

npT_EX の開発（方針）について



北川 弘典

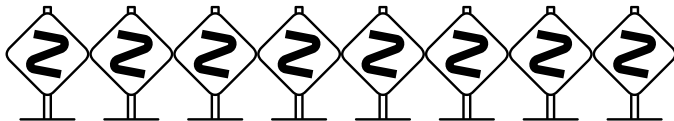
2023-11-11

T_EXConf 2023

- はじめに
- 過去を見るか．未来を見るか
- Unicode の直接入力
- 開発者への説明
- X_YTeX ベースにする場合の懸案事項



このスライドのように、
タイトルが白抜きになっているものは
発表では原則飛ばします。



スライドはどこかにアップロードの**予定**



L^AT_EX 2_ε 2020-02-02 以降，expl3 が
フォーマット内に読み込まれるように

- L^AT_EX 2_ε が徐々に expl3 の機能を用いて書き直される

L^AT_EX 2_ε 2020-02-02 以降，expl3 が
フォーマット内に読み込まれるように

- L^AT_EX 2_ε が徐々に expl3 の機能を用いて書き直される

例 フック機能の統一化 (2020-10-01)

everyshi, filehook, atbegshi, ... が「ほぼ不要」になった

L^AT_EX 2_ε 2020-02-02 以降，expl3 が
フォーマット内に読み込まれるように

- L^AT_EX 2_ε が徐々に expl3 の機能を用いて書き直される

例 フック機能の統一化 (2020-10-01)

everyshi, filehook, atbegshi, ... が「ほぼ不要」になった

- Tagged PDF project

Frank Mittelbach and Chris Rowley. “L^AT_EX Tagged PDF—
A blueprint for a large project”, TUGboat **41**:3, 2020.

ところが……

```
\RequirePackage{pdfmanagement-testphase}  
\DocumentMetadata{backend = dvipdfmx}  
\documentclass{article}  
\usepackage{hyperref,pxjahyper}  
\hypersetup{pdftitle={日本語}}
```

```
\RequirePackage{pdfmanagement-testphase}  
\DocumentMetadata{backend = dvipdfmx}  
\documentclass{article}  
\usepackage{hyperref,pxjahyper}  
\hypersetup{pdftitle={日本語}}
```

結果 ((u)pL^AT_EX <2023-02-14u04>+1)

! LaTeX Error: Invalid UTF-8 string: missing continuation
byte (x3).


```
\RequirePackage{pdfmanagement-testphase}  
\DocumentMetadata{backend = dvipdfmx}  
\documentclass{article}  
\usepackage{hyperref,pxjahyper}  
\hypersetup{pdftitle={日本語}}
```

結果 ((u)pL^AT_EX <2023-02-14u04>+1)

! LaTeX Error: Invalid UTF-8 string: missing continuation
byte (x3).

“Japanese characters in metadata on (u)pLaTeX”, aminophen,
[latex3/pdfresources/#18](#), 2021-06-14

```
\RequirePackage{pdfmanagement-testphase}  
\DocumentMetadata{backend = dvipdfmx}  
\documentclass{article}  
\usepackage{hyperref,pxjahyper}  
\hypersetup{pdftitle={日本語}}
```

2年以上前
今も未解決

結果 ((u)pL^AT_EX <2023-02-14u04>+1)

! LaTeX Error: Invalid UTF-8 string: missing continuation
byte (x3).

“Japanese characters in metadata on (u)pL^AT_EX”, aminophen,
[latex3/pdfresources/#18](#), 2021-06-14

(u)pT_EX の特殊性：和文文字トークン

1/3

(u)pT_EX は**欧文部分は 8 bit エンジン**だが
和文文字トークンという「異質なもの」がある：

(u)pT_EX は**欧文部分は 8 bit エンジン**だが
和文文字トークンという「異質なもの」がある：

- 0-15 の範囲の `\catcode` を持たない
- 文字コードが 256 以上になりうる
- 「`\lowercase` トリック」等が不可能

(u)pT_EX は欧文部分は 8 bit エンジンだが
和文文字トークンという「異質なもの」がある：

- 0-15 の範囲の `\catcode` を 持たない
- 文字コードが 256 以上になりうる
- 「`\lowercase` トリック」等が不可能

→ 和文文字トークンを考慮しないパッケージは、

(u)pL_AT_EX で動かないことがある

とくに expl3 の一部

前スライドの通り，比重がますます増している expl3 では「和文文字トークンを満足にサポート」とは言えない．

前スライドの通り，比重がますます増している expl3 では「和文文字トークンを満足にサポート」とは言えない．

しかし，今の日本語 T_EX 開発コミュニティには，和文文字トークンについて

- パッチを自前で書く
- L^AT_EX チーム他に事情説明する

だけの余力はない．

話は聞かせてもらったぞ！
(u)pL^AT_EX は滅亡する！



- l3regex において，和文文字トークン「±」(U+00B1) と UTF-8 列中のバイト B_1 が区別されない

「[l3regex] [upTeX] 和文文字と欧文バイトがマッチしてしまう」，h20y6m, [h20y6m/plexpl3/#2](https://h20y6m.github.io/plexpl3/#2), 2022-09-02

- T_EXConf 2019 に話す予定だったネタも無関係ではなかった

p_TE_X 4.0.0 (TL2022) で解決

Hironori Kitagawa, “Distinguishing 8-bit characters and Japanese characters in (u)p_TE_X”, TUGboat **41**:3, 2020.



次を目標とした新エンジンの開発を計画中：

「[npTeX] 欧文トークンに似せた CJK トークン」，t-tk，
tex-jp-build/#150, 2022-10-16.

- L^AT_EX 本体の変更に伴う保守が「楽に」行える
- 従来の upL^AT_EX 用クラス・パッケージが
(ほぼ) そのまま処理できる

次を目標とした新エンジンの開発を計画中：

「[npT_EX] 欧文トークンに似せた CJK トークン」, t-tk,
tex-jp-build/#150, 2022-10-16.

- L^AT_EX 本体の変更に伴う保守が「楽に」行える
- 従来の upL^AT_EX 用クラス・パッケージが
(ほぼ) そのまま処理できる

名称だけは npT_EX と決まっている.

→ ttk さんの案. new pT_EX, next pT_EX, novel pT_EX,
notable pT_EX, nippon (p)T_EX, nihongo pT_EX

名称は決まったが……

1/2

npT_EX については、名称以外は

- 内部コードは Unicode

入力は EUC/Shift_JIS も……

- 1 文字について 1 つの文字トークンを生成

- **文字トークン段階**では、欧文・和文の
区別は必須でない

LuaT_EX-j_a はそれでやっている

以外はほとんど決まっていない。

名称は決まったが……



2/2

npT_EX については、名称以外は

- 内部コードは Unicode

入力は EUC/Shift_JIS も……

- 1 文字について 1 つの文字トークンを生成

- **文字トークン段階**では、欧文・和文の
区別は必須でない

LuaT_EX-ja はそれでやっている

以外はほとんど決まっていない。

以下 「何を悩んでいるのか」 を説明。

トークン以外でどう欧文・和文を区別？



文字ノードを1つずつ作っていく際に、
その文字コードの `\cjkkxcode` (仮称) の値を見て判定
従来の `\kcatcode` で一度に制御された属性を
`\cjkkxcode` ではビットで分割してもよい：

- 欧文文字ノードか和文文字ノードか
- `\jchrwidowpenalty` の挿入対象か
- 直後の改行は半角空白を発行するか

LuaTeX-jā では

attribute

kcatcode 最下位

欧文/和文に連動



大きな開発方針は 2 案

- i ϵ -up $\text{T}_{\text{E}}\text{X}$ から文字トークンの扱いを変える

2022.10 に途中まで試して放置

- ii $\text{X}_{\text{E}}\text{T}_{\text{E}}\text{X}$ に $\text{pT}_{\text{E}}\text{X}$ の和文組版機能を載せる

2023.8 に試作品→そこから放置

- $\text{X}_{\text{E}}\text{T}_{\text{E}}\text{X}$ の「pdf 出力部」は `xdvipdfmx`
- 「`fontspec`, `unicode-math` は欧文のみ」
「和文フォントは JFM のみ」までが現実的

どこまで「過去を見る」か？

1/3

前提 既存の T_EX ソースから組版を変えたくないのなら
T_EX 環境 (T_EX Live xxxx) ごと保存すべき.

どこまで「過去を見る」か？

2/3

前提 既存の $\text{T}_{\text{E}}\text{X}$ ソースから組版を変えたくないのなら
 $\text{T}_{\text{E}}\text{X}$ 環境 ($\text{T}_{\text{E}}\text{X}$ Live xxxx) ごと保存すべき.

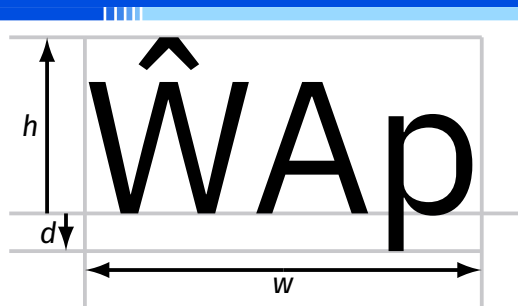
今新しく文書を作成する際に、
(数式以外で) 8 bit フォントを用いる必要性はどの程度？

前提 既存の $\text{T}_{\text{E}}\text{X}$ ソースから組版を変えたくないのなら
 $\text{T}_{\text{E}}\text{X}$ 環境 ($\text{T}_{\text{E}}\text{X}$ Live xxxx) ごと保存すべき.

今新しく文書を作成する際に、
(数式以外で) 8 bit フォントを用いる必要性はどの程度？

- Type 1 フォントのサポート終了記事：[Adobe](#), [Microsoft](#)
- dvips による PostScript 出力の必要性？
- フォントによる組版の差異をどこまで許容？

フォントによる組版の差異の例



	enc	h	d	w
cmr10	OT1	9.47221	1.94444	23.33339
ecrm1000	T1	8.9978	1.94397	23.32764
lmroman10-regular	TU	8.69	1.94	22.23



LuaT_EX-j_a への移行の 3 ハードル

LuaT_EX-j_a への移行の 3 ハードル

1 fontspec など, Unicode L^AT_EX の流儀

LuaT_EX-j_a への移行の 3 ハードル

- 1 fontspec など, Unicode L^AT_EX の流儀
- 2 マクロパッケージであることによる制限
 - `\ltjsetParameter{xkanjiskip=...}` など
 - 単位 Q, H, zw, zh なし

LuaT_EX-j_a への移行の 3 ハードル

- 1 fontspec など, Unicode L^AT_EX の流儀
- 2 マクロパッケージであることによる制限
 - `\ltjsetParameter{xkanjiskip=...}` など
 - 単位 Q, H, zw, zh なし
- 3 メモリ食い・遅い

LuaT_EX-j_a への移行の 3 ハードル

- 1 fontspec など, Unicode L^AT_EX の流儀
- 2 マクロパッケージであることによる制限
 - `\ltjsetParameter{xkanjiskip=...}` など
 - 単位 Q, H, zw, zh なし
- 3 メモリ食い・遅い

「X₃T_EX ベースの npT_EX」では 2, 3 はない

→ Unicode T_EX への移行を促せるか？

a		b	c
トークン	バイト単位	Unicode 文字	Unicode 文字
フォント	8 bit	8 bit	Unicode
直接入力 実現法	アクティブ文字 (inputenc)		—

- a 8 bit エンジン上 (pdfL^AT_EX, (u)pL^AT_EX)
- b Unicode L^AT_EX + T1
- c X_ƎL^AT_EX, LuaL^AT_EX (Unicode L^AT_EX) 標準

a		b	c
トークン	バイト単位	Unicode 文字	Unicode 文字
フォント	8 bit	8 bit	Unicode
直接入力 実現法	アクティブ文字 (inputenc)		—

- a 8 bit エンジン上 (pdfL^AT_EX, (u)pL^AT_EX)
- b Unicode L^AT_EX + T1, ε-upT_EX ベースの npT_EX
- c X_YL^AT_EX, LuaL^AT_EX (Unicode L^AT_EX) 標準

a		b	c
トークン	バイト単位	Unicode 文字	Unicode 文字
フォント	8 bit	8 bit	Unicode
直接入力 実現法	アクティブ文字 (inputenc)	<u>未考慮</u>	—

a 8 bit エンジン上 (pdfL^AT_EX, (u)pL^AT_EX)

b Unicode L^AT_EX + T1, ε-upT_EX ベースの npT_EX

c X_YL^AT_EX, LuaL^AT_EX (Unicode L^AT_EX) 標準

8 bit エンジンでの Unicode 直接入力



\P Q でバイト $0xPQ$ を, \hoge で名称が hoge の制御綴を表す

1 UTF-8 上位バイト (\C 8 – \F F) をアクティブ化 (latex.ltx)

2 各上位バイトに「UTF-8 列を読み取る命令」を割り当て (utf8.def)

\C 2 \F 9 $\longrightarrow \text{\UTFviii@two@octets}$ \C 2 \F 9 $\longrightarrow \dots \longrightarrow \text{\u8:}\text{\C 2}\text{\F 9}$

3 Unicode 文字の実際の出力命令を定義
(\DeclareUnicodeCharacter, utf8enc.dfu)

$\text{\u8:}\text{\C 2}\text{\F 9} \longrightarrow \text{\IeC}\{\text{\ss}\} \longrightarrow \text{\ss}$



`\DeclareUnicodeCharacter` などを再定義し,
「各 Unicode 文字を個別にアクティブ化, 出力命令割当」は？

$$\beta \longrightarrow \dots \longrightarrow \text{\IeC{\ss}} \longrightarrow \text{\ss}$$

`\DeclareUnicodeCharacter` などを再定義し,
「各 Unicode 文字を個別にアクティブ化, 出力命令割当」は？

$$\beta \longrightarrow \dots \longrightarrow \text{\IeC{\ss}} \longrightarrow \text{\ss}$$

本家 $\text{\LaTeX} 2_{\epsilon}$ で対応されるのが望ましい

`\DeclareUnicodeCharacter` などを再定義し,
「各 Unicode 文字を個別にアクティブ化, 出力命令割当」は?

$$\beta \longrightarrow \dots \longrightarrow \text{\IeC{\ss}} \longrightarrow \text{\ss}$$

本家 $\text{\LaTeX} 2_{\epsilon}$ で対応されるのが望ましい

懸案事項 「`\Gödel`」 「`\straße`」 など

アクティブ化以外の実装？



`\nptexnoderecipe c = <token list>`

<token list> が空でない場合，水平モード中の
「文字コード *c*，`\catcode 11` の文字トークン」は
「文字コード *c* の文字ノード」を作るのではなく，
<token list> のように振る舞う．

- 「`\nptexnoderecipe c`」が非空かどうか見分ける手段は？
- `\accent` との兼ね合い？



L^AT_EX や多くのパッケージでは、現状次のような認識：

`\sys_if_engine_luatex:p v \sys_if_engine_xetex:p`

⇔ (`\Umathcode` 定義済)

⇔ (Unicode 1 文字でいつでも 1 トークン)

⇔ (本文の欧文フォントは Unicode)

L^AT_EX や多くのパッケージでは、現状次のような認識：

`\sys_if_engine_luatex:p v \sys_if_engine_xetex:p`

⇔ (`\Umathcode` 定義済)

⇔ (Unicode 1 文字でいつでも 1 トークン)

⇔ (本文の欧文フォントは Unicode)

「 ϵ -upT_EX ベースの npT_EX」や、「X₃T_EX ベースの npT_EX,
本文 8 bit フォント」では上記の認識は誤りとなる.

X_YTeX 互換 (`\sys_if_engine_xetex:p`) を謳うなら
L^ATeX やパッケージの改変は少なくて済むが、

- X_YTeX 拡張と pTeX 拡張の干渉への対応
 - 入力の文字コード (次スライド)
 - `\XeTeXinterchartoken`

X_YT_EX 互換 (`\sys_if_engine_xetex:p`) を謳うなら
L^AT_EX やパッケージの改変は少なくて済むが、

- X_YT_EX 拡張と pT_EX 拡張の干渉への対応
 - 入力の文字コード (次スライド)
 - `\XeTeXinterchartoken`
- dvips を想定するなら、「dvi 互換モード」必要
→ 「ε-upT_EX ベースの npT_EX」実装と実質的に同じ作業

■ pT_EX 系列 (ptexenc)

文字コード EUC, Shift_JIS, ISO-2022-JP, UTF-8

ただし, ISO-2022-JP は常に認識

Unicode 正規化 濁点・半濁点付き仮名を合成

■ X_YT_EX

文字コード UTF-8, UTF-16, バイト単位, ICU

Unicode 正規化 無変換, NFC, NFD から選択

両者のマージは大変そう → $\backslash\left\{\begin{smallmatrix} \text{epT} \\ \text{XeT} \end{smallmatrix} \text{eX}\right\} \text{inputencoding}$ の値で二者択一？

余談：“native font”にしか効かない

「和文フォントは従来通り JFM のみ」とすれば、
これらは pT_EX 由来の和文組版処理とは干渉せず

- $\backslash\text{XeTeXdashbreakstate}$, $\backslash\text{XeTeXlinebreak}$ $\left\{ \begin{array}{l} \text{locale} \\ \text{skip} \\ \text{penalty} \end{array} \right\}$
- Unicode の合字，異体字シーケンスなど

← upT_EX に実装計画あり．npT_EX では？

「upT_EX+dvipdfmx」異体字セクタ、Unicode 合成文字」，
t-tk, [tex-jp-build/#46](https://t-tk.org/tex-jp-build/#46), 2018-01-28.



- はじめに
- 過去を見るか．未来を見るか
- Unicode の直接入力
- 開発者への説明
- X₃TeX ベースにする場合の懸案事項

typeset by npTeX $\frac{3.141592653}{\text{T}_{\text{E}}\text{X}}$ $\frac{2.6}{\varepsilon\text{-T}_{\text{E}}\text{X}}$ $\frac{0.999995}{\text{X}_{\text{3}}\text{T}_{\text{E}}\text{X}}$ $\frac{0.0000}{\text{npT}_{\text{E}}\text{X}}$

