

SERVICII DE STRUCTURARE A CONȚINUTULUI PAGINILOR WEB

Viorel CĂRBUNE, Andrei NICU, Dorin TURTUREAN, Irina MOROI, Liudmila CAZAC

Universitatea Tehnică a Moldovei

Abstract: În acest articol s-a descris o metodologie de structurare a conținutului paginilor Web care oferă o gamă largă de posibilități pentru Web developeri de a grupa conținutul din pagină în dependență de tipul acestuia, cu scopul de a optimiza indexarea paginilor Web de către motoarele de căutare. Pe lângă aceasta s-au identificat potențialele avantaje ale aplicării acestei metodologii la trecerea de la Web-ul clasic la cel semantic. Deasemeni s-a argumentat creșterea productivității procesului de comunicare dintre server și motoarele de căutare.

Cuvinte cheie: pagină Web, motor de căutare, Web semantic, Web developer, crawler, html, server.

Odată cu apariția Internetului a apărut și posibilitatea schimbului de date dintre utilizatorii acestuia. Apariția limbajului HTML a fost ca urmare a necesității prezentării adecvate a acelui-ași conținut pe diferite mașini de calcul. Astfel informația postată de către autor putea fi disponibilă cititorilor din întreaga lume indiferent de configurația mașinilor de calcul de care dispuneau aceștia. Limbajul HTML – Hyper Text Markup Language drept răspuns la cerințele clienților începe să evolueze și pe lângă text apare posibilitatea de a plasa pe pagini Web imagini, video, audio, forme și multe altele, astfel abătându-se de la sarcina inițială pentru care a fost conceput – de a marca textul, și capătă ca urmare tot mai multe funcționalități secundare. Din aceste considerente s-a ajuns la etapa în care conținutul fișierelor HTML a devenit neomogen. Lucrul Web developerilor cu astfel de fișiere a devenit mult mai anevoios. O tentativă de separare funcțională a fost divizarea stilurilor de restul conținutului paginii prin introducerea fișierelor CSS în arhitectura paginilor Web. Această structură arhitecturală oferă posibilitatea de a schimba modul de afișare a conținutului fără a modifica însuși conținutul. Necătfînd la aceste posibilități limbajul HTML își păstrează imunitatea în ceea ce privește conținutul.

Odată cu dezvoltarea tehnologiilor Web au apărut și alte tipuri de utilizatori decît cei clasici și anume roboții de căutare – utilizatori intermediari. Necătfînd la faptul că fiecare tip de utilizator are nevoie de tipuri de date diferite, Web serverele nu țin cont de acest lucru și asigură transferul integral de date prin protocolul HTTP livrînd tot conținutul paginii Web către client. Din aceste considerente ar fi mai eficient ca Web serverul să fie capabil de a oferi doar informația necesară sau în cel mai rău caz – întreaga categorie din care face parte informația cerută de către client. Pentru îndeplinirea acestui scop se recomandă gruparea conținutului paginii Web după careva categorii **Figura 1**.

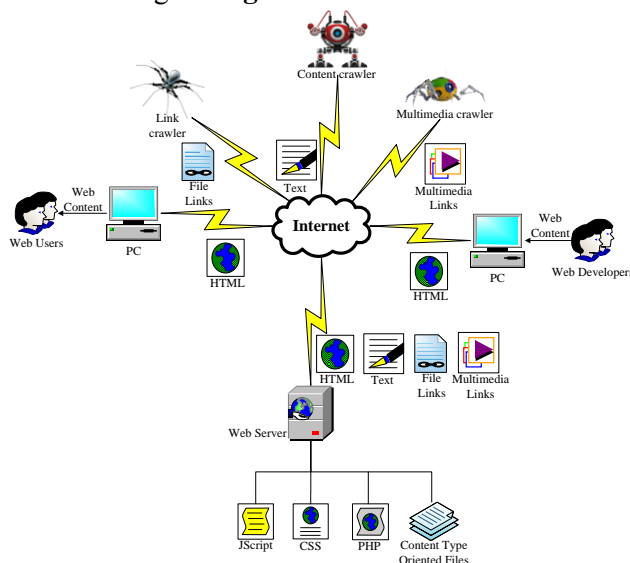


Fig. 1: Traficul tipurilor de conținut în internet.

Din **Figura 1** se observă faptul că Web Users sau utilizatorii clasici sunt deserviți în continuare prin aceeași metodă, adică la cererile acestora serverul răspunde prin a oferi informația deplină conținută pe pagina Web. Spre deosebire de utilizatorii clasici, utilizatorii de tip Link crawler sau roboții de analiză a referințelor nu au nevoie de conținutul textual de pe pagină și nici de codul HTML integral ci doar de referințe. Pe utilizatorii de tip Content crawler i-ar interesa doar conținutul textual. Grupul de utilizatorii Multimedia crawler înglobează în sine mai multe tipuri de roboți: de căutare a imaginilor, a clipurilor video, a fișierelor audio etc.

Procesul de structurare a conținutului este lăsată pe seama Web developerilor, însă se presupune că informația din aceeași categorie este amplasată în același fișier ce poartă denumirea paginii din care este extrasă iar extensia este definită de categoria din care face parte. De exemplu dacă din fișierul index.html se extrage întreg conținutul textual, acesta este plasat în fișierul real sau virtual index.cnt. Analogic index.ref, index.img și index.vid vor indica conținutul de tip adrese de referințe, căi către imagini, căi către fișierele video respectiv. Astfel rămîne la discreția Web developerilor de a indica apartenența conținutului la diverse categorii. Evident devine necesară adoptarea unui standard după care vor fi servite grupurile de Web clienți atipici. Aceștia vor avea posibilitatea de a accesa direct sau indirect fișierele cu conținutul necesar fără a încărca și altă informație adăugătoare.

Pentru a separa conținutul textual de codul HTML s-a apelat la facilitatea extinsă a browserelor moderne de a interpreta taguri introduse de utilizator, ceea ce ne dă posibilitatea de a marca locul unde va fi situat careva tip de content pe pagină. Pentru a ține cont de ordinea de amplasare a conținutului s-a introdus restricția ca fiecare linie din fișierele orientate pe conținut să corespundă unei perechi de taguri din fișierul HTML. Prin urmare fișierul html în rezultat nu va mai conține conținut ci doar marcajul componentelor HCML – Hyper Component Markup Language .

HTML: 212 B	HCML: 197 B	index.cnt: 40 B
<html>	<html>	HTML
<head>	<head>	Text1
<title>HTML</title>	<title><ci></ci></title>	Text2
</head>	</head>	LinkText1
		LinkText2
<body>	<body>	
<p>Text1</p> 	<p><ci></ci></p> 	index.ref: 12 B
<p>Text2</p> 	<p><ci></ci></p> 	
<a	<a	Link1
href="Link1">LinkText1 	href=""><ci></ci> 	Link2
<a	LinkText2 	href=""><ci></ci> 	index.img: 7 B
	 	
</body>	</body>	Source1

Implementarea metodei de structurare a conținutului paginilor Web devine posibilă odată cu implicarea directă a Web developerilor prin respectarea unor cerințe și standarde comune în procesul de dezvoltarea a paginilor Web. Implementarea acestei metodologii la etapa actuală ar presupune modificări introduse la nivelul de servere Web și anume implementarea posibilității de clasificare conform categoriilor a conținutului paginii Web clasice pentru clienții atipici prin adaugarea pe server a unui script care ar implementa această funcționalitate. În momentul în care spre server se adresează clienți atipici, aceștea pot cere acces către fișierele reale sau virtuale cu contentul necesar prin intermediul acestui script care va oferi acces direct către fișierul real sau îl va genera în cazul fișierului virtual.

Variația capacității conținutului interogărilor la trecerea de la tehnologia Web clasică la cea orientată pe tipologia utilizatorului este prezentată în **Tabelul 1**.

Tabelul 1. Capacitatea conținutului interogărilor.

HTML	HCML	Content	Reference	Image
212	212	212	212	212
212	212	212	212	212
256	197	40	12	7
256	197	40	12	7

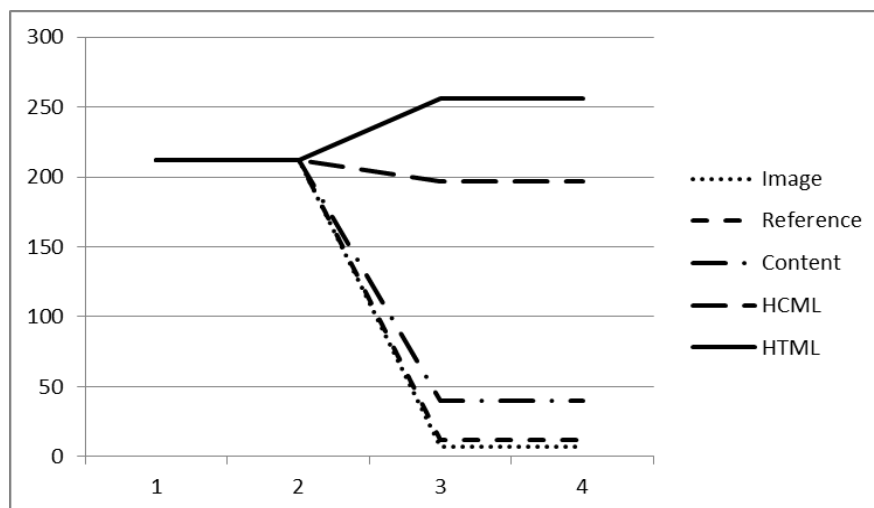


Fig. 2: Optimizarea traficului tipurilor de conținut în internet.

Graficul din **Figura 2** marchează o descreștere semnificativă a traficului de date pentru utilizatorii atipici, indiferent de categoria de conținut și o creștere nesemnificativă a traficului de date pentru utilizatorii clasici. Această concluzie însă nu poate fi univocă deoarece în cazul paginilor Web orientate preponderent spre utilizatorii clasici, creșterea nesemnificativă a traficului este proporțională în timp cu numărul de utilizatori ce accesează pagina și aplicarea acestei metodologii ar putea fi nerezonabilă din acest punct de vedere. Dacă însă prin aplicarea acestei metodologii se urmărește separarea conținutului textual al paginii Web de restul conținutului inclusiv de codul HTML sau gruparea conținutului conform cărorva criterii specifice, atunci metoda propusă poate contribui cu succes la rezolvarea acestor probleme. În cazul unor servicii Web cum ar fi serviciile de noutăți, care sunt orientate inclusiv în mare parte spre utilizatorii atipici reprezentați de alte pagini Web, implementarea acestei metodologii și-ar găsi rostul. Astfel înainte de a implementa această metodă, mai întâi de toate ar fi recomandabil de a analiza categoriile de utilizatori țintă și de a evalua raportul dintre aceștia pentru a putea estima avantajele și dezavantajele implementării metodologiei expuse.

Implementarea metodologiei expuse în lucrare nu presupune neapărat utilizarea metodelor de implementare descrise anterior. Metodele de implementare utilizând fișiere reale sau virtuale propuse anterior au ca scop final de a argumenta și de a exemplifica doar posibilitatea reală de utilizare a metodologiei de structurare a conținutului paginilor Web, și pot fi substituite în dependență de specificul problemei prin metode cât mai optimal posibil.

Bibliografie

1. <http://www.w3schools.com/>
2. <http://www.w3.org/>
3. <http://www.google.md/intl/ro/goodtoknow/web/101/>
4. <http://www.portalroman.com/articole/Internetul-140.html>
5. <http://invatza.info/paginiweb/introducere-paginiweb/5-ce-este-internetul.html>