

# When Typography Meets Programming\*

Jean-Michel HUFFLEN

LIFC (EA CNRS 4157)

University of Franche-Comté

16, route de Gray

25030 BESANÇON CEDEX

FRANCE

[jmhufflen@lifc.univ-fcomte.fr](mailto:jmhufflen@lifc.univ-fcomte.fr)

<http://lifc.univ-fcomte.fr/home/~jmhufflen>

## Abstract

$\TeX$ 's language is a wonderful tool to customise  $(\LaTeX)$ 's behaviour. However, some applications related to 'pure' programming are difficult to put into action, as shown by  $\text{Lua}\TeX$ , where some tasks can be delegated to Lua procedures. We show some examples directly related to typesetting, and demonstrate some solutions using 'more classical' programming languages.

**Keywords**  $(\LaTeX)$ , Scheme, programming workarounds.

## Streszczenie

Makrojęzyk  $\TeX$ a jest cudownym narzędziem do modyfikowania zachowania  $(\LaTeX)$ a. Jednak niektóre zastosowania związane z „czystym“ programowaniem są w nim trudne do zrealizowania, co demonstruje  $\text{Lua}\TeX$ , w którym niektóre zadania mogą być delegowane do procedur Lua. Pokażemy kilka przykładów bezpośrednio związanych z typografią i zaprezentujemy ich rozwiązania przy użyciu „bardziej klasycznych“ języków programowania.

**Słowa kluczowe**  $(\LaTeX)$ , Scheme, rozwiązanie zastępcze.

## 0 Introduction

$\TeX$ 's language is known as a wonderful tool to typeset texts nicely. However there are some examples related to typography, where using a 'more classical' programming language can be judicious. The purpose of this talk is to show some examples of that. Hereafter, a particular case is given *in extenso*, other examples will be demonstrated at the  $\text{Bacho}\TeX$  conference. Reading this article only requires basic knowledge of  $\TeX$ , the commands we mention are documented in [14].

## 1 Dealing with 'false' small capitals

Let us assume that you want to typeset 'Bachotek' using small capitals, that is, 'BACHOTEK'. Unluckily you are using an 'exotic' font that does not provide small capitals directly<sup>1</sup>, so a workaround is to use smaller-sized capital letters. That is—even if the result may seem to be ugly—you can replace ' $\text{\textsc{Bachotek}}$ ' by ' $\text{B}\text{\fake{ACHOTEK}}$ ', using

the  $\text{\fake}$  command defined in Fig. 1. Let us notice that putting:

```
 $\text{\fake{BACHOTEK}}$ 
```

—or ' $\text{\fake{\uppercase{Bachotek}}}$ '—is incorrect, because the initial 'B' must be an 'actual' capital letter, in normal size. In addition, a command such as  $\text{\TeX}$  would remain untouched in:

```
 $\text{\textsc{Bacho}\TeX}$ 
```

(' $\text{BACHO}\TeX$ '), so our workaround should apply to this last example as ' $\text{B}\text{\fake{ACHO}\TeX}$ '.

Some solutions to this problem have already been proposed—e.g., in [13, § 103]—but as far as we know, all these solutions are not robust: they yield incorrect results if they are applied to texts containing uppercase letters or nested commands, that is, cases such as:

```
 $\text{\textsc{\textbf{Bacho}\TeX}}$ 2
```

If we would like to define an alternative command for  $\text{\textsc}$ , the solution is to examine *each* letter of

\* Title in Polish: *Gdy typografia spotyka się z programowaniem.*

<sup>1</sup> ... or your publisher forces you to use such a font.

<sup>2</sup> In this last case, a font substitution may also be ordered by means of the  $\text{\DeclareFontShape}$  command [16]. We do not go thoroughly into this way.

```

\def\fake#1{\small#1}

\def\alttextsc#1{%
  \def\iter{%
    \afterassignment\testcase\let\next= %
  }%
  \def\testcase{%
    \expandafter\ifx\next\enditer\else%
    \ifcat\noexpand\next\relax\noexpand\next%
    \else%
      \if\next a\fake{A}\else%
      \if\next b\fake{B}\else%
        ...
      \fi%
      \fi%
    \fi\let\next=\relax%
  \fi\next\iter}%
  \def\enditer##1{\relax}%
  \iter#1\enditer%
}

```

**Figure 1:** Using ‘false’ small capitals.

the argument of this command, so the result would look like:

```

\textbf{%
  B\fake{A}\fake{C}\fake{H}\fake{O}\TeX%
}

```

for our last example. It is not difficult for a  $\TeX$ pert to write a loop iterating on a string, but quite tedious. We need a marker for the end of this string, and a macro using the `\afterassignment` command puts this *modus operandi* into action<sup>3</sup>, as shown in Fig. 1 by our `\alttextsc` command.

However, we can observe the huge number of ‘`\if... \fi`’ constructs. Unfortunately,  $\TeX$  does not provide any construct to check if a letter is upper-case (or down-case). First, the tables used by the commands `\uppercase` and `\downcase` are not directly readable; second, changing these tables by means of the commands `\lccode` or `\uccode` may affect other commands’ behaviour. The second remark holds on for the table controlling character category (changeable by means of the `\catcode` command). In fact, in this last case,  $\TeX$  provides the construct ‘`\ifcat... \fi`’ in order to check a character’s category... but uppercase and lowercase letters default to the same category. In addition, let us notice that even if the best solution induced changes in a table (by `\catcode`, `\lccode`, or `\uccode`) for 26 letters — uppercase or lowercase letters — such a

<sup>3</sup> We gave a similar example in [5]: a macro computing the numbers of characters of a string.

way would result in a lengthy source text for our new `\alttextsc` command.

If we decide to use a ‘more classical’ programming language in order to implement the replacement of the string ‘`\textsc{Bachotek}`’ by a quite equivalent string using our `\fake` command — that is, doing ‘macro’-replacement before applying  $\TeX$ ’s macros — the advantage is that most modern languages provide efficient constructs to iterate a function along a string. More precisely, they allow a function to be iterated on each character of a string, and then all these subresults can be collected into a structure, most often into another string<sup>4</sup>. Unfortunately, processing each character separately in the case of the `\textsc` command’s argument fails for ‘`\textsc{BachO\TeX}`’. Another solution using a ‘mini’  $\TeX$  parser is given in Fig. 2, using the Scheme programming language<sup>5</sup> [11]. The comments we put in this figure should make clear the broad outlines of our method. In addition to some predefined functions and special forms of Scheme, our `alttextsc` function uses some additional definitions provided by `SRFIS`<sup>6</sup>: the `let*-values` macro is defined in [3], the `cons*` function in [17], the `string-concatenate-reverse/shared` function in [18]. This `alttextsc` function also uses definitions belonging to `MIBIB $\TeX$` <sup>7</sup> source files [4]:

- (`mk-r-string-thunk s`) may be viewed as a character-generator: this expression returns a *thunk*<sup>8</sup>; successive calls of this *thunk* return the successive characters of the *s* string, in turn; when the end of *s* is reached, all the subsequent calls return the false value ‘`#f`’;
- (`t-next-token thunk char x`) where *thunk* is a character-generator, *char* the last character read within the string or the false value ‘`#f`’ if the end has been reached, and *x* an escape function to be called in case of an error, returns two values: the next token parsed according to  $\TeX$ ’s conventions, and the character after; the resulted token may be a single character, or a command’s name given as a string if this token

<sup>4</sup> Good examples are Scheme functions `string-map` and `string-for-each` [18]. So are the *iterators* of the `STL` (Standard Template Library), part of the C++ language [10].

<sup>5</sup> Readers interested in an introductory book to this functional programming language can refer to [19].

<sup>6</sup> Scheme Requests For Implementation. That is an effort to coordinate libraries and other additions to the Scheme language between implementations.

<sup>7</sup> MultiLingual `BIB $\TeX$` . The functions we mention hereafter are used by `MIBIB $\TeX$`  to parse a  $\LaTeX$  document’s preamble, in order to know which multilingual packages are used [7].

<sup>8</sup> That is, a zero-argument function, w.r.t. Scheme’s terminology.

```

(define alttextsc
  (let ((backslash-as-string "\\") ; This escape character must be escaped itself.
        (fake-open-string "\\fake{")
        (fake-close-string "}"))
    (lambda (string-0)
      (let ((r-thunk (mk-r-string-thunk string-0)) ; Builds a character-generator along string-0.
            (x (lambda (y) #f))) ; This function will be called by t-next-token in case
                                ; of an error. y is supposed to be an error label.
          (let loop ((char (r-thunk)) ; Launches the parsing of the string.
                    (acc '()))
            (let*-values (((ttoken char-0) (t-next-token r-thunk char x))) ; 2 values returned: the next
                          ; token — #f, a character, or a string — and the character after.
                          (if ttoken ; W.r.t. MIBBTEX's terminology, 'ttoken' is for 'TEX token.
                              (loop char-0
                                   (if (char? ttoken)
                                       (if (char-lower-case? ttoken)
                                           (cons* fake-close-string (string (char-upcase+ ttoken))
                                                 fake-open-string acc)
                                           (cons (string ttoken) acc))
                                       ; ; Otherwise, ttoken is a command's name, given as a string. We restore the '\
                                       ; ; character before this name.
                                       (cons* ttoken backslash-as-string acc))))
                              ; ; If ttoken is the 'false' value #f, this means that we have reached the string's end. To get the
                              ; ; final result, all the strings of the accumulator acc are concatenated by using reverse order.
                              (string-concatenate-reverse/shared acc))))))))))

```

**Figure 2:** Generating a source text typeset using ‘false’ small capitals.

begins with the ‘\’ character; if the end of the input or string port has been reached, the two values are equal to #f;

- the `char-lower-case+?` (resp. `char-upcase+`) function extends the `char-lower-case?` (resp. `char-upcase`) function — dealing to the ‘basic’ ASCII<sup>9</sup> encoding — to the Latin 1 encoding<sup>10</sup>; they should be replaced by Unicode-compliant functions in Scheme’s future versions.

This implementation of the `alttextsc` command — as *pre-processing* in Scheme, the result being processed by L<sup>A</sup>T<sub>E</sub>X — could be improved by putting as few occurrences of the `\fake` as possible. Of course, it requires some knowledge of Scheme, but it seems to us to be not as heavy as TEX’s. Moreover, we think this point holds on with most modern programming languages, provided that they can use a TEX-like parser.

## 2 Requirements for an interface

Roughly speaking, if we are interested in doing some pre-processing about a fragment of a TEX document, we do not have to pay attention to some subtleties related to the application of macros and commands

<sup>9</sup> American Standard Code for Information Interchange.

<sup>10</sup> In [6], we explain why this encoding is internally used by MIB<sub>B</sub>TEX’s present version.

by TEX, we do not have to be careful with commands related to *control features* like `\expandafter` [14]. We do not have to be careful with the space characters ending a command’s name and gobbled by this command, either<sup>11</sup>. But let us recall that such a ‘mini’ TEX parser must rule out comments. More precisely, it must skip the characters located between a ‘%’ sign and the end of the current line, these two delimiter characters belonging to the sequence to be skipped. Space characters at the beginning of a line must be skipped, too.

Values returned by a TEX-like tokeniser may be three types: the character type<sup>12</sup> for ‘simple’ characters — in particular, contents of paragraphs to be typeset — the type implementing commands’ names, and a one-valued type for the parsing process’ end. To do that, our ‘mini’ TEX parser respectively uses characters in Scheme, strings in Scheme, and the false value ‘#f’. As an alternative implementation,

<sup>11</sup> So does the ‘mini’ TEX parser developed as part of MIB<sub>B</sub>TEX (cf. Footnote 7, p. 1002). On the contrary, the parser of bibliography data (.bib) files, built by means of analogous methods, *expands* some L<sup>A</sup>T<sub>E</sub>X commands. Examples are given by accent commands whose result belongs to the Latin 1 encoding: ‘\’{e}’ is expanded into the ‘é’ character within internal forms used by MIB<sub>B</sub>TEX.

<sup>12</sup> Possibly Latin 1, Latin 2, . . . or Unicode-compliant, according to what you would like to parse.

let us mention that we have written another T<sub>E</sub>X-like tokeniser, only experimental, using the object-oriented Ruby programming language [15].

### 3 Conclusion

In the past, some programming languages were designed to be universal, that is, to serve all purposes. All of these languages PL/1 [8], Ada [1] have failed to be accepted as filling this role. Nowadays, only the C programming language [12] is still used for a very range of applications. But this language can be viewed as a high-level assembly language. In particular, handling strings in C may be quite tedious. On the contrary, T<sub>E</sub>X's language, like most specialised languages, is very efficient inside its domain, quite unsuitable outside. It offers efficient primitives for advanced typesetting (`\ifcat`, `\ifmode`, ...) [14] but programming a sort procedure is nightmare.

The cooperation among procedures written using different programming languages is not new: for example, Ada's second version [20] provides interfaces with other programming languages (e.g., C). The new typesetting engine LuaT<sub>E</sub>X [2] belongs to this framework. Work related to typesetting is done by T<sub>E</sub>X, whereas other tasks may be *delegated* to procedures written using Lua [9]. But this cooperation is unidirectional: LuaT<sub>E</sub>X can run procedures in Lua, but these procedures cannot call T<sub>E</sub>X's commands. T<sub>E</sub>X's major drawback is that it recognises only its own formats. We hope that the mini-parsers we wrote will indirectly make T<sub>E</sub>X's features open to other environments.

### 4 Acknowledgements

Many thanks to Jerzy B. Ludwichowski, who has translated the abstract and keywords in Polish.

### References

- [1] ANSI: *The Programming Language Ada<sup>®</sup> Reference Manual*. Technical Report ANSI/MIL-STD-1815A-1983, American National Standard Institute, Inc. LNCS No. 155, Springer-Verlag. 1983.
- [2] Hans HAGEN: "The Luafication of T<sub>E</sub>X and ConT<sub>E</sub>Xt". In: *Proc. BachoT<sub>E</sub>X 2008 Conference*, pp. 114–123. April 2008.
- [3] Lars T. HANSEN: *Syntax for Receiving Multiple Values*. March 2000. <http://srfi.schemers.org/srfi-11/>.
- [4] Jean-Michel HUFFLEN: "MIBIBT<sub>E</sub>X in Scheme (First Part)". *Biuletyn GUST*, Vol. 22, pp. 17–22. In *BachoT<sub>E</sub>X 2005 conference*. April 2005.
- [5] Jean-Michel HUFFLEN: "T<sub>E</sub>X's Language within the History of Programming Languages". *Biuletyn GUST*, Vol. 22, pp. 23–32. In *BachoT<sub>E</sub>X 2005 conference*. April 2005.
- [6] Jean-Michel HUFFLEN: "Managing Order Relations in MIBIBT<sub>E</sub>X". *TUGboat*, Vol. 29, no. 1, pp. 101–108. EuroBachoT<sub>E</sub>X 2007 proceedings. 2007.
- [7] Jean-Michel HUFFLEN: "Managing Languages within MIBIBT<sub>E</sub>X". *TUGboat*, Vol. 30, no. 1, pp. 49–57. July 2009.
- [8] IBM SYSTEM 360: *PL/1 Reference Manual*. March 1968.
- [9] Roberto IERUSALIMSKY: *Programming in Lua*. 2nd edition. Lua.org. March 2006.
- [10] ISO: *Programming Languages—C++*. Technical Report ISO/IEC 14882:2003, ISO/IEC. 2003.
- [11] Richard KELSEY, William D. CLINGER, Jonathan A. REES, Harold ABELSON, Norman I. ADAMS IV, David H. BARTLEY, Gary BROOKS, R. Kent DYBVIK, Daniel P. FRIEDMAN, Robert HALSTEAD, Chris HANSON, Christopher T. HAYNES, Eugene Edmund KOHLBECKER, JR, Donald OXLEY, Kent M. PITMAN, Guillermo J. ROZAS, Guy Lewis STEELE, JR, Gerald Jay SUSSMAN and Mitchell WAND: "Revised<sup>5</sup> Report on the Algorithmic Language Scheme". *HOSC*, Vol. 11, no. 1, pp. 7–105. August 1998.
- [12] Brian W. KERNIGHAN and Denis M. RITCHIE: *The C Programming Language*. 2nd edition. Prentice Hall. 1988.
- [13] Marie-Paule KLUTH : *FAQ L<sup>A</sup>T<sub>E</sub>X française pour débutants et confirmés*. Vuibert Informatique, Paris. Également disponible sur [CTAN:help/LaTeX-FAQ-francaise/](http://CTAN:help/LaTeX-FAQ-francaise/). Janvier 1999.
- [14] Donald Ervin KNUTH: *Computers & Typesetting. Vol. A: The T<sub>E</sub>Xbook*. Addison-Wesley Publishing Company, Reading, Massachusetts. 1984.
- [15] Yukihiro MATSUMOTO: *Ruby in a Nutshell*. O'Reilly. English translation by David L. Reynolds, Jr. November 2001.
- [16] Frank MITTELBAACH and Michel GOOSSENS, with Johannes BRAAMS, David CARLISLE, Chris A. ROWLEY, Christine DETIG and Joachim SCHROD: *The L<sup>A</sup>T<sub>E</sub>X Companion*. 2nd edition. Addison-Wesley Publishing Company, Reading, Massachusetts. August 2004.
- [17] Olin SHIVERS: *List Library*. October 1999. <http://srfi.schemers.org/srfi-1/>.
- [18] Olin SHIVERS: *String Library*. December 2000. <http://srfi.schemers.org/srfi-13/>.
- [19] George SPRINGER and Daniel P. FRIEDMAN: *Scheme and the Art of Programming*. The MIT Press, McGraw-Hill Book Company. 1989.
- [20] S. Tucker TAFT and Robert A. DUFF, eds.: *Ada 95 Reference Manual. Language and Standard Libraries*. No. 1246 in LNCS. Springer-Verlag. International Standard ISO/IEC 8652:1995(E). 1995.