

Zpravodaj Československého sdružení uživatelů TeXu

Jiří Kosek
Sazba XML

Zpravodaj Československého sdružení uživatelů TeXu, Vol. 13 (2003), No. 1, 3–6

Persistent URL: <http://dml.cz/dmlcz/149910>

Terms of use:

© Československé sdružení uživatelů TeXu, 2003

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

V posledních měsících proběhlo v texové diskusní skupině i jinde několik diskusí o tom, zda je lepší při přípravě dokumentů používat $\text{T}_{\text{E}}\text{X}$ nebo XML. Pokusím se příliš nezneužít možnosti zaplnit stránky Zpravodaje názory, které jsem v diskuzích zastával já. Konec konců celá diskuse neměla moc smysl, protože $\text{T}_{\text{E}}\text{X}$ a XML spolu mohou žít bok po boku bez nějakých velkých problémů.

Asi se všichni shodneme, že $\text{T}_{\text{E}}\text{X}$ je výtečný nástroj pro sazbu dokumentů. Po dlouhá léta se dokumenty pro sazbu připravovaly přímo v $\text{T}_{\text{E}}\text{X}$ u. $\text{T}_{\text{E}}\text{X}$ tak sloužil jako jazyk pro zápis dokumentu a také jako systém pro jeho sazbu a zlom. Pokud je naším cílem jen dokument připravit, pěkně vysázet a vytisknout, je potom $\text{T}_{\text{E}}\text{X}$ vhodným a postačujícím nástrojem.

Pokud však chceme s připraveným dokumentem naložit i jinak, než jej jen vytisknout, třeba převést do HTML nebo z něj dotazem vybrat nějaká data, není už $\text{T}_{\text{E}}\text{X}$ jako formát pro primární pořízení dokumentu nejvhodnějším kandidátem. Vhodnější pro tyto účely jsou značkovací jazyky, v dnešní době především XML. S tím by možná někteří diskutéři nesouhlasili, ale zkusme to brát jako fakt. XML se zkrátka hodí pro ukládání informací, protože je lze dále poměrně snadno zpracovávat. $\text{T}_{\text{E}}\text{X}$ se hodí pro sazbu dokumentů, protože je rychlý a má implementovány kvalitní typografické algoritmy. Nic nám tedy nebrání ve využití výhod obou technologií – dokumenty ukládat do XML a pro potřeby tisku je zpracovávat $\text{T}_{\text{E}}\text{X}$ em. Získáme tak výhody obou světů.

Možností, jak sazbu dokumentu XML provést pomocí $\text{T}_{\text{E}}\text{X}$ u, je několik a my se je pokusíme přiblížit. V navazujících článcích jsou pak ty nejpoužívanější metody rozebrány více do hloubky včetně příkladů použití.

Přímá sazba XML

Asi nejjednodušší variantou je přímá sazba XML pomocí $\text{T}_{\text{E}}\text{X}$ u. Nepotřebujeme žádné přídatné programy, XML se zpracovává čistě prostředky $\text{T}_{\text{E}}\text{X}$ u. Tento přístup má však dva poměrně nepříjemné nedostatky.

Napsat parser XML přímo v $\text{T}_{\text{E}}\text{X}$ u je velmi obtížný úkol. Existující parsery jsou schopné načítat jen podmnožinu XML, obvykle nepodporují entity, které umožňují jeden dokument XML rozložit do více fyzických souborů.

Druhým nedostatkem je samotný princip, kdy musíme mapovat výskyty značek v XML na texové sekvence. V mnoha případech toto mapování můžeme poměrně jednoduše vytvořit, ale jsou případy, kdybychom v místě výskytu značky

potřebovali znát informace z jiné části dokumentu XML. Lze to samozřejmě obejít pomocí ukládání stavových informací nebo víceprůchodovým zpracováním dokumentu apod. Ale v porovnání se speciálními nástroji, které během zpracování umožňují přístup k celému dokumentu, je to velmi nepohodlné.

Pro přímou sazbu XML pomocí $\text{T}_{\text{E}}\text{X}$ u se dnes používají dva systémy. Prvním z nich je `xmltex`. Jedná se o jednoduchý parser implementovaný v $\text{T}_{\text{E}}\text{X}$ u, kterému můžeme v konfiguračním souboru určit, na jaké $\text{T}_{\text{E}}\text{X}$ ové sekvence se mají převádět jednotlivé značky v dokumentu XML. Podrobnější popis `xmltex`u naleznete v samostatném článku *Použití parseru XML v $\text{T}_{\text{E}}\text{X}$ u* v tomto čísle *Zpravodaje*.

Druhým produktem, který si s XML poradí přímo, je `Con $\text{T}_{\text{E}}\text{X}$ t`. Tato uživatelsky příjemná makronadstavba nad $\text{T}_{\text{E}}\text{X}$ em obsahuje i nástroje na zpracování XML.

Konverze XML do $\text{T}_{\text{E}}\text{X}$ u

Přímou sazbu XML $\text{T}_{\text{E}}\text{X}$ em lze použít jen v jednoduchých situacích. I při složitějších požadavcích však nemusíme na použití $\text{T}_{\text{E}}\text{X}$ u rezignovat. Pro převod XML do podoby $\text{T}_{\text{E}}\text{X}$ ového kódu můžeme použít sofistikovanější nástroje, které nám umožní pohodlné načtení XML a jeho převedení.

Asi nejvhodnějším nástrojem pro tuto úlohu je jazyk XSLT. Tento jazyk primárně slouží k popsání transformaci dokumentu XML do podoby HTML stránky, případně XML dokumentu s jinou strukturou. XSLT však s jistými omezeními umí generovat i obyčejné textové soubory, a tedy i $\text{T}_{\text{E}}\text{X}$ ové dokumenty. XSLT pomocí dotazovacího jazyka XPath nabízí přístup k libovolné části dokumentu XML kdykoliv během zpracování a řeší tak pohodlně problémy jako je generování obsahu, křížových odkazů apod.

Napsání stylu XSLT pro převod XML do $\text{T}_{\text{E}}\text{X}$ u je poměrně jednoduché. Musíme si dát pozor jen na správné ošetření a převod znaků, které mají v $\text{T}_{\text{E}}\text{X}$ u speciální význam (např. #, \$, &).

Pěknou ukázkou použití XSLT při sazbě knihy z podkladů v XML jste mohli nalézt v minulém čísle *Zpravodaje* v článku Zdeňka Wagnera *Využití XML a $\text{L}_{\text{A}}\text{T}_{\text{E}}\text{X}$ u při sazbě odborných knih*.

Využití standardních stylových jazyků

Výše popsané metody nám umožňují z dokumentů XML získat kvalitní tištěný výstup pomocí $\text{T}_{\text{E}}\text{X}$ u. Skalní příznivce značkovacích technologií by však tomuto přístupu přece jen jedno vytkl: formátovací předpis je úzce svázán se systémem $\text{T}_{\text{E}}\text{X}$ – nelze jej využít s jiným sázecím systémem a nemůžeme tak dostat jiný

výstup než formáty podporované (pdf) \TeX em – DVI, PostScript a PDF. Co kdybychom však chtěli výstup v jiném formátu, jako je třeba RTF nebo HTML? Vytvořená konverze do \TeX u by nám byla k ničemu.

Tento problém ve světě XML řeší stylové jazyky. Ty umožňují popsat vzhled dokumentu poměrně abstraktně a nezávisle na výstupním formátu a médiu. Takovýmto obecným předpisem – stylem – se může řídit převod do mnoha formátů v různých aplikacích. Stačí, když se standardizuje jazyk, kterým se styl popisuje. Je to jen přirozené pokračování snahy vytvořit nezávislý formát dat – tím je XML. Vzhled dokumentu by měl být také definován nezávisle, a to jak na dokumentu, tak na programu, který se použije pro výsledné formátování.

Stylových jazyků, které se hodí pro formátování XML dokumentů, existuje mnoho, ale mezi dva nejpoužívanější a „nejmocnější“ patří DSSSL a XSL. Shodou okolností pro oba dva jazyky existují i jejich implementace, které pro samotné formátování používají \TeX . V případě DSSSL se jedná o Jade \TeX a v případě XSL pak o Passive \TeX . I jejich základním principům a použití se věnujeme v samostatných článcích (články o Jade \TeX u i o Passive \TeX u jsou v tomto čísle Zpravodaje).

Použití stylových jazyků přináší kromě nezávislosti na použitém formátovacím nástroji i další výhody. Většinu parametrů sazby lze nastavit jednoduše jako hodnoty vlastností a není potřeba pronikat do makrojazyka \TeX u. DSSSL i XSL obsahují dotazovací jazyk, který umožňuje kdykoliv pracovat s libovolnou částí vstupního dokumentu XML – nejsme tak omezeni jako u \TeX u, který vstupní soubor čte sekvenčně. Své výhody má i \TeX – lze jej plně programovat, a jeho možnosti záleží jen na schopnostech uživatele, ne na jednou pevně daných vyjadřovacích možnostech standardizovaného stylového jazyka.

Další vývoj

Sazba XML pomocí \TeX u je možná a v mnoha případech můžeme získat výborné výsledky. Nicméně jsou dvě oblasti, které by potřebovaly výrazně zlepšit. \TeX je sám o sobě 8bitový a zpracování dokumentů používajících jako znakovou sadu Unicode mu činí potíže. Jde je sice obejít, ale postup je krkolomný a výsledek není vždy zcela robustní. Řešením může být přechod na systém Omega, který do \TeX u přidává právě podporu Unicode.

Druhým omezením je samotný vstupní procesor \TeX u. Jde s ním sice dělat kouzla, ale napsat v něm plnohodnotný parser XML nebo jiného formátu není příliš jednoduché. Vhodnější by byl modulární sázeč systém, kde by šlo použít uživatelské vstupní moduly speciálně připravené na čtení XML. Něco podobného by mělo umožňovat NTS.

Vše záleží na tom, jak rychle se budou „nové verze \TeX u“ jako NTS a Omega vyvíjet, a jak rychle se rozšíří mezi uživateli. Bez jejich rozšíření se asi nikdo nebude obtěžovat s přepisováním stávajících maker pro tyto nové nadějně projekty.

Summary: Typesetting XML

This article summarizes methods suitable for processing XML documents by the \TeX system – direct typesetting (`xmlltex`, `ConTeXt`), conversion to \TeX (`XSLT`) and \TeX based stylesheet language implementations (`XSL`, `DSSSL`). The article acts as an introduction for more detailed articles about processing XML with \TeX .

Jiří Kosek
jirka@kosek.cz

Použití parseru XML v \TeX u

JIŘÍ KOSEK

Chceme-li XML dokumenty sázet čistě prostředky \TeX u bez nutnosti používání dalších pomocných programů, je asi nejvhodnějším nástrojem `xmlltex`. Jedná se o makro nadstavbu na \LaTeX em, která implementuje parser XML. Ten umožňuje načítání dokumentů XML a umí podle pravidel definovaných v konfiguračním souboru převádět jednotlivé elementy na sekvence texového kódu, které pak řídí sazbu. Autorem `xmlltex`u je David Carlisle, který je uznávaným odborníkem jak na oblast \TeX u, tak i XML a zejména pak stylového jazyka XSL.

Většina moderních distribucí \TeX u již `xmlltex` obsahuje jako samostatný formát. Vysázení dokumentu XML je pak otázkou spuštění jediného příkazu:

```
xmlltex dokument.xml
```

resp.

```
pdfxmlltex dokument.xml
```

v případě, že chceme získat výstup v PDF pomocí `pdf \TeX u`. Nemáme-li `xmlltex` připravený jako samostatný formát, jde využít jednoduchý pomocný dokument v \TeX u, který nadefinuje jméno dokumentu ke zpracování a poté aktivuje `xmlltex`:

```
\def\xmlfile{dokument.xml}  
\input xmlltex.tex  
\end{document}
```

Tento dokument se pak zpracuje pomocí klasického \LaTeX u.