

XML jako wejście Omegi

Paweł Grams

ENST-Bretagne BP 832

29-285 Brest CEDEX, France

pawel.grams@enst-bretagne.fr

Pracę zgłosił: Andrzej Borzyszkowski

Streszczenie

Obecne publikacje opracowane za pomocą \TeX a lub \LaTeX a rzadko są przekształcane na postać hipertekstową. Ponadto komplikacja języka \TeX , a w szczególności skomplikowana i zależna od kontekstu składnia, sprawiają iż trudno jest stworzyć oprogramowanie do automatycznej obróbki danych w tym formacie. Co prawda możliwa jest zmiana formatu \LaTeX na HTML i odwrotnie, lecz wymaga to odpowiedniego sformatowania pliku \LaTeX , co jednocześnie wiąże się z ograniczeniem swobody autora podczas tworzenia struktury dokumentu.

Wprowadzenie

\TeX jest językiem programowania dającym doskonałe możliwości formatowania tekstu. Jednakże stopień jego komplikacji czyni go niedostępnym dla początkującego użytkownika. Głównie z tego powodu powstał \LaTeX , udostępniający użytkownikowi zestaw prostych makr, pozwalających w pełni wykorzystać możliwości \TeX a przy minimum wiedzy o składaniu tekstu. System Omega również daje takie możliwości, jednocześnie pozwalając na obsługę wszystkich znaków Unicode. Kierunkiem dalszej ewolucji jest bardziej zunifikowana wersja \LaTeX 3.

Tematem mojego stażu, prowadzonego na École Nationale Supérieure des Télécommunications de Bretagne pod nadzorem prof. Yannisa Haralambousa, jest zaadaptowanie Omegi do środowiska XML poprzez nadanie jej możliwości czytania plików w tym formacie. Docelowo Omega będzie miała również zdolność zapisywania w plikach wyjściowych nie tylko jak dotychczas danych DVI, lecz również dodatkowych, redundantnych danych w formatach takich jak SVG, MathXML lub inne. Pozwoli to na efektywne przetwarzanie danych wejściowych za pomocą istniejących narzędzi ze świata XML.

Realizacja projektu

Główną przeszkodą utrudniającą stworzenie spójnego zestawu znaczników XML na podstawie komend \LaTeX a jest mieszanie w tym ostatnim, zawartości merytorycznej dokumentu z danymi służącymi do formatowania tekstu. Jako przykład niech posłuży komenda:

```
\begin{tabular}{c|c|c}
```

Wyraźnie widać, że w nawiasach klamrowych mamy tylko i wyłącznie informacje potrzebne do po-

prawnego ułożenia tekstu na stronie. Natomiast w komendzie `\emph{tekst}` tekst w nawiasach zawiera treści dokumentu. Najtrudniejsze do obsługi są komendy gdzie tekst dokumentu miesza się z parametrami niezbędnymi do składu:

```
\begin{tabular}{r@{równie}l}
```

Taka konstrukcja komend \LaTeX bardzo utrudnia automatyczną konwersję danych wejściowych na inny format, zautomatyzowane sprawdzanie poprawności plików, oraz konstrukcję narzędzi do wstępnej obróbki danych.

Przy tworzeniu $X\LaTeX$ starałem się stworzyć znaczniki XML podobne w konstrukcji do odpowiadających im komend \LaTeX a. Miało to na celu ułatwienie korzystania z nowego narzędzia osobom które już używały \LaTeX a. Pozostawiłem również wiele środowisk w ich pierwotnej formie. Można się do nich odwołać za pomocą znaczników typu: `<?math tekst ?>`. Wewnątrz takich znaczników używać można tylko komend (\LaTeX) dozwolonych dla danego środowiska. Miało to na celu zmniejszenie liczby znaczników XML do niezbędnego minimum, które pozwoli na wprowadzenie \LaTeX a do świata XML, a jednocześnie nie będzie wymagało od użytkowników czasochłonnego zgłębiania listy wszystkich dostępnych znaczników.

Strona techniczna

Ze strony technicznej, pierwsza część przedsięwzięcia została zaprojektowana tak, aby pozwolić na łatwą adaptację do późniejszych zmian wprowadzonych w strukturze $X\LaTeX$.

Mianowicie jedyna zmiana w kodzie Omegi polega na wbudowaniu możliwości korzystania z zewnętrznego filtra wstępnie przetwarzającego dane wejściowe. Filtr ten, pobierał będzie dane z pliku,

konwertował je na format \LaTeX i na bieżąco przekazywał do Omegi. Ze strony użytkownika wymagane będzie tylko podanie w linii komend Omegi nazwy pliku zawierającego filtr niezbędny do poprawnego odczytania pliku wejściowego. Użytkownik korzystający z formatu $X\LaTeX$, nie będzie miał więc w ogóle do czynienia z plikami w formacie \LaTeX .

Zastosowanie filtru pozwoli zaawansowanym użytkownikom zainteresowanym zmianami w formacie $X\LaTeX$ lub wprowadzeniem nowego formatu, na łatwą adaptację Omegi do swoich potrzeb.

Struktura $X\LaTeX$

Podczas tworzenia struktury znaczników $X\LaTeX$ starałem się rozdzielić treść dokumentu od parametrów komend \LaTeX , definiując te ostatnie jako wartości argumentów elementów XML. Jedynym wyjątkiem, gdzie nie udało się uniknąć pomieszania tych dwóch typów danych jest środowisko `tabular`. Jak wiadomo, w definicji tabeli, za pomocą komendy `@{}` można wstawić dowolny tekst, który następnie będzie systematycznie umieszczany w tym miejscu, w każdym wierszu. Aby uniknąć komplikacji struktury $X\LaTeX$ zaproponowałem dwie wersje znacznika `<tabular>`. W pierwszej, format tabeli podawany jest bezpośrednio jako wartość argumentu, zaś w drugiej, jest on objęty znacznikami `<format>`.

Jeśli chodzi o wartości argumentów, powinny one dokładnie odpowiadać wartościom jakie przyjmują parametry odpowiednich komend języka \LaTeX . Nie mogą one zawierać znaczników $X\LaTeX$. Trzeba również pamiętać o tym, iż argumenty obowiązkowe dla komend \LaTeX w $X\LaTeX$ również pozostają argumentami, których obecność jest obowiązkowa.

Jedynym dodanym sztucznie argumentem jest `star`. Został on wprowadzony w celu zmniejszenia ilości znaczników $X\LaTeX$. Jeśli jego wartością jest `yes`, otrzymujemy środowisko „z gwiazdką”. Jeśli argument jest nieobecny, lub różny od `yes`, wówczas otrzymamy środowisko bez gwiazdki.

Należy zauważyć iż podczas zmiany formatu plików z \TeX na $X\LaTeX$, zmienił się również zestaw znaków, które należy chronić. Obecnie należy chronić znaki `<` `>` `&` w całym pliku, oraz `"` `'` wewnątrz wartości argumentów.

Środowiska

W $X\LaTeX$ lista znaczników XML obejmuje wiele podstawowych komend niezbędnych do napisania poprawnego dokumentu. Pozostałe komendy można uzyskać korzystając bezpośrednio z możliwości wpisywania kodu \TeX . Taką możliwość daje instrukcja zmiany trybu: `<?tex >`. Postanowiliśmy jednak

wprowadzić dodatkowo również inne tryby, odpowiadające często używanym środowiskom.

W większości przypadków przejście wymaga od filtru jedynie użycia istniejącego środowiska. Jednakże w przypadku środowisk `verbatim` i `verbatim*` zdecydowaliśmy się na stworzenie nowego środowiska, które będzie mogło być używane bez ograniczeń wewnątrz innych środowisk. W tym celu, filtr dodaje¹ automatycznie na początek pliku definicję nowego środowiska: `\newenvironment{my}{\begin{alltt}}{\end{alltt}}`. Następnie wszystkie znaki `\` obecne wewnątrz, są zamieniane na `\textbackslash`, oraz wprowadzana jest ochrona nawiasów klamrowych. W efekcie otrzymujemy tekst identyczny z tym produkowanym przez środowisko `verbatim`, lecz możliwy do uzyskania w każdych warunkach.

Środowisko `verbatim*` jest możliwe do uzyskania poprzez użycie trybu:

```
<?verbatim star="yes" >
```

który daje identyczny efekt. Zasada działania jest taka sama jak dla `verbatim`, z tym iż dodatkowo spacje zamieniane są na `\vs` zdefiniowane jako:

```
\def\vs{\leavevmode\hbox{\tt\char'\ }} i odpowiadające dokładnie znakowi używanemu do zaznaczania spacji w środowisku verbatim* Definicja ta jest dodawana1 automatycznie przez filtr na początku pliku.
```

Zakończenie

Celem tej publikacji, jak i wystąpienia na konferencji `BachTeX 2003`, jest uzyskanie odzewu ze strony środowiska \TeX -owego, na temat poprawności ustalonego przez nas zestawu znaczników $X\LaTeX$. Chcielibyśmy w niedalekiej przyszłości ogłosić ostateczny zestaw znaczników $X\LaTeX$, w postaci użytecznej i niewymagającej poważnych zmian. Kolejnym krokiem będzie nadanie Omedze możliwości zapisywania dodatkowych danych w plikach wyjściowych, również w postaci znaczników XML (ewolucja formatu DVI do DVX). Być może uda się uzyskać zapis dodatkowej porcji danych w formatach SVG, MathXML lub innych podobnych. Taki kierunek ewolucji Omegi pozwoli na zwiększenie możliwości obróbki tekstu, zarówno w formie jeszcze nie przetworzonej, jak i już gotowej do druku.

¹ dodaje virtualnie, gdyż Omega nigdzie nie zapisuje pliku stworzonego przez wyjście filtru.

Lista markerów XML zdefiniowana do dnia 02-04-2003

documentclass, usepackage	<documentclass options="options" name="name"/> => \documentclass[options]{class}
include, includeonly, input	<include file="file"/> => \include{file}
pagestyle, thispagestyle	<pagestyle name="name"/> => \pagestyle{name}
table	<table pos="pos" star="yes"> text </table> => \begin{table*}[pos] text \end{table*}
letter	<letter><recipient> text </recipient> maintext </letter> => \begin{letter}{text} maintext \end{letter}
tabular	<tabular format="format"> text </tabular> lub: <tabular><format> format </format> text </tabular> => \begin{tabular}{format} text \end{tabular}
Pomiędzy znacznikami format można stosować inne znaczniki \LaTeX	
chapter, paragraph, part, subparagraph, section, subsection, subsubsection	<section star="yes"> text </section> => \section*{text} lub: <section><header> text1 </header><toc> text2 </toc> text </section> => \section[text1]{text}\tocsection{text2}
linebreak, nolinebreak, nopagebreak, pagebreak	<linebreak weight="weight"/> => \linebreak[weight]
address, date, signature	<address> text <address/> => \address{text}
glossary, index	<index entry="entry/>" => \index{entry}
Argument output powinien zawierać komendy \LaTeX niezbędne do poprawnego wyświetlenia tekstu	lub: <index output="output"> text <index output="output1"> text1 </index></index> => \index{text@output!text1@output1}
multicolumn	<multicolumn cols="colsnumber" pos="pos"> text </multicolumn> => \multicolumn{colsnumber}{pos}{text}
cline	<cline range="range"/> => \cline{range}
item	<item mark="mark"> text </item> => \item[mark] text
hyphenation	<hyphenation words="words"/> => \hyphenation{words}
group	<group> text </group> => {text}
abstract, center, description, document, enumerate, flushleft, flushright, itemize, quotation, quote, tabbing, titlepage, trivlist, verse	<document> text </document> => \begin{document} text \end{document}
caption	<caption star="yes"> text </caption> => \caption*{text} lub: <caption><header> text1 </header> text </caption> => \caption[text1]{text}
author, cc, closing, date, displaylines, encl, fbox, hbox, mbox, opening, ps, thanks, title, vbox	<mbox> text </mbox> => \mbox{text}
and, appendix, bigskip, cleardoublepage, clearpage, dotfill, fussy, hline, hrulefill, indent, kill, ldots, listoffigures, listoftables, makeglossary, makeindex, maketitle, medskip, newline, newpage, noindent, sloppy, smallskip, tableofcontents, vline	<newpage/> => \newpage

Paweł Grams

bf, emph, footnotesize, huge, Huge, it, large, Large, LARGE, md, normal, normalsize, rm, sc, scriptsize, sf, sl, small, tiny, tt, up

	<code><rm> text </rm> ⇒ \begin{rm} text \end{rm}</code>
footnote	<code><footnote id="id"> text </footnote> ⇒ \footnote[id]{text}</code>
footnotemark	<code><footnotemark id="id"/> ⇒ \footnotemark[id]</code>
footnotetext	<code><footnotetext id="id"> text </footnotetext> ⇒ \footnotetext[id]{text}</code>
figure	<code><figure pos="pos"> text </figure> ⇒ \begin{figure}[pos] text \end{figure}</code>
rule	<code><rule voffset="off" width="width" height="height"/> ⇒ \rule[off]{width}{height}</code>
marginpar	<code><marginpar><left> textleft </left> textright </marginpar> ⇒ \marginpar[textleft]{textright}</code>
picture	<code><picture size="size" offset="offset"> <?tex text ?></picture> ⇒ \begin{picture}(size)(offset) text \end{picture}</code>
Wewnątrz znaczników przechodzimy bezpośrednio do trybu T _E X	
minipage	<code><minipage pos="pos" width="width"> text </minipage> ⇒ \begin{minipage}[pos]{width} text \end{minipage}</code>
framebox, makebox	<code><makebox pos="pos" width="width"> text </makebox> ⇒ \makebox[width][pos]{text}</code>
parbox	<code><parbox pos="pos" width="width"> text </box> ⇒ \parbox[pos]{width}{text}</code>
addvspace	<code><addvspace length="length"/> ⇒ \addvspace{length}</code>
thebibliography	<code><thebibliography offset="offset"> text </thebibliography> ⇒ \begin{thebibliography}{offset} text \end{thebibliography}</code>
bibliographystyle	<code><bibliographystyle name="name"/> ⇒ \bibliographystyle{name}</code>
bibitem, cite	<code><bibitem text="text" id="id"/> ⇒ \bibitem[text]{id}</code>
pagenumbering	<code><pagenumbering name="name"/> ⇒ \pagenumbering{name}</code>
shortstack	<code><shortstack pos="pos"> text </shortstack> ⇒ \shortstack[pos]{text}</code>
label, pageref, ref	<code><label id="id"/> ⇒ \label{id}</code>
nowa linia	<code><br star="yes" skip="skip"/> ⇒ *[skip]</code>
	<code><TeX/> ⇒ \TeX{}</code>
	<code><LaTeX/> ⇒ \LaTeX{}</code>
	<code><cr/> ⇒ \cr</code>
	<code><italcor/> ⇒ \/</code>
	<code><ndash/> ⇒ ~~~~~2012</code>
inne	<code><mdash/> ⇒ ~~~~~2014</code>
	<code><nbspsp/> ⇒ ~</code>
	<code><par/> ⇒ \par</code>
	<code><spc/> ⇒ \u</code>
	<code><tab/> ⇒ & (w środowisku tabular) lub: \>(w środowisku tabbing)</code>
	<code><tabdef/> ⇒ \=</code>
	<code><?tex text ?> ⇒ text</code>
	<code><?math text ?> ⇒ \$text\$</code>
	<code><?displaymath text ?> ⇒ \$\$text\$\$</code>
	<code><?equation text ?> ⇒ \begin{equation}text \end{equation}</code>
środowiska nie przekształcane na X _E L _A T _E X	<code><?verbatim text ?> ⇒ \begin{myverb}text \end{myverb}</code>
	<code><?verbatimstar text ?> ⇒ \begin{myverb*}text \end{myverb*}</code>
	<code><?alltt text ?> ⇒ \begin{alltt}text\end{alltt}</code>