

XML versus T_EX – jemné porovnání

Martin Tůma

31. srpna 2004

Obsah

1	Úvod	3
1.1	TEX	3
1.2	XML	4
2	Srovnání v jednotlivých oblastech použití	4
2.1	Dokumenty všedního dne	4
2.2	Sazba matematiky a vědecké literatury	7
2.3	Elektronické publikování	7
3	Závěr	8

1 Úvod

Srovnání T_EXu, robusního a časem prověřeného typografického systému s poměrně mladou, ale o to více se rozvíjející technologií XML by jistě vydalo na celou knihu, tento článek je proto pouhým shrnutím možností obou technologií při řešení nejzákladnějších uživatelských činností. Především jde o tvorbu klasických článků a knih, vědeckých článků s matematickými výrazy a elektronické publikování. Na úvod krátké shrnutí obou technologií.

1.1 T_EX

T_EX je sázecí systém, který umožňuje vytvářet dokumenty složité struktury i vysoké typografické kvality. Jeho autorem je Donald Ervin Knuth, profesor počítačových věd na Standfordské univerzitě, a jeho vznik se datuje již do roku 1978.

Názvem T_EX se často označuje kompletní instalace T_EXového prostředí. Ta se skládá z nástrojů pro manipulaci s písmem a virtuálními písmem, nástrojů pro generování bitmapových instancí písem, nástrojů pro manipulaci s metrikami písem, ze systému pro popis písem METAFONT, z vlastního programu T_EX, i maker se základními nadstavbami T_EXu – plain T_EX, L^AT_EX či ConT_EXt a z volně šiřitelných písem. V distribucích T_EXu, jako je například teT_EX¹ či T_EXLive², dále bývají zahrnuta i národní rozšíření, pro české dokumenty jsou to například *csplain* a *cslatex*.

T_EX však je nejenom sázecí systém, ale i programovací jazyk. Patří do kategorie makrojazyků a makra jsou základní programovací strukturou T_EXu. Makra v T_EXu začínají znakem `\` a končí před prvním nepísmenným znakem. Závorky `{` a `}` pak ohraničují oblast platnosti lokálních definic, definice a argumenty makra, znak `#` označuje číslo argumentu. Speciální význam mají také znaky `%` a `$` – `%` uvozuje komentář a znak `$` přepíná sazbu do matematického módu. Příklad makra pro nadpis dokumentu pak může vypadat například následovně:

```
\def\nadpis#1{                % definice nadpisu
  \removeataskip\bigskip      % odmaže poslední vert. mezeru a přidá vlastní
  \indent{\titulfont #1}      % odsazený text nadpisu větším fontem
  \par\nobreak\medskip}      % konec řádku, zakázaný zlom, menší mezera
```

Jeho použití v textu pak vypadá takto:

```
\nadpis{Toto je nadpis}
```

Základních příkazů T_EXu, tzv. primitiv, ze kterých se makra skládají, je pouze něco přes 300. Vše ostatní jsou již složená makra, která ale mohou být díky možnostem jazyka velice mocná a umožňují tak poměrně dobře rozlišit v dokumentu obsah a formu.

¹<http://www.tug.org/teTeX>

²<http://www.tug.org/texlive>

Výstupním formátem \TeX u je obvykle `dvi`, ale existují i jiné možnosti – výstup do `pdf` pomocí programu `pdf \TeX` ³ či do HTML pomocí `L \TeX 2HTML`⁴.

1.2 XML

Na rozdíl od \TeX u není XML (eXtensible Markup Language) typografickým systémem, ale jedná se o značkovací jazyk. XML vzniklo v roce 1998 výběrem nejužitečnějších vlastností SGML (Standard Generalized Markup Language) a jeho doplněním o některé prvky. Důležitou vlastností XML převzatou právě z SGML je možnost tvorby vlastních *definic typů dokumentů* – DTD. Důsledkem toho je, že si každý může vytvořit svůj vlastní formát dokumentu zcela vyhovující jeho potřebám.

Základní věci jsou však pro všechny XML dokumenty stejné. Každý XML dokument se skládá z *elementů*, které jsou do sebe navzájem vnořené. Elementy se v textu vyznačují pomocí tzv. *tagů*. Názvy tagů se zapisují mezi znaky `<` a `>`, v případě, že je element párový, tzn. uzavírá nějaký text, má ukončovací tag před svým názvem ještě znak `/`, který musí být v XML uveden i na konci každého nepárového tagu. Každý element může mít dále své *atributy*. Atributy se zapisují dovnitř tagu ve formátu `atribut="hodnota"`. Zápis nadpisu by v XML mohl vypadat třeba takto:

```
<nadpis jazyk="cz">Toto je nadpis</nadpis>
```

Otázkou zůstává, jak bude takovýto element vlastně v konečném dokumentu vypadat, samotné XML tento problém nijak neřeší. K tomuto účelu vždy slouží některý ze stylových jazyků, nejčastěji CSS (Cascading Style Sheets) nebo XSL (eXtensible Stylesheet Language) a příslušný stylový procesor. Jako příklad lze uvést procesor formátovacích objektů FOP⁵ či XSLT procesor Saxon⁶.

Výstupní formát záleží pouze na zvoleném stylovém procesoru a jeho možnostech. Nejčastěji používaným formátem pro tištěný výstup je formát *pdf* ale výstup je obvykle možné získat i v dalších formátech jako RTF, HTML či PostScript.

2 Srovnání v jednotlivých oblastech použití

2.1 Dokumenty všedního dne

Pod tímto poněkud netypickým názvem jsou myšleny klasické články, dopisy, a jiné „kancelářské“ záležitosti. Zde je možností jaký nástroj zvolit k tvorbě takového dokumentu opravdu nepřehledně a to jak v případě \TeX u, tak v případě XML.

V obou případech lze při tvorbě takového dokumentu jít v dané technologii velice do hloubky. V \TeX u to znamená kompletní nadefinování všech maker v `plain \TeX` u, v případě

³<http://www.tug.org/applications/pdftex>

⁴<http://www.latex2html.org>

⁵<http://xml.apache.org/fop>

⁶<http://saxon.sourceforge.net>

XML by pak mohlo jít o návrh vlastního DTD i XLS stylu. Tento přístup je však vhodný pouze pro opravdu rozsáhlá díla a to ještě ne vždy. Navíc je ještě nutné dodat, že pro opravdu profesionální typografii se nehodí ani jeden ze systémů a v současné době se téměř výhradně používají různé proprietární DTP systémy.

Je možné ale zvolit i opačný „extrém“ a jít cestou WYSIWYG (What You See Is What You Get – Co vidíš to dostaneš) editorů, a nutno přiznat, že takto v současnosti vzniká většina dokumentů. \TeX nabízí pod GPL šířený LyX⁷ a TeXmacs⁸, či komerční Scientific Word⁹. Na straně XML jsou „nejtěžším kalibrem“ bezpochyby OpenOffice.org¹⁰, jejichž nativním datovým formátem je právě XML. Výhodou takového přístupu je potřeba minimálních či spíše žádných znalostí \TeX u či XML. Nevýhodou je pak horší kvalita výstupu a prakticky nulové oddělení obsahu dokumentu od jeho formy, což podstatně ztěžuje případné vyhledávání relevantních informací v dokumentu.

Rozumným kompromisem, kdy za určitou snahu získáme téměř profesionální výstup je pak použití některé makronástavby \TeX u, například \LaTeX u. Na straně XML je pak srovnatelným řešením použití DocBooku¹¹. Jak dokumenty v \LaTeX u tak dokumenty v DocBooku jdou psát velmi dobře „ručně“, tzn. bez nutnosti použití programu doplňujícího formátování, pouze v textovém editoru. Následující fragmenty ukazují typickou strukturu dokumentů v DocBooku a \LaTeX u.

DocBook:

```
<?xml version='1.0' encoding='iso-8859-2'?>
<!DOCTYPE book PUBLIC '-//OASIS//DTD DocBook XML V4.2//EN'
'http://www.oasis-open.org/docbook/xml/4.2/docbookx.dtd'>

<article lang="cs">
  <artheader>
    <title>Můj článek</title>
    <author>
      <firstname>Martin</firstname>
      <surname>Tůma</surname>
    </author>
  </artheader>

  <sect1>
    <title>První část</title>
    <para>První odstavec</para>
    <para> ... </para>
```

⁷<http://www.lyx.org>

⁸<http://www.texmacs.org>

⁹<http://www.mackichan.com>

¹⁰<http://www.openoffice.org>

¹¹<http://www.docbook.org>

```
<sect2>
  <title>První podčást </title>
  <para> ... </para>
</sect2>
</sect1>
```

```
<sect1>
  ...
</sect1>
```

```
<bibliography> ... </bibliography>
```

```
</article>
```

L^AT_EX:

```
\documentclass[a4paper, 12pt]{article}
\usepackage{czech}
\usepackage{a4wide}
```

```
\title{Můj článek}
\author{Martin Tůma}
```

```
\begin{document}
\maketitle
```

```
\section{První část}
První odstavec
```

```
...
```

```
\subsection{První podčást}
```

```
...
```

```
\section{ ... }
```

```
\begin{thebibliography}
```

```
...
```

```
\end{thebibliography}
```

```
\end{document}
```

Struktura tedy vypadá v obou případech podobně, jen způsob zápisu je odlišný. Výhoda \LaTeX u spočívá v dalším zpracování takového dokumentu. Takto zapsaný text se jednoduše „prožene“ interpretem a výsledkem je perfektní výstup ve formátu *dvi*. Podrobný postup viz například [1]

U DocBooku je situace o trochu složitější. Před finální transformací do konečné podoby je potřeba ještě, jak již bylo zmíněno, vytvořit odpovídající styl. Naštěstí existují již vytvořené kvalitní styly, které lze rovnou na dokument aplikovat a získat tak požadovaný výstup. Více o problematice aplikace stylů lze nalézt v [4].

Výhoda XML v tomto případě spočívá v rozmanitosti možných výstupních formátů, které lze aplikováním různých stylů z jednoho zdroje získat. Toho lze v omezené míře dosáhnout i u \LaTeX u, ale již ne tak efektivně jako v případě XML – v mnoha případech je nutné upravovat i původní dokument.

2.2 Sazba matematiky a vědecké literatury

Jestliže v případě tvorby klasických dokumentů lze bez větších problémů použít jak \TeX u tak XML, v případě tvorby matematických či odborných publikací s vysokým podílem matematických výrazů, je situace zcela odlišná – zde má \TeX zatím navrch.

XML sice nabízí MathML, ale jeho využití je spíše ve výměně dat mezi matematickými programy. MathML podporuje například Maple¹² či Mathematica¹³. Je však potřeba přiznat, že v poslední době množství software využívající MathML pro sazbu matematiky roste, obsáhlý seznam je možné najít na stránkách organizace W3C¹⁴.

\TeX má však oproti všem řešením založeným na MathML dvě obrovské výhody. Tou první je perfektní kombinace sazby matematiky se sazbou ostatního textu. Druhou, a možná ještě mnohem důležitější výhodou, je pak zápis matematických vzorců v \TeX u. Zápis matematiky je velice efektivní a přirozený a hlavně podobný „ručnímu“ zápisu. Oproti tomu zápis v MathML je velice nepřehledný a většina editorů MathML proto nabízí „klikací“ zápis, tedy stejný druh zápisu jako má například editor rovnic v MS Wordu. Ten je však ve většině případů mnohem zdlouhavější, než přímý zápis ve formě \TeX ového kódu.

Že je tato forma zápisu výhodná pak potvrzuje i to, že jí převzal i kancelářský balík OpenOffice.org, který umožňuje zápis matematiky se syntaxí velice podobnou \TeX u. OpenOffice.org však jednak používají upravené MathML, což výrazně ztěžuje přenositelnost takto vytvořených vzorců, a navíc je typografická kvalita matematického výstupu velice nízká.

2.3 Elektronické publikování

V dnešní době je stále více potřeba vytvořený dokument exportovat do několika různých formátů, aby jej bylo možné nejenom vytisknout, ale také publikovat na internetu. Te-

¹²<http://www.maplesoft.com>

¹³<http://www.wolfram.com>

¹⁴<http://www.w3.org/Math/implementations.html>

oreticky je sice možné dokument na internetu vystavit ve formátu PDF, chceme-li však dokument zpřístupnit opravdu „on-line“, je potřeba zvolit HTML, či ještě lépe XHTML.

Pokud pro tvorbu dokumentu použijeme \TeX , je potřeba pro převod do HTML použít nějaký externí převaděč, ostatně stejně jako pro převod do jakéhokoliv jiného formátu než dvi. Vyčerpávající seznam HTML/XML převaděčů lze nalézt na internetové stránce:

<http://www.cse.ohio-state.edu/~gurari/TeX4ht/mn.html>, nejznámější a nejpoužívanější je však již zmíněný \LaTeX2HTML pro \LaTeX . Za zmínku stojí $\text{Html}\text{\TeX}$ ¹⁵, jehož autorem je Sergey Brin, jeden ze zakladatelů Googlu. Při použití \LaTeX2HTML získáme poměrně dobrý HTML výstup, i když pravda, standartní vzhled se mě osobně příliš nelíbí. Dokument vyhovuje zvolenému DTD, které však lze zvolit pouze z množiny HTML, na modernější XHTML \LaTeX2HTML transformovat zatím neumí. Přetransformovat lze velká většina \LaTeX ových dokumentů, problémy jsou hlavně s matematikou a vlastními makry, přičemž u matematiky to není vinou \LaTeX2HTML , ale bídou podporou MathML v současných prohlížečích.

U XML je situace poněkud odlišná – XHTML je podmnožinou XML, takže dokumenty lze teoreticky tvořit rovnou v HTML. Toto však není dobrý nápad, pokud chceme dokumenty i tisknout. Kaskádové styly používané pro HTML na kvalitní tiskový výstup většinou nestačí. Pokud tedy dokument netvoříme pouze pro web, je nejefektivnější cestou zvolit obecnější XML formát, například DocBook, a pomocí XSL či DSSSL stylů jej konvertovat na HTML/XHTML.

XML má však v oblasti elektronického publikování jednu obrovskou výhodu – může sloužit i jako univerzální datový formát. To znamená, že z XML dokumentu lze velice účinně extrahovat potřebné informace. Existuje dokonce „dotazovací“ jazyk pro XML – XQuery¹⁶. A právě možnost efektivně pracovat s daty v XML dokumentech je jednou z nejdůležitějších vlastností XML.

3 Závěr

V předcházejících částech byly naznačeny možnosti obou technologií v jednotlivých oblastech. Kde je tedy výhodné použít \TeX a kde XML?

\TeX je bezesporu vynikající na sazbu matematiky a tím pádem také vědeckých článků. Zde je jeho pozice stále prakticky neotřesitelná. Velmi dobrým nástrojem je \TeX také na psaní článků a referátů v případě, že výsledný dokument bude určen primárně pro tisk. Zde je však nutno říci, že \TeX em je zde myšlen především \LaTeX , použití „čistého“ \TeX u se vyplatí až u rozsáhlejších děl, nebo pokud vyžadujeme specifický výstup.

Problémem \TeX u může být dostupnost dokumentace. Poměrně dobře je zdokumentován \LaTeX , viz například [2] či *\LaTeX navigátor* na stránkách [cstugu](http://www.cstug.cz)¹⁷. Horší to již ale je s dokumentací pro plain \TeX . Zde dle mého názoru chybí středně odborná literatura. K dispozici jsou buď publikace pro úplné začátečníky, nebo velmi odborné publikace (\TeX book,

¹⁵<http://www-db.stanford.edu/~sergey/htmltex>

¹⁶<http://www.w3.org/TR/xquery>

¹⁷<http://www.cstug.cz/latex/lm/frames.html>

TeXbook naruby). Naopak často zmiňovaná složitá instalace TeXu je dle mého názoru mýtus, tedy alespoň v případě LINUXu. Prakticky všechny moderní distribuce obsahují TeX připravený k okamžitému použití a to většinou i s českými rozšířeními. Pro Windows pak existuje MikTeX¹⁸, který nabízí velice dobře zpracovaný instalátor. Složitější situace je s rozšířeními TeXu jako L^ATeX2HTML, které potřebují nějaký interpret, v tomto případě Perl, což na Windows představuje další složitou instalaci.

Pokud dokument potřebujeme ve více výstupních formátech s důrazem na elektronické publikování je výhodnější zvolit XML. Instalace potřebného software (XML parser, XSLT procesor, . . .) sice také není triviální záležitost, je však společná pro všechny výstupní formáty. Volba XML navíc přináší oproti TeXu mnohem lepší možnost zpracování dat v dokumentech. Jako výhodu XML je ještě nutné zmínit nativní podporu UNICODE kódování, které TeX vůbec nepodporuje a v neposlední řadě kvalitní dostupnou dokumentaci, především díky konsorciu W3C.

V současnosti osobně využívám jak TeX(L^ATeX), tak XML, měl-li bych se ale rozhodnout pouze pro jedno z nich, zvolil bych s výhledem do budoucnosti mocně se rozvíjející XML.

¹⁸<http://www.miktex.org>

Literatura

- [1] Olšák: *První setkání s T_EXem*
<ftp://math.feld.cvut.cz/pub/cstex/doc/prvni.pdf>
- [2] Oetiker, Partl, Hyna, Schlegl: *The Not So Short Introduction to L^AT_EX2_ε*
www.ctan.org/tex-archive/info/lshort/english/lshort.pdf
- [3] Kosek: *XML pro každého* Grada 2000
- [4] Kosek: *DocBook, stručný úvod do tvorby a zpracování dokumentů*
<http://www.kosek.cz/xml/db>
- [5] Diskuzní fórum `cz.comp.cstex`
<http://usenet.jyxo.cz/cz.comp.cstex>