

## Typographers' Inn

Peter Flynn

### Page numbering revisited

In my last column [2] I mentioned retrofitting page numbers from the PDF back into a web version of a document, and said it was relatively trivial with T<sub>E</sub>X. Daniel Nemenyi of KCL emailed me to ask how, so I had to dig out the code and see.

I know others have done this, but I don't know if anyone has documented it anywhere. What we implemented was for a client using XML, generating one transformation to X<sub>Q</sub>L<sup>A</sup>T<sub>E</sub>X for creating PDF, and another to HTML for their in-house web site. The code is not proprietary but I can't extract it directly without exposing a lot of their in-house naming, so I rewrote a short version for Daniel.

The implementation was done using the `fwlw` package (which makes catchwords available for each page: the first and last words on the page plus the first word of the next page). With this, we modified the `\pagestyle` provided to typeset the catchwords at the bottom of the page, in white, so they were not visible, and separated them by an otherwise unused delimiter so we could extract them reliably: we used the ASCII decimal 172 (0xAC) character (–) or NOT symbol (see Figure 1).

```
\documentclass{article}
\usepackage{lipsum,fwlw}
\usepackage{xcolor}
\makeatletter
\def\ps@pagerange{\let\@mkboth\@gobbletwo
\let\@oddhead\@empty\let\@evenhead\@oddhead
\def\@oddfoot{\rlap{\color{white}%
Page=\thepage–First=\usebox\FirstWordBox–%
Last=\usebox\LastWordBox–%
Next=\usebox\NextWordBox–}%
\hfil\thepage\hfil}%
\let\@evenfoot\@empty
\let\chaptermark\@gobble
\let\sectionmark\@gobble
\let\subsectionmark\@gobble
}
\makeatother
\pagestyle{pagerange}
\begin{document}
\lipsum[1-100]
\end{document}
```

**Figure 1:** Minimum worked example to expose catchwords for retrieval.

For the extraction we used *pdftotext*, a freely-available utility which creates a plaintext version of a PDF document. In this, page-breaks are signalled

with an ASCII decimal 12 (0x0C) character, which is the Control-L or FF (FormFeed). In this example, the few lines immediately above each of the page-breaks contains the page number preceded by the delimited string we defined in `\ps@pagerange` in Figure 1.

In Figure 2 you can see two fragments of the output, the first from page 1 and the second from page 16 showing a problem where the page number occasionally gets imbrangled in the ‘Next’ catchword. This has not been resolved.

```
lorem lorem, interdum eu, tincidunt sit amet,
Page=1–First=–Last=amet,–Next=laoreet–
1
~Llaoreet vitae, arcu. Aenean faucibus pede eu
ante, Praesent enim elit, rutrum

tempus magna. Aliquam ut purus. Proin tellus.
Page=16–First=amet,–Last=tellus.–Next=
16
Vestibulum–
~LVestibulum ante ipsum primis in faucibus
```

**Figure 2:** Text fragments of output from Figure 1 at pages 1 and 16.

Now that the data is plaintext, you can use the standard *grep* and *awk* text utilities (or Perl, or Python, or Lua, or whatever is your favourite scripting language *du jour*) to pull out the lines with the delimited page number, first, last, and next words. You can then programmatically step through each page number and locate the span of text delimited by the First and Last words, using the Next word as a cross-check.

The tricky bit is application-dependent: you then need to be able to reliably read your source text programmatically, find the first word on a page, scan forward to the last, check the following word is the next value, and then do whatever is needed to insert the page number at whatever point is appropriate for your document.<sup>1</sup>

In the case in point, the production text was stored as XML, so the delimiters they used for the line of data embedded in the PDF were actually < and > characters, so the extracted fragments were already XML. That way the lines extracted from the text file were used in XSLT to identify each location in the XML source, push the page numbers into

<sup>1</sup> Daniel did suggest it might be more tractable to write the page-break data to a separate external file rather than embedding it: I'd be interested to hear from anyone implementing this.

