

pT_EX 系列の文字列化における 和文欧文の区別

北川 弘典

2019-10-12

~~T_EXConf 2019~~ 中止

本発表の話題は次の現象に関連した諸々

入力

```
\documentclass{minimal}%           UTF-8 入力  
\begin{document} f \end{document}% f: U+017F
```

エラーメッセージ

```
! Package inputenc Error: Unicode character 顛 (U+C4CF)  
(inputenc) not set up for use with LaTeX.
```

“Inconsistent error message”, JulienPalard,
2019-05-27, platex/#84

- 入出力は UTF-8 と想定
- 主に内部コード EUC の pT_EX を扱う

記法

hoge hoge という名前の制御綴

XY 8 ビット値 0xXY

もしくはその文字コードの欧文文字トークン

AB
CD 16 ビット値 0xABCD

もしくはその文字コードの和文文字トークン

■導入：T_EXとトークン

■和文・欧文の区別が消えるとき

■現行の実装と改善案

(p)TeX は入力をトークン単位で処理する。

入力例 `\hbox{ β 漢 \hskip3pt}`

$\begin{matrix} \boxed{C} & \boxed{9} & \boxed{E} & \boxed{B} & \boxed{A} \\ \boxed{3} & \boxed{F} & \boxed{6} & \boxed{C} & \boxed{2} \end{matrix}$

前処理 ptexenc により，内部コードへ変換

`\hbox{^^c3^^9f $\frac{B}{4}\frac{C}{1}$ \hskip3pt}`

トークン列

`\hbox` { $\frac{C}{3}$ $\frac{9}{F}$ 漢 `\hskip` 3 p t }

■：和文文字トークン 無：欧文文字トークン

トークンの補足

^^ 記法 間接的な入力法，常に欧文扱い¹

- ^^pq (p, q: 0-9, a-f)
- ^^X (X: ASCII で 0-9, a-f 以外)

upTeX \kcatcode で欧文扱いと和文扱い切り替え

<code>hbox</code>	{	<table border="1"><tr><td>C</td><td>9</td></tr><tr><td>3</td><td>F</td></tr></table>	C	9	3	F		<table border="1"><tr><td>E</td><td>B</td><td>A</td></tr><tr><td>6</td><td>C</td><td>2</td></tr></table>	E	B	A	6	C	2		<code>hskip</code>	3 p t }
C	9																
3	F																
E	B	A															
6	C	2															
		β		漢													

¹pTeX p3.1.4 以降. p3.1.8 で制御綴名内の ^^ 記法の扱い改善.

```
\font\TT=ec-lmtt10 \TT % 256文字あるフォント
\def\fuga{^^c5^^bf顛 ƒ β }
```

- \fuga の展開結果のトークン列

C5 BF 顛 C5 BF C3 GF

- \fuga の結果 (DVI)

Åƒ顛ÅƒÃ§ (*)

- ソース中で直に ^^c5^^bf顛ƒβ としても結果は (*)

■ 導入： $\text{T}_{\text{E}}\text{X}$ とトークン

■ 和文・欧文の区別が消えるとき

- plain $\text{pT}_{\text{E}}\text{X}$ の例
- $\text{pL}^{\text{A}}\text{T}_{\text{E}}\text{X}$ での例
- $\text{upT}_{\text{E}}\text{X}$ での状況

■ 現行の実装と改善案

\meaning 他の挙動

```
\font\TT=ec-lmtt10 \TT  
\def\fuga{^^c5^^bf顛𐀀}
```

	C5BF
UTF-8	「
EUC-JP	顛

- `\writeN{\fuga}` 顛顛顛 𐀀³
- `\meaning\fuga` macro:-> 顛顛顛 𐀀
- `\detokenize\expandafter{\fuga}` 顛顛顛 𐀀
- `\message{\fuga}` 顛顛顛 𐀀³

後ろの3命令は「文字列化」を行う命令

制御綴名と \string

```
\catcode"C5=11 \catcode"BF=11
```

```
\def\顛{P} \def\l{Q}
```

X	\X	\csname X\endcsname	\string\X
顛	Q	Q	\顛
l	Q	Q	\顛
^^c5^^bf	Q	Q	\顛

“Inconsistent error message” 再掲

JulienPalard, 2019-05-27, latex/#84

入力

```
\documentclass{minimal}% UTF-8 入力  
\begin{document} f \end{document}
```

エラーメッセージ

```
! Package inputenc Error: Unicode character 顛 (U+C4CF)  
(inputenc) not set up for use with LaTeX.
```

`u8:f` (= `u8:C5BFBF`) の `\string` の結果 から作られる

実際には `\u8:顛`

“Non-ASCII symbols in labels”

aminophen, 2019-07-13, platex/#86

入力

元々は L^AT_EX 2_ε 2019-10-01 関係

```
\documentclass{article}
```

```
\UseRawInputEncoding%
```

"raw" encoding

```
\begin{document}
```

```
\section{a}\label{abβ}\ref{abβ}
```

```
\end{document}
```

警告 (消えない)

LaTeX Warning: Reference `ab^β` on page 1 undefined
on input line 4.

aminophen, 2019-07-13, platex/#86

aux ファイル 3 行目

```
\newlabel{ab $\overset{\text{C}}{\underset{3}{9}}\text{f}$ }{1}{1}
```

→読み込み時に ptexenc は次のように変換：

```
\newlabel{ab $\overset{\text{A2}}{\underset{\text{AF}}{9}}\text{f}$ }{1}{1}
```

→ $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ はラベル「ab $\overset{\text{A2}}{\underset{\text{AF}}{9}}\text{f}$ 」が定義されたと認識

亜種：違うラベル名だが……？

入力

```
\documentclass{article}
\UseRawInputEncoding % `raw' encoding
\begin{document}
\section{\label{f}\label{顛}}
\end{document}
```

警告メッセージ（消えない）

LaTeX Warning: Label `顛' multiply defined.

理由

aux ファイル 3, 4 行目が同一 ($\newlabel{顛}{{1}{1}}$)

- (実質的には) ptexenc によるコード変換なし
→ latex/#84, #86 は upL^AT_EX では影響なし
- 文字列化後の再トークン化は `\kcatcode` 依存

```
\font\TT=ec-lmtt10 \TT
```

```
\def\fuga{^^e3^^81^^82あβ}
```

```
\kcatcode"3042=16
```

```
\fuga%
```

ãÁćあÃ§

```
\meaning\fuga%
```

macro:-> ああÃ§

```
\expandafter\string\csname\fuga\endcsname
```

```
%
```

ああÃ§

- (実質的には) ptexenc によるコード変換なし
→ latex/#84, #86 は upL^AT_EX では影響なし
- 文字列化後の再トークン化は `\kcatcode` 依存

```
\font\TT=ec-lmtt10 \TT
```

```
\def\fuga{^^e3^^81^^82あβ}
```

```
\kcatcode"3042=15
```

```
\fuga%
```

```
ãあÃ§
```

```
\meaning\fuga%
```

```
macro:->ããÃ§
```

```
\expandafter\string\csname\fuga\endcsname
```

```
%
```

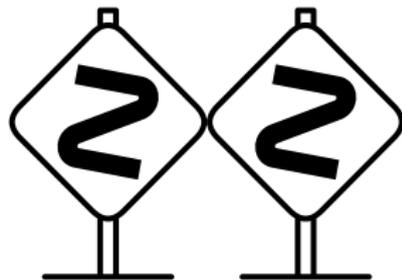
```
\ããÃ§
```

■導入：T_EX とトークン

■和文・欧文の区別が消えるとき

■現行の実装と改善案

- 文字列化
- 端末・ファイルへの出力
- 端末・ファイルからの入力



改善案 tex-jp-build 個人フォーク²で実験中

²https://github.com/h-kitagawa/tex-jp-build/tree/printkanji_16bit

(途中の) 結果を一旦配列 *str_pool* に格納する処理

- *str_pool* にはバイト列として格納される
→和文由来か欧文由来かの情報は消える
- プリミティブ達
 - `\the`, `\detokenize`, `\number`, `\string`, `\meaning`, `\euc` 他
→結果を文字列化したものを再トークン化
 - `\message`, `\errmessage`, `\special`, `\write18`
→文字列化したものを出力 (実行)
 - `\scantokens`
→文字列化したものを疑似ファイルと扱う

「\the, \detokenize, ..., \euc 他」の一覧

\the, \detokenize 以外は command code が *convert*

\TeX 82 \the, \number, \romannumeral, \string, \meaning,
\fontname, \jobname

ϵ - \TeX \detokenize, \eTeXrevision

(u) $\text{p}\TeX$ \kansuji, \euc, \sjis, \jis, \kuten, \ucs,
\ptexrevision, \uptexrevision

$\text{pdf}\TeX$ \pdfstrcmp, \pdfcreationdate, \pdffilemoddate,
\pdffilesize, \pdfmdfivesum, \pdffiledump,
\pdfuniformdeviate, \pdfnormaldeviate

(途中の) 結果を一旦配列 *str_pool* に格納する処理

現行 *str_pool* の各要素は 0 .. 255 を格納

改善案 *str_pool* の各要素は 0 .. 511 を格納³

0-255 欧文由来バイト

256-511 和文由来バイト

 改造量を考え、「1 要素 1 文字」にはしなかった

 \pdf... などで X_YTEX のコード流用可能
(X_YTEX では *str_pool* の各要素は 0 .. 65535 を格納)

³WEB ソース中では 0 .. 32768 (unsigned short[] にするため).

余談：fmt ファイルの容量

pLaTeX2e <2019-10-01> (based on LaTeX2e <2019-10-01>)

[KiB] (切り上げ)	英語のみ			全部		
	現行	改善案		現行	改善案	
pL ^A T _E X	898	994	+96	4449	4584	+136
upL ^A T _E X	894	990	+96	4445	4580	+136

- 多数のハイフネーションパターンの下では影響薄
- 容量が気になれば X_YT_EX, LuaT_EX のように gzip 圧縮？
(T_EX Live の pT_EX 系列は SyncT_EX のために zlib リンク済)

```
\def\fuga{^^c5^^bf顛β\cr}\meaning\fuga
```

- 1 \fuga の中身を表すトークン列

C₅ B_F 顛 C₃ 9_F cr

- 2 文字列化の結果 (*str_pool* の中身)

macro : - > C₅ B_F C₃ 9_F \ c r _

- 3 再トークン化

macro : - > 顛 顛 C₃ 9_F \ c r _

和文文字と解釈できるバイト列は和文文字トークンに

```
\def\fuga{^^c5^^bf顛β\cr}\meaning\fuga
```

- 1 \fuga の中身を表すトークン列

```
Ⓞ5 ⓄF 顛 Ⓞ3 ⓄF Ⓞcr
```

- 2 文字列化の結果 (*str_pool* の中身)

```
macro : - > Ⓞ5 ⓄF Ⓞ01C5 Ⓞ01BF Ⓞ3 ⓄF \ cr ␣
```

- 3 再トークン化

```
macro : - > Ⓞ5 ⓄF 顛 Ⓞ3 ⓄF \ cr ␣
```

フラグあり (Ⓞ₀₁₀₀ - Ⓞ₀₁_{FF}) の列からのみ和文文字トークンに

```
\def\fuga{^^e3^^81^^82あβ\cr}\meaning\fuga
```

- 1 \fuga の中身を表すトークン列

$\boxed{\text{E}_3}$ $\boxed{\text{8}_1}$ $\boxed{\text{8}_2}$ 顛 $\boxed{\text{C}_3}$ $\boxed{\text{9}_F}$ $\boxed{\text{cr}}$

- 2 文字列化の結果 (*str_pool* の中身)

macro : - > $\boxed{\text{E}_3}$ $\boxed{\text{8}_1}$ $\boxed{\text{8}_2}$ $\boxed{\text{E}_3}$ $\boxed{\text{8}_1}$ $\boxed{\text{8}_2}$ $\boxed{\text{C}_3}$ $\boxed{\text{9}_F}$ \ c r _

- 3 再トークン化

macro : - > あ あ $\boxed{\text{C}_3}$ $\boxed{\text{9}_F}$ \ c r _

和文文字トークンにするかは \kcatcode 依存

```
\def\fuga{^^e3^^81^^82あβ\cr}\meaning\fuga
```

- 1 \fuga の中身を表すトークン列

```

E3  81  82  顛 C3  9F  cr

```

- 2 文字列化の結果 (*str_pool* の中身)

```

macro : - > E3  81  82  01E3  0181  0182  C3  9F  \ c r _

```

- 3 再トークン化

```

macro : - > E3  81  82  あ C3  9F  \ c r _

```

フラグあり (0100-01FF) の列からのみ和文文字トークンに

制御綴名も配列 *str_pool* に格納される

→ここでも和文由来・欧文由来の区別はない

例

X	\X \csname X\endcsname	\string\X
---	------------------------	-----------

顛	内部では <table border="1"><tr><td>C</td><td>B</td></tr><tr><td>5</td><td>F</td></tr></table>	C	B	5	F	\顛
C		B				
5		F				

𐀀		\顛
---	--	----

^^c5^^bf		\顛
----------	--	----

制御綴名も配列 *str_pool* に格納される

→ここでも和文由来・欧文由来の区別をつける

例

X	\X	\csname X\endcsname	\string\X
---	----	---------------------	-----------

顛	内部では	<table border="1"><tr><td>01</td><td>01</td></tr><tr><td>C5</td><td>BF</td></tr></table>	01	01	C5	BF	\顛
01	01						
C5	BF						

ŀ	内部では	<table border="1"><tr><td>C</td><td>B</td></tr><tr><td>5</td><td>F</td></tr></table>	C	B	5	F	\C B
C			B				
5	F						

^^c5^^bf			\C B
----------	--	--	------

制御綴名と \string (upTeX)

\kcatcode 依存なのは**制御綴名の再トークン化時**か？
制御綴名の内部表現か？

% 現行

改善案

```
\font\TT=ec-lm tt10 \TT
\def\あ{hoge}\def\B{\あ}
\catcode"E3=11 \catcode"81=11 \catcode"82=11
{\kcatcode"3042=15
  \あ                % hoge                ! Undefined
  \meaning\B        % macro:->\ãŒ        macro:->\あ
\bye
```

pTEX 系列における 1 バイトの出力は複数段階

$print \xrightarrow{\text{(ptexenc)}} print_char \xrightarrow{\text{(ptexenc)}} putc2 \xrightarrow{\text{(ptexenc)}} \text{実際の出力}$

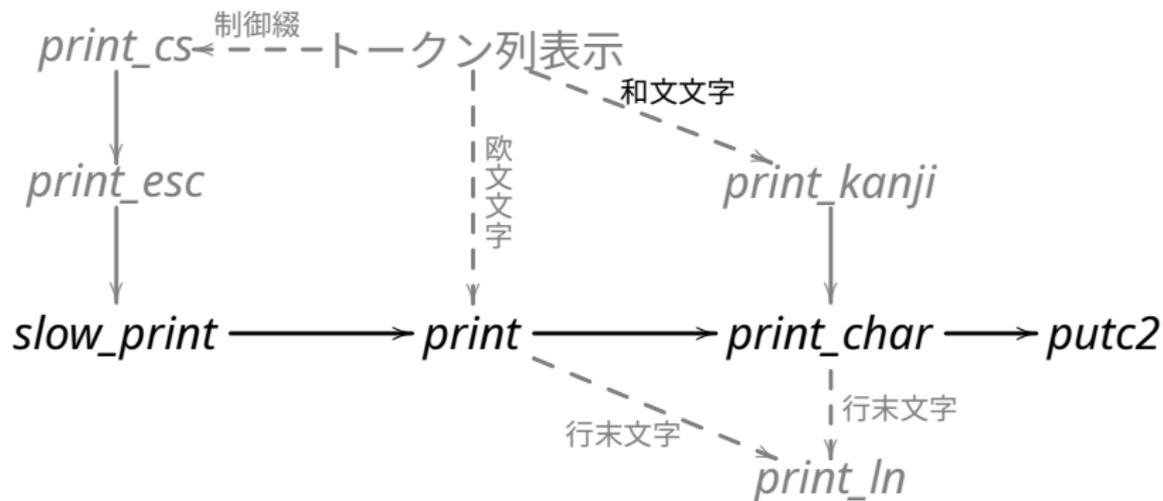
$\begin{matrix} \text{9} \\ \text{F} \end{matrix}$	$\begin{matrix} \text{9} \\ \text{F} \end{matrix}$	$\wedge \wedge 9 f$	$\wedge \wedge 9 f$	下記以外
a	a	a	a	ASCII
$\begin{matrix} \text{C} \\ \text{5} \end{matrix}$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix}$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix}$	(次を待つ)	和文先頭
$\begin{matrix} \text{C} \\ \text{5} \end{matrix} A^4$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} A$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} A$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} A$	
$(\begin{matrix} \text{C} \\ \text{5} \end{matrix} \begin{matrix} \text{B} \\ \text{F} \end{matrix})^5$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} \begin{matrix} \text{B} \\ \text{F} \end{matrix}$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} \begin{matrix} \text{B} \\ \text{F} \end{matrix}$	$\begin{matrix} \text{C} \\ \text{5} \end{matrix} A$	顛

putc2 はバイトが和文由来か欧文由来か知らない

⁴実際は *print(0xc5); print(0x41)* という 2 回の呼び出し。

⁵和文文字トークンの出力時には *print* は経由しない。

print 系プロシージャ



	文字列	1 バイト	
汎用	<code>slow_print</code>	<code>print</code>	(必要なら ^^ 記法)
高速	<code>print</code>	<code>print_char</code>	

str_pool で 0100 - 01FF を和文由来にしたことと合わせる

(ptexenc)

print → *print_char* → *putc2* → 実際の出力

P	Q	{	P	Q	出力可能 出力不能
^	^		p	q	

無理	01C5	01C5	(次を待つ)	}	ほぼ現行
	01C5 0141	01C5 A	C A		
	01C5 01BF	01C5 01BF	顛		

1 文字 *c* を出力する目的の *print(c)* は変更が必要
(非常に素直だが)

1 文字を出力する *print*

$c \geq 256$ のとき *print(c)* は「 c 番の文字列」を出力
→ **改善案** では場合わけが必要

例：関数 *slow_print(s)* 現行（と同じ動作のコード）
出力不能な文字を含みうる， s 番の文字列を出力

```
while  $j < str\_start[s + 1]$  do begin  
   $c \leftarrow so(str\_pool[j]);$   
  print(c);  
  incr(j);  
end;
```

1 文字を出力する *print*

$c \geq 256$ のとき *print(c)* は「 c 番の文字列」を出力
→ **改善案** では場合わけが必要

例：関数 *slow_print(s)* **改善案**

出力不能な文字を含みうる， s 番の文字列を出力

```
while  $j < str\_start[s + 1]$  do begin  
   $c \leftarrow so(str\_pool[j]);$   
  if  $c \geq 256$  then print_char(c) else print(c);  
  incr(j);  
end;
```

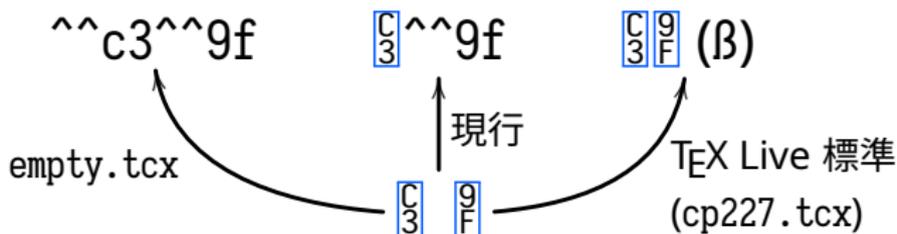
`0` (`_`)-`7E` (`~`) 以外に……

現行

和文の内部コードで現れるバイトは必ず出力可

改善案

内部コードと離れて再考可能. cf. TCX ファイル



TCX ファイル 1

TeX Character Translation

- 8ビットエンコーディング用の入出力変換機構
- 出力可能性では「何も変換しない」TCX利用

欧文 TeX での例

```
\message{\string\^^A.\string\^^K.^^9f.^^c1.^^f7}
```

empty.tcx \^^A.\^^K.^^9f.^^c1.^^f7

cp8bit.tcx \^^A.\^^K.

9	F	1	7
---	---	---	---

8	F
---	---

 も出力可

cp227.tcx \^^A.\

0	9	F	1	7
---	---	---	---	---

 ↑ と HT, LF, VT

natural.tcx \

9	0	F	1	7
---	---	---	---	---

0	F
---	---

 出力可

TCX ファイル 2

p_TE_X 系列の TCX サポートは限定的

- そもそも「変換する」TCX の利用価値なし
→ここでは出力可能性のみ話題に
- 出力可能性がフォーマット作成時 (inip_TE_X) に
固定されてしまう

```
eptex -ini -etex -translate-file=HOGE platex.ini
```

```
% →出力可能性について効力あり
```

```
platex -translate-file=FUGA ... % 効力なし
```

```
\write1{^^c5^^bf顛abβ}
```

1 出力すべきトークン列

C5 BF 顛 a b C3 9F

2 ptexenc へ渡るもの

C5 BF C5 BF a b C3 ^ ^ 9 f

3 出力

顛顛abC3^^9f

\message は文字列化を挟むがほぼ同じ

```
\write1{^^c5^^bf顛abβ}
```

```
TCX: empty.tcx
```

1 出力すべきトークン列

C₅ B_F 顛 a b C₃ F₉

2 ptexenc へ渡るもの

C₅ B_F 01_{C5} 01_{BF} a b ^ ^ c 3 ^ ^ 9 f

3 出力

```
^^c5^^bf顛ab^^c3^^9f
```

\message は文字列化を挟むがほぼ同じ

```
\write1{^^c5^^bf顛abβ}
```

```
TCX: cp227.tcx
```

1 出力すべきトークン列

```

 $\begin{matrix} \text{C} & \text{B} & & & \text{C} & \text{F} \\ \text{5} & \text{F} & \text{顛} & \text{a} & \text{b} & \text{3} & \text{9} \end{matrix}$ 

```

2 ptexenc へ渡るもの

```

 $\begin{matrix} \text{C} & \text{B} & \text{01} & \text{01} & & \text{C} & \text{F} \\ \text{5} & \text{F} & \text{C5} & \text{BF} & \text{a} & \text{b} & \text{3} & \text{9} \end{matrix}$ 

```

3 出力

```

 $\begin{matrix} \text{C} & \text{B} & & \text{C} & \text{F} \\ \text{5} & \text{F} & \text{顛ab} & \text{3} & \text{9} \end{matrix}$ 
  (= [顛abβ])

```

\message は文字列化を挟むがほぼ同じ

各行は ptexenc がコード変換し，配列 *buffer* に格納

- JIS X 0208 範囲外の Unicode 文字
→ ^^ 記法に変換
- 1 バイト目に現れないバイト ($\text{\C}\text{-}\text{\F}$, $\text{\F}\text{-}\text{\F}$) 単独
→変換せずそのまま
- $\text{\C}\text{-}\text{\F}$ が決まった数だけ続かないもの

$\text{\C}\text{-}\text{\F}$ $\text{\E}\text{-}\text{\F}$ $\text{\C}\text{-}\text{\D}$ $\text{\F}\text{-}\text{\F}$ $\text{\C}\text{-}\text{\D}$ $\text{\C}\text{-}\text{\D}$

→ TL2019 以前では $\text{\A}\text{\F}$ (EUC), $\text{\C}\text{\D}$ (SJIS) に変換

これが platex/#86 の一要因

各行は ptexenc がコード変換し，配列 *buffer* に格納

- JIS X 0208 範囲外の Unicode 文字
→ ^^ 記法に変換
- 1 バイト目に現れないバイト ($\text{\80-}\text{\91}$, $\text{\95-}\text{\9F}$) 単独
→変換せずそのまま
- $\text{\80-}\text{\9F}$ が決まった数だけ続かないもの

$\text{\92-}\text{\94}$ $\text{\9E-}\text{\9F}$ $\text{\80-}\text{\9F}$ $\text{\90-}\text{\94}$ $\text{\80-}\text{\9F}$ $\text{\80-}\text{\9F}$

→ TL2019 以前では $\text{\A2}\text{\AF}$ (EUC), $\text{\81}\text{\AD}$ (SJIS) に変換
TL2020 では ^^ 記法に変換 (r52071)

各行は ptexenc がコード変換し，配列 *buffer* に格納

- 1 バイト目に現れないバイト ($\text{\textcircled{0}}-\text{\textcircled{1}}$, $\text{\textcircled{F5}}-\text{\textcircled{FE}}$) 単独
→変換せずそのまま
- $\text{\textcircled{8}}-\text{\textcircled{BF}}$ が決まった数だけ続かないもの
→ ^^ 記法に変換 (安全)

UTF-8 入力 (補足)

- $\begin{matrix} A2 \\ AF \end{matrix}$ (EUC), $\begin{matrix} 81 \\ AD \end{matrix}$ (SJIS) は何者？

→ JIS X 0208 で最初に未割り当て」なコード

UTF-8 列 → スカラー値 → JIS X 0208 → 内部コード
 $\begin{matrix} C \\ 9 \end{matrix}$ 「0」 「2区15点」 A2AF (EUC)

- 冗長な UTF-8 バイト列の扱い

$\text{p}\TeX$ 特に問題にしない

$\text{u}\text{p}\TeX$ 最短のバイト数での表現に変換

但し U+0000-U+007F については

$\begin{matrix} 0 \\ 0 \end{matrix}$ $\begin{matrix} 1 \\ F \end{matrix}$ のように前に $\begin{matrix} 0 \\ 0 \end{matrix}$ が付加

入力の各行を格納する配列

制御綴名読み取りの際、*buffer* 内で ^^ 記法を変換

前 a^^c3^^9f.\  ^^c5^^bf^^c3^^9f.
顛

後 a^^c3^^9f.\  .
顛

仮定：-の catcode は 11

入力の各行を格納する配列

制御綴名読み取りの際、*buffer* 内で ^^ 記法を変換

前 a^^c3^^9f. \

C	B
S	F

 ^^c5^^bf^^c3^^9f.

顛

後 a^^c3^^9f. \

C	B	C	B	C	9
F	F	F	F	F	F

 .

1 1 0 0 0 0

仮定：

8
F

 の catcode は 11

制御綴名の和文・欧文の区別も簡単に実装したい

→和文由来バイトのフラグを別配列に格納

 *buffer* は 8 bit のまま (ptexenc は他プログラムでも使う)

■ 導入：T_EX とトークン

■ 和文・欧文の区別が消えるとき

- plain pT_EX の例
- pL^AT_EX での例
- upT_EX での状況

■ 現行の実装と改善案

- 文字列化
- 端末・ファイルへの出力
- 端末・ファイルからの入力