

How to handle compound and bounded words

just another hyphenation method

Hans Hagen

April 1 1996

One of T_EX's strong points in building paragraphs is the way hyphenations are handled. Although for real good hyphenation of non-english languages some extensions to the program are needed, fairly good results can be reached with the standard mechanisms and an additional macro, at least in Dutch.

1 \unprotect

ConT_EXt originates in the wish to typeset educational materials, especially in a technical environment. In production oriented environments, a lot of compound words are used. Because the Dutch language poses no limits on combining words, we often favor putting dashes between those words, because it facilitates reading, at least for those who are not that accustomed to it.

In T_EX compound words, separated by a hyphen, are not hyphenated at all. In spite of the multiple pass paragraph typesetting this can lead to parts of words sticking into the margin. The solution lays in saying `spoelwater||terugwinunit` instead of `spoelwater-terugwinunit`. By using a one character command like `|`, delimited by the same character `|`, we get ourselves both a decent visualization (in T_EXedit and colored verbatim we color these commands yellow) and an efficient way of combining words.

The sequence `||` simply leads to two words connected by a hyphen. Because we want to distinguish such a hyphen from the one inserted when T_EX hyphenates a word, we use a bit longer one.

`spoelwater||terugwinunit` `spoel- wa- ter- te- rug- win- unit` `spoelwater-terugwinunit`

As we already said, the `|` is a command. This commands accepts an optional argument before it's delimiter, which is also a `|`.

`polymeer|*|chemie` `po- |y- meer* che- mie` `polymeer*chemie`

Arguments like `*` are not interpreted and inserted directly, in contrary to arguments like:

`polymeer|~|chemie` `po- |y- meers che- mie` `polymeerchemie`
`|(|polymeer|)|chemie` `(po- |y- meer-)| che- mie` `(polymeer)chemie`
`polymeer|(|chemie)|` `po- |y- meer (c-che- mie-)` `polymeer(chemie)`

Although such situations seldom occur —we typeset thousands of pages before we encountered one that forced us to enhance this mechanism— we also have to take care of comma's.

`op||, in|| en uitstellen` `op- in- en uit stel len` `op-, in- en uitstellen`

The next special case (concerning quotes) was brought to my attention by Piet Tutelaers, one of the driving forces behind rebuilding hyphenation patterns for the dutch language.¹ We'll also take care of this case.

`AOW|'|er` `AOW' en` `AOW'er`
`cd|'|tje` `cd tje` `cd'tje`
`ex|-|PTT|'|er` `ex- PTT' en` `ex-PTT'er`
`rock|-|'n|-|roller` `rock' 'n' roller` `rock-'n-roller`

The mechanism described here is one of the older inner parts of ConT_EXt. The most recent extensions concerns some special cases as well as the possibility to install other characters as delimiters. The preferred way of specifying compound words is using `||`, which is installed by:

`\installdiscretionaries || -`

Some alternative definitions are:

¹In 1996 the spelling of the dutch language has been slightly reformed which made this topic actual again.

```

\installdiscretionaries ** -
\installdiscretionaries ++ -
\installdiscretionaries // -
\installdiscretionaries ~~ -

```

after which we can say:

```

test**test**test   test test test   test-test-test
test++test++test   test test test   test-test-test
test//test//test   test test test   test-test-test
test~~test~~test   test test test   test-test-test

```

Now let's go to the macros. First we define some variables. In the main ConTeXt modules these can be tuned by a setup command. Watch the (maybe) better looking compound hyphen.

```

2 \def\compoundhyphen    {-}\kern-.25ex{-}
\def\beginofsubsentence {---}
\def\endofsubsentence   {---}

```

The last two variables are needed for subsentences —like this one— which we did not yet mention.

We want to enable breaking but at the same time don't want compound characters like - or – to be separated from the words. T_EX hackers will recognise the next two macro's:

```

3 \def\prewordbreak  {\penalty10000\hskip0pt\relax}
\def\postwordbreak  {\penalty0\prewordbreak}

```

We first show the original implementation, which only supports | as command and delimiter. Before activating | we save it's value:

```

\edef\domathmodediscretionary{\string|}

```

after which we're ready to define it's meaning to:

```

\catcode'\|= \@@active

\protected\def|{%
  {\ifmmode
    \expandafter\domathmodediscretionary
  \else
    \expandafter\dotextmodediscretionary
  \fi}

```

We need a two stage \futurelet because we want to look ahead for both the compound character definition and the (optional) comma that follows it, and because we want to prevent that T_EX puts this comma on the next line. We use \next for easy and fast checking of the argument, we save this argument (which can consist of more tokens) and also save the character following the |#1| in \nextnext.

```

\def\dotextmodediscretionary%
  {\bgroup
   \futurelet\next\dodotextmodediscretionary}

\def\dodotextmodediscretionary#1|%
  {\def\betweendiscretionaries{#1}%
   \futurelet\nextnext\dododotextmodediscretionary}

```

The main macro consists of quite some \ifx tests while \checkafterdiscretionary handles the commas. We show the simplified version here:

```

\def\dododotextmodediscretionary%
  {\let\nextnextnext=\egroup
   \ifx | \next
     \checkafterdiscretionary
     \prewordbreak\hbox{\compoundhyphen\nextnext}\postwordbreak
   \else\ifx= \next
     \prewordbreak\compoundhyphen
   \else\ifx~ \next

```

```

\discretionary{-}{\thinspace}\postwordbreak
\else\ifx(\next
\prewordbreak\discretionary}{(-){(\prewordbreak
\else\ifx)\next
\prewordbreak\discretionary{-})}{})\postwordbreak
\else\ifx'\next
\prewordbreak\discretionary{-}{}'\postwordbreak
\else
\checkafterdiscretionary
\prewordbreak\hbox{\betweendiscretionaries\nextnext}\postwordbreak
\fi\fi\fi\fi\fi\fi
\nextnextnext}

```

```

\def\checkafterdiscretionary%
{\ifx,\nextnext
\def\nextnextnext{afterassignment\egroup\let\next=}%
\else
\let\nextnext=\relax
\fi}

```

The most recent implementation is more advanced. As demonstrated we can install delimiters, like:

```
\installdiscretionaries || \compoundhyphen
```

This time we have to use a bit more clever way of saving the math mode specification of the character we're going to make active. We also save the user supplied compound hyphen. We show the a bit more traditional implementation first.

```

\def\installdiscretionaries#1%
{\catcode'#1\@@other
\expandafter\doinstalldiscretionaries\string#1}

\def\doinstalldiscretionaries#1%
{\setvalue{mathmodediscretionary#1}{#1}%
\catcode'#1\@@active
\dodoinstalldiscretionaries}

\def\dodoinstalldiscretionaries#1#2%
{\setvalue{textmodediscretionary\string#1}{#2}%
\protected\def#1{\discretionarycommand#1}}

```

A bit more *catcode* and character trickery enables us to discard the two intermediate steps. This trick originates on page 394 of the *T_EXbook*, in the appendix full of dirty tricks. The second argument has now become redundant, but I decided to reserve it for future use. At least it remembers us of the symmetry.

```

4 \def\installdiscretionaries#1#2#3%
  {\setvalue{mathmodediscretionary\string#1}{\char'#1}%
  \setvalue{textmodediscretionary\string#1}{#3}%
  \catcode'#1=\@@active
  \scratchcounter=\the\uccode'~
  \uccode'~='#1
  \uppercase{\protected\def~{\discretionarycommand~}}%
  \uccode'~=\scratchcounter}

5 \def\dohandlemathmodebar#1%
  {\getvalue{mathmodediscretionary\string#1}}

6 \def\discretionarycommand%
  {\ifmmode
  \expandafter\dohandlemathmodebar
  \else
  \expandafter\dotextmodediscretionary
  \fi}

```

Although adapting character codes and making characters active can interfere with other features of macropackages, normally there should be no problems with things like:

```
\installdiscretionary || +
\installdiscretionary ++ =
```

The real work is done by the next set of macros. We have to use a double `\futurelet` because we have to take following characters into account.

```
7 \def\dotextmodediscretionary#1%
  {\bgroup
  \def\dodotextmodediscretionary##1#1%
    {\def\betweendiscretionary{##1}%
    \futurelet\nextnext\dododotextmodediscretionary}%
  \let\discretionarycommand=#1%
  \def\textmodediscretionary{\getvalue{textmodediscretionary}\string#1}%
  \futurelet\next\dodotextmodediscretionary}

8 \def\dododotextmodediscretionary%
  {\let\nextnextnext=\egroup
  \ifx\discretionarycommand\next
    \checkafterdiscretionary
    \prewordbreak\hbox{\textmodediscretionary\nextnext}\postwordbreak
  \else\ifx=\next
    \prewordbreak\textmodediscretionary
  \else\ifx~\next
    \prewordbreak\discretionary{-}{\thinspace}\postwordbreak
  \else\ifx(\next
    \ifdim\lastskip>\!zeropoint\relax
      (\prewordbreak
      \else
        \prewordbreak\discretionary}{(-){}\prewordbreak
      \fi
    \else\ifx)\next
      \ifx\nextnext\blankspace
        \prewordbreak)\relax
      \else
        \prewordbreak\discretionary{-)}{)}\postwordbreak
      \fi
    \else\ifx'\next
      \prewordbreak\discretionary{-}{'}\postwordbreak
    \else\ifx<\next
      \hbox{\beginofsubsentence}\prewordbreak
    \else\ifx>\next
      \prewordbreak\endofsubsentence
    \else
      \checkafterdiscretionary
      \prewordbreak\hbox{\betweendiscretionary\nextnext}\postwordbreak
    \fi\fi\fi\fi\fi\fi\fi\fi
  \nextnextnext}

9 \def\checkafterdiscretionary%
  {\ifx,\nextnext
  \def\nextnextnext{\afterassignment\egroup\let\next=}%
  \else
  \let\nextnext=\relax
  \fi}
```

Before we show some more tricky alternative, we first install the mechanism:

```
10 \installdiscretionaries || \compoundhyphen
```

One of the drawbacks of this mechanism is that characters can be made active afterwards. The next alternative can be used in such situations. This time we don't compare the arguments directly but use the `\uccode`'s instead. `TEX` initializes these codes of the alphabetic glyphs to their uppercase counterparts. Normally the other characters remain zero. If so, we can use the `\uccode` as a signal.

The more advanced mechanism is activated by calling:

```
\enableactivediscretionaries
```

which is defined as:

```
11 \def\enableactivediscretionaries%
    {\uccode'=(\relax \uccode')=\relax \uccode'=='\relax
     \uccode'<=<\relax \uccode'>=>\relax
     \uccode'=''\relax \uccode'~='~\relax
     \let\dottextmodediscretionary = \activedottextmodediscretionary
     \let\dododottextmodediscretionary = \activedododottextmodediscretionary}
```

We only have to redefine two macros. While saving the \uccode in a macro we have to take care of empty arguments, like in ||.

```
12 \def\activedottextmodediscretionary#1%
    {\bgroup
     \def\dodottextmodediscretionary##1#1%
       {\def\betweendiscretionary{##1}%
        \def\nextuccode####1####2\relax%
         {\ifcat\noexpand####1\noexpand\relax
          \edef\nextuccode{0}%
          \else
           \edef\nextuccode{\the\uccode'####1}%
          \fi}%
         \nextuccode##1@\relax
         \futurelet\nextnext\dododottextmodediscretionary}%
     \let\discretionarycommand=#1%
     \def\textmodediscretionary{\getvalue{textmodediscretionary}\string#1}}%
     \futurelet\next\dodottextmodediscretionary}
```

This time we use \ifnum:

```
13 \def\activedododottextmodediscretionary%
    {\let\nextnextnext=\egroup
     \ifx\discretionarycommand\next
       \checkafterdiscretionary
       \prewordbreak\hbox{\textmodediscretionary\nextnext}\postwordbreak
     \else\ifnum\uccode'==\nextuccode
       \prewordbreak\textmodediscretionary
     \else\ifnum\uccode'~=\nextuccode
       \prewordbreak\discretionary{-}{\thinspace}\postwordbreak
     \else\ifnum\uccode'(\=\nextuccode
       \ifdim\lastskip>!!zeropoint\relax
         (\prewordbreak
         \else
          \prewordbreak\discretionary{-}{(-){}\prewordbreak
         \fi
       \else\ifnum\uccode')=\nextuccode
         \ifx\nextnext\blankspace
           \prewordbreak)\relax
         \else
          \prewordbreak\discretionary{-}{-}{})\postwordbreak
         \fi
       \else\ifnum\uccode'='\nextuccode
         \prewordbreak\discretionary{-}{-}'\postwordbreak
       \else\ifnum\uccode'<=\nextuccode
         \hbox{\beginofsubsentence}\prewordbreak
       \else\ifnum\uccode'>=\nextuccode
         \prewordbreak\endofsubsentence
       \else
         \checkafterdiscretionary
```

```

\prewordbreak\hbox{\betweendiscretionary\nextnext}\postwordbreak
\fi\fi\fi\fi\fi\fi\fi\fi
\nextnextnext}

```

Now we can safely do things like:

```

\catcode'<=\@active \def<{hello there}
\catcode'>=\@active \def>{hello there}
\catcode'(\@active \def({hello there}
\catcode')=\@active \def){hello there}

```

In normal day-to-day production of texts this kind of activation is seldom used.² If so, we have to take care of the math mode explicitly, just like we did when making | active. It can be confusing too, especially when we load macropackages afterwards that make use of < in \ifnum or \ifdim statements.

14 \protect

²In the ConT_EXt manual the < and > are made active and used for some cross-reference trickery.