71_Man-Made_Thing
 - crm:E28_Conceptual_Object

³⁴ <http://data.gov.ro/>

- crm:E55_Type
 - crm:E56_Language
- crm:E90_Symbolic_Object
 - crm:E41_Appellation
 - crm:E82_Actor_Appellation
 - *ride:Anthroponym*
 - crm:E42_Identifier
 - *ride:URI*
 - crm:E89_Propositional_Object
 - crm:E73_Information_Object
 - crm:E33_Linguistic_Object
 - *ride:Annotation*
 - crm:E31_Document
 - crm:E32_Authority_Document
- crm:E2_Temporal_Entity
 - crm:E4_Period
 - crm:E5_Event
 - crm:E7_Activity
 - *ride:Recording*
- rdfs:Literal
 - rdf:XMLLiteral
 - crm:E59_Primitive_Value
 - crm:E61_Time_Primitive
 - *ride:Date_Time*

Taxonomia proprietăților

edm:isRelatedTo

- dc:subject
 - crm:P67i_is_referred_to_by
 - edm:isAnnotationOf
 - *ride:has_qualifier*
 - crm:P129_is_about

- dc:relation
 - crm:P12_occurred_in_the_presence_of
 - crm:P11_had_participant
 - crm:P14-carried_out_by
 - crm:P148_has_component
 - crm:P4_has_time-span
 - *ride:has_timing*
- crm:P2_has_type
 - *ride:has_gender*
- ride:has_identifier*
 - crm:P1_is_identified_by
 - *ride:has_URL*
 - crm:P131_is_identified_by
 - *ride:has_appellation*
 - *ride:has_anthroponym*
 - *ride:has_preferred_anthroponym*
 - *ride:has_preferred_identifier*
 - *ride:has_preferred_anthroponym*
 - *ride:has_URI*
- dc:language
 - *ride:has_language*
 - crm:P72_has_language
- dc:description
 - rdf:value
 - crm:P3_has_note
 - dc:date
 - rdfs:label
 - *ride:Recording*

La dezvoltarea acestei ontologii (prin extinderea ontologiei CIDOC-

CRM), am folosit intens maparea lui Martin Doerr.³⁵ Adaosurile specifice nouă sunt evidențiate prin prefixul „ride.”

Anexa B.

Exemplu de balizare XML a unui antroponim

În acest exemplu se poate observa cum se asociază segmentele formei de afișare cu cheile de sortare/indexare.

```
<rml:anthroponym xmlns="http://ride.culturalia.ro"
  xmlns:rml="http://ride.culturalia.ro" ... rml:language="ro">
  <rml:displayForm>
    Sfântul
    <rml:segment rml:id="1">Albert</rml:segment>
    cel Mare, zis
    <rml:segment rml:id="2">Albertus</rml:segment>
    Magnus
  </rml:displayForm>
  <rml:indexEntry rml:relatedSegmentId="1" rml:isMainEntry="true"
rml:language="ro">    <rml:textKey-1>ALBERT CEL
MARE</rml:textKey-1>
  </rml:indexEntry>
  <rml:indexEntry rml:relatedSegmentId="2" rml:language="la">
    <rml:textKey-1>ALBERTUS MAGNUS</rml:textKey-1>
  </rml:indexEntry>
</rml:anthroponym>
```

NB. Fiecărei intrări de index i se poate asocia o limbă, iar limba intrării principale dă și limba formei de afișare.

³⁵ http://www.cidoc-crm.org/docs/EDM-DC-ORE-CRM-FRBR_Integration_ORE_fix.ppt