

Hyphenation Patterns

Pattern files

\TeX has two mysterious commands that the average user will never or seldom meet:

```
\hyphenation{as-so-ciates}
\patterns  {.ach4}
```

Both commands can take multiple strings, so in fact both commands should be plural. The first command can be given any time and can be used to tell \TeX that a word should be hyphenated in a certain way. The second command can only be issued when \TeX is in virgin mode, i.e. starting with a clean slate. Normally this only happens when a format is generated.

The second command is more mysterious than the first one and its entries are a compact way to tell \TeX between what character sequences it may hyphenate words. The numbers represent weights and the (often long) lists of such entries are generated with a special program called `pat.gen`. Since making patterns is work for specialists, we will not go into the nasty details here.

In the early stage of Con \TeX t development it came with its own pattern files. Their names started with `lang-` and their suffixes were `pat` and `hyp`.

However, when Con \TeX t went public, I was convinced to drop those files and use the files already available in distributions. This was achieved by using the Con \TeX t filename remapping mechanism. Although those files are supposed to be generic, this is not always the case, and it remains a gamble if they work with Con \TeX t. Even worse, their names are not consistent and the names of some files as well as locations in the tree keep changing. The price Con \TeX t users pay for this is lack of hyphenation until such changes are noticed and taken care of. Because constructing the files is an uncoordinated effort, all pattern files have their own characteristics, most noticeably their encoding.

After the need to adapt the name mapping once again, I decided to get back to providing Con \TeX t specific pattern files. Pattern cooking is a special craft and \TeX users may count themselves lucky that it's taken care of. So, let's start with thanking all those \TeX experts who dedicate their time and effort to get their language hyphenated. It's their work we will build (and keep building) upon.

In the process of specific Con \TeX t support, we will take care of:

- consistent naming, i.e. using language codes when possible as a prelude to a more sophisticated naming scheme, taking versions into account
- consistent splitting of patterns and hyphenation exceptions in files that can be recognized by their suffix
- making the files encoding independent using named glyphs
- providing a way to use those patterns in plain \TeX as well

Instead of using a control sequence for the named glyphs, we use a different notation:

```
[ssharp] [zcaron] [idiaeresis]
```

The advantage of this notation is that we don't have to mess with spacing so that parsing and cleanup with scripts becomes more robust. The names conform to the Con \TeX t way of naming glyphs and the names and reverse mappings are taken

from the encoding files in the ConT_EXt distribution, so you need to have ConT_EXt installed.

The ConT_EXt pattern files are generated by a ruby script. Although the converting is rather straightforward, some languages need special treatment, but a script is easily adapted. If you want a whole bunch of pattern files, just say:

```
ctxtools --patterns all
```

or, if you want one language:

```
ctxtools --patterns nl
```

If for some reason this program does not start, try:

```
texmfstart ctxtools --patterns nl
```

When things run well, this will give you four files:

lang-nl.pat	the patterns in an encoding independent format
lang-nl.hyp	the hyphenation exceptions
lang-nl.log	the conversion log (can be deleted afterwards)
lang-nl.rme	the preambles of the files used (copyright notices and such)

If you redistribute the files, it makes sense to bundle the rme files as well, unless the originals are already in the distribution. It makes no sense to keep the log files on your system. When the file lang-all.xml is present, the info from that file will be used and added to the pattern and hyphenation files. In that case no rme and log file will be generated, unless --log is provided.

In the Dutch pattern file you will notice entries like the following:

```
e[ediaeresis]n3
```

So, instead of those funny (encoding specific) \^e or (format specific) \e we use names. Although this looks ConT_EXt dependent it is rather easy to map those names back to characters, especially when one takes into account that most languages only have a few of those special characters and we only have to deal with lower case instances.

The ConT_EXt support module supp-pat.tex is quite generic and contains only a few lines of code. Actually, most of the code is dedicated to the simple XML handler. Loading a pattern meant for EC encoded fonts in another system than ConT_EXt is done as follows:

```
\bgroup
\input supp-pat
\lccode"E4="E4 \definepatterntoken adiaeresis \^e4
\lccode"F6="F6 \definepatterntoken odiaeresis \^f6
\lccode"FC="FC \definepatterntoken ediaeresis \^fc
\lccode"FF="FF \definepatterntoken ssharp \^ff
\enablepatterntokens
\enablepatternxml
\input lang-de.pat
\input lang-de.hyp
\egroup
```

In addition to this one may want to set additional lower and uppercase codes. In ϵ -T_EX these are stored with the language.

Just for completeness we provide the magic command to generate the XML variants:

```
ctxtools --patterns --xml all
```

This will give you files like:

```
<?xml version='1.0' standalone='yes'?>
<!-- some comment -->
<patterns>
... e&ediaeresis;n3 ...
</patterns>
```

This is also accepted as input but for our purpose it's probably best to stick to the normal method. The pattern language is a T_EX specific one anyway.

Installing languages

Installing a language in ConT_EXt should not take too much effort assuming that the language is supported. Language specific labels are grouped in lang-* files, like lang-ger.tex for the germanic languages.

Patterns will be loaded from the files in the general T_EX distribution unless lang-nl.pat is found, in which case ConT_EXt assumes that you prefer the ConT_EXt patterns. In that case, run

```
ctxtools --patterns all
```

You need to move the files to the ConT_EXt base path that you can locate with:

```
textools --find context.tex
```

You can also use kpswhich, but the above method does an extensive search. Of course you can also generate the files on a temporary location. Now it's time to generate the formats:

```
texexec --make --all
```

Since X_YT_EX needs patterns in utf-8 encoding, we provide a switch for achieving that:

```
texexec --make --all --utf8
```

Beware: you need to load patterns for each language and encoding combination you are going to use. You can configure your local cont-usr file to take care of this. When an encoding does not have the characters that are needed, you will get an error. When using the non ConT_EXt versions of the patterns this may go unnoticed because the encoding is hard coded in the file. Of course it will eventually get noticed when the hyphenations come out wrong.

The ConT_EXt distribution has a file lang-all.xml that holds the copyright and other notes of the patterns. A description looks like:

```
<description language='nl'>
  <sourcefile>nehyph96.tex</sourcefile>
  <title>TeX hyphenation patterns for the Dutch language</title>
  <copyright>
    <year>1996</year>
    <owner> Piet Tutelaers (P.T.H.Tutelaers@tue.nl)</owner>
    <comment>8-bit hyphenation patterns for TeX based upon the
      new Dutch spelling, officially since 1 August 1996.
      These patterns follow the new hyphenation rules in the
      'Woordenlijst Nederlandse Taal, SDU Uitgevers, Den Haag
```

```

1995' (the so called 'Groene Boekje') described in
section 5.2 (Het afbreekteken)</comment>
</copyright>
</description>

```

This file is 'work in process': more details will be added and comments will be enriched.

Commands

You can at any moment add additional hyphenation exceptions to the language specific dictionaries. For instance:

```
\language[nl] \hyphenation{pa-tiën-ten}
```

Switching to another language is done with the `\language` command. The document language is set with `\mainlanguage`.

If you want to let \TeX know that a word should be hyphenated in a special way, you use the `\-` command, for instance:

```
Con\-\TeXt
```

Compound words are not recognized by the hyphenation engine, so there you need to add directives, like:

```
the ConTeXt|-|system
```

If you are using XML as input format, you need to load the hyphenation filter module. Here we assume that utf encoding is used:

```
\useXMLfilter[utf,hyp]
```

In your XML file you can now add:

```

<hyphenations language='nl' regime='utf'>
  <hyphenation>pa-tiën-ten</hyphenation>
  <hyphenation>pa-tiën-ten-or-ga-ni-sa-tie</hyphenation>
  <hyphenation>pa-tiën-ten-plat-form</hyphenation>
</hyphenations>

```

This filter also defines some auxiliary elements. Explicit hyphenation points can be inserted as follows:

```
Zullen we hier af<hyphenate/>bre<hyphenate/>ken of niet?
```

The compound token can be anything, but keep in mind that some tokens are treated special (see other manuals).

```
Wat is eigenlijk een patiënten<compound token="-"/>platform?
```

A language is set with:

```
nederlands <language code="en">english</language> nederlands
```

If you set attribute `scope` to `global`, labels (as used for figure captions and such) adapt to the language switch. This option actually invokes `\mainlanguage`.

Languages

When users in a specific language area use more than one font encoding, patterns need to be loaded multiple times. In theory this means that one can end up with more instances than \TeX can host. However, the number of sensible font encodings is limited as is the number of languages that need hyphenation. Now that memory is cheap and machines are fast, preloading a lot of pattern files is no problem. The

following table shows the patterns that are preloaded in the version of ConT_EXt that is used to process this file.

language	encoding	mapping	number	left _{min}	right _{min}
nl	texnansi	texnansi	1	2	2
nl	ec	ec	2	2	2
fr	texnansi	texnansi	3	2	2
fr	ec	ec	4	2	2
de	texnansi	texnansi	5	2	2
de	ec	ec	6	2	2
it	texnansi	texnansi	7	2	2
it	ec	ec	8	2	2
pt	texnansi	texnansi	9	2	2
pt	ec	ec	10	2	2
hr	ec	ec	11	2	2
pl	p10	p10	12	2	2
pl	ec	ec	13	2	2
pl	qx	qx	14	2	2
cz	il2	il2	15	2	2
cz	ec	ec	16	2	2
sk	il2	il2	17	2	2
sk	ec	ec	18	2	2
sl	il2	il2	19	2	2
sl	ec	ec	20	2	2
en	ec	ec	22	2	2
da	ec	ec	23	2	2
sv	ec	ec	24	2	2
af	ec	ec	25	2	2
no	ec	ec	26	2	2
deo	ec	ec	27	2	2
uk	ec	ec	28	2	2
us	ec	ec	29	2	2
es	ec	ec	30	2	2
ca	ec	ec	31	2	2
la	ec	ec	32	2	2
ro	ec	ec	33	2	2
tr	ec	ec	34	2	2
fi	ec	ec	36	2	2
hu	ec	ec	37	2	2

In the (near) future the somewhat arcane p10 and il2 encodings will go away since they are only used for Polish and Czech/Slovak computer modern fonts, which can be replaced by Latin Modern alternatives. Also, a new dense encoding may find its way into this list.

Hans Hagen