

Nuweb Version 1.1.1
A Simple Literate Programming Tool

Preston Briggs¹

preston@tera.com

HTML scrap generator by John D. Ramsdell

ramsdell@mitre.org

scrap formatting and continuing maintenance by Marc W. Mengel

mengel@fnal.gov

¹This work has been supported by ARPA, through ONR grant N00014-91-J-1989.

Contents

1	Introduction	1
1.1	Nuweb	1
1.1.1	Nuweb and HTML	2
1.2	Writing Nuweb	3
1.2.1	The Major Commands	3
1.2.2	The Minor Commands	4
1.3	Running Nuweb	5
1.4	Generating HTML	6
1.5	Restrictions	6
1.6	Acknowledgements	7

Chapter 1

Introduction

In 1984, Knuth introduced the idea of *literate programming* and described a pair of tools to support the practise [2]. His approach was to combine Pascal code with T_EX documentation to produce a new language, WEB, that offered programmers a superior approach to programming. He wrote several programs in WEB, including **weave** and **tangle**, the programs used to support literate programming. The idea was that a programmer wrote one document, the web file, that combined documentation (written in T_EX [3]) with code (written in Pascal).

Running **tangle** on the web file would produce a complete Pascal program, ready for compilation by an ordinary Pascal compiler. The primary function of **tangle** is to allow the programmer to present elements of the program in any desired order, regardless of the restrictions imposed by the programming language. Thus, the programmer is free to present his program in a top-down fashion, bottom-up fashion, or whatever seems best in terms of promoting understanding and maintenance.

Running **weave** on the web file would produce a T_EX file, ready to be processed by T_EX. The resulting document included a variety of automatically generated indices and cross-references that made it much easier to navigate the code. Additionally, all of the code sections were automatically prettyprinted, resulting in a quite impressive document.

Knuth also wrote the programs for T_EX and METAFONT entirely in WEB, eventually publishing them in book form [4, 5]. These are probably the largest programs ever published in a readable form.

Inspired by Knuth's example, many people have experimented with WEB. Some people have even built web-like tools for their own favorite combinations of programming language and typesetting language. For example, CWEB, Knuth's current system of choice, works with a combination of C (or C++) and T_EX [7]. Another system, FunnelWeb, is independent of any programming language and only mildly dependent on T_EX [9]. Inspired by the versatility of FunnelWeb and by the daunting size of its documentation, I decided to write my own, very simple, tool for literate programming.¹

1.1 Nuweb

Nuweb works with any programming language and L^AT_EX [6]. I wanted to use L^AT_EX because it supports a multi-level sectioning scheme and has facilities for drawing figures. I wanted to be able to work with arbitrary programming languages because my friends and I write programs in many languages (and sometimes combinations of several languages), *e.g.*, C, Fortran, C++, yacc, lex, Scheme, assembly, Postscript, and so forth. The need to support arbitrary programming languages has many consequences:

No prettyprinting Both WEB and CWEB are able to prettyprint the code sections of their documents because they understand the language well enough to parse it. Since we want to use *any* language, we've got to abandon this feature. However, we do allow particular individual formulas or fragments of L^AT_EX

¹There is another system similar to mine, written by Norman Ramsey, called *noweb* [8]. It perhaps suffers from being overly Unix-dependent and requiring several programs to use. On the other hand, its command syntax is very nice. In any case, nuweb certainly owes its name and a number of features to his inspiration.

code to be formatted and still be parts of output files. Also, keywords in scraps can be surrounded by `@_` to have them be bold in the output.

No index of identifiers Because WEB knows about Pascal, it is able to construct an index of all the identifiers occurring in the code sections (filtering out keywords and the standard type identifiers). Unfortunately, this isn't as easy in our case. We don't know what an identifier looks like in each language and we certainly don't know all the keywords. (On the other hand, see the end of Section 1.2.2)

Of course, we've got to have some compensation for our losses or the whole idea would be a waste. Here are the advantages I can see:

Simplicity The majority of the commands in WEB are concerned with control of the automatic prettyprinting. Since we don't prettyprint, many commands are eliminated. A further set of commands is subsumed by L^AT_EX and may also be eliminated. As a result, our set of commands is reduced to only four members (explained in the next section). This simplicity is also reflected in the size of this tool, which is quite a bit smaller than the tools used with other approaches.

No prettyprinting Everyone disagrees about how their code should look, so automatic formatting annoys many people. One approach is to provide ways to control the formatting. Our approach is simpler—we perform no automatic formatting and therefore allow the programmer complete control of code layout. We do allow individual scraps to be presented in either verbatim, math, or paragraph mode in the T_EX output.

Control We also offer the programmer complete control of the layout of his output files (the files generated during tangling). Of course, this is essential for languages that are sensitive to layout; but it is also important in many practical situations, *e.g.*, debugging.

Speed Since nuweb doesn't do too much, the nuweb tool runs quickly. I combine the functions of **tangle** and **weave** into a single program that performs both functions at once.

Page numbers Inspired by the example of noweb, nuweb refers to all scraps by page number to simplify navigation. If there are multiple scraps on a page (say, page 17), they are distinguished by lower-case letters (*e.g.*, 17a, 17b, and so forth).

Multiple file output The programmer may specify more than one output file in a single nuweb file. This is required when constructing programs in a combination of languages (say, Fortran and C). It's also an advantage when constructing very large programs that would require a lot of compile time.

This last point is very important. By allowing the creation of multiple output files, we avoid the need for monolithic programs. Thus we support the creation of very large programs by groups of people.

A further reduction in compilation time is achieved by first writing each output file to a temporary location, then comparing the temporary file with the old version of the file. If there is no difference, the temporary file can be deleted. If the files differ, the old version is deleted and the temporary file renamed. This approach works well in combination with **make** (or similar tools), since **make** will avoid recompiling untouched output files.

1.1.1 Nuweb and HTML

In addition to producing L^AT_EX source, nuweb can be used to generate HyperText Markup Language (HTML), the markup language used by the World Wide Web. HTML provides hypertext links. When an HTML document is viewed online, a user can navigate within the document by activating the links. The tools which generate HTML automatically produce hypertext links from a nuweb source.

1.2 Writing Nuweb

The bulk of a nuweb file will be ordinary L^AT_EX. In fact, any L^AT_EX file can serve as input to nuweb and will be simply copied through, unchanged, to the documentation file—unless a nuweb command is discovered. All nuweb commands begin with an “at-sign” (@). Therefore, a file without at-signs will be copied unchanged. Nuweb commands are used to specify *output files*, define *macros*, and delimit *scraps*. These are the basic features of interest to the nuweb tool—all else is simply text to be copied to the documentation file.

1.2.1 The Major Commands

Files and macros are defined with the following commands:

@o *file-name flags scrap* Output a file. The file name is terminated by whitespace.

@d *macro-name scrap* Define a macro. The macro name is terminated by a return or the beginning of a scrap.

A specific file may be specified several times, with each definition being written out, one after the other, in the order they appear. The definitions of macros may be similarly specified piecemeal.

The above are typeset without page-breaks. To allow page breaks, use the variants @O and @D.

Scraps

Scraps have specific begin markers and end markers to allow precise control over the contents and layout. Note that any amount of whitespace (including carriage returns) may appear between a name and the beginning of a scrap.

@{*anything*@} where the scrap body includes every character in *anything*—all the blanks, all the tabs, all the carriage returns. This scrap will be typeset in verbatim mode.

@[*anything*@] where the scrap body includes every character in *anything*—all the blanks, all the tabs, all the carriage returns. This scrap will be typeset in paragraph mode, allowing sections of T_EX documents to be scraps, but still be prettyprinted in the document.

@(*anything*@) where the scrap body includes every character in *anything*—all the blanks, all the tabs, all the carriage returns. This scrap will be typeset in math mode. This allows this scrap to contain a formula which will be typeset nicely.

Inside a scrap, we may invoke a macro.

@<*macro-name*@> Causes the macro *macro-name* to be expanded inline as the code is written out to a file. It is an error to specify recursive macro invocations.

@<*macro-name*@(*a1* @, *a2* @) @> Causes the macro *macro-name* to be expanded inline with the parameters *a1*, *a2*, etc. Up to 9 parameters may be given.

@1, @2, ..., @9 In a macro causes the n'th macro parameter to be substituted into the scrap. If the parameter is not passed, a null string is substituted.

Note that macro names may be abbreviated, either during invocation or definition. For example, it would be very tedious to have to type, repeatedly, the macro name

```
@d Check for terminating at-sequence and return name if found
```

Therefore, we provide a mechanism (stolen from Knuth) of indicating abbreviated names.

```
@d Check for terminating...
```

Basically, the programmer need only type enough characters to identify the macro name uniquely, followed by three periods. An abbreviation may even occur before the full version; nuweb simply preserves the longest version of a macro name. Note also that blanks and tabs are insignificant within a macro name; each string of them is replaced by a single blank.

Sometimes, for instance during program testing, it is convenient to comment out a few lines of code. In C or Fortran placing `/* ... */` around the relevant code is not a robust solution, as the code itself may contain comments. Nuweb provides a command, to be used only inside scraps, whose effect is similar to L^AT_EX's "`%`":

`@%`

Note that each use of "`@%`" creates a newline within the scrap, which L^AT_EX's "`%`" does not.

When scraps are written to a program file or a documentation file, tabs are expanded into spaces by default. Currently, I assume tab stops are set every eight characters. Furthermore, when a macro is expanded in a scrap, the body of the macro is indented to match the indentation of the macro invocation. Therefore, care must be taken with languages (*e.g.*, Fortran) that are sensitive to indentation. These default behaviors may be changed for each output file (see below).

Flags

When defining an output file, the programmer has the option of using flags to control output of a particular file. The flags are intended to make life a little easier for programmers using certain languages. They introduce little language dependences; however, they do so only for a particular file. Thus it is still easy to mix languages within a single document. There are three "per-file" flags:

- d Forces the creation of `#line` directives in the output file. These are useful with C (and sometimes C++ and Fortran) on many Unix systems since they cause the compiler's error messages to refer to the web file rather than to the output file. Similarly, they allow source debugging in terms of the web file.
- i Suppresses the indentation of macros. That is, when a macro is expanded within a scrap, it will *not* be indented to match the indentation of the macro invocation. This flag would seem most useful for Fortran programmers.
- t Suppresses expansion of tabs in the output file. This feature seems important when generating `make` files.

1.2.2 The Minor Commands

We have two very low-level utility commands that may appear anywhere in the web file.

`@@` Causes a single "at sign" to be copied into the output.

`@_` Causes the text between it and the next `@_` to be made bold (for keywords, etc.)

`@i file-name` Includes a file. Includes may be nested, though there is currently a limit of 10 levels. The file name should be complete (no extension will be appended) and should be terminated by a carriage return.

`@rx` Changes the escape character '`@`' to '`x`'. This must appear before any scrap definitions.

Finally, there are three commands used to create indices to the macro names, file definitions, and user-specified identifiers.

`@f` Create an index of file names.

`@m` Create an index of macro name.

`@u` Create an index of user-specified identifiers.

I usually put these in their own section in the L^AT_EX document; for example, see Chapter ??.

Identifiers must be explicitly specified for inclusion in the @u index. By convention, each identifier is marked at the point of its definition; all references to each identifier (inside scraps) will be discovered automatically. To “mark” an identifier for inclusion in the index, we must mention it at the end of a scrap. For example,

```
@d a scrap @{
Let's pretend we're declaring the variables FOO and BAR
inside this scrap.
@| FOO BAR @}

```

I've used alphabetic identifiers in this example, but any string of characters (not including whitespace or @ characters) will do. Therefore, it's possible to add index entries for things like <<= if desired. An identifier may be declared in more than one scrap.

In the generated index, each identifier appears with a list of all the scraps using and defining it, where the defining scraps are distinguished by underlining. Note that the identifier doesn't actually have to appear in the defining scrap; it just has to be in the list of definitions at the end of a scrap.

1.3 Running Nuweb

Nuweb is invoked using the following command:

```
nuweb flags file-name...
```

One or more files may be processed at a time. If a file name has no extension, .w will be appended. L^AT_EX suitable for translation into HTML by L^AT_EX2HTML will be produced from files whose name ends with .hw, otherwise, ordinary L^AT_EX will be produced. While a file name may specify a file in another directory, the resulting documentation file will always be created in the current directory. For example,

```
nuweb /foo/bar/quux
```

will take as input the file /foo/bar/quux.w and will create the file quux.tex in the current directory.

If it should happen that a source file extension *already* is .tex, then we need to prevent nuweb's documentation file from over-writing the source file, so the documentation filename is given the special suffix _nuweb.tex. The same suffix can be forced for all files via the -m command-line flag (see below).

By default, nuweb performs both tangling and weaving at the same time. Normally, this is not a bottleneck in the compilation process; however, it's possible to achieve slightly faster throughput by avoiding one or another of the default functions using command-line flags. There are currently three flags that significantly limit output:

- c Avoid testing output files for change before updating them.
- o Suppress generation of the output files.
- t Suppress generation of the documentation file.

Thus, the command

```
nuweb -to /foo/bar/quux
```

would simply scan the input and produce no output at all.

There are several additional command-line flags:

- d Print "dangling" identifiers – user identifiers which are never referenced, in indices, etc.
- m Forces documentation filenames to end with _nuweb.tex
- n Forces scraps to be numbered sequentially from 1 (instead of using page numbers). This form is perhaps more desirable for small webs.

- p Suppress prepending some default \LaTeX macros
- s Doesn't print list of scraps making up each file following each scrap.
- v For "verbose," causes nuweb to write information about its progress to `stderr`.

1.4 Generating HTML

Nikos Drakos' $\text{\LaTeX}2\text{HTML}$ Version 0.5.3 [1] can be used to translate \LaTeX with embedded HTML scraps into HTML. Be sure to include the document-style option `html` so that \LaTeX will understand the hypertext commands. When translating into HTML, do not allow a document to be split by specifying "`--split 0`". You need not generate navigation links, so also specify "`--no_navigation`".

While preparing a web, you may want to view the program's scraps without taking the time to run $\text{\LaTeX}2\text{HTML}$. Simply rename the generated \LaTeX source so that its file name ends with `.html`, and view that file. The documentations section will be jumbled, but the scraps will be clear.

1.5 Restrictions

Because nuweb is intended to be a simple tool, I've established a few restrictions. Over time, some of these may be eliminated; others seem fundamental.

- The handling of errors is not completely ideal. In some cases, I simply warn of a problem and continue; in other cases I halt immediately. This behavior should be regularized.
- I warn about references to macros that haven't been defined, but don't halt. This seems most convenient for development, but may change in the future.
- File names and index entries should not contain any `@` signs.
- Macro names may be (almost) any well-formed \TeX string. It makes sense to change fonts or use math mode; however, care should be taken to ensure matching braces, brackets, and dollar signs. When producing HTML, macros are displayed in a preformatted element (`PRE`), so macros may contain one or more A, B, I, U, or P elements or data characters.
- Anything is allowed in the body of a scrap; however, very long scraps (horizontally or vertically) may not typeset well.
- Temporary files (created for comparison to the eventual output files) are placed in the current directory. Since they may be renamed to an output file name, all the output files should be on the same file system as the current directory.
- Because page numbers cannot be determined until the document has been typeset, we have to rerun nuweb after \LaTeX to obtain a clean version of the document (very similar to the way we sometimes have to rerun \LaTeX to obtain an up-to-date table of contents after significant edits). Nuweb will warn (in most cases) when this needs to be done; in the remaining cases, \LaTeX will warn that labels may have changed.

Very long scraps may be allowed to break across a page if declared with `@O` or `@D` (instead of `@o` and `@d`). This doesn't work very well as a default, since far too many short scraps will be broken across pages; however, as a user-controlled option, it seems very useful. No distinction is made between the upper case and lower case forms of these commands when generating HTML.

1.6 Acknowledgements

Several people have contributed their times, ideas, and debugging skills. In particular, I'd like to acknowledge the contributions of Osman Buyukisik, Manuel Carriba, Adrian Clarke, Tim Harvey, Michael Lewis, Walter Ravenek, Rob Shillingsburg, Kayvan Sylvan, Dominique de Waleffe, and Scott Warren. Of course, most of these people would never have heard or nuweb (or many other tools) without the efforts of George Greenwade.

Since maintenance has been taken over by Marc Mengel, online contributions have been made by:

- Walter Brown <wb@fnal.gov>
- Nicky van Foreest <n.d.vanforeest@math.utwente.nl>
- Javier Goizueta <jgoizueta@jazzfree.com>
- Alan Karp <karp@hp.com>

Bibliography

- [1] Nikos Drakos, *From text to hypertext: A post-hoc rationalisation of latex2html*, Computer Networks and ISDN Systems **27** (1994), 215–224.
- [2] Donald E. Knuth, *Literate programming*, The Computer Journal **27** (1984), no. 2, 97–111.
- [3] ———, *The T_EXbook*, Computers and Typesetting, vol. 1986aA, Addison-Wesley, Reading, MA, USA, 1986.
- [4] ———, *T_EX: The program*, Computers and Typesetting, vol. B, Addison-Wesley, Reading, MA, USA, 1986b1986.
- [5] ———, *Metafont: The program*, Computers and Typesetting, vol. D, Addison-Wesley, Reading, MA, USA, 1986d1986.
- [6] Leslie Lamport, *L^AT_EXa document preparation system user's guide and reference manual*, Addison-Wesley, Reading, MA, USA, 1985.
- [7] Silvio Levy, *WEB adapted to C, another approach*, TUB **8** (1987), no. 1, 12–13.
- [8] N. Ramsey, *Literate programming simplified*, IEEE Software **11** (1994), no. 5, 97–105.
- [9] Ross Williams, *Funnelweb user's manual*, `ftp.adelaide.edu.au` in `/pub/compression` and `/pub/funnelweb`, University of Adelaide, Adelaide, South Australia, Australia, 1992.