

# EncT<sub>E</sub>X

možnost překódování vstupu v T<sub>E</sub>Xu

6. 9. 1997, 3. 1. 2003, 13. 1. 2003

Petr Olšák

EncTeX je volné programové vybavení; můžete jej dále šířit a modifikovat podle podmínek „GNU General Public License“, kterou publikovala Free Software Foundation; použijte verzi 2 této licence nebo (podle Vaší volby) libovolnou pozdější verzi.

Balíček najdete na Internetu na

`ftp://math.feld.cvut.cz/pub/olsak/encTeX/`.

Tento balíček je rozšiřován v naději, že bude užitečný, avšak BEZJAKÉKOLI ZÁRUKY; neposkytují se ani odvozené záruky PRODEJNOSTI anebo VHODNOSTI PRO URČITÝ ÚČEL. Další podrobnosti hledejte v Obecné veřejné licenci GNU.

Kopii „GNU General Public License“ jste měl obdržet spolu s tímto programem; pokud se tak nestalo, napište o ni Free Software Foundation, Inc., 675 Mass Ave, Cambridge, MA 02139, USA. Český překlad této licence najdete na <http://www.gnu.cz/gplcz.html>.

© 1997, 2002, 2003 RNDr. Petr Olšák

TeX je ochranná známka American Mathematical Society.

Autor TeXu je profesor Donald Knuth. TeX je volné programové vybavení se specifickou licencí, viz dokumentaci k tomuto programu.

## 1. Základní informace

Balík `encTeX` je jednoduché rozšíření `TeX`u pro takové implementace, ve kterých se `TeX` instaluje ze zdrojového kódu `tex.web`. Tuto podmínku například splňuje implementace `web2c`, určená pro UNIXy a jiné operační systémy s kvalitním překladačem jazyka C.

Rozšíření `encTeX` je zpětně kompatibilní s originálním `TeX`em. Přidává osm nových primitivů, kterými lze číst nebo naplňovat vnitřní kódovací tabulky, podle nichž je znak transformován na úrovni vstupního procesoru `TeX`u nebo při výstupu na terminál, do `log` souboru a `\write` souborů. Tyto tabulky se ukládají do formátových souborů, takže po načtení formátu se inicializují ve stejném stavu, v jakém byly v okamžiku příkazu `\dump`.

Změna `TeX`u je důkladně testovaná a prošla též testem TRIP s těmito dvěma odlišnostmi:

- Odlišný banner
- Počet „multiletter control sequences“ je o osm větší.

### 1.1. Instalace

Viz soubor `INSTALL`.

### 1.2. Verze

V roce 1997 byla zveřejněna první verze `encTeX`u, která umožňovala konverze pouze v režimu „byte na byte“ a nastavovala tisknutelnost znaků (primitivy `\xordcode`, `\xchrcode`, `\xprncode`).

V prosinci roku 2002 a v lednu 2003 jsem zapracoval možnost překódování více bytů na jeden byte nebo na kontrolní sekvenci. Nová verze má označení `Jan. 2003` a přidává dalších pět primitivů `\mubyte`, `\endmubyte`, `\mubytein`, `\mubyteout` a `\specialout`. To umožní definovat vstup pro UTF-8 kódované soubory.

### 1.3. Konflikt s TCX tabulkami

Protože jak `encTeX` tak TCX tabulky (přepínač `-translate-file`) pracují se stejnými kódovacími vektory `xord` a `xchr`, konfliktu při současném používání se nevyhneme. Proto jsem také původně stáhnul `encTeX` a přestal ho po omlazení TCX tabulek v roce 1998 prosazovat. Nicméně tyto tabulky neumějí UTF-8 kódované soubory a jsou podle mého názoru podstatně méně flexibilní. Z principiálních důvodů TCX tabulky nebudou nikdy umět deklarovat konverzi z více bytů na kontrolní sekvenci. V roce 2003 jsem tedy přidal podporu UTF-8 do nové verze `encTeX`u a začal jej znovu prosazovat.

Upozornění: Pokud máte `TeX` pozměněný `encTeX`em, pak jsou TXC tabulky neaktivní. I když napíšete na příkazovém řádku `-translate-file`, nestane se nic a soubor se nenačte. Pro hodnoty `xord` a `xchr` vektorů je rozhodující pouze to, co s nimi udělá `encTeX` prostřednictvím primitivů `\xordcode` a `\xchrcode`. Na možné koordinaci `encTeX`u s TCX tabulkami se pracuje.

### 1.4. Problém s licenci `TeX`u

`EncTeX` rozšiřuje `TeX` o nové primitivy, takže bychom neměli tomuto programu říkat `TeX`. Na druhé straně ale Knuth samotný předpokládá, že vnitřnosti `TeX`u budou odstíněny od prostředí operačního systému. Proto implementoval `xord` a `xchr` vektory. V `encTeX`u můžeme nastavit podle zvyklostí operačního systému vstupní a výstupní překódovací tabulky a pak nastavit všem novým primitivům význam `\undefined`. Dále se bude `TeX` modifikovaný `encTeX`em chovat stejně, jako originální `TeX`. Navíc můžeme překódovací tabulky nastavit při generování formátu a v produkční verzi `TeX`u zakázat přístup k primitivům. Produkční verze `TeX`u se pak chová zcela stejně jako originální `TeX`. Knuth předpokládá, že odstínění od prostředí operačního systému se provede vždy při kompilaci zdrojového kódu `TeX`u, zatímco `encTeX` umožňuje tuto otázku řešit později, například v době generování formátu. Umožnění úpravy některých (například paměťových) parametrů až za běhu `TeX`u také není nic nového a známe to skoro u všech distribucí `TeX`u.

Domnívám se, že druhý řádek na terminálu a v logu dostatečně informuje o tom, že se jedná o modifikovanou verzi `TeX`u. Také se domnívám, že pokud se velmi rozšíří kódování UTF-8, pak není zbytné a takové konverze jsou v 8 bitové verzi `TeX`u nezbytné.

Je důležité rovněž připomenout, že implicitní chování `encTeX`u je takové, že pokud se nepoužijí rozšířené primitivy `encTeX`u, pak se chová naprosto stejně jako originální `TeX`.

Podle mého názoru novější implementace web2c  $\TeX$ u taky není v přísném slova smyslu  $\TeX$ . Umožňuje totiž změnu chování programu, pokud na prvním řádku dokumentu za znaky `%&` je cosi specifického napsáno. To je podle mého názoru větší přestupek oproti licenci  $\TeX$ u, než `encTeX`ovými primitivami nastavit ve formátu prostředí systému a pak tyto primitivy v produkční verzi  $\TeX$ u zakázat.

## 2. Překódování byte na byte pomocí vektorů `xord`, `xchr`

### 2.1. Vektory `xord` a `xchr`

Vektory `xord` a `xchr` mají velikost 255 bytů a obsahují informaci o překódování znaku vstupujícího do  $\TeX$ u nebo vystupujícího na terminál a do textových souborů. Jedná se o pole vestavená do programu, přes která jsou filtrovány veškeré textové vstupní a výstupní informace. Má-li znak na vstupu kód  $x$  a chceme, aby měl uvnitř  $\TeX$ u kód  $y$ , pak musí být nastaven vektor `xord` tak, aby `xord[x]=y`. Při zpětném výstupu znaku na terminál, do logu a do souborů zpracovávaných pomocí `\write` platí tato pravidla: Není-li znak s kódem  $y$  označen jako „tisknutelný“, pak vystupuje pomocí přepisu `^^kód y`. Je-li tisknutelný, pak vystupuje s kódem  $z=xchr[y]$ .

Standardně bývají v systémech s kódem ASCII nastaveny hodnoty těchto vektorů tak, že

`xord[i]=xchr[i]=i` pro všechna  $i$  v rozsahu 0 až 255.

Na systémech, které nepoužívají ASCII, se může mapovat 94 tisknutelných ASCII znaků jinam. Mimoto je deklarovaná vlastnost „tisknutelnosti“ znaku v ASCII takto: Znak je tisknutelný, pokud má kód  $y$  v rozsahu 32 až 126. Ostatní znaky se považují za netisknutelné a  $\TeX$  je standardně přepisuje pomocí dvojité stříšky.

Po instalaci balíčku `encTeX` je možno přímo nastavovat a číst obsahy vektorů `xord` a `xchr` prostřednictvím primitiv `\xordcode` a `\xchrcode` a dále nastavovat vlastnost „tisknutelnosti“ znaku pomocí primitivu `\xprncode`. Syntaxe všech tří nových primitiv je naprosto stejná, jakou známe například u primitiv `\lccode` a `\uccode`. Například:

```
\xordcode"AB="CD \xchrcode\xordcode"AB="AB \the\xchrcode200
```

nastavuje `xord[0xAB]=0xCD`; `xchr[xord[0xAB]]=0xAB` a dále vytiskne hodnotu `xchr[200]`.

Na rozdíl od podobných primitiv `\catcode`, `\lccode`, `\sfcode` a dalších však nově zavedené primitivy mají jednu podstatnou výjimku. Reprezentují interní registry  $\TeX$ u, které vždy mají globální platnost. Proto je nastavení `\xordcode` a `\xchrcode` uvnitř skupiny za všech okolností globální, ačkoli to explicitně nepíšeme. Ústupem z požadavku na možnost lokálního deklarování hodnot jsem dosáhl podstatně větší efektivity výsledného kódu programu.

### 2.2. Tisknutelnost znaků nastavená pomocí `\xprncode`

Primitiv `\xprncode` umožňuje nastavovat vlastnost „tisknutelnosti“ znaku takto: Znak s kódem  $y$  je tisknutelný právě tehdy, když je  $y$  v rozsahu 32 až 126 nebo je `\xprncode y > 0`. Napíšeme-li například `\xprncode255=1`, bude tisknutelný znak s kódem 255. Na druhé straně, nastavení `\xprncode'a` třeba na nulu nemá na chování programu žádný vliv, protože kód znaku `a` je v rozsahu 32 až 126. Tímto opatřením program vykazuje určitý pud sebezáchovy, protože zlý uživatel by mu mohl nastavit všechny znaky jako netisknutelné a program by ztratil schopnost se vyjadřovat. Hodnoty `\xprncode` lze nastavit jako u ostatních nových primitiv v rozsahu nula až 255, ovšem otázka tisknutelnosti je totožná s otázkou na kladnou hodnotu bez ohledu na to, jak velká tato hodnota je.

Výchozí hodnoty pro kódování v době `iniTeX`u jsou následující:

- `\xordcode i = i` pro všechna  $i$  v rozsahu 0...255,
- `\xchrcode i = i` pro všechna  $i$  v rozsahu 0...255,
- `\xprncode i = 0` pro  $i$  v rozsahu 0...31, 127...255,
- `\xprncode i = 1` pro  $i$  v rozsahu 32...126.

První dva řádky jsou pravdivé jen na operačních systémech, které přijaly kódování anglické abecedy podle ASCII. Pokud tomu tak není, pak jsou výchozí hodnoty vektorů `xord` a `xchr` pozměněny tak, aby mapovaly tisknutelné znaky podle systému do ASCII uvnitř  $\TeX$ u. Taková změna se týká jen 95 základních tisknutelných znaků, které jsou v ASCII na pozicích 32 až 126.

### 3. Konverze více bytů na jeden byte nebo kontrolní sekvenci

Od verze Dec 2002 enc $\TeX$  umí také konvertovat na úrovni vstupního procesoru více bytů na jeden byte nebo kontrolní sekvenci. Při výstupu do logu a  $\backslash\text{write}$  souborů je pak tento objekt zpětně převeden na původních více bytů. Tato vlastnost nechce nahradit chybějící interpret regulárních výrazů ve vstupním procesoru  $\TeX$ u. Byla implementována pouze z důvodu umožnit pracovat s UTF-8 kódovanými soubory v běžném 8 bitovém  $\TeX$ u tak, že znaky z UTF-8 z nejčastěji používané abecedy mohou být mapovány na jeden znak, který může mít svůj  $\backslash\text{catcode}$ ,  $\backslash\text{uccode}$  atd. Jiné znaky z UTF-8 mohou být mapovány na libovolné kontrolní sekvence.

Pro nastavení takové konverze jsou do  $\TeX$ u přidány nové čtyři primitivy:  $\backslash\text{mubytein}$ ,  $\backslash\text{mubyteout}$ ,  $\backslash\text{mubyte}$  a  $\backslash\text{endmubyte}$ . Primitivy  $\backslash\text{mubytein}$  a  $\backslash\text{mubyteout}$  jsou celočíselné registry implicitně s nulovou hodnotou, tj. konverze vstupu a výstupu podle konverzní tabulky se neprovádějí. Je-li  $\backslash\text{mubytein}$  nastaveno na kladnou hodnotu,  $\TeX$  okamžitě zahájí konverze vstupního řádku podle konverzní tabulky. Je-li  $\backslash\text{mubyteout}$  nastaveno na kladnou hodnotu,  $\TeX$  začne konvertovat do logu a výstupních souborů podle stejné konverzní tabulky. Implicitně je konverzní tabulka prázdná a jednotlivé řádky se do ní přidávají pomocí dvojice primitivů  $\backslash\text{mubyte}$ ,  $\backslash\text{endmubyte}$  s touto syntaxí:

```
 $\backslash\text{mubyte}$   $\langle\text{first\_token}\rangle\langle\text{one\_optional\_space}\rangle\langle\text{optional\_prefix}\rangle\langle\text{byte\_sequence}\rangle\backslash\text{endmubyte}$ 
```

Každá  $\langle\text{byte\_sequence}\rangle$  bude převedena ve vstupním procesoru na  $\langle\text{first\_token}\rangle$ . Je-li  $\langle\text{first\_token}\rangle$  znakem (tj. není to kontrolní sekvence), pak se ignoruje jeho kategorie, protože konverze je prováděna v input procesoru podle schématu:  $\langle\text{byte\_sequence}\rangle$  na jeden  $\langle\text{byte}\rangle$ . Při výstupu do logu a  $\backslash\text{write}$  souborů se pak každý takový  $\langle\text{byte}\rangle$  znovu převede na  $\langle\text{byte\_sequence}\rangle$ .

Pokud je  $\langle\text{first\_token}\rangle$  kontrolní sekvence, pak se na úrovni vstupního procesoru promění každá  $\langle\text{byte\_sequence}\rangle$  na tuto kontrolní sekvenci implementovanou ve formě neměnitelného tokenu. Token procesor tuto sekvenci tedy znovu neinterpretuje a zůstává za ní ve stavu neignorování mezer. Při  $\backslash\text{mubyteout}<2$  není při výstupu do  $\backslash\text{write}$  souborů tato kontrolní sekvence zpětně převáděna na původní  $\langle\text{byte\_sequence}\rangle$ , ale podléhá jen běžné expanzi, jako ostatní kontrolní sekvence. Při  $\backslash\text{mubyteout}>=2$  se i tyto kontrolní sekvence převádějí do  $\backslash\text{write}$  souborů na původní  $\langle\text{byte\_sequence}\rangle$ . Aby se ale mohly převést, nesmějí před tím expandovat, tj. musejí mít v době expanze význam neexpandovatelné kontrolní sekvence nebo musejí být označeny pomocí  $\backslash\text{noexpand}$ . Při  $\backslash\text{mubyteout}>=3$  enc $\TeX$  potlačí expanzi kontrolních sekvencí deklarovaných v  $\backslash\text{mubyte}$  automaticky (podrobněji viz sekci 3.7). Výstup do logu a na terminál, který není produktem příkazu  $\backslash\text{write}$ , ponechává kontrolní sekvence nezměněny.

Význam  $\langle\text{optional\_prefix}\rangle$  je vysvětlen v sekci 3.4.

#### 3.1. Zanášení údajů do konverzní tabulky

Záznamy do konverzní tabulky jsou pomocí primitivů  $\backslash\text{mubyte}$ ,  $\backslash\text{endmubyte}$  zanášeny globálně, zatímco hodnoty v registerch  $\backslash\text{mubytein}$  a  $\backslash\text{mubyteout}$  mají obvykle lokální platnost.

Dvojice primitivů  $\backslash\text{mubyte}$ ,  $\backslash\text{endmubyte}$  pracuje analogicky, jako dvojice  $\backslash\text{csname}$ ,  $\backslash\text{endcsname}$ . Rozdíl je pouze v tom, že první token  $\langle\text{first\_byte}\rangle$  se neexpanduje a že za ním může (po expanzi) následovat  $\langle\text{one\_optional\_space}\rangle$ . Při skenování  $\langle\text{optional\_prefix}\rangle$  a  $\langle\text{byte\_sequence}\rangle$  již probíhá úplná expanze a při ní se nesmí objevit na vstupu do hlavního procesoru token typu kontrolní sekvence, jinak nastane chyba, kterou už známe z používání  $\backslash\text{csname}$ ,  $\backslash\text{endcsname}$ :

```
! Missing  $\backslash\text{endmubyte}$  inserted.
```

Primitiv  $\backslash\text{mubyte}$  na rozdíl od  $\backslash\text{csname}$  neprovádí činnost na úrovni expand procesoru, ale jedná se o přiřazovací primitiv zpracovaný na úrovni hlavního procesoru. Takže po

```
 $\backslash\text{edef}\backslash\text{a}\{\backslash\text{mubyte X ABC}\backslash\text{endmubyte}\}$ 
```

bude makro  $\backslash\text{a}$  obsahovat tokeny:  $\backslash\text{mubyte X ABC}\backslash\text{endmubyte}$ .

Příklady:

```
 $\backslash\text{mubyte}$   $\text{~}^{\text{c}}1$   $\text{~}^{\text{c}}3\text{~}^{\text{c}}81\backslash\text{endmubyte}$  % Ā  
 $\backslash\text{mubyte}$   $\text{~}^{\text{e}}1$   $\text{~}^{\text{c}}3\text{~}^{\text{a}}1\backslash\text{endmubyte}$  % á  
% atd. -- implementace UTF8
```

```
 $\backslash\text{mubyte}$   $\backslash\text{endash}$   $\text{~}^{\text{c}}4\text{~}^{\text{f}}6\backslash\text{endmubyte}$  % příklad na kontrolní sekvenci  
 $\backslash\text{mubyte}$   $\backslash\text{integral}$  INT $\backslash\text{endmubyte}$  % příklad pro ilustraci, viz dále.
```

```
\mubytein=1 \mubyteout=1 % od této chvíle je překódování aktivní
```

```
\def\endash {--}  
\def\integral {\ifmmode \int\else $\int$\fi}
```

V tomto příkladě je v místě *<one optional space>* více mezer a tabulátorů. Protože tabulátory mají kategorii mezery, jsou všechny tyto znaky přeměněny token procesorem na jedinou mezeru požadovanou v syntaktickém pravidle pro `\mubyte`, `\endmubyte`.

Po použití definic z příkladu se slovo INTEGRAL promění v token `\integral` okamžitě následovaný písmeny „EGRAL“. V textu INT EGRAL bude za tokenem `\integral` mezera a teprve pak písmena „EGRAL“. Také jsou možné konstrukce typu `\defINT{něco}` apod. Pokud je `\integral` nedefinovaná kontrolní sekvence, pak si při použití slova INTEGRAL budeme muset zvyknout na poněkud podivnou chybovou hlášku:

```
! Undefined control sequence.  
1.13 tady je slovo INT  
          EGRAL.
```

Když napíšeme `\show INT`, dostaneme odpověď:

```
> \integral=undefined.  
1.13 \show INT
```

a `\string INT` se expanduje na text: `\integral`.

Po deklaraci INT podle předchozího příkladu se může stát, že někdo napíše: `\INT`. Správně by to mělo vést na prázdnou kontrolní sekvenci (`\csname\endcsname`) následovanou kontrolní sekvencí `\integral`. Protože se ale s prázdnými kontrolními sekvencemi v  $\TeX$ u moc často nepracuje a pro uživatele by to mohlo být matoucí, rozhodl jsem se tuto situaci ošetřit tak, že `\INT` je převedeno pouze na token `\integral`. Pozor na skutečnost, že za sekvencí `\INT` není  $\TeX$  ve stavu ignorování mezer a navíc může za ní okamžitě následovat písmeno.

### 3.2. Vlastnosti konverze

Multibytové sekvence jsou převedeny ze vstupu pouze tehdy, pokud jsou celé obsaženy v jediném řádku. Přesah do dalšího řádku není možný. Připojený `\endlinechar` na konci řádku se může stát předmětem konverze podle konverzní tabulky.

Sekvence `^^c3^81` se nepromění ani po použití definic z příkladu na byte „Á“, protože převod dvojitých zobáků na jednotlivé byty probíhá v token procesoru, tj. později, než převody více bytů na jeden podle `\mubyte`.

Převod více bytů na jeden byte nebo kontrolní sekvenci probíhá později než konverze podle `\xordcode` a při výstupu do `\write` a log souborů pak převod podle `\mubyte` probíhá dříve než konverze podle `\xchrcode`. *<byte\_sequence>* tedy musí obsahovat sekvenci bytů tak, jak jsou tyto byty konvertovány ze vstupního souboru pomocí `\xordcode`.

Postupné procesy na vstupu a výstupu si můžeme naznačit takto:

```
vstupní text -> \xordcode -> připojení \endlinechar ->  
                \mubyte -> token procesor -> expanze ...  
argument \write -> expanze -> \mubyte -> \xchrcode -> výstup
```

Při výstupu do `\write` souborů a logů zpětně konvertované *<byte\_sequence>* už nepodléhají další konverzi na formát typu `^^c3^81` ani opakované konverzi podle `\mubyte`. Tyto *<byte\_sequence>* se pouze převedou podle hodnot `\xchrcode`.

Výstup do logu a na terminál má výjimku z pravidla, že je tento výstup modifikován podle hodnot konverzní tabulky a primitivu `\mubyteout`. Jedná se o případy nezměněného přepisu vstupních řádků do výstupu. Je-li `\mubytein` kladný, pak je vstupní řádek přepsán do logu nebo na terminál bez `\mubyte` konverze tam ani zpět. Aktivní zůstávají jen konverze podle `\xchrcode` a `\xordcode`. Byl-li například vstupní řádek kódován v UTF-8, bude tento řádek byte po byte shodně vypadat v logu a na terminálu. Přepisy na tvar `^^aa` jsou v takovém případě také potlačeny bez závislosti na hodnotě `\xprncode`. Tato výjimka se týká např. přepisu vstupního řádku při chybě, kdy  $\TeX$  roztržením tohoto řádku na terminálu

a v logu dává na jevo místo, kde došlo k chybě. Výjimka se samozřejmě nedotýká výstupů příkazu `\write` a `\message`, které jsou vždy podmíněny stavem konverzní tabulky a hodnotou primitivu `\mubyteout`.

### 3.3. Co se stane při sporných údajích v konverzní tabulce

Pokud existují v konverzní tabulce dvě *byte\_sequence* se stejným začátkem tak, že jedna je případně podsekvencí druhé, pak má přednost později zadaná hodnota v tabulce a dříve zadaná hodnota je zcela ignorována. Příklad:

```
\mubyte X ABC\endmubyte
{\mubytein=1 nyní se ABC konvertuje na X}
\mubyte Y ABCDE\endmubyte
\mubyte W ABFG\endmubyte
{\mubytein=1 nyní se ABC nemění a ABCDE se konvertuje na Y
 a ABFG se konvertuje na W}
\mubyte Z AB\endmubyte
{\mubytein=1 nyní se ABCDE promění na ZCDE}
```

Tato konvence umožňuje vymazat z tabulky všechny řádky, kde *byte\_sequence* mají společné první písmeno tak, že napíšeme jednoznakovou *byte\_sequenci*, která se má konvertovat na stejný znak. Například:

```
\mubyte A A\endmubyte
```

odstraní z konverzní tabulky všechny *byte\_sequence* začínající písmenem A. Je-li *first\_token* roven *byte\_sequenci*, pak primitiv `\mubyte` opravdu maže z konverzní tabulky řádky začínající na *first\_token* a tím uvolňuje hlavní paměť T<sub>E</sub>Xu, kde jsou tato data uložena. Ve všech ostatních případech primitiv `\mubyte` pouze zakládá další řádek do konverzní tabulky s tím, že některé předchozí řádky tam mohou zůstat neaktivní.

Následující kód promaže celou tabulku:

```
{\catcode'\^^@=12
\gdef\clearmubytes{\bgroup \count255=1
 \loop \uccode'X=\count255
 \uppercase{\mubyte XX\endmubyte}%
 \advance\count255 by1
 \ifnum\count255<256 \repeat
 \mubyte ^^@^^@\endmubyte
 \egroup}
}
```

### 3.4. Vstupní a výstupní část konverzní tabulky

Konverzní tabulka konstruovaná pomocí `\mubyte`, `\endmubyte` má dvě nezávislé části: vstupní, se kterou pracuje input procesor a výstupní, která se používá při zpětných konverzích. Údaje je možné zanést nezávisle do každé části při použití neprázdného *optional\_prefix* (viz syntaktické pravidlo `\mubyte` na začátku této kapitoly). Je-li *optional\_prefix* prázdný, pak se požadavek na konverzi zanese dvojmo do vstupní i výstupní části. Je-li ale *optional\_prefix* znak kategorie 8 (obvykle znak `_`), pak se údaj zanese jen do vstupní části tabulky. Obsahuje-li *optional\_prefix* dvojici znaků kategorie 8 (obvykle tedy `--`), pak se údaj zanese jen do výstupní části tabulky.

Při *optional\_prefix* `--` (výstupní část tabulky) je dovoleno mít prázdnou *byte\_sequenci*. V takovém případě se původní údaje z výstupní části tabulky, které odpovídají *first\_token*, vymažou. Pokud tam žádné takové údaje nebyly, nestane se nic.

Vraťme se ke kódu na vymazání tabulky z předchozí sekce. Tento kód vymaže vše ze vstupní části tabulky a z výstupní jen ty údaje, které jsou vázány na *first\_token* ve tvaru *byte*. Údaje z výstupní části vázané na kontrolní sekvence nejsou tímto způsobem promazány. Promazání jednoho údaje uděláme pomocí: `\mubyte \foo --\endmubyte`.

### 3.5. Vkládání dalších kontrolních sekvencí

Je-li  $\langle first\_token \rangle$  ve tvaru kontrolní sekvence a navíc  $\langle optional\_prefix \rangle$  je token kategorie 6 (obvykle znak #), pak  $\langle byte\_sequence \rangle$  zůstane zachována, jen před ní vloží input procesor deklarovanou kontrolní sekvenci. Údaj se zanese jen do vstupní části konverzní tabulky. Příklad použití:

```
\mubyte \warntwobytes #^^c3^^80\endmubyte
\mubyte \warntwobytes #^^c3^^82\endmubyte
\mubyte \warntwobytes #^^c3^^83\endmubyte
% atd...
\def\warntwobytes #1#2{\bgroup\mubyteout=0
  \message{WARNING: the UTF8 code: #1#2 is not defined i my macros.}
  \egroup}
```

Při  $\mubytein=1$  a při vkládání kontrolních sekvencí je zachování  $\langle byte\_sequence \rangle$  absolutní, tj. žádná část  $\langle byte\_sequence \rangle$  nepodléhá další konverzi. Na druhé straně při  $\mubytein>1$  je možná další konverze některé části  $\langle byte\_sequence \rangle$ .

```
\mubyte \foo #ABC\endmubyte \mubyte X BC\endmubyte
\mubytein=1 Nyní ABC přechází na \foo ABC
\mubytein=2 Nyní ABC přechází na \foo AX
```

### 3.6. Rozpoznání začátku řádku

Existují-li v konverzní tabulce  $\langle byte\_sequence \rangle$  s prvním znakem shodným s aktuálním  $\endlinechar$ , tj.  $\langle byte\_sequence \rangle$  jsou ve tvaru  $\langle endlinechar \rangle \langle zbytek \rangle$ , pak input procesor navíc ověřuje, zda je  $\langle zbytek \rangle$  shodný se začátkem každého řádku. Pokud ano, provede požadovanou konverzi. Příklad použití:

```
\bgroup \uccode'X=\endlinechar \uppercase{\gdef\echar{X}}\egroup
\mubyte \fooB \echar ABC\endmubyte % vyhovuje ABC na začátku řádku
\mubyte \fooE ABC\echar \endmubyte % vyhovuje ABC na konci řádku
\mubyte \fooW \spce\space ABC\space \endmubyte
% vyhovuje ABC jako slovo s mezerami vpředu i vzadu
\mubyte \foo #\echar ABC\endmubyte %
% je-li ABC na začátku řádku, vloží před něj \foo
```

### 3.7. Potlačení expanze v parametrech write

Chceme-li převádět kontrolní sekvence zpětně na  $\langle byte\_sequence \rangle$  při zápisu do  $\write$  souborů, musíme potlačit případnou expanzi těchto kontrolních sekvencí například pomocí  $\let\macro=\relax$ . Protože ale  $\write$  často pracuje asynchronně a kontrolních sekvencí mapujících UTF-8 můžeme mít stovky nebo tisíce, umožňuje encTeX nastavit v době expanze parametrů  $\write$  příslušným kontrolním sekvencím význam  $\relax$  automaticky. Dělá to při  $\mubyteout>=3$  a význam  $\relax$  přiřadí právě těm kontrolním sekvencím, které mají ve výstupní části konverzní tabulky neprázdnou  $\langle byte\_sequenci \rangle$ . Jakékoli jiné expanze mimo parametr  $\write$  probíhají normálním způsobem. Příklad:

```
\mubyte \foo ABC\endmubyte \def\foo{macro body}
\mubyteout=2
\immediate\write16{testwrite: \foo} % zapíše "testwrite: macro body"
\immediate\write16{testwrite: \noexpand\foo} % zapíše "testwrite: ABC"
\mubyteout=3
\immediate\write16{testwrite: \foo} % zapíše "testwrite: ABC"
\message{testmessage: \foo} % zapíše "testmessage: macro body"
\message{testmessage: \noexpand\foo} % zapíše "testmessage: \foo"
\edef\a{testedef: \foo} % expanduje na macro body
\foo % expanduje na macro body
\immediate\write16{\meaning\foo} % zapíše "\relax"
\message{\meaning\foo} % zapíše "macro:->macro body"
```

Pomocí zápisu  $\mubyte \langle control\_sequence \rangle \relax \endmubyte$  je možno přidělit kontrolní sekvenci příznak, aby se neexpandovala v parametrech  $\write$  při  $\mubyteout>=3$ , ale na druhé straně nebude

konvertována do žádné  $\langle byte\_sequence \rangle$ , ale vypíše se jako obvykle. Uvedený zápis má tedy stejný význam jako  $\backslash mubyte \langle control\_sequence \rangle \_ \_ \backslash string \langle control\_sequence \rangle \backslash space \backslash endmubyte$ , ale navíc šetří paměť  $\TeX$ u, neboť  $\TeX$  není nucen ukládat string  $\langle byte\_sequence \rangle$  do poolu.

### 3.8. Asynchronní zpracování příkazu write

Je známo, že pokud nepoužijeme  $\backslash immediate$ , pak se argument příkazu  $\backslash write$  expanduje až později: ne v okamžiku výskytu příkazu. Příkaz  $\backslash write$  si proto uloží do své paměti aktuální hodnotu registru  $\backslash mubyteout$  v době prvního zpracování a pak při expanzi a zápisu do souboru tuto hodnotu použije.

Díky této vlastnosti můžeme třeba pro soubor s obsahem zapisovat s hodnotou  $\backslash mubyteout=3$  a současně při zápisu do jiného souboru ponecháme hodnotu  $\backslash mubyteout=0$ . To může být žádoucí například proto, že soubor je určen ke zpracování programem, který nemá implementovanou schopnost práce s UTF-8 kódováním. Vyzkoušejte:

```
\newwrite\tocfile \newwrite\indexfile
\immediate\openout\tocfile=\jobname.toc
\immediate\openout\indexfile=\jobname.idx
\mubyteout=3
\write\tocfile{parametr se bude později konvertovat do UTF-8}
{\mubyteout=0 \write\indexfile{parametr zůstane nezměněný bez konverze}}
\write\tocfile{zde se znovu provede konverze}
\end % a teprve v tento okamžik se všechny tři zápisy provedou
```

### 3.9. Hodnoty registru mubyteout

Kromě již zmíněných hodnot 0, 1, 2 a 3 registru  $\backslash mubyteout$  může být někdy užitečné nastavit tento registr na hodnoty  $-1$ ,  $-2$  a  $-3$ . Význam těchto hodnot je vysvětlen v následující tabulce:

$\backslash mubyteout$	$\langle byte \rangle \rightarrow \langle byte\_sequence \rangle$	$\langle cs\_name \rangle \rightarrow \langle byte\_sequence \rangle$	potlačení expanze
0	ne	ne	ne
1	ano	ne	ne
2	ano	ano	ne
3	ano	ano	ano
-1	ano	ne	ano
-2	ne	ne	ano

Je-li zapnutá konverze  $\langle byte \rangle \rightarrow \langle byte\_sequence \rangle$ , pak se tato konverze provádí i do logu a na terminál, zatímco konverze  $\langle cs\_name \rangle \rightarrow \langle byte\_sequence \rangle$  a potlačení expanze se týkají jen argumentů  $\backslash write$  a  $\backslash special$ .

## 4. Argumenty primitivu special

V argumentech  $\backslash special$  se často objevují texty v přirozeném jazyce (například texty pro záložky do PDF dokumentu). Při stále častějším používání UTF-8 je žádoucí, aby tyto texty byly kódovány v tomto kódování. Enc $\TeX$  tuto možnost nabízí.

Argument primitivu  $\backslash special$  je zpracován podle hodnoty celočíselného registru  $\backslash specialout$ , který má implicitní hodnotu 0.

- $\backslash specialout=0$  – žádná konverze argumentu se neprovede.
- $\backslash specialout=1$  – provede se konverze jen podle vektoru  $xchr$ .
- $\backslash specialout=2$  – provede se konverze jen podle hodnoty  $\backslash mubyteout$ .
- $\backslash specialout=3$  – provede se konverze podle hodnoty  $\backslash mubyteout$  následovaná konverzí podle  $xchr$ .

Primitiv  $\backslash special$  expanduje svůj argument okamžitě. Při  $\backslash specialout$  2 nebo 3 se expanze provede podle hodnoty  $\backslash mubyteout$  stejně jako u primitivu  $\backslash write$ . Pak si příkaz  $\backslash special$  uloží do paměti aktuální hodnoty  $\backslash specialout$  a  $\backslash mubyteout$  a tyto hodnoty použije ještě jednou při skutečném výstupu argumentu do dvi souboru.

## 5. Dokumentace k přiloženým souborům maker

Tato část dokumentace nebyla ve verzi Dec. 2002 revidovaná a je ponechána ve stavu z roku 1997 s výjimkou následujícího odstavce.

### 5.1. Kódování UTF-8

Pro vstupní kódování UTF-8 jsou připraveny soubory `utf8-csf.tex` a `utf8-t1.tex`. V tomto případě je překódování implementováno pomocí `\mubyte` a vektory `xord`, `xchr` jsou nastaveny tak, že na jejich úrovni je zachováno identické zobrazení.

### 5.2. Formáty typu plain-x-y

V balíčku jsou připraveny inicializační soubory pro vygenerování formátu podobnému standardnímu formátu `plain`. Například příkazem

```
$ tex -i plain-1250-cs
```

vygenerujeme formát analogický `plainu`, který čte vstupní soubory v kódování CP1250 a pracuje s CS-fonty.

V balíku jsou k dispozici tyto inicializační soubory pro `plain`:

```
plain-il2-cs ... vstup podle ISO8859-2, textové fonty v TeXu: CS-font
plain-kam-cs ... vstup podle Kamenických, textové fonty v TeXu: CS-font
plain-1250-cs ... vstup podle CP1250, textové fonty v TeXu: CS-font
plain-852-cs ... vstup podle CP852, textové fonty v TeXu: CS-font
plain-il2-dc ... vstup podle ISO8859-2, textové fonty v TeXu: DC
plain-kam-dc ... vstup podle Kamenických, textové fonty v TeXu: DC
plain-1250-dc ... vstup podle CP1250, textové fonty v TeXu: DC
plain-852-dc ... vstup podle CP852, textové fonty v TeXu: DC
```

### 5.3. Poznámka k dlouhým názvům souborů

Všechny soubory `*.tex` v balíčku splňují DOSové omezení na délku názvu 8+3. Výjimkou z tohoto pravidla jsou pouze soubory `plain-x-y` popsané výše a analogické inicializační soubory pro `LaTeX`. Pokud používáte systém, který je omezen na 8+3, doporučuji pro každé kódování zvolit jedno písmeno (například `c=cs`, `d=dc`, `i=il2`, `w=1250`, `p=852`, `k=kam`, `o=koi8`, `m=mac`) a nahradit názvy souborů v distribuci těmito názvy:

```
plain-il2-cs.tex    plain-ic.tex
plain-kam-cs.tex    plain-kc.tex
plain-1250-cs.tex   plain-wc.tex
plain-852-cs.tex    plain-pc.tex
plain-il2-dc.tex    plain-id.tex
plain-kam-dc.tex    plain-kd.tex
plain-1250-dc.tex   plain-wd.tex
plain-852-dc.tex    plain-pd.tex
kam-latex.tex       latex-ki.tex
852-latex.tex       latex-pi.tex
```

Obsah `\message` v souborech `plain-x-y` neměňte. Například formát `plain-wc` se po spuštění představí svým plným jménem

```
The format: plain-1250-cs <Sep. 1997>.
```

### 5.4. Kódovací tabulky

Protože změna vektorů `xord` a `xchr` může totálně rozhodit chování `TeXu` zcela k nepoznání, doporučuji používat určité soubory, které nastaví požadované kódování, a dále s primitivy `\xordcode`, `\xchrcode` a `\xprncode` za běhu `TeXu` moc nelaškovat. V balíčku `encTeX` jsou k dispozici soubory, které změnu vektorů pro běžná kódování definují. Tyto soubory mají obvyklou příponu `tex`. Říkáme jim kódovací tabulky. Rozlišujeme dva typy kódovacích tabulek.

## 5.5. První typ kódovacích tabulek

První typ tabulek deklaruje vnitřní kódování  $\TeX$ u ve vztahu ke kódování, které je běžně používané v hostitelském operačním systému. Máme-li například v systému kódování ISO-8859-2 a vnitřní kódování  $\TeX$ u volíme podle Corku (kódování je označováno jako T1), pak tabulka musí předefinovat xord vektor tak, aby mapoval znaky z ISO-8859-2 do T1 a vektor xchr musí převádět zpátky z T1 do kódování systému.

Tento typ tabulek je použit v inicializačních souborech `plain-*.tex` a obsahuje v názvu souboru vstupní i cílové vnitřní kódování  $\TeX$ u. Podívejte se, jak vypadá například tabulka `il2-t1.tex`, která definuje vnitřní kódování  $\TeX$ u podle Corku a vstupní kódování ISO8859-2.

Každá tabulka prvního typu čte soubor `encmacro.tex` s definicemi maker `\setcharcode`, `\expandto`, `\texaccent`, `\texmacro` a `\redefaccent`.

- `\setcharcode #1 #2 #3 #4 #5 #6 #7` deklaruje  $\TeX$ ové kódy pro jeden znak. Nastaví `xord[#1]=#2`, `xchr[#2]=#1`, `\xprncode#2=#7` a postupně nastaví `\lccode`, `\uccode`, `\sfcode` a `\catcode` znaku s kódem #2 na hodnoty #3, #4, #5 a #6. Je-li #1 otazník, pak se xord a xchr nenastaví.
- `\expandto {<definice>}` definuje aktivní podobu znaku #2 z posledního `\setcharcode` tak, že tento token expanduje na `<definici>`. Podrobněji: je-li v `\setcharcode` uvedeno #6=13, pak bude každý výskyt znaku #2 expandovat na `<definici>`. Není-li v `\setcharcode` řečeno #6=13, pak k expanzi znaku #2 na `<definici>` dojde teprve tehdy, když bude (třeba později) nastaveno `\catcode` znaku #2 na 13.
- `\texaccent` uvzápis akcentu připravuje expanzi „zápisu akcentu“ na znak s kódem #2 z naposledy použitého `\setcharcode`. Například zápis `\v C` bude po načtení souboru `il2-t1.tex` expandovat na znak s kódem "83. Pokud zápis pro akcent není v tabulce uveden, zůstává v původním významu, tj. třeba `\v g` expanduje na primitiv `\accent`, který usadí háček nad písmeno g. K aktivaci všech „zápisů akcentu“ dojde až po použití makra `\redefaccent` (viz níže).
- `\texmacro #1` deklaruje makro #1 tak, že bude expandovat na znak s kódem #2 z naposledy použitého `\setcharcode`. K předefinování makra #1 dojde (na rozdíl od `\texaccent`) okamžitě. Například makro `\S` bude po načtení souboru `il2-t1.tex` expandovat na znak s kódem 9F, protože na této pozici je podle Corku znak paragraf.
- `\redefaccent #1` aktivuje expanzi zápisů podle `\texaccent` pro jeden konkrétní akcent #1.

Kromě toho je na začátku tabulky čten soubor definic závislých na kódování textového fontu  $\TeX$ u. V naší ukázce jde například o soubor `t1macro.tex`. Definují se tam sekvence `\promile`, `\clqq` a další.

Může se stát, že nechceme uvedená makra použít, ale hodnoty z tabulky načíst chceme. Pak můžeme přistoupit k následujícímu triku: Definujeme si makra `\setcharcode` až `\redefaccent` sami a dále provedeme načtení tabulky takto:

```
\let\originput=\input \def\input #1 \originput il2-t1
\let\input=\originput
```

V balíčku jsou připraveny tyto tabulky prvního druhu:

Název souboru	vstupní kódování	vnitřní kódování $\TeX$ u
<code>il2-csf.tex</code>	ISO8859-2	CS-font
<code>kam-csf.tex</code>	Kamenických	CS-font
<code>1250-csf.tex</code>	CP1250, MS-Windows	CS-font
<code>852-csf.tex</code>	CP852, PC Latin2	CS-font
<code>il2-t1.tex</code>	ISO8859-2	T1 alias Cork
<code>kam-t1.tex</code>	Kamenických	T1 alias Cork
<code>1250-t1.tex</code>	CP1250, MS-Windows	T1 alias Cork
<code>852-t1.tex</code>	CP852, PC Latin2	T1 alias Cork

Za zmínku stojí první uvedená tabulka `il2-csf.tex`, protože ta jediná ponechává vektory xord a xchr beze změny. Tuto tabulku je tedy možné použít i v  $\TeX$ u, který neobsahuje rozšíření `enc $\TeX$` . Všechny ostatní tabulky `enc $\TeX$`  explicitně vyžadují.

## 5.6. Druhý typ kódovacích tabulek

Druhý typ tabulek provádí překódování pouze na vstupní straně  $\TeX$ u. Poznáme je podle toho, že nemají na konci názvu značku pro vnitřní kódování  $\TeX$ u (tj. `t1` nebo `csf`), ale značku používanou

pro kódování operačního systému (např. `il2`, `kam`). Třeba tabulka `kam-il2.tex` provádí na vstupní straně konverzi z kódování kamenických do kódování ISO8859-2. Tento typ tabulek pozměňuje pouze vektor `xchr`, ale výstupní vektor `xord` ponechává beze změny. Takovou tabulku použijeme, pokud  $\TeX$ em načítáme soubor, který je v jiném kódování, než běžně používáme na našem operačním systému. Přitom výstup do `log`, `aux` apod. ponecháme v kódování podle našeho systému. Tyto změny kódování je možné provádět i v průběhu zpracování jediného dokumentu.

Druhý typ tabulek navazuje na vstupní kódování deklarované dříve tabulkou prvního typu. Nastavení vnitřního kódování  $\TeX$ u není vůbec druhým typem tabulek měněno. Uvedeme příklad. Při generování formátu jsme použili tabulku prvního typu `il2-t1.tex`, takže vnitřní kódování máme podle Corku. Nyní můžeme při zpracování dokumentu na přechodnou dobu vybrat některou z tabulek `*-il2.tex`, třeba:

```
\input kam-il2
\input dokument
\restoreinputencoding
nyní mohu pracovat v původním kódování...
\end
```

V době, kdy probíhá načítání souboru `dokument.tex` se provádí překódování z Kamenických do T1, uvnitř  $\TeX$ u se vše zpracovává v T1 a výstup na terminál a do logu máme v ISO8859-2. V tomto kódování je také zapsán další text pod `\restoreinputencoding`. Tabulka totiž deklaruje toto makro, aby byl možný návrat k původnímu nastavení vektoru `xord`.

Při použití tabulek druhého typu musíme dát velký pozor, abychom něco neudělali špatně. V našem příkladě jsou všechny výstupy do souborů typu `aux` v ISO-8859-2, takže je při opakovaném spuštění  $\TeX$ u nesmíme načítat v okamžiku, kdy máme nastaven vstupní kód podle Kamenických. To je také důvod, proč nedoporučuji generovat formát příkazem `\dump` v situaci, kdy máme načtenou tabulku druhého typu.