

de la HTML la PDF (cu Perl și LaTeX)

Vlad Bazon, <http://sitsco.com>

Dec 2007

Obținem varianta tipărită a unui document HTML, folosind modulul *perl* HTML::Latex și limbajul de marcare/programare LaTeX; o anumită subrutină *perl* va asigura *download*-area documentului PDF obținut.

- Din nou, despre *documente*
- Transformarea documentelor cu <http://docs.google.com/>
- Utilizarea modulului HTML::Latex
- Elemente conjuncturale de limbaj LaTeX
- Subrutină *perl* pentru *download*
- Bibliografie

Din nou, despre *documente*

document provine din latinescul **documentum** care se trage de la **docere** care înseamnă **a învăța**; “document” este asociat adesea cu “a informa”, dar “a informa” — și-atât — este o *marginalizare* (obișnuită astăzi) pentru “a învăța”.

În principiu, Web-ul este o infinitate de documente “electronice” care sunt interconectate între ele și sunt accesibile de oriunde (datorită unor scheme universale de adresare, de identificare și de comunicare). Documentele de pe Web *nu* au în vedere “formatul A4” (nu vizează operația de tipărire pe hârtie), pentru motivul evident că de pe hârtie nu este posibilă interconectarea cu alte documente.

Dar operarea cu documente Web poate să aibă ca scop final chiar tipărirea pe hârtie (la edituri, instituții, etc.); astfel că au devenit necesare instrumente de conversie de la reprezentarea specifică documentelor Web (bazată pe limbaje de marcare, în principal HTML și XML), la “formatul A4”.

Pentru reprezentarea pe hârtie au fost produse două categorii de instrumente: așa numitele “procesoare de text” WYSIWYG, de exemplu aplicația Word din pachetul comercial Microsoft-Office, sau aplicația necomercială OpenOffice.org-Writer; respectiv, limbaje de marcare/programare (TeX, cu varianta LaTeX) precum și formate de reprezentare corespunzătoare, “independente de dispozitiv” (deci și pentru pagina de hârtie), precum PDF.

Instituțiile (funcționarii acestora) folosesc de regulă produsele de tip *point-and-click*, în principal Word; în general, ele nu produc documente *interconectate*, ci doar documente tipărite (folosind Web-ul doar pentru *transmiterea* fișierelor .DOC către instituțiile subordonate sau către beneficiari - urmând ca fișierele descărcate să fie tipărite local).

În cadrul sistemului de învățământ se studiază în mod extensiv, produsele *point-and-click* (și anume, exact acelea din varianta *comercială*, de la Microsoft). Motivul ar fi *ușurința* cu care se poate învăța să faci *click* pe diversele meniuri și iconițe, vizualizând *imediat* efectul (în principal, efectul de *formatare* asupra porțiunii selectate). Văzând această preferință exclusivă pentru *point-and-click*, poate fi îndreptățită bănuiala că scopul învățământului ar fi acela de a forma funcționari. Adevărul este însă acesta:

- este mult mai ușor de folosit un editor de text *obișnuit*, decât un “editor Word”; în orice caz, programatorii folosesc editoare de text (nicidecum, Word) atât pentru scrierea programelor, cât și pentru realizarea unor documente oricât de complexe;

- ca să folosești un editor de text *obișnuit* nu ai nevoie de o *licență* costisitoare, precum în cazul folosirii Microsoft-Office;
- editoarele de text *obișnuit* sunt independente de platformă (sau au versiuni pentru Windows, Linux, Mac, etc.);
- dacă ai învățat să folosești un *limbaj de marcare* (de exemplu, HTML), atunci va fi ușor să înveți și *alte* limbaje de marcare (CSS, XML; LaTeX);
- cunoscând două-trei limbaje de marcare, poți folosi un editor de text *obișnuit* pentru a realiza *orice* tip de document Web.

Desigur, însușirea unui limbaj de marcare poate să nu fie la fel de ușoară precum învățarea mecanismului *point-and-click*, mai ales dacă ai neglijat experiența lucrului cu fișiere text și cu linia de comandă și te-ai lăsat păcălit de ”valențele” acestui mecanism. Marele avantaj al mecanismului *point-and-click* este faptul că vezi imediat efectul fiecărei operații de editare; păcăleala vine de la obișnuința pe care o capeți de a opera pe bucățele, scăpând din vedere aspectele globale de structură și de formatare.

În cazul folosirii limbajelor de marcare, nu mai poți vizualiza *imediat* efectul. Mai întâi trebuie să descrii în fișierul text respectiv, structura și caracteristicile globale documentului (folosind marcajele și comenzile specifice acelui limbaj de marcare); apoi, trebuie să apelezi la un compilator sau la un program care poate interpreta marcajele respective și care fie produce ca rezultat un alt fișier text care va putea fi vizualizat cu o anumită aplicație existentă, fie chiar vizualizează documentul respectiv.

De exemplu, pentru a realiza un document HTML folosind un editor de text obișnuit, de obicei trebuie să completezi întâi secțiunea <head> precizând diverse caracteristici globale (de exemplu, *content*=”text/html; charset=utf-8”; stilurile de bază folosite; bibliotecile sau funcțiile angajate; etc.); apoi desigur, se adaugă conținutul propriu-zis, sub tag-ul <body>; vizualizarea se poate realiza prin încărcarea sursei respective într-un browser. De asemenea, pentru a realiza un document folosind LaTeX trebuie specificate de la bun început caracteristicile globale și pachetele de stiluri folosite; fișierul text respectiv trebuie compilat (apelând *latex*, sau *pdflatex*, etc.) obținând un fișier text care poate fi vizualizat cu anumite aplicații (”Acrobate Reader”, etc.).

Prin urmare, pentru a folosi limbaje de marcare ești obligat să deprinzi obiceiul (care și este foarte firesc) de a prevedea din start structura documentului și cerințele principale de formatare pentru secțiuni, sau pentru diverse categorii de elemente de conținut. Este adevărat însă că poți și să ajungi la *degradarea* limbajelor de marcare, lucrând ca și în Word (doar că, în loc să faci click-uri, scrii un simbol de marcare); de asemenea, în principiu se poate opera și în Word în maniera specifică limbajelor de marcare (optând de la început pentru un anumit *șablon* de document).

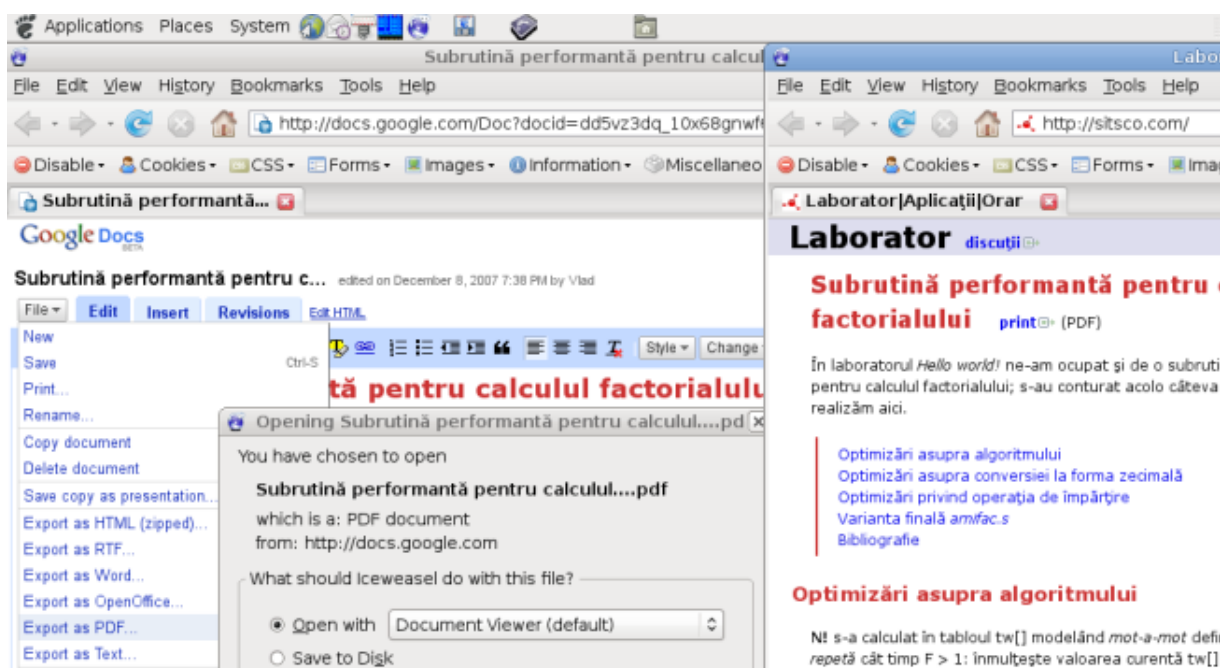
Limbajele de marcare (mai ales TeX) se bazează pe câteva aspecte care sunt binecunoscute de către oricine a răsfoit mai cu atenție vreo două cărți obișnuite și vreo două ”articole științifice”: un document are un *titlu* și un *autor*; are apoi, o mică ”introducere” sau un *sumar* (sau ”abstract”, sau poate o ”prefață”); are una sau mai multe secțiuni sau capitole, are un ”cuprins”, poate să aibă un glosar sau un index de termeni, are o ”bibliografie”; fiecare secțiune sau subsecțiune are un titlu și este formată din *paragrafe*, eventual *tabele*, *figuri*, sau imagini. Un limbaj de marcare instituie anumite secvențe de caractere (denumite *tag-uri*, sau *simboluri de marcare*) cu ajutorul cărora se pot marca elementele de conținut specificate mai sus; de obicei se prevăd și diverse simboluri pentru eventuala formatare *ad-hoc* a unei secvențe cât de mici, a documentului.

Documentul —fișier-text, nimic altceva— marcat astfel este ”consumat” de obicei de un browser (Internet Explorer, Firefox, etc.) care știe să interpreteze marcajele respective, redând documentul sub forma indicată prin marcajele fixate în text de către autor. Unele limbaje de marcare (de exemplu, LaTeX) sunt de fapt, chiar limbaje *de programare* (textul sursă trebuind să fie compilat) și oferă mari facilități: constituirea automată a cuprinsului, generarea indecșilor, gestionarea referințelor la pagini, tabele, la note de subsol sau la bibliografie, facilități de împărțire automată a unor porțiuni de text pe două sau mai multe coloane, și altele; oferă de asemenea, marcaje corespunzătoare pentru scrierea de formule matematice și pot produce în final (în funcție desigur, de ce a prevăzut autorul) documente de înaltă calitate tipografică.

Transformarea documentelor cu

Știm desigur, că un fișier .DOC poate fi exportat în format HTML direct de sub Word, utilizând opțiunea “Save As...” din meniul “File”. De obicei fișierul HTML rezultat arată destul de rău (sursa HTML obținută trebuie “curățată” de tag-urile inutile și eventual trebuie efectuate multe “reparații” asupra unor tag-uri și atribute care astăzi sunt declarate “deprecated”); de vină poate fi autorul .DOC-ului, fiindcă nu s-a preocupat câtuși de puțin de structurarea coerentă a documentului și nu a gândit unitar asupra formatării; dar chiar dacă s-ar proceda unitar, de exemplu formatând din capul locului toate paragrafele la fel - HTML-ul rezultat tot va conține atributele (aceleași pentru toate paragrafele) în cadrul fiecărui tag de paragraf...

Pentru editarea documentelor și pentru transformarea unui document dintr-un format în altul - și tot pe bază de *point-and-click* - Google a pus la dispoziție tuturor aplicația Docs. Descriem aici o modalitate de folosire — anume, ne propunem să obținem formatul PDF pentru un articol existent pe <http://sitsco.com/>.



Am accesat <http://sitsco.com/> și apoi (într-o a doua instanță Firefox) <http://docs.google.com/>; am selectat cu mouse-ul textul articolului și l-am copiat cu CTRL-C; în Google Docs am făcut click pe *New* și am făcut “paste” (prin CTRL-V), iar ca urmare am obținut articolul respectiv în fereastra de editare oferită; apoi, am deschis meniul *File* și am făcut click pe opțiunea *Export as PDF...*, obținând fereastra de dialog prin care se poate opta pentru deschiderea documentului PDF cu *Document Viewer*, sau pentru salvarea lui pe disk

Desigur, înainte de a exporta ca PDF, ar fi fost nevoie de unele retușuri în fereastra de editare oferită — de exemplu trebuie eliminate secvențele “^ TOP ^”. Rezultatul PDF se obține foarte ușor și este de bună calitate; dar poate fi și nemulțumitor, de vreme ce nu ai posibilitatea de a opta pentru un anumit format (de exemplu pagină cu două coloane, sau format de “articol științific”). Google Docs este adresat desigur utilizatorilor obișnuiți; programatorul ar vrea mai degrabă un fișier text intermediar (cu marcajele corespunzătoare) pe care să-l poată eventual modifica, înainte de a genera PDF-ul.

Utilizarea modului HTML::Latex

Avem nevoie de un program prin care tag-urile HTML să fie înlocuite corespunzător prin simboluri de marcarea LaTeX; de exemplu, `<h2>Titlul secțiunii</h2>` să fie înlocuit prin `\section*{Titlul secțiunii}` (sau de exemplu, prin `\chapter{Titlul secțiunii}` - după cum dorim, sau în funcție de genul de document LaTeX care se potrivește). Documentul HTML inițial poate fi reprezentat printr-un arbore (sau rețea) de obiecte de memorie, în care fiecare nod corespunde unui element HTML (înregistrând cel puțin tag-ul respectiv, lista atributelor din acel element, un pointer la “parent”-ul elementului și unul la primul “child”); având acest arbore, rămâne ca el să fie traversat (recursiv) pentru a aplica fiecărui nod funcția de înlocuire.

Pe CPAN am găsit un singur modul, **HTML::Latex**, destinat să creeze un fișier LaTeX dintr-unul HTML; se folosește HTML::TreeBuilder pentru a obține arborele de memorie corespunzător documentului HTML (fiecare nod este un obiect HTML::Element) și XML::Simple pentru a reconstrui arborele corespunzător noilor marcaje (și desigur, acea funcție care în asociere cu diverse alte subrutine, asigură corespondența dintre tag-uri HTML și marcaje LaTeX). Modulul respectiv este cam vechi (din anul 2000) și cu siguranță că el s-ar putea rescrie acum mai bine (cel puțin, în privința considerării “caracterelor internaționale”); folosim funcționalitatea *minimă* oferită de el, intenționând (și preferând deocamdată) să operăm manual modificările dorite sau necesare asupra fișierului LaTeX rezultat.

```
#!/usr/bin/perl -w
# ./html-latex.pl htmllatex.html htmllatex.tex
use HTML::Latex;
my ($fi, $fo) = @ARGV; # preia de pe linia de comandă, numele celor două fișiere
my ($in, $out);
open $in, ':encoding(utf8)', $fo or die $!;
my $parser = new HTML::Latex(); # obiect capabil să analizeze HTML (din $in) și
$parser->html2latex($in, $out); # să "traducă" în LaTeX (în $out)
```

Pe acest site avem `<meta http-equiv="Content-Type" content="text/html; charset=utf-8">` și desigur, putem să includem această specificație și în `<head>`-ul fișierului pe care-l transformăm; dar HTML::Latex nu ține cont de specificațiile din `<head>`-ul fișierului (deși unele dintre acestea vizează chiar conținutul propriu-zis și deci nu pot fi considerate “colaterale”; este drept însă că există posibilitatea “legării” unor funcții proprii, care să rezolve diverse aspecte particulare).

De aceea am specificat *utf8* pentru deschiderea celor două fișiere — în absența acestor specificații, caracterele românești “dispar”; mai precis, un “caracter internațional” este reprezentat prin 2-6 octeți de memorie care sunt interpretați la redare ca fiind un *singur* caracter *numai* în prezența anumitor specificații de codare (cum ar fi pentru browser, “Encoding: UTF-8”) — altfel, sunt redade eventual “caracterele” aferente celor 2-6 octeți constituienți și nu caracterul reprezentat unitar de aceștia.

Elemente conjuncturale de limbaj LaTeX

Redăm porțiunile semnificative din fișierul *htmllatex.tex* obținut prin executarea programului redat mai sus și completat apoi manual (în *preambul*—partea dinaintea `\begin{document}`); am adăugat un “abstract”); secvențele marcate prin %’ sunt “comentarii” (și vor fi ignorate de compilator).

```
\documentclass[10pt,a4paper]{article}
\usepackage{fullpage, url, ucs, graphicx} % pachete LaTeX
\usepackage[utf8x]{inputenc}
\renewcommand{\abstractname}{ } % elimină denumirea implicită "Abstract"
\setlength{\parskip}{1ex} % spațiul dintre paragrafe
\setlength{\parindent}{0ex}
\graphicspath{{/home/vb/lar/root/static/}} % path-ul pentru imagini (.PNG)

\begin{document}
```

```

\title{de la HTML la PDF (cu Perl și LaTeX)}
\author{Vlad Bazon, \texttt{http://sitsco.com}}
\date{Dec 2007}
\maketitle % comanda de includere a titlului

\begin{abstract}
Obținem varianta tipărită a unui document HTML, folosind modulul ...
\begin{itemize} % enumerare (similar cu <ul> și <li> din HTML)
\item Din nou, despre \emph{documente}
\item Transformarea documentelor cu \emph{\url{http://docs.google.com/}}
...
\end{itemize}
\end{abstract}

\section*{Din nou, despre \emph{documente}}
\textbf{document} provine din latinescul ...

\^In principiu, Web-ul este o infinitate de documente ...
...

\section*{ Transformarea documentelor cu \includegraphics[scale=0.6]{docslogo.png} }
Știm desigur, că un fișier .DOC poate fi exportat \^{\i}n format HTML:.
...

\begin{center}\includegraphics[scale=0.7]{sitdocs.png}\end{center}

Am accesat \emph{\url{http://sitsco.com/}} și apoi (\^{\i}ntr-o ...
...

\section*{ Utilizarea modulului HTML::Latex}
Avem nevoie de un program prin care tag-urile HTML să fie \^{\i}nlocuite...
...

\begin{verbatim} % similar cu <pre> din HTML
#!/usr/bin/perl -w

# ./html-latex.pl htmllatex.html htmllatex.tex

use HTML::Latex;
...
\end{verbatim}

...

\end{document}

```

Preambulul conține comenzi care afectează întregul document. Comanda `\documentclass[options]{class}` precizează compilatorului ce tip de document alege să creeze autorul (în cazul de față - document din clasa "article", de tipărit pe hârtie cu formatul A4, cu dimensiunea de bază a fontului de 10 puncte). După executarea comenzilor din preambol (cu efectul setării corespunzătoare a diversilor parametri globali), compilatorul va identifica textul documentului (conținutul propriu-zis) pe baza marcajelor `\begin{document}` și `\end{document}` și va aplica asupra lui comenzile (și marcajele de formatare) întâlnite¹.

LaTeX prevede mai multe clase de document, de exemplu: *article* pentru articole destinate revistelor științifice, pentru rapoarte de mică întindere, pentru documentații de programe, etc.; *report* pentru cărți

¹asemenea aserțiuni reducționiste nu trebuie luate "mot-a-mot" ...; de obicei lucrurile sunt mult mai complicate

mici sau rapoarte care conțin mai multe capitole, pentru teze de doctorat; *book* pentru cărți; *slides* pentru prezentări (analog celor produse cu PowerPoint); etc.

Pentru fiecare tip de document sunt prevăzute mai multe nivele imbricate de secționare a documentului; de exemplu, un *article* poate cuprinde `\part{Numele-Părții}`, `\section{Nume-Secțiune}`, `\subsection{...}`, `\subsubsection{...}`, `\paragraph{...}`, `\subparagraph{...}`. Secțiunile (la fel, tabelele, figurile, notele de subsol) sunt numerotate automat, într-o manieră care ține cont de ierarhia lor; titlurile de secțiuni pot fi adăugate automat într-un tabel de cuprins.

Imaginile pot fi incluse în document prevăzând în preambul includerea pachetului *graphics* și folosind apoi comanda `\includegraphics[opțiuni]{cale-nume-fișier-imagine}`; anumite conversii sunt efectuate automat (de la GIF la PNG); *graphics* oferă numeroase comenzi pentru integrarea imaginii, de exemplu scalare, “clip”, rotire, “wrap”, alipire de imagini, etc.

Diverse alte pachete LaTeX permit gestionarea aspectelor privitoare la tabele (*tabular*, *table*, etc.), figuri, redarea secvențelor de program sau de “pseudocod” (*listings*, *algorithm*, etc.), împărțirea textului în coloane (*multicol*, etc.), folosirea culorilor, folosirea de formule matematice oricât de complexe (*math*, *displaymath*, *amsmath*, etc.), folosirea caracterelor internaționale (*ucs*, *inputenc*), etc.

Pentru a obține direct un document PDF corespunzător fișierului *htmlatex.tex* creat, putem apela din linia de comandă compilatorul **pdflatex** (vb@debian: \$ pdflatex htmlatex.tex). Eventual, comanda trebuie executată de mai multe ori consecutiv, dacă există referiri (citări) la note de subsol, la bibliografie, la figuri, etc. Fișierul PDF obținut astfel ocupă mai mult spațiu decât fișierul DVI (“Device independent file format”) care s-ar fi obținut prin compilare cu **latex**, dar include toate fonturile necesare și de asemenea, imaginile (DVI doar le referă, urmând ca fișierele cu imagini să fie “incluse” doar în momentul tipăririi documentului); în plus, *pdflatex* recunoaște formatele de imagine JPG și PNG, în timp ce *latex* recunoaște numai formatul “Encapsulated PostScript” (EPS).

Fișierul PDF (“Portable Document Format”) poate fi tipărit ca atare (în general, imprimanta va “ști” să producă documentul conform indicațiilor și comenzilor existente în fișier), sau poate fi vizualizat (și tipărit) folosind aplicații precum Acrobat Reader, Xpdf, Evince, etc.

Subrutină perl pentru *download*

Articolele de pe acest site sunt fișiere HTML, dar nu sunt “documente HTML” în sensul standard; ele *nu* conțin partea “obligatorie” de `<head>`, fiind în principal doar o succesiune de `<div>`-uri (aferele secțiunilor din articol) și de paragrafe (desigur și de alte elemente HTML, precum `<h2>` folosit pentru titlurile secțiunilor, `<a>` pentru diverse referințe interne sau externe, `` etc.). Am putut “evita” să includ câte o parte de `<head>` în fiecare articol, prin existența unei “pagini de bază” *index.html* al cărei `<head>` prevede toate cele necesare (“content-type”, fișiere de stiluri, biblioteci *javascript* de inclus, etc.) și de asemenea, prin implicarea mecanismelor *Ajax* prin care fișierul respectiv este transferat (la click pe link-ul corespunzător, din cadrul paginii de bază) și este redat de browser în locul indicat în pagina de bază pentru aceasta. Cu alte cuvinte, articolele respective sunt *integrate la cerere* paginii de bază și *nu* sunt accesibile (sau redade) în mod independent (dezavantajul major fiind acela că articolele nu sunt “văzute” de motoarele de căutare; iar soluția acestei probleme pare a fi una singură...).

Așa stând lucrurile, pare dificil și să copiezi (altfel decât ca text neformat) și să printezi vreunul dintre articole (mai ales dacă e vorba să faci aceasta de sub Windows și cu Internet Explorer; altfel, am arătat mai sus că se poate foarte ușor cu Google Docs). În loc de a renunța la Ajax (soluție nedorită, dar... de loc rea) - sau poate și mai bine... de a abandona totul - ne-am pus această chestiunea practică: a asocia unora dintre articolele de aici o “variantă PDF”. Am arătat mai sus pe larg, demersurile și instrumentele necesare pentru aceasta (plus diverse conexiuni *specifice...*) și ar rămâne acum de arătat cum anume producem “download/save” pentru PDF-ul respectiv.

Ca interfață, am adăugat link-ul “print” (alături de titlul articolului) care (la *click*) accesează în mod explicit serverul pentru a executa funcția necesară (transmițându-i ca parametru numele fișierului PDF care trebuie returnat). Funcția respectivă este încorporată în modulul Controller::Qnote:

```

package lar::Controller::Qnote;
use strict;
use warnings;
use base 'Catalyst::Controller'; # vezi http://catalyst.perl.org (și [5])

# alte module incluse, alte subrutine...

sub down_exe : Local {
    my ($self, $c, $fin) = @_;          # $fin preia numele fișierului de downloadat
    my $fi = $c->config->{home}.'/root/'.$fin;
    my $size = (stat($fi))[7];          # dimensiunea fișierului
    open(INX, $fi) || die "nu se poate accesa $fi, $!";
    local $/;
    my ($data) = <INX>;
    close(INX);
    $c->res->content_length($size);     # transmite întâi header-ele necesare
    $c->res->headers->header('Content-Type' => 'application/pdf');
    $c->res->headers->header('Content-Disposition' => "attachement;filename=$fin");
    $c->res->body($data);                # transmite fișierul respectiv
}

```

Până să poată fi considerat ca "acceptabil" pentru a fi oferit spre *download*, PDF-ul respectiv necesită de obicei, unele retușuri; de exemplu pentru cazul articolului de față, a apărut o problemă chiar nebanală: textul din secțiunea anterioară care descrie structura de principiu a unui fișier LaTeX, conține *comenzi* LaTeX și trebuie încorporat și el în fișierul "htmlatex.tex" — se pune problema de a *evita* ca aceste comenzi să fie și executate de către *pdflatex*; inițial, am inclus textul respectiv sub "environment"-ul *verbatim*, care are o comportare asemănătoare cu <pre> din HTML (mai precis, programul *html-latex.pl* care folosea HTML::Latex a transformat <pre> în *verbatim*, dar fără să "escapeze" *corect* comenzile Latex din interiorul elementului <pre>). Numai că în cadrul *verbatim* poate apărea orice secvență de caractere (și va fi redată ca atare), cu excepția oricărei *comenzi* LaTeX (adică nu este permis să apară de exemplu "\begin{document}"); soluția este desigur "escaparea" caracterului special "\", prin "\textbackslash{" (deci comenzile din cadrul *verbatim* trebuie scrise sub forma: \textbackslash{ }begin{document}).

Bibliografie

- [1] http://en.wikibooks.org/wiki/LaTeX_guide_to_the_LaTeX
(From Wikibooks, the open-content textbooks collection)
- [2] <http://www.latex-project.org/>
LaTeX – A document preparation system (LaTeX project site)
- [3] <http://ctan.org/>
CTAN (Comprehensive TeX Archive Network)
- [4] <http://cpan.org/>
CPAN (Comprehensive Perl Archive Network)
- [5] Vlad Bazon - *Încorporarea unui editor WYSIWYG într-o pagină Web*
Gazeta de Informatică, mai 2007