

upTeX – Unicode version of pTeX with CJK extensions

Takuji Tanaka 田中 琢爾

upTeX project

Oct 26, 2013

Outline / 概要

- (1) Introduction
- (2) Unicodization / Unicode 化
 - ▶ Japanese / 日本語
 - ▶ CJK / 中韓 / 中・日・
 - ▶ with European languages / 欧文との親和性
 - ▶ world languages / 世界の言語
- (3) Imprementation / 実装
 - ▶ Unicodization / Unicode 化
 - ▶ `\kcatcode`
 - ▶ `set3`
- (4) \upTeX vs. Ω , $X_{\exists}\TeX$, ...
- (5) Present & future / 現在と今後

Part I

Introduction

ASCII pTeX/pLaTeX

It's **great**:

High quality Japanese typesetting

incl. vertical writing, Japanese hyphenation, . . .

Japanese **standard** T_EX/L_AT_EX

Strong support by environment

—DVIware, packages, macros, softwares, books, . . .

but has **weakness**:

Japanese local

— 8bit Latin/Chinese/Korean are not available

Limited character set

by legacy encodings (Shift_JIS, EUC-JP)

Motivation

Support **wider character set** of Japanese
by Unicode

Support **babel**
by switching Latin–CJK tokens

Support **Chinese/Korean**

Keep quality & environment of pT_EX

Feature of upTeX/upL^AT_EX

- (1) **High quality** CJK typesetting
based on pTeX/pL^AT_EX
- (2) **Compatible** with pTeX/pL^AT_EX
- (3) **Unicode / UTF-8**
- (4) Switching Latin (12bit) / **CJK (29bit) tokens**
- (5) **CJK with Babel** (Latin/Cyrillic/Greek. . .)
- (6) **Over BMP** — incl. SIP (U+2xxxx)

Part II

Unicodization / Unicode 化

Unicodization / Unicode 化

Strategies of Unicodization

- (1) Unicodize only IO
Ex: `\usepackage[utf8]{inputenc}`
- (2) Imprement Unicode functions
Ex: `X3TEX`
- (3) Comromise
up \TeX : Intenal: Unicodize only CJK,
IO: Fully Unicodize

Partial Unicodization / 折衷的 Unicode 化

		TEX	pTEX	upTEX
Latin	7bit Latin	azAZ	azAZ	azAZ
	8bit Latin	æœÆŒ		æœÆŒ
	inputenc	гдГД		гдГД
Japanese	JIS X 0208		あア亜	あア亜
	Unicode			高
CK	Unicode			汉字 漢字

pTEX, upTEX consists of two parts

- (1) As same as original TEX
- (2) pTeX–JIS X 0208, upTeX–Unicode

New JIS : JIS X 0213

upTeX treats new JIS X 0213 (over JIS X 0208)

焔 燂焮

づかけ

燿燒

滯濃濘燻燹(株)(有)

燻燈燉

鄧小平 李承燁 里見弴 草彌剛 朴璐美 森鷗外 森雞二

王銘琬 宮崎あおい 蔣介石 你好 深圳 東日本旅

客鐵道株式会社 尾骿骨 生酖仕込 夙月堂 屯寿 今寿

圓塙函數

啞然 火焰 嚙む 任俠 長身瘦軀 石鹼 屢 刺繡 醬油

蟬時雨 隔靴搔痒 奧飛驒 箆笥 搦む 充填 顛末 祈禱

瀆職 土囊 潑刺 醱酵 頰紅 素麵 麴町 蓬萊 蠟燭 攢竹

Characters out of JIS / JIS 外字

over JIS X 0213 (new JIS)

高島屋、内田百間、
柿落とし、安全ホ一、吉野家

source

高島屋、内田百間、柿落
とし、安全ホ一、吉野家

output

Platform dependent characters are now in Unicode

ミ、キ、センメーグ、ラト、アーヘクリツワッカロド、センパーミリペー
リ、ロチ、トルム、ンル、タルトルト、リー、ルト、セントバルジ
mmcmkmmgkgccm²穢 “ „ K.K. ①②③④⑤(株)(有)(代)明治大正昭和
高間塚徳豊崎彌淳燁珉鄧

Chinese/Japanese/Korean

中 · 日 ·

```
\schrm 简体中文: 你好  
\tchrm 繁體中文: 早晨  
\jpnrm 日本語: こんにちは  
\korrm :
```

source

```
简体中文: 你好  
繁體中文: 早晨  
日本語: こんにちは  
:
```

output

Difference of glyphs among CJK / CJKのグリフの違い

Simplified Chinese	骨練	平直。	神祀	才次	.
Traditional Chinese	骨練	平直。	神祀	才次	.
Japanese	骨練	平直。	神祀	才次	.
Korean	骨練	平直。	神祀	才次	.

end-of-line

Please give↓
me beer.

请给我↓
啤酒。

ビールを私に↓
下さい。



.

Please give me beer.
(treated as space)

请给我啤酒。
(ignored)

ビールを私に下さい。
(ignored)



.

(treated as space)

Control word by CJK characters

```
\def\    {%  
\number\year    %  
\number\month    %  
\number\day    %  
}  
Today: 《\    》
```

Today: 《2013 10 26
》

Japanese-OTF package

```
\usepackage[uplatex,...]{otf}
```

```
...
```

```
Adobe-Korea1-1:\\
\CIDK{8322}\CIDK{8588}
```

```
...
```

```
Adobe-Japan1-5:\\
\問\答\ajRecycle{10}%
\ajLig{学校法人}%
\ajPICT{野球}\\
\ajMaru{1}...
```

```
Adobe-Korea1-1:
  ㅁ ㅂㅅㅈㅊ
```

```
Adobe-Japan1-5:
  問答爬慶
    3 4    (七)
```

Japanese-OTF package also supports CK.

Unification / 統合

	standard	full-width
Cyrillic	Ⓚ U+0416	Ⓚ U+0416
Latin	Ⓜ U+0057	Ⓜ U+FF37

No “full-width” code in Greek, Cyrillic in Unicode.
 It is a barrier to Unicodize Japanese softs.
 upT_EX can treat full-width Greek, Cyrillic by markup.

inputenc & UTF-8

```
\usepackage[utf8]{inputenc}
\usepackage[T1]{fontenc}
\kcatcode'ç=15
```

...

“¿But aren't Kafka's Schloß and Æsop's Œuvres often naïve vis-à-vis the dæmonic phoenix's official rôle in fluffy soufflés?”

“¿But aren't Kafka's Schloß and Æsop's Œuvres often naïve vis-à-vis the dæmonic phoenix's official rôle in fluffy soufflés?”

Babel

```
\usepackage[french,...]%  
{babel}  
...  
\selectlanguage{english}  
English ... \today  
...  
\selectlanguage{russian}  
Русский ... \today  
...  
\selectlanguage{japanese}  
日本語 ... \today
```

English
October 26, 2013

Français
26 octobre 2013

Deutsch
26. Oktober 2013

Czech
26. října 2013

Русский
26 октября 2013 г.

日本語
2013年10月26日

It's a small world

up $\text{T}_{\text{E}}\text{X}$ can treat CJK, Latin, Cyrillic and Greek.
up $\text{T}_{\text{E}}\text{X}$ cannot directly treat Arabic, Brahmic, . . .

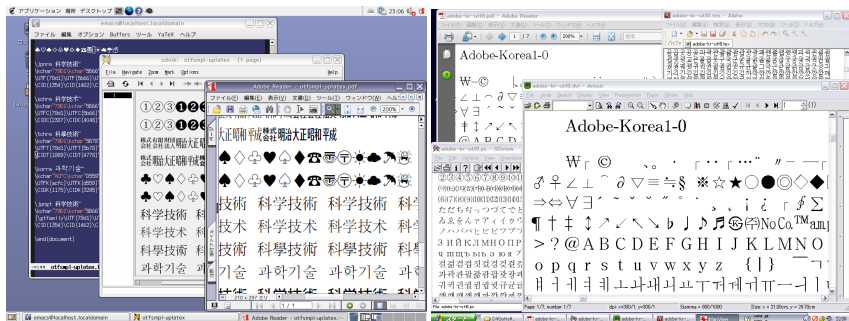
Part III

Implementation / 実装

Unicodization / Unicode 化

- (1) IO: EUC/SJIS in p $\text{T}_{\text{E}}\text{X}$ UTF8 in up $\text{T}_{\text{E}}\text{X}$
(ptexenc library)
- (2) Internal buffer: 16bit in p $\text{T}_{\text{E}}\text{X}$ 29bit in up $\text{T}_{\text{E}}\text{X}$
(Ref. Omega)
- (3) Unicodize standard macros, libraries
- (4) up $\text{T}_{\text{E}}\text{X}$ support of DVIWARE

DVIware



ptetex3+ / Linux

W32TeX / Windows

dvipdfmx, dvips, xdvi, dvi2tty & DVIOUT are available

\kcatcode

kcat code	cat code	kind	e.g.	control word	end of line
			
	10	space	␣		
15	11	char	azAZ	yes	as space
	12	other char	(.!?	no	as space
			
16		Kanji	汉漢	yes	ignore
17		Kana	かナ	yes	ignore
18		CJK symbol	《 ． ｡ 』	no	ignore
19		Hangul		yes	as space

If \kcatcode is 15, the character is treat as Latin and upT_EX works as same as original T_EX.

Part IV

upTeX vs. XeTeX, ...

upTeX vs. XeTeX, ...

		TeX	pTeX	upTeX	XeTeX
Compatibility	Latin				
	Japanese	—			×
Advancedness		×	×	×	×
Multilingual	Latin				
	Japanese	—			
	CK	—	—		
	others	—	—	—	
Integrity	(Japanese)				
Popularity	Japan				
	World				

> > > ×

Part V

Present & Future / 現在と今後

History

Year	
1995	ASCII pTeX ver.2, pLaTeX2e
2007	upTeX first release, alpha version
2007	upTeX is in W32TeX
2008	e-upTeX by Kitagawa-san
2012	upTeX 1.00
2012	upTeX is in TeX Live
2013	upTeX presentation in TUG2013

Future / 今後

Currently, up $\text{T}_{\text{E}}\text{X}$ has capability of multilingual (CJK, Latin, Cyrillic, Greek) typesetting.

Possible items in the future are:

- (1) **Document classes** for Chinese/Korean
(Any volunteer?)
- (2) **Babel options** for Chinese/Korean
(It will be useful in ko.TeX etc. Any volunteer?)
- (3) Does up $\text{T}_{\text{E}}\text{X}$ have a potential
to be a **useful CJK $\text{T}_{\text{E}}\text{X}$** ?

Part VI

Appendix / おまけ

Latin/CJK tokens

		T_EX	pT_EX	upT_EX	
Latin	I/O	8bit (multibytes) [†]	7bit 1byte	8bit (multibytes) [†]	
	token	charcode	8bit	8bit	8bit
		catcode	4bit	4bit	4bit
CJK	I/O	—	EUC etc. 8bit 2bytes	UTF-8 8bit 2–4bytes	
	token	charcode	—	16bit	24bit
		kcatcode	—	—	5bit
	Latin/CJK classification		—	fixed	customizable
inputenc		OK	NG	OK	
Babel		full	partial	full	

[†]: with inputenc

Character encoding in upTeX

	Latin	CJK		comment
	TeX compatible <256	upTeX extended BMP	over BMP	
.tex / .aux I/O buffer		UTF8		
token	1byte	2–3bytes	4bytes	
	12bit		29bit	with (k)catcode
.dvi / .vf	set1 T1 etc. 8bit	set2 UCS2 16bit	set3 UTF32 24bit	
.tfm	T1 etc. 8bit	UCS2 16bit	—†	†treated as Kanji 'jfm' for CJK
.ps / CMap	T1 etc. 8bit	UCS2 16bit	UTF16 2×16bit	

kcatcode

kcat code	cat code	kind	e.g.	control word	end of line
			
	10	space	␣		
15	11	char	azAZ	yes	as space
	12	other char	(.!?)	no	as space
			
16		Kanji	漢漢	yes	ignore
17		Kana	かナ	yes	ignore
18		CJK symbol	《・。』	no	ignore
19		Hangul		yes	as space