

# Gestiunea datelor personale bazată pe microformate

Marius-Gabriel Butuc

Sabin-Corneliu Buraga

Facultatea de Informatică, Universitatea „A.I.Cuza” din Iași

Str. Berthelot, nr. 16, Iași 700483

{mbutuc, busaco}@infoiasi.ro

## REZUMAT

În stadiul actual al Web-ului resimțim o adevărată înflorire a laturii sociale prin posibilitatea comunicării dintre tot mai multe entități (aparent) separate între ele. În scopul înlesnirii unei comunicări facile și a înțelegerii între entități, avem nevoie de standarde care să descrie ce spunem, ce ascultăm și ce stocăm în datele noastre. Astfel, dezvoltarea marcajelor semantice nu vizează doar îmbogățirea înțeleșurilor pentru oameni, dar și la înțelegerea de către mașină a semanticilor definite uman. Scopul acestui proiect este acela de a descrie o aplicație Web – Microwler – care să fie capabilă să ajute utilizatorii să își îmbogățească experiența socială pe Web. Microwler oferă utilizatorului posibilitatea să caute și să obțină datele personale de contact ale unei persoane ce le-a făcut disponibile, indiferent de restul conținutului paginilor. Utilizatorul are și posibilitatea de a vizualiza relațiile pe care acea persoană le are în rețeaua sa de prieteni XHTML (XFN).

## Cuvinte cheie

Microformate, regăsire de informații, gestiune date personale, interacțiune Web.

## Clasificare ACM

H5.4. Hypertext/Hypermedia: Navigation, User Issues.

## PREAMBUL

Odată cu creșterea constantă a popularității Web-ului, barierele între autor și cititor se estompează din ce în ce mai mult. Nivelul responsabilității și al încrederii crește permițând cititorului să devină chiar autor și să contribuie la cantitatea comună de cunoștințe. Tindem să devenim din ce în ce mai sociabili *on-line*, să ne publicăm datele de contact într-un mod care să poată fi ușor accesat de către oricine, să creăm grupuri, rețele de prieteni *on-line*, să împărtășim informații despre evenimentele la care participăm etc. Astfel, au apărut diverse modele de structurare informațională a resurselor publice și, de asemenea, în concordanță cu așteptările utilizatorilor: există numeroase situri focalizate pe *blog*-uri, grupuri sociale, semne de carte colaborative, cunoștințe colaborative, prezentări și recenzii de produse sau companii, portaluri de știri – toate aliniate Web-ului social (a se consulta [2]).

Web-ul Semantic se regăsește într-o etapă în care din ce în ce mai multe informații devin disponibile *on-line* și în care avem la dispoziție șabloane de proiectare ce oferă semantică acestor informații. În acest context, grupul de inițiativă focalizat asupra microformatelor [7] propune specificarea unor metode de lucru cu aceste șabloane, adică instituirea unor standarde de publicare a informațiilor incluzând semantici repetitive – apar premisele marcării semantice a datelor, conducând Web-ul

către un nou stadiu în care informația este accesibilă în egală măsură oamenilor și mașinilor.

Cu o istorie ce a început odată cu articolul lui Vanevar Bush [3] din 1945, regăsirea de informații – *information retrieval* (IR) – a devenit un domeniu vast. Un proces de recuperare de informații începe în momentul în care un utilizator introduce o interogare în sistem. Un motor de căutare este un exemplu de sistem de IR, scurtând atât timpul necesar găsirii informației, cât și cantitatea de date ce trebuie parcursă. Cele mai vizibile aplicații de tip IR sunt motoarele de căutare Web – ca Google, Yahoo Search sau Live Search – care caută informații în World Wide Web (WWW).

Microwler reprezintă un motor de căutare Web specializat pe microformate ce oferă utilizatorului posibilitatea de a recupera informații de contact din pagini Web îmbogățite semantic prin hCard-uri și îi facilitează vizualizarea relațiilor între entități marcate cu XFN.

După cum notează Toby Segaran [5], algoritmi pentru căutarea textelor sunt printre cei mai importanți algoritmi de inteligență colaborativă și ideile inovatoare în acest domeniu au un succes imediat. Se consideră că rapida metamorfoză a Google dintr-un proiect academic în cel mai popular motor de căutare a avut loc grație algoritmului PageRank. Microwler a fost construit pentru parcurgerea informațiilor, indexare și căutare în datele salvate, păstrând clasificarea în funcție de rang pentru o dezvoltare viitoare.

Lucrarea de față este structurată după cum urmează: mai întâi, în secțiunea 2, vom defini noțiunea de microformat și vom descrie posibilele sale utilizări în contextul vizat. Apoi, în cea de-a treia secțiune vom prezenta algoritmi ce stau la baza dezvoltării proiectului, împreună cu modelul folosit pentru date. Secțiune patru se concentrează pe metodele de prezentare ale aplicației: interfața cu utilizatorul, ca design, dar mai ales ca interacțiune. Lucrarea se termină cu o sinteză a problematicilor abordate și prezintă direcțiile ulterioare de dezvoltare.

## MICROFORMATE

### Definiții

Microformatele, o abordare evoluționară ce adaugă semantică marcajelor bazate pe *HyperText Markup Language* (HTML) au fost adoptate cu succes de tot mai mulți editori Web și dezvoltatori de servicii în ultimii 2-3 ani. Actuala definiție a acestora, conform sitului oficial al inițiativei dedicate microformatelor [7], este:

*Proiectate în primul rând pentru oameni și în al doilea rând pentru mașini, microformatele reprezintă un set de formate de date simple și deschise, construite pe baza standardelor existente și larg adoptate.*

Secțiunea de discuții a aceluiași sit furnizează o definiție mult mai accesibilă:

*Microformatele sunt convenții simple pentru încapsularea de semantică în cadrul marcajelor HTML, în scopul facilitării dezvoltării descentralizate.*

Încă mai precis de atât, microformatele reprezintă anumite convenții de marcare a documentelor Web, pe baza elementelor *eXtensible HyperText Markup Language* (XHTML), via nume și valori de atribute, având o semantică precisă predefinită.

### Caracteristici importante

Principiile de bază în proiectarea microformatelor sunt *simplitatea* – proiectate pentru a rezolva o problemă specifică – și *interconectivitatea slabă* – reprezintă blocuri mici, specializate, prin a căror compunere se pot obține blocuri din ce în ce mai largi și se pot exprima semantici din ce în ce mai complexe, însă a căror interconectare este opțională, nediminuând puterea de expresivitate la nivel semantic.

Microformatele își ating scopurile fie prin adăugiri la marcajele XHTML (microformatele elementare), fie prin specificarea de valori pentru atributele asociate unor elemente ale limbajului XHTML și imbricări de astfel de elemente, construindu-se structuri pentru conținutul ce trebuie marcat semantic.

Specificarea (definirea) microformatelor este încă în desfășurare: o parte dintre acestea sunt complet specificate, iar altele sunt în stadiu de lucru. Indiferent de acest aspect, microformatele sunt larg răspândite, fie explicit, fie implicit. Pentru mai multe detalii se poate consulta secțiunea de discuții a sitului oficial [7] care prezintă, pentru fiecare microformat, și un index al implementărilor actuale.

### Microformate de interes

Lista actuală a microformatelor oficiale este compusă din următoarele: *hCalendar*, *hCard*, *rel-license*, *rel-nofollow*, *rel-tag*, *VoteLinks*, *XFN*, *XMDP*, *XOXO*, *adr*, *geo*, *hAtom*, *hResume*, *hReview*, *rel-directory*, *rel-enclosure*, *rel-home*, *rel-payment*, *Robots Exclusion* și *xFolk*.

Microformatele de interes pentru Microwler sunt cele care încapsulează date personale și relațiile dintre contacte:

- *hCard*-urile conțin descrieri de entități: persoane, companii etc.;
- *XFN* descrie tipurile de relații cu persoanele din lista de contacte salvate în format *hCard*.

### Exemplificări

Un exemplu modular de marcaj XHTML ce poate descrie autorul unui blog, recurge la *hCard* pentru a modela datele personale după cum urmează:

```
<div class="vcard">
  
  <a class="url fn n"
    href="http://purl.oclc.org/net/mariusb">
    <span class="given-name">Marius</span>
    <span class="family-name">Butuc</span>
```

```
</a>
<div class="adr">
  <span class="street-address">
    no. 9, Toma Cozma street
  </span>
  <span class="locality">Ia&#x219;i</span>,
  <abbr class="region"
    title="Ia&#x219;i">IS</abbr>,
  <span class="postal-code">700554</span>
</div>
<div class="telecom">
  <span class="tel">+40742754488</span>
  <a class="email"
    ref="mailto:mbutuc@infoiasi.ro">
    mbutuc@infoiasi.ro
  </a>
</div>
</div>
```

Relațiile sociale dintre autor și persoanele din lista de contacte se vor marca în format XFN astfel:

```
<a
  href="http://fenrir.info.uaic.ro/~evalica/"
  rel="friend met colleague">Ecaterina</a>
<a href="http://ppuuffaann.blogspot.com/"
  rel="friend met co-worker">Andi</a>
```

Pe baza microformatelor, putem extrage în mod automat atât datele de contact, cât și relațiile sociale.

### ALGORITMI ȘI ARHITECTURĂ

Primul pas în dezvoltarea unui motor de căutare este dezvoltarea componentei pentru adunarea informației de care avem nevoie. Dat fiind caracterul ubicuu al datelor de contact în Web-ul actual – cu precădere în blogosferă – singura soluție viabilă constă în *crawling*: pornim cu un set redus de documente și mergem din legătură în legătură către altele. Deși standardul *hCard* prevede numeroase proprietăți moștenite din standardul *vCard*, acestea se supun totuși principiului lui Pareto. Din acest motiv alegem doar un subset dintre aceste proprietăți, cele care se întâlnesc cel mai des în viața de zi cu zi. Cât despre vocabularul XFN, acesta va fi considerat în întregime.

Funcția ce realizează *crawling*-ul în documente poate fi definită recursiv, astfel încât fiecare legătură să apeleze din nou funcția, dar aplicând algoritmul BFS putem modifica ulterior codul mult mai ușor dacă vrem să căutăm continuu sau să salvăm o listă cu pagini neindexate pentru un *crawling* ulterior.

După ce am adunat informațiile, le stoca și le vom indexa. Aceasta implică una sau mai multe tabele care să conțină informația și sursa acesteia. Deși ne-am limitat la un subset din proprietățile *hCard*, putem întâlni alte excepții. Una dintre cele mai des întâlnite probleme apare în momentul în care avem *hCard*-uri diferite care au în comun unele câmpuri. E drept că mai multe persoane pot lucra la aceeași organizație, sau mai multe organizații pot avea aceeași poziționare geografică, dat fiind gradul actual de abstractizare. Totuși, problemele apar în clipa în care apar duplicate în câmpurile ce se presupun a fi unice, cum ar fi email, URL sau numele formatat. În acest caz vom

considera fiecare câmp ca având o podere, egală cu a celorlalte câmpuri și adunăm numărul câmpurilor prezente, comparând sumele între cele două *hCard*-uri. Dacă rata de asemănare este peste 80%, vom suprascris informațiile vechi cu cele noi. Dacă asemănarea este mai mică, vom combina datele dintre cele două contacte.

Din multitudinea de pagini Web ce conțin date personale de contact ce pot prezenta interes din punct de vedere social, nu toate se bucură de un cod XHTML valid. Pentru a preveni și rezolva această problemă vom folosi în proiect *Beautiful Soup* [8] – un *parser* HTML/XML ce analizează și marcaj invalid XHTML (pentru că a fost construit având în vedere marcajele SGML).

O altă problemă ce poate să apară este descoperirea unei versiuni noi al aceluiași *hCard* la un URL dat. Pentru a trata această excepție, nu salvăm în baza de date doar data la care am adăugat URL-ul procesat, ci și data ultimei modificări. Astfel avem un control mult mai bun asupra versiunilor informațiilor personale pe care le adunăm.

Ca pas final, informația trebuie returnată utilizatorului și afișată într-o manieră cât mai prietenoasă, clasificând lista de contacte în funcție de tipul relațiilor dintre cele două entități. Secțiune următoare dezvoltă metodele de prezentare ale aplicației: interfața cu utilizatorul, ca design, dar mai ales ca interacțiune.

#### INTERFAȚA CU UTILIZATORUL

În continuare, vom discuta interfața cu utilizatorul a aplicației și mecanismele de interacțiune dintre utilizator și aplicație. Încă din etapa de proiectare, Microwler a fost gândit avâns în minte utilizatorul final, fie el utilizator uman sau mașină. URL-urile au fost proiectate spre a fi mai curate, simple, mult mai ușor de citit sau chiar de intuit. Spre exemplu:

```
/hcard.php?hcard_id=4
```

reprezintă poate primul mod în care ne-am gândi să accesăm *hCard*ul cu ID-ul 4. Același efect îl obținem prin:

```
/hcard/4/
```

Fiind însuși un motor de căutare, Microwler a fost dezvoltat urmărind tehnicile *Search Engine Optimization* (SEO) [4]. Optimizarea merge mai departe de URL-urile prietenoase amintite anterior. Încă de la nivel de template, în timpul implementării s-a urmărit respectarea standardelor și ca utilizatorii ce folosesc browsere sau platforme diferite să nu trăiască experiențe diferite.

Întreaga interfață (a se consulta Figura 2) a fost împărțită în diviziuni logice:

- Zona superioară – antetul – găzduiește logo-ul, o scurtă descriere și sistemul navigațional principal reprezentat de un meniu orizontal;
- Sub antet regăsim trei coloane:
  - în stânga avem fie câmpurile pe care se efectuează interogarea, fie rădăcina arborelui de căutare – entitatea (persoana/organizația) după care utilizatorul a efectuat căutarea. Avem la dispoziție informațiile de contact culese de pe situl original. Pentru un plus de interactivitate, utilizatorul are

posibilitatea de a salva aceste date în format *vCard*, pentru a le folosi cu orice aplicație de gestiune a datelor personale ce suportă acest standard;

- coloana mijlocie conține informațiile de contact ale tuturor entităților care au o legătură directă cu persoana/organizația menționată anterior. Relațiile sociale sunt reprezentate grafic folosind setul de simboluri din Figura 1, create de Wolfgang Bartelme și Chris Messina [6];
  - în coloana din stânga avem cele mai noi zece date de contact adăugate sau modificate din baza de date;
- Dedesubt, utilizatorul are la îndemână un meniu orizontal ce include cele mai importante acțiuni.

	normal	met
me		n/a
friend		
sweetheart		
colleague		

Figura 1. Reprezentarea grafică a relațiilor sociale dintre persoane

Pentru a facilita accesul la cele mai noi zece date de contact adăugate sau modificate din baza de date cu un dispozitiv de tip *mouse*, s-a recurs la poziționarea acestei unități logice în zona medie pe verticală și extremă stângă pe plan orizontal. Am luat aceasta decizie pentru a scurta timpul de acces, pentru că se știe că majoritatea utilizatorilor sunt dreptaci.

Tot din punct de vedere al accesibilității, pentru a îmbunătăți experiența atât a utilizatorilor cu deficiențe de vedere cât și a acelor ce folosesc unelte alternative de vizualizare, în implementarea interfeței aplicației Web se va folosi un design relativ, *em-based* [1]. În acest mod, toate zonele logice ale aplicației au dimensiuni relative la mărimea caracterelor. Dacă un utilizator va dori să mărească dimensiunea caracterelor, toate diviziunile logice ale aplicației se vor redimensiona automat direct proporțional cu noile dimensiuni. De asemenea, suntem datori să luăm în calcul și scenariul în care un utilizator va folosi un dispozitiv mobil sau alt tip de dispozitiv cu afișaj de dimensiuni reduse ori rezoluții mici pentru a accesa aplicația. Acesta va putea să beneficieze în continuare de experiența completă, datorită faptului că întreaga interfață va fi scalată relativ la dimensiunea implicită, pentru acel dispozitiv, a caracterelor.

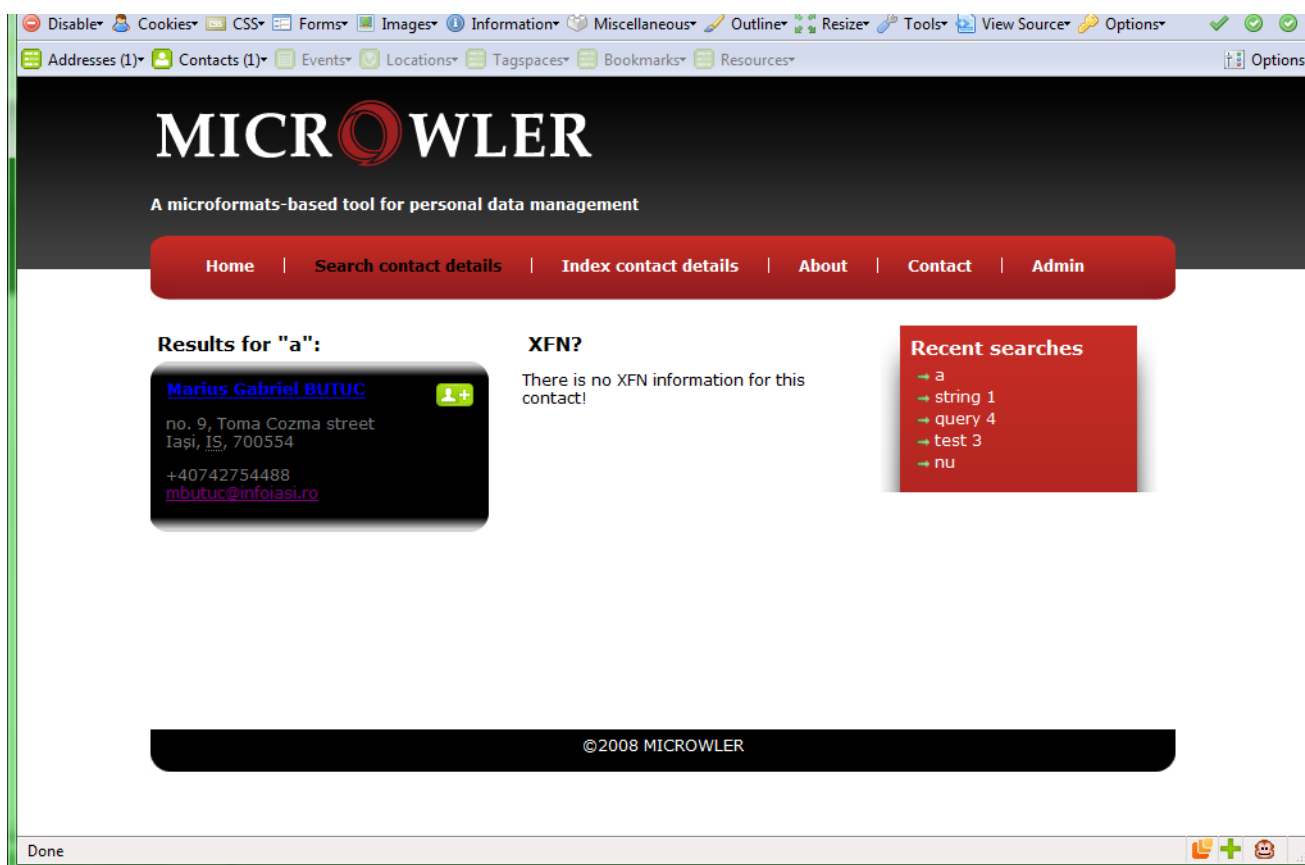


Figura 2. Interfața aplicației în momentul căutării unui contact

## CONCLUZII ȘI DIRECȚII VIITOARE DE DEZVOLTARE

Acest articol descrie un pas proaspăt în evoluția unui Web Social, în care utilizatorul devine autor: *microformatele*.

Când ritmul de viață tot mai accelerat ne determină să socializăm cu persoane/organizații uneori chiar de pe alte continente, datele personale de contact devin un lucru din ce în ce mai de preț.

O direcție importantă în care se poate avansa este aceea a recomandărilor: aplicația să ofere posibilitatea regăsirii de entități deja cunoscute de utilizator din cercul de cunoștințe apropiate ale contactelor sale.

Conceput și dezvoltat ca o aplicație extensibilă, Microwler este pregătit pentru permanenta evoluție a microformatelor. Cum propunerea unor microformate noi este deschisă larg comunității Web, apariția unor specificații care să încapsuleze noi semantici reprezintă o posibilitate de a îmbunătăți aplicațiile dedicate prin facilitarea accesului la date care sunt momentan „ascunse” de lipsa unui marcaj corespunzător. În clipa de față, pentru a suporta un microformat nou este nevoie doar să scriem un nou modul pentru Microwler pentru a putea parcurge structura corespunzătoare.

O altă oportunitate de dezvoltare o prezintă integrarea coordonatelor geografice regăsite în *hCard*-uri cu hărți Google, pentru a spori experiența utilizatorului.

Astăzi, când instalăm o aplicație software, există riscul ca versiunea noastră să fie deja învechită. Planurile pe termen lung includ transformarea Microwler într-un serviciu Web complet funcțional.

## REFERINȚE

- [1] J. Allsopp. *Microformats: Empowering Your Markup for Web 2.0*. Friends of Ed, 2007.
- [2] S. Buraga (coord.), *Tendențe actuale în proiectarea și dezvoltarea aplicațiilor Web* (în limba română), Matrix Rom, 2006.
- [3] V. Bush. *As we may think*. O'Reilly Network, 1945.
- [4] J. Grappone and G. Couzin. *Search Engine Optimization: An Hour a Day*. Wiley Publishing, Inc., 2006.
- [5] T. Segaran. *Programming Collective Intelligence: Building Smart Web 2.0 Applications*. O'Reilly Media, 2007.
- [6] C. Messina and W. Bartelme. *Microformats Icons*. <http://factorycity.net/projects/microformats-icons/>
- [7] \* \* \*, *Situl oficial dedicat microformatelor*: <http://microformats.org/>
- [8] \* \* \*, *Beautiful soup*. <http://crummy.com/software/BeautifulSoup>