

# *Cahiers* **GUT**enberg

☞ LE TRAITEMENT DES TEXTES POLONAIS  
AVEC LE LOGICIEL T<sub>E</sub>X

☞ Hanna KOŁODZIEJSKA

*Cahiers GUTenberg*, n° 0 (1988), p. 3-10.

<[http://cahiers.gutenberg.eu.org/fitem?id=CG\\_1988\\_\\_0\\_3\\_0](http://cahiers.gutenberg.eu.org/fitem?id=CG_1988__0_3_0)>

© Association GUTenberg, 1988, tous droits réservés.

L'accès aux articles des *Cahiers GUTenberg*

(<http://cahiers.gutenberg.eu.org/>),

implique l'accord avec les conditions générales

d'utilisation (<http://cahiers.gutenberg.eu.org/legal.html>).

Toute utilisation commerciale ou impression systématique  
est constitutive d'une infraction pénale. Toute copie ou impression  
de ce fichier doit contenir la présente mention de copyright.



## LE TRAITEMENT DES TEXTES POLONAIS AVEC LE LOGICIEL T<sub>E</sub>X

Hanna KOŁODZIEJSKA

Instytut Informatyki  
Uniwersytet Warszawski  
PKiN p.850  
00-901 Warszawa, Pologne

### Résumé

*Je présente dans ce document quelques difficultés d'utilisation du logiciel T<sub>E</sub>X pour les textes polonais, liées à l'insertion des lettres avec des signes diacritiques, à la division des mots et à certaines particularités de la typographie polonaise.*

### Introduction

Etant une langue slave, le polonais appartient au groupe des langues indo-européennes. Il est parlé par environ 50 millions de personnes en Pologne et à l'étranger. Le polonais n'est pas facile à apprendre à cause de sa grammaire assez compliquée, son orthographe et sa prononciation difficile, tellement différente de celle des langues romanes par exemple. Il est rarement connu par les étrangers, c'est pourquoi la plupart des documents scientifiques, ceux concernant les mathématiques en particulier, écrits en Pologne, sont publiés en anglais.

### L'alphabet polonais

L'alphabet polonais est basé sur l'alphabet latin. Il est composé de 32 lettres suivantes:

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| a | A | f | F | m | M | ś | Ś |
| ą | Ą | g | G | n | N | t | T |
| b | B | h | H | ń | Ń | u | U |
| c | C | i | I | o | O | w | W |
| ć | Ć | j | J | ó | Ó | y | Y |
| d | D | k | K | p | P | z | Z |
| e | E | l | L | r | R | ź | Ź |
| ę | Ę | ł | Ł | s | S | ż | Ż |

Il y a neuf lettres avec de différents signes diacritiques: ą, ć, ę, ł, ń, ó, ś, ź, ż. Chacune d'elles est prononcée de manière tout à fait différente que la correspondante lettre latine. Par exemple: le mot "róża" (*rose*) est prononcé en polonais de la même façon que le mot

*rouja* serait prononcé en français. Sans signes diacritiques nous prononcerions ce mot comme *rosa* en français.

Parfois l'omission des signes diacritiques change le sens du mot. Par exemple les mots "pólka" et "polka" signifient respectivement *rayon* et *polka* (une danse d'origine tchèque). Ainsi le manque des signes diacritiques dans le texte polonais change essentiellement la prononciation des mots, parfois même leur sens et écorche le polonais.

### Insertion des lettres polonaises

Rares sont en Pologne les claviers comportant tout l'alphabet polonais. Le plus souvent nous nous servons des claviers n'ayant que 26 lettres latines ce qui nécessite l'utilisation des commandes  $\TeX$  pour accentuer certaines lettres. Suivant le signe diacritique nous distinguons quatre groupes des lettres accentuées:

- (1.)  $\acute{c}$   $\acute{n}$   $\acute{o}$   $\acute{s}$   $\acute{z}$  (et les correspondantes lettres majuscules)—nous pouvons les obtenir par surimpression de l'accent aigu sur la lettre latine appropriée, par exemple:

$$\backslash'c \rightarrow \acute{c}$$

- (2.)  $\dot{z}$  ( $\dot{Z}$ )—nous l'obtenons de la manière suivante:

$$\backslash.z \rightarrow \dot{z}$$

$$\backslash.Z \rightarrow \dot{Z}$$

- (3.)  $\l$  ( $\text{Ł}$ )—cette lettre, qui apparaît dans l'alphabet polonais et lituanien, est appelée dans [K]: "Polish suppressed-L" et peut être obtenue à l'aide des commandes suivantes:

$$\backslash l \rightarrow \l$$

$$\backslash L \rightarrow \text{Ł}$$

- (4.)  $\a$   $\e$  ( $\text{Ą}$   $\text{Ę}$ )—ce signe diacritique, appelé "ogonek" dans [ISO], manque dans  $\TeX$ . Il diffère de la cédille par sa forme et avant tout par la courbure inversée. En utilisant la cédille nous obtenons les caractères suivants:

$$\backslash c a \rightarrow \a$$

$$\backslash c e \rightarrow \e$$

Leur qualité typographique n'est pas bonne. "Ogonek" qui apparaît dans ce document, symétrique à la cédille, a été obtenu par une simple modification du fichier `amr10.px1` (le fichier `amr10.tfm` n'a suivi aucun changement).

Les lettres accentuées apparaissent très souvent dans les mots polonais si bien que l'utilisation des commandes  $\TeX$  n'est pas commode. Par exemple pour obtenir la phrase suivante:

"*Żółty liść spadł z dębu*"

(*Une feuille jaune est tombée du chêne*)

il faut écrire:

$$\backslash.Z\backslash'o\backslash l \text{ ty li}\backslash's\backslash'c \text{ spad}\backslash\backslash \text{ z d}\backslash c \text{ ebu}$$

ce qui est pratiquement illisible.

D'autres exemples: après l'insertion de la suite des caractères:

$$\text{wzi}\backslash c \text{ a}\backslash'c \text{ si}\backslash c \text{ e w gar}\backslash's\backslash'c$$

nous obtiendrons l'expression:

“wziąć się w garść”

(*reprendre son sang froid, se remettre*)

et après l'insertion de `\L\'od\'z` c'est le mot “Łódź”, le nom d'une grande ville polonaise, qui sera imprimé.

Pour faciliter l'insertion des lettres polonaises nous pouvons adopter la solution qui ressemble à celle trouvée dans notre Institut pour d'autres logiciels. Nous distinguons un signe particulier (`@` — à *commercial*) informant que la lettre qui le suit doit être accentuée.

Nous insérons:

`@a`  
`@c`  
`@e`  
`@l`  
`@n`  
`@o`  
`@s`  
`@z`  
`@r`

pour obtenir:

ą  
 ć  
 ę  
 ł  
 ń  
 ó  
 ś  
 ź  
 ż

Nous faisons de même pour les lettres majuscules. Nous écrivons `@r` pour la lettre ‘ż’ afin de la distinguer du ‘z’ décrit par `@z` (il y a deux signes diacritiques différents au-dessus de la même lettre). La désignation `@r` n'est pas fortuite. Nous l'avons choisie, parce que la lettre ‘ż’ et la consonne composée ‘rz’ sont prononcées en polonais de la même façon (comme *j* dans le mot français *jeudi*). Les expressions présentées ci-dessus en tant qu'exemples seraient insérées de la manière suivante:

```
@R@o@lty li@s@c spad@l z d@ebu
wzi@a@c si@e w gar@s@c
@L@od@z
```

Cette solution facilite l'insertion du texte polonais et le rend un peu plus lisible en entrée. Voilà le texte complet des définitions `TeX` ajoutées:

```
1 \catcode'@\active
2 \def@{\leavevmode\futurelet\next\lettertest}
3 \def\lettertest
4   {\ifx\next a\def\diacr{6}\else
5    \ifx\next A\def\diacr{7}\else
6    \ifx\next c\def\diacr{1}\else
7    \ifx\next C\def\diacr{1}\else
8    \ifx\next e\def\diacr{0}\else
9    \ifx\next E\def\diacr{0}\else
10   \ifx\next l\def\diacr{4}\else
11   \ifx\next L\def\diacr{5}\else
12   \ifx\next n\def\diacr{1}\else
13   \ifx\next N\def\diacr{1}\else
14   \ifx\next o\def\diacr{1}\else
15   \ifx\next O\def\diacr{1}\else
16   \ifx\next r\def\diacr{2}\else
```

```

17 \ifx\next R\def\diacr{3}\else
18 \ifx\next s\def\diacr{1}\else
19 \ifx\next S\def\diacr{1}\else
20 \ifx\next z\def\diacr{1}\else
21 \ifx\next Z\def\diacr{1}\else
22 \fi\fi\fi\fi\fi
23 \fi\fi\fi\fi\fi
24 \fi\fi\fi\fi\fi
25 \fi\fi\fi
26 \ifcase\diacr
27 {\c \next}\or
28 {\accent"13\next}\or
29 {\accent"5Fz}\or
30 {\accent"5FZ}\or
31 {\char'40l}\or
32 {\setbox0=\hbox{L}\hbox to\wd0{\hss\char'40L}}\or
33 {\setbox0=\hbox{a}a\kern-\wd0\char24\kern-.23em}\or
34 {\setbox0=\hbox{A}A\kern-\wd0\kern.16em\char24\kern-.12em}\fi
35 \let\next}

```

Le signe © défini comme *actif* et la lettre qui le suit seront remplacés par des commandes respectives. Il faut remarquer que dans ce cas-là le contrôle de macro-instructions à l'aide de `\tracingmacros` devient plus compliqué à cause des très grandes dimensions du fichier .log.

Outre une certaine commodité d'insertion des lettres polonaises cette solution a encore un autre avantage. Même après avoir créé avec METAFONT des polices des caractères polonais, la façon d'insérer les textes resterait invariable. Après avoir effectué de simples modifications des fichiers .t<sub>fm</sub> chaque lettre avec un signe diacritique serait traitée comme ligature. Par exemple ©c constituerait la ligature 'ć'. Le traitement des lettres accentuées comme ligatures permettrait à T<sub>E</sub>X de diviser les mots polonais (voir [A], [K]). Il suffit seulement, conformément à ce que W. Appelt a proposé dans [A], de changer la catégorie du signe © pour 11:

```
\catcode'\©=11 \lccode'\©='©
```

Il est évidemment possible de passer aux ligatures de ce type même dans le cas où nous insérons le texte en utilisant les commandes du PLAIN T<sub>E</sub>X pour créer les lettres accentuées. Il suffit uniquement de modifier leur définition:

```

1 \def\.#1{\if#1z{©r}\else
2 \if#1Z{©R}\else
3 {\accent 95 #1}\fi\fi}
4 \def\'#1{\if#1c{©c}\else
5 \if#1C{©C}\else
6 \if#1n{©n}\else
7 \if#1N{©N}\else
8 \if#1o{©o}\else
9 \if#1O{©O}\else

```

```

10      \if#1s{0s}\else
11      \if#1S{0S}\else
12      \if#1z{0z}\else
13      \if#1Z{0Z}\else
14      {\accent 19 #1}\fi\fi\fi\fi\fi\fi\fi\fi\fi\fi}
15 \let\ccc=\c
16 \def\c#1{\if#1a{0a}\else
17      \if#1A{0A}\else
18      \if#1e{0e}\else
19      \if#1E{0E}\else{\ccc#1}\fi\fi\fi\fi}
20 \def\l{0l}
21 \def\L{0L}

```

Remarquons qu'après cette dernière modification les lettres non polonaises ayant des signes diacritiques seront traitées par T<sub>E</sub>X comme jusqu'à présent.

### Polices des caractères polonais

Les problèmes concernant la qualité typographique des lettres polonaises avec des signes diacritiques et la division des mots nécessitent la création de nouvelles polices de caractères avec le logiciel METAFONT.

Toutes les lettres polonaises devraient se trouver dans ces polices, les minuscules ainsi que les majuscules. Même en adoptant, en tant que solution temporaire, la modification des polices déjà existantes il faut que toutes les 18 lettres accentuées polonaises (9 minuscules et 9 majuscules) s'y trouvent.

Remarquons que dans la police Tél<sub>é</sub>type les signes diacritiques des lettres 'Ź' et 'ż' avaient été remplacés par d'autres caractères. Par conséquent cette police n'est pas utile pour le polonais.

### La division des mots polonais

Les problèmes concernant la division des mots avec des lettres accentuées ont été présentés dans les documents: [A], [D1], [D2], [RS] consacrés à l'adaptation du système T<sub>E</sub>X à des langues autres que l'anglais. La plupart de ces problèmes sont aussi valables pour la langue polonaise.

Aucun dictionnaire polonais n'indique les divisions des mots ce qui empêche l'utilisation du programme PATGEN ([L]). Il reste à écrire "à la main" les motifs requis par T<sub>E</sub>X, de même que J. Désarménien l'a fait pour le français.

Nous n'avons pas encore fini notre travail aux motifs pour le polonais, nous pouvons toutefois indiquer les difficultés essentielles.

Voilà les règles de division du polonais (d'après [S]):

- (1.) **Critères.** La division des mots s'accomplit conformément à deux critères: phonétique (entre les syllabes) et morphologique (entre le préfixe et la racine ou, dans le cas des mots composés, entre les éléments de composition). Le critère morphologique est supérieur au phonétique.

(2.) **Il est interdit de diviser:**

- les mots composés d'une seule syllabe;
- les groupes de deux consonnes désignant un son: **ch, cz, dz, dź, dż, rz, sz**, par exemple: **mu-cha, o-cho-czy**;
- les groupes composés d'une consonne suivie de la lettre 'i' et d'une voyelle, par exemple: **cio-tka, ciot-ka**;
- les groupes composés des lettres 'dzi' suivies d'une voyelle, par exemple: **sie-dzieć**.

(3.) **Groupes de consonnes.** Ils se divisent de façon arbitraire, par exemple: **wa-rstwa, war-stwa, wars-twa, warst-wa**. La coupure juste après tout le groupe est la seule qui est interdite. Il faut toutefois diviser les groupes composés de deux consonnes identiques, par exemple: **pan-na, las-so**.

(4.) **Division entre le préfixe et la racine.** Conformément au supérieur critère morphologique la coupure comme **wez-brać**, par exemple, est la seule autorisée (préfixe "wez-", racine "brać"). D'autres coupures: **we-zbrać** ou **wezb-rać** sont interdites.

(5.) **Mots composés.** Il faut les diviser entre les éléments de composition, par exemple: **pół-noc (mi-nuit)**. La division: **pół-moc** n'est pas correcte.

Après l'analyse des règles présentées ci-dessus nous remarquons deux problèmes essentiels: la division des groupes de consonnes et la division du mot entre le préfixe et la racine.

Le nombre de différents groupes de consonnes qui apparaissent dans les mots polonais est très grand. Il est pratiquement impossible de construire les motifs qui permettraient de trouver tous les points de division d'autant plus que personne, jusqu'à présent, n'a établi la liste complète de suites de consonnes possibles. La solution la plus simple est de diviser les mots juste après la voyelle, avant tout le groupe de consonnes. Selon la troisième règle présentée ci-dessus ce point de division est correcte. Il faut toutefois remarquer que les règles mentionnées, observées en Pologne depuis 1936, sont plus libérales que les précédentes formulées en 1918. L'essentielle différence concerne la division des groupes de consonnes (ici: la 3<sup>ème</sup> règle). La règle précédente avait désigné exactement les points de division et exclu, en particulier, la coupure entre la voyelle et le groupe de consonnes. N'étant plus valide, elle est profondément enracinée dans les esprits des Polonais pour qui la coupure comme **wa-rstwa** reste toujours "bizarre" et "à éviter".

L'autre problème concerne la division du mot entre le préfixe et la racine. Il y a environs trente préfixes d'origine slave: **bez-, na-, nad-, naj-, o-, ob-, od-, po-, pod-, prze-, przed-, we-, wes-, wez-**, etc. Il n'y a aucune règle pour distinguer automatiquement le préfixe d'un mot, par exemple il y a deux façons de décomposer le mot "nadrobić": **na-drobić (émiéter)** et **nad-robić (faire qqch en plus)**. Puisque le critère morphologique est supérieur au phonétique, la solution proposée ci-dessus de diviser le mot juste après une voyelle n'est pas satisfaisante. Il faut absolument exclure la division à l'intérieur d'un préfixe éventuel.

En adoptant la solution proposée nous éviterons les coupures fautives, mais beaucoup



de points de division correctes ne seront jamais trouvés. En outre il faut ajouter des motifs pour empêcher la coupure avant le groupe de consonnes qui termine le mot, par exemple: -szcz, -rstw, -ństw, -ńcz, -rsz.

Remarquons qu'une liste d'exception n'est pas utile pour la langue polonaise dans laquelle non seulement les verbes se conjuguent mais aussi les noms, les adjectifs et les pronoms se déclinent.

## La typographie polonaise

Il n'existe pas de difficultés essentielles pour adapter T<sub>E</sub>X à la typographie polonaise. Une des différences concerne les guillemets qui ont la forme suivante: „ ”. Les guillemets < > ne sont utilisés qu'à l'intérieur d'une autre citation, comme dans la phrase suivante:

Autor pisze: „przeczytałem < Hamleta > Szekspira w wieku 10 lat”

(L'auteur dit: < j'ai lu < Hamlet > de Shakespeare à l'âge de 10 ans >)

Les guillemets qui se trouvent au début d'une citation peuvent être définis à l'aide de la virgule:

```
\def\"{\leavevmode\hbox to.5em{,\hss,}}
```

Une autre différence concerne la façon de distinguer certains mots du texte. Assez souvent nous imprimons séparément les lettres d'un tel mot, par exemple:

“Wyraz tekst jest wyróżniony”

(Le mot “tekst” est différencié).

Mentionnons que dans la typographie polonaise, basée sur la tradition française, on mesure en points didot.

## Conclusions

L'adaptation de T<sub>E</sub>X pour le polonais est possible bien qu'elle exige un certain effort. Parmi les problèmes présentés dans ce document celui concernant la division des mots est sûrement le plus difficile à résoudre.

Remarquons pour terminer que pour accomplir ce travail une profonde connaissance linguistique est nécessaire.

## Remerciements

Je tiens à remercier Anna Borkowska de m'avoir aidée à rédiger en français ce document ainsi que mes collègues, Janusz S. Bień et Krzysztof Szafran, pour leurs conseils et leur amical intérêt.

## Bibliographie

- [A] Appelt (W.), *The Hyphenation of Non-English Words with T<sub>E</sub>X*, in: T<sub>E</sub>X for Scientific Documentation, Proc. of the First European Conference, Addison-Wesley, 1984, pp.61-65.
- [D1] Désarménien (J.), *How to run T<sub>E</sub>X in French*, Stanford University, Report No. STAN-CS-1013, 1984.

- [D2] Désarménien (J.), *The use of T<sub>E</sub>X in French: hyphenation and typographie*, in: T<sub>E</sub>X for Scientific Documentation, Proc. of the First European Conference, Addison-Wesley, 1984, pp.41–59.
- [ISO] ISO, International Organisation for Standardization, *Information processing—Coded character sets for text communication*, Standard No. ISO 6937, 1983.
- [K] Knuth (D.E.), *The T<sub>E</sub>Xbook*, Addison-Wesley, Reading, Mass., 1984.
- [L] Liang (F.M.), *Word Hy-phen-a-tion by Com-put-er*, Stanford University, Report No. STAN-CS-83-977, 1983.
- [RS] Romberger (S.), Sundblad (Y.), *Adapting T<sub>E</sub>X to Languages that use Latin Alphabetic Characters*, in: T<sub>E</sub>X for Scientific Documentation, Proc. of the First European Conference, Addison-Wesley, 1984, pp.27–40.
- [S] Szymczak (M.), *Słownik ortograficzny języka polskiego wraz z zasadami pisowni i interpunkcji*, PWN, Warszawa, 1981.