

Zpravodaj Československého sdružení uživatelů TeXu

Marcel Svitalský

Google Summer of Code 2009 a TUG

Zpravodaj Československého sdružení uživatelů TeXu, Vol. 19 (2009), No. 1-2, 94–101

Persistent URL: <http://dml.cz/dmlcz/150036>

Terms of use:

© Československé sdružení uživatelů TeXu, 2009

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>

Abstrakt

Zpráva přináší překlad webové stránky s idejemi TUGu [9] pro Google Summer of Code 2009 [7].

Klíčová slova: Google, Google Summer of Code 2009, TUG, T_EX Users Group.

Ročník 2009

TUG se letos nezúčastní programu Google Summer of Code [7] v roli poradní organizace [16]; naše žádost byla pro letošní rok zamítnuta, avšak kterýkoli zájemce si může vyhledat seznam ideí, jež jsme společně vytvořili, a nechat se jím inspirovat. V T_EXovém světě zbývá stále ještě víc než dost vývojářské práce, již je třeba udělat.

V roce 2008 se pod záštitou TUGu zúčastnili GSoC tři studenti. Jejich projekty [18] a kódy [19], které vytvořili, stejně jako oznámení TUGu pro rok 2008 [8], si můžete v angličtině pročíst a prohlédnout na Internetu.

Všechny organizační diskuse TUGu týkající se GSoC se odehrávají na mailing listu `summer-of-code@tug.org` [10]; bez obav se do něj můžete přihlásit nebo jeho archivy [21] pročíst.

Nápady a myšlenky projektů byly:

- Přístupné PDF s pomocí T_EXu 94
- Dublin Core metadata a T_EX 96
- Rozpoznání rukou psaných L^AT_EXových symbolů 96
- Zvýraznění syntaxe s hyperlinkem pro T_EX 97
- Nové šablony dokumentů pro L^AT_EX 98
- Implementace `fontspec` pro LuaT_EX 98
- Mikrokernl pro L^AT_EX3 100

1. Přístupné PDF s pomocí T_EXu

Cílem tohoto projektu je obohacení pdfT_EXu o schopnost produkovat označované PDF ve shodě se specifikacemi PDF/A, PDF/UA, ISO 32000 a (dosud nepublikovanou) ISO-32000-2. To zahrnuje jak strukturní, tak obsahové značkování pro užití se čtečkami obrazovky (pro zrakově postižené) a softwarovými pluginy schopnými zobrazit a vylepšit strukturu i obsah matematických vzorců.

Úspěšné završení tohoto projektu musí zahrnovat následující úlohy:

- a) Vyhledejte všechna místa ve standardním L^AT_EXovém kódu, kde je tagování struktury nebo obsahu vhodné.
- b) Vytvořte soubor(y) obsahující alternativní definice maker zahrnující programové „přípojky“ umožňující umístění tagů na patřičná místa ve výstupním streamu, tj. do T_EXového „vertikálního seznamu“.
- c) Vyhledejte ta místa, kde upravená makra nebudou postačující ke vložení strukturních nebo obsahových tagů optimálním způsobem. To bude indikovat potřebu přidat pod nimi běžícímu procesoru dodatečné schopnosti.
- d) Dodejte soubor definic (La)T_EXových maker, poskytujících všechny strukturní, obsahové a matematické tagy, podporované již ve specifikaci PDF 1.7 (nebo vyšší) a jež budou podporovány normou ISO-32000-2.
- e) Dodejte „podpurné“ soubory umožňující propojení „přípojek“ z bodu b) s definicemi maker z bodu d). Takových podpurných souborů může být zapotřebí více, aby poskytovaly různé úrovně podpory vyhovující požadavkům různých úrovní tagování, kupř. pro PDF/X, PDF/A (různých verzí), PDF/UA a pro plné tagování matematických vzorců.
- f) Dodejte soubory/ovladače, řídicí způsob, jímž „přípojky“ z bodu e) vkládají tagy do výstupu vytvářeného určitým procesorem. Zpočátku bude třeba podporovat pouze jeden procesor, např. pdfT_EX, avšak kód by měl být natolik modulární, aby podporu dalších procesorů, např. X_YT_EXu a LuaT_EXu, bylo možné zajistit pouhou výměnou souboru/ovladače. Všechny další moduly by mělo být možné užít bez jakýchkoli úprav.
- g) Vytvořte patřičnou dokumentaci popisující, co mají jednotlivé moduly dělat a jak jich užívat.

Naplnění těchto zadání bude společným úsilím zahrnujícím všechny účastníky projektu, vizte níže, nikoliv pouze studenta. Student by měl nepochybně být dobře obeznámen s postupy programování T_EXových maker, jakož i se způsoby, jimiž jsou tato užívána ve vnitřních makrech L^AT_EXu.

Znalost dalších aspektů programování, jako PDF/PostScript, Tangle & Weave či programování pro pluginy Adobe, bude považována za bonus, který se může ukázat užitečným, není však nezbytným. Diskuse k této problematice hostí T_EXový mailing list [11] v River Valley Technologies.

Personální obsazení (poradci):

- * Hàn Thé Thành, River Valley Technologies, Německo. Vývojář pdfT_EXu.
- * Ross Moore, Senior Lecturer in Mathematics, Macquarie University, Sydney, Australia. T_EXové/L^AT_EXové programování, předsednictvo TUGu.
- * Neil Soiffer, Senior Scientist, Design Science Inc., California and St. Paul, Minnesota, USA. Plug-in software k podpoře matematického obsahu.
- * CV Radhakrishnan, River Valley Technologies, India, Trivandrum, Kerala, India. T_EXové/L^AT_EXové programování, publikování vědeckých prací.

2. Dublin Core metadata a \TeX

Celý text návrhu [5] je dostupný na Internetu a v samostatném překladu v tomto Zpravodaji, zde uvádíme shrnutí. S projektem Dublin Core se lze v češtině podrobně seznámit na stránkách [4].

Dublin Core Metadata Initiative je otevřená organizace zapojená do vývoje standardů pro interoperabilní online metadata podporujících široké rozpětí účelů a aplikačních modelů. Vyvinula abstraktní framework pro metadata a několik strojově čitelných reprezentací metadatových výrazů, mimo jiné i ty v Resource Description Frameworku (RDF).

Jedním z významných uživatelů RDF metadat je Adobe, tvůrce formátu PDF. Jejich eXtensible Metadata Platform (XMP) umožňuje tvůrcům PDF dokumentů vkládat do PDF libovolná metadata. Ta jsou viditelná jak aplikacím Adobe, tak i rostoucímu množství dalších vyhledávacích a archivačních nástrojů, včetně Spotlight z Mac OS X. XMP je implementováno v XML reprezentaci RDF.

Klíčovými výstupy tohoto projektu by měly být:

1. \TeX ová implementace Dublin Core Abstract modelu;
2. metody pro export metadat z abstraktního modelu do vnějších souborů v různých formátech, zejména RDF+XML, popř. také DC-TEXT a N3;
3. automatické vkládání XMP paketů do vytvářeného PDF souboru v případě pdf \LaTeX u, s defaultním minimem XMP výrazů Z39.88 OpenURL COinS polí, jak pro vlastní metadata dokumentu, tak i pro všechny citované odkazy a vnější hyperlinky;
4. uživatelsky přívětivé rozhraní pro vytváření metadatových výrazů;
5. chybějí-li autorské deklarace v pdf \LaTeX ovém dokumentu, měla by být vložena všechna metadata, jež lze detekovat automaticky;
6. metody pro autory balíků, umožňující deklarovat nové množiny metadatových elementů a slovníků, aby autoři mohli zapisovat metadata příslušné oblasti jejich zájmu. Osobně uvažuji o Learning Object Metadata, avšak mapování LOM do Dublin Core je problematické.

Poradci projektu budou Peter Flynn a Matthew Leingang.

3. Rozpoznání rukou psaných \LaTeX ových symbolů

Tento projekt dosud nemá přiřazeného poradce. Kdokoliv, kdo by měl zájem jako poradce působit, je srdečně zván ke kontaktu prostřednictvím mailing listu [10].

Sázecí systém \LaTeX [15] poskytuje příkazy pro sazbu tisíců rozličných symbolů potřebných pro přípravu dokumentů z oblastí lingvistiky, matematiky, hudby, techniky, fyziky a mnoha dalších. Najít \LaTeX ový název pro daný glyf může být

pro tvůrce dokumentu obtížným úkolem. V současnosti je nejlepším řešením užít Comprehensive L^AT_EX Symbol List [3] (úplný seznam L^AT_EXových symbolů), souhrn tabulek symbolů organizovaných do ad hoc kategorií a indexovaných dle L^AT_EXového názvu symbolu. Problémem tohoto přístupu je, že různí uživatelé spojují s týmž glyfem odlišné názvy, čímž se vyhledávání stává obtížným.

Uvažte pro příklad, že se pokoušíte vyhledat L^AT_EXový název pro kroužek s tečkou uprostřed. Astronom bude možná hledat „slunce“, matematik „mez“, lingvista „mlaskavou souhlásku“, tvůrce map „městské centrum“, někdo písíci o alchymii „zlato“. Ve skutečnosti je výpisu významů tohoto symbolu věnována celá stránka Wikipedie [1]. Neanglicky mluvící budou dále znevýhodněni tím, že většina L^AT_EXových symbolů byla pojmenována anglicky.

Domníváme se, že uživatelům L^AT_EXu by byl velkou nápomocí webový nástroj vyhledávající symboly založený na rozpoznávání textu. To jest, představujeme si webovou stránku, na níž by uživatel mohl nakreslit svůj symbol a následně by mu byl zobrazen seznam L^AT_EXových symbolů (příkazů a jejich výstupů) nejvíce odpovídajících uživateli kresbě.

Student by musel vyhodnotit četné možnosti při rozpoznávání ručně vykreslených symbolů, najít vhodnou vnitřní reprezentaci pro tisíce L^AT_EXových symbolů a asociovaná metadata, vytvořit vhodné uživatelské rozhraní pro práci s rozpoznávacím symbolů a zajistit, že výsledný software bude spravovatelný, zejména s přihlédnutím k četnosti přidávání nových symbolů do L^AT_EXu.

Taková idea je nepochybně na implementaci náročná. Zároveň to však bude zcela jistě vzrušující a prospěšná zkušenost, vzhledem k přehrášli technologií, jichž se bude týkat: T_EX/L^AT_EX, rozpoznávání textu, rozličné webové technologie, a patrně řada programovacích jazyků. Lze se zde mnoho naučit a TUG velmi rád bude pomáhat studentovi se zájmem a schopností uskutečnit projekt převodu rukopisných symbolů do L^AT_EXu.

Návrh podal Scott Pakin, je však třeba poradce.

4. Zvýraznění syntaxe s hyperlinkem pro T_EX

Nyní je běžné, že zobrazuje-li se kód na webové stránce, je užito zvýraznění syntaxe. Děje se tak např. na stránce Google Code [6].

Tento projekt má poskytnout zvýraznění syntaxe pro T_EXový kód – jak pro dokumenty, tak pro makra – s další přidanou schopností. Každý zvýrazněný příkaz bude zároveň hyperlinkem nabízejícím bublinovou nápovědu a, pokud naň bude kliknuto, vedoucím k další dokumentaci.

Dvěma z nejlepších zvýrazňovačů syntaxe jsou:

- Pygments [20] v jazyku Python, a
- Chili [17] pro JavaScript, založeném na jQuery.

Tento projekt sestává ze tří částí. Prvou je poskytnutí pokročilého zvýrazňování syntaxe pro $\text{T}_{\text{E}}\text{X}$ ový kód. Druhou je vytvoření databáze příkazů. Třetí je vzájemné spojení prvé a druhé části.

V závislosti na obtížích, jež se vyskytnou, se tento projekt může ukázat jako příliš rozsáhlý. Pokud by to nastalo, předpokládáme, že student uskuteční pouze jeho část.

Poradcem projektu bude Jonathan Fine.

5. Nové šablony dokumentů pro $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$

Lamport vytvořil pro $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ řadu šablon dokumentů, jako jsou třídy `book`, `article` atd. Dokonce i dnes užívá velké množství $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ ových dokumentů tyto šablony, což ústí ve zřetelný „ $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ ový vzhled“.

Autorům jsou rovněž k dispozici třídy pro specializovaná užití, např. pro časopisy, konkrétní konference nebo diplomové práce té či oné univerzity. Avšak tyto třídy jsou velmi často jen uzpůsobenými Lamportovými šablonami a udržují si svůj vzhled. Z opačné strany vzato by bylo vhodnou součástí tohoto projektu prozkoumání existujících tříd poskytujících výrazně odlišný vzhled, jako jsou koma-scripty `scr*` [14] či `smfart` [2] Francouzské matematické společnosti.

Tento projekt by měl poskytnout alternativní šablony pro široké užití. Tyto šablony mohou být vhodné i pro knihy a články, popř. pro další účely. Měly by pokud možno být doprovázeny návodem pro možné uživatele s ohledem na jejich vhodné použití, např. pro knihy převážně bez matematiky či pro automaticky generované texty.

Poradcem projektu bude Jim Hefferon.

6. Implementace `fontspec` pro $\text{LuaT}_{\text{E}}\text{X}$

$\text{X}_{\text{E}}\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ je populárním rozšířením $\text{T}_{\text{E}}\text{X}$ u umožňujícím užití TrueType a OpenType písem, aniž by musely být upraveny do Type 1 nebo TFM souborů. Za část svého úspěchu vděčí balíku `fontspec` Willa Robertsona, poskytující $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ ové rozhraní primitivům $\text{X}_{\text{E}}\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ u pro nahrání písem. Cílem tohoto projektu je implementovat balík `fontspec` pro $\text{LuaT}_{\text{E}}\text{X}$.

Jelikož $\text{X}_{\text{E}}\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ a $\text{LuaT}_{\text{E}}\text{X}$ představují velmi odlišná paradigmaty v přístupu k rozšiřování $\text{T}_{\text{E}}\text{X}$ u, $\text{LuaT}_{\text{E}}\text{X}$ ová část `fontspec` bude výrazně odlišná od nyní existujícího kódu pro $\text{X}_{\text{E}}\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$: zatímco tento využívá řady dodatečných systémových knihoven, aby $\text{T}_{\text{E}}\text{X}$ u dodal pokročilou podporu písem, $\text{LuaT}_{\text{E}}\text{X}$ poskytuje spojení mezi $\text{T}_{\text{E}}\text{X}$ em a jazykem Lua a umožňuje zapojit kód Lua do $\text{T}_{\text{E}}\text{X}$ ového sázecího procesoru. Tudíž převážná většina práce bude spočívat v implementaci pokročilých rysů písem v Lua. Tedy toho, co již dělá `CONTEX`T verze „Mark IV“.

Je třeba věnovat pozornost těmto aspektům:

Vyhledávání písem: Aby LuaTeX dokázal užívat systémová písma stejně jako XeTeX, je třeba pro ně implementovat vyhledávací mechanismus. Nabízí se několik možností:

- XeTeX je slinkován s knihovnou `fontconfig`, což ale není případem LuaTeXu; avšak na systémech, kde je `fontconfig` k dispozici také jako samostatný program, jej lze spustit a použít jeho výstup.
- Mohlo by být možné přeložit `fontconfig` jako sdílenou knihovnu a tu užívat z LuaTeXu. Lua totiž disponuje vydařeným rozhraním pro jazyk C umožňujícím mu přímo užívat knihovny tohoto jazyka. Zřejmým problémem takového přístupu by však byly značné potíže s přenositelností mezi různými operačními systémy.
- Bylo by také možné implementovat vlastní jednoduchý vyhledávací mechanismus, prohledávající adresáře standardní pro jednotlivé platformy, konkrétně C:\Windows\Fonts pod Microsoft Windows, /System/Library/Fonts, /Library/Fonts a ~/Library/Fonts v Mac OS, /usr/share/fonts v Linuxu atd. Takto to dělá ConTeXt Mark IV: spouští skript kešující všechny soubory a jména písem.

Tyto schopnosti by samozřejmě byly doplněním knihovny `kpathsea`, kterou lze také volat z Lua. Měli bychom tudíž duální situaci velmi podobnou XeTeXu.

Nahrávání písma: LuaTeXový primitiv `\font` se na úrovni procesoru chová přibližně stejně jako týž primitiv v pdfTeXu, a nemůže tedy nahrát soubory TrueType nebo OpenType písem. Očekává defaultní typ souboru TFM. Potřebujeme proto „přetížít“ tento primitiv tak, aby emuloval chování stejné jako v XeTeXu.

Rozložení OpenType: I zde je několik možností:

- LuaTeX disponuje nástroji k implementaci kompletního procesoru rozložení OpenType. Vyžadovalo by to značné úsilí, avšak základ kódu je již nyní dostupný pro ConTeXt a bylo by možné ho s výhodou využít.
- XeTeX používá externí knihovny jako je International Components for Unicode (ICU) [12] od IBM. Podobně jako ve druhé možnosti prvního bodu bychom mohli zkompileovat sdílenou knihovnu obsahující wrapper v jazyce C poskytující rozhraní pro Lua.

LaTeXové rozhraní: V ideálním případě by LuaTeXová implementace `fontspec` měla mít stejné vysokoúrovňové rozhraní jako současný balík.

Poradci projektu budou Will Robertson a Arthur Reutenauer.

7. Výchozí mikrokernel \LaTeX 3

Sázecí systém \LaTeX po mnoho let znamenal $\text{\LaTeX} 2_{\epsilon}$. Poslední vývoj následníka, $\text{\LaTeX} 3$ [13], byl soustředěn převážně na nový nízkourovňový programovací systém pro \TeX . Jelikož tato nízkourovňová práce již dozrává, lze zkoušet nové programovací postupy pro práci na vyšší úrovni.

Cílem projektu mikrokernelu $\text{\LaTeX} 3$ je započít zkoumání, jak lze uplatnit nízkourovňový systém k vytvoření systému, který bude možno použít k vysázení jednoduchého dokumentu ve stylu $\text{\LaTeX} 2_{\epsilon}$, aniž by bylo zapotřebí nahrát návrh současného $\text{\LaTeX} 2_{\epsilon}$ kernel. Jako úvodní cíl bude použit k otestování základní dokument:

```
\documentclass{minimal}
\begin{document}
\emph{Hello World!}
\end{document}
```

Zde popsaný mikrokernel nemá být úplnou implementací kernelu $\text{\LaTeX} 2_{\epsilon}$, tj. `latex.ltx`. Existuje řada dosud nezodpovězených otázek ohledně uživatelského rozhraní $\text{\LaTeX} 3$. Cílem projektu je vytvořit základní samostatný kernel, jenž poté bude možno pomalu rozšiřovat implementováním dalších vlastností. Nejpravděpodobněji z `latex.ltx`, možná i z přídatných \LaTeX ových balíků.

Dalším logickým krokem po vytvoření systému schopného pracovat s testovacím dokumentem budou dodatky jako základní strukturní příkazy a prostředí (seznamy, zarovnání atp.). Jisté oblasti mohou být také vyloučeny: New Font Selection Scheme, komplexní výstupní rutiny či plovoucí obsah jsou všechny mimo rámec tohoto projektu.

Poradcem projektu bude Joseph Wright.

Seznam literatury [on-line 14. 5. 2009]

- [1] Circled dot – Wikipedia, the free encyclopedia [on-line]. URL: http://en.wikipedia.org/wiki/Circle_with_a_point_at_its_centre
- [2] Classes for Société mathématique de France publications: smflatex package. <http://www.ctan.org/tex-archive/macros/latex/contrib/smflatex/>
- [3] CTAN: The Comprehensive \LaTeX Symbol List [on-line]. URL: <http://www.ctan.org/tex-archive/info/symbols/comprehensive/>
- [4] Dublin Core – Czech homepage [on-line]. URL: http://www.ics.muni.cz/dublin_core/
- [5] Dublin Core metadata project [on-line]. URL: <http://tug.org/gsoc/dublincore.html>
- [6] Google Code [on-line]. URL: <http://code.google.com/>

- [7] Google Open Source Programs [on-line].
URL: <http://socghop.appspot.com/>
- [8] Google Summer of Code 2008 and TUG [on-line].
URL: <http://tug.org/gsoc/2008.html>
- [9] Google Summer of Code and TUG [on-line].
URL: <http://tug.org/gsoc/>
- [10] Google Summer of Code discussions for TUG [on-line].
URL: <http://lists.tug.org/summer-of-code>
- [11] T_EX Info Page [on-line]. URL:
<http://lists.river-valley.com/cgi-bin/mailman/listinfo/tex>
- [12] International Components for Unicode [on-line]. Home Page.
URL: <http://www.icu-project.org/>
- [13] L^AT_EX project: The L^AT_EX3 project [on-line].
URL: <http://www.latex-project.org/latex3.html>
- [14] KOMA-Script: A versatile bundle of L^AT_EX document classes and packages.
<http://ctan.org/tex-archive/macros/latex/contrib/koma-script/>
- [15] L^AT_EX project: L^AT_EX – A document preparation system [on-line].
URL: <http://www.latex-project.org/>
- [16] List participating Organizations [on-line]. URL: http://socghop.appspot.com/program/accepted_orgs/google/gsoc2009
- [17] Notes Log >> Chili [on-line].
URL: <http://noteslog.com/category/chili/>
- [18] Organization Information: T_EX Users Group [on-line].
URL: <http://code.google.com/soc/2008/tex/about.html>
- [19] Project Home: Code samples from students working with T_EX Users Group for Google Summer of Code 2008 [on-line].
URL: <http://code.google.com/p/google-summer-of-code-2008-tex/>
- [20] Pygments – Python syntax highlighter [on-line].
URL: <http://pygments.org/>
- [21] The summer-of-code Archives [on-line].
URL: <http://tug.org/pipermail/summer-of-code/>

Summary: Google Summer of Code 2009 and TUG

The report is a Czech translation of the web page presenting ideas of the T_EX Users Group for Google Summer of Code 2009.

Key words: Google, Google Summer of Code 2009, TUG, T_EX Users Group.

Přeložil: Marcel Svitalský
marcel.svitalsky@centrum.cz