

# Interpreter

v.1.0, July 2011

Paul Isambert

zappathustra AT free DOT fr

## Introduction

*Interpreter* preprocesses input files before their contents is fed to  $\TeX$ . It is meant to write document with whatever markup one wishes to define while using normal  $\TeX$  macros in the background. As a simple example, suppose you have a macro `\bold` to put text in boldface; then *Interpreter* lets you map `*text*`, or `<strong>text</strong>`, or simply `!text`, or anything else, to `\bold{text}`. *Interpreter* doesn't perform any trickery with active characters; instead, it manipulates the strings representing the lines of a file and search for patterns.

There are two main advantages: first,  $\TeX$  documents can be typeset with a completely non- $\TeX$  syntax; second, if one uses some lightweight markup language, the source file is much easier to read and might even be more useful than the typeset PDF file, e.g. for some technical documentation you want to read directly in your text editor while writing code (powerful editors generally have their own documentation in such a format, for a good reason). A third advantage, not explored in this documentation, is that while feeding modified lines to  $\TeX$  you can also translate the original lines into, say, HTML, and write them to an external file, thus creating both PDF and HTML output at once.

## Input files

Once *Interpreter* is loaded with

```
\input interpreter
```

in plain  $\TeX$  or

```
\usepackage{interpreter}
```

in  $\LaTeX$ , files to be processed are input as follows:

```
\interpretfile{<language>}{<file>}
```

There should exist a file `i-<language>.lua` containing the language used in `<file>`. For instance, the source of this documentation is `interpreter-doc.txt`, input in the master file `interpreter-doc.tex` with

```
\interpretfile{doc}{interpreter-doc.txt}
```

and the interpretation to be used is defined in `i-doc.lua`. The contents of such an interpretation file is the object of the rest of this documentation.

## Paragraphs

*Interpreter* doesn't process lines one by one. Instead, it gathers an entire paragraph and then processes the lines. It is important because you can manipulate an entire paragraph when a given pattern is detected, and modify several lines according to what happens in only one. A paragraph in *Interpreter* has nothing to do with what  $\TeX$  considers a paragraph; instead, it is defined by the following string.

### `interpreter.paragraph`

(Default: *blank line with spaces ignored*)

A string to be interpreted as a paragraph boundary when *Interpreter* collects lines before processing them. The string actually represents a pattern, so magic characters are obeyed. The default is `%*`, i.e. a blank line is considered a paragraph boundary, spaces notwithstanding. Of course, the end of the file itself is a paragraph boundary.

## Declaring patterns

Once the lines of a paragraph have been collected, *Interpreter* searches them trying to match declared patterns, but it doesn't do so indiscriminately: patterns are searched in a given order, as explained below.

Patterns are searched for in each line only, i.e. no match can occur across lines. However, since you can manipulate entire paragraphs based on a match in one line, the limitation easily vanishes.

### `interpreter.add_pattern(<table>)`

This is the basic function used to defined patterns. The `<table>` may contain the following entries, along other entries *Interpreter* won't use but which can be useful to you, especially with `call` below. The function returns a table.

#### `class` (Default: *intepreter.default\_class*)

The class of the pattern. See the section on classes.

#### `pattern`

The pattern to match. Lua's magic characters are in force and

should be escaped with % if necessary, unless `nomagic` is true (or the pattern itself is the result of `interpreter.nomagic`).

### **nomagic** (Default: *false*)

A boolean deciding whether the pattern should be transformed with `interpreter.nomagic`.

### **replace**

The replacement for the pattern, applied only if there is no `call` entry. This may be a string, a table or a function. *Interpreter* simply executes something similar to `string.gsub()`, hence the replacement follows this function's ordinary syntax. More precisely, if `replace` is a string, the pattern is replaced with it; in this string, %n may be used to denote the *n*th capture in the pattern. If `replace` is a table, the first capture or the entire match (if there is no capture) is used as the key, and the associated value is used as the replacement. If `replace` is a function, it is called with the captures passed as arguments, or the entire match if there is no capture. For instance, the following pattern will replace all `*text*` with `\bold{text}`:

```
interpreter.add_pattern{
  pattern = "%*(.-)%*",
  replace = [[\bold{%1}]
}
```

### **offset** (Default: *0*)

The number of positions *Interpreter* should shift to the right after a match has occurred. Normally, *Interpreter* starts searching for another occurrence of the current pattern at the same position where it found the last one. However, loops might easily occur: the replacement for a pattern may very well contain another match for the same pattern, so *Interpreter* will get stuck. Suppose for instance you want to replace `TeX` with `\TeX`. The first match will do that, but then *Interpreter* will start searching again at the backslash, producing `\\TeX`, then `\\\TeX`, etc. In this case, if you set `offset` to 2 in the pattern, then search will start again at the e and no new match will occur.

### **call**

This entry shall contain a function to be called if there is a match (if this entry exists, `replace` isn't applied). It is

meant to perform complex tasks that aren't amenable to simple string replacement. The function will be executed as follows:

```
function (paragraph, line, index, pattern)
```

`paragraph` is a table representing the current paragraph; lines are stored at successive indices. The last line of this paragraph is always the paragraph boundary (see `interpreter.paragraph`), unless the paragraph stopped at the end of the file. The second argument, `line`, is a number representing the index in paragraph containing the line where the pattern was found; `index` is the position in this line where the match occurred. Finally, `pattern` is the entire table declared with `interpreter.add_pattern` and containing all the entries discussed here.

The function may return zero, one, or two numbers. If it returns none, the search for the next occurrence of the pattern will start again on the same line (rather, on the line with the same position in the paragraph), at `index`. If it returns one number, the search will resume at the same line but at position *n*, with *n* the returned number. Finally, if two numbers are returned, the search will resume at line *m* at position *n*, *m* and *n* being the returned values. Specifying which line should be examined when the search resumes might be necessary if the function adds new lines in the paragraph *before* the current line, since *Interpreter* only keeps count of line numbers.

The entire paragraph can thus be modified if necessary. For instance, suppose you want to declare comments in your source file with only `!Comment` in the first line, i.e.  $\TeX$  should ignore a paragraph such as:

```
!Comment
This should be ignored
by TeX
```

Then the following pattern will do (where the function requires only the first argument):

```
local function comment (paragraph)
  for n, l in ipairs(paragraph) do
    paragraph[n] = "%" .. l
  end
end
```

```

interpreter.add_pattern{
  pattern = "^!Comment",
  call    = comment
}

```

### interpreter.nomagic (string)

A function which reverses the usual Lua magic for patterns: ordinary magic characters are normal characters here, unless they are prefixed with %, in which case they are magic again. For instance, a pattern like `.+` is normally interpreted as “one or more characters”. If passed to this function, a pattern is returned meaning “a dot followed by a plus sign”. On the contrary, `%.%+` normally has the second interpretation, while with `interpreter.nomagic` it has the first one. The function makes another transformation: `...` is used to denote a capture `(.-)`. Thus `interpreter.nomagic('*...*')` returns a pattern matching any number of characters surrounded by stars and capturing those characters; this would be expressed in ordinary Lua magic as `%(.-)%*`.

### Classes

As already alluded to, the search for patterns isn't done at random. Instead, patterns are organized in classes, which are applied one after the other. More precisely, the process is as follows: *Interpreter* searches the entire paragraph for the first pattern in class 1, then for the second pattern in the same class, then for the third, etc., then when there is no pattern left in class 1 it does the same with class 2, up to class  $n$ , where  $n$  is the highest class number such that there exists a class  $n - 1$  (in other words, classes should be numbered consecutively). Finally, the same goes for the patterns in class 0 (which always exists, even if it contains no pattern).

Inside a class, patterns are ordered by length from long to short, or alphabetically if two patterns have the same length. This means that if you use e.g. `/text/` for italics and `//text//` for bold, you don't need to put the second pattern in a class before the first to avoid `//text//` being interpreted as two empty arguments in italics surrounding a text in roman. Since the way the bold-pattern will be declared, e.g. `/(.-)//`, is probably longer than for the italic-pattern, e.g. `/(.-)/`, it will always match first.

That said, the sorting isn't very clever and simply relies on the number of symbols, no matter what they mean; in

the patterns above, the parentheses denote a capture but they still count in the pattern's length as understood by *Interpreter*. Alternatively, while `.*` denotes “zero or more character” and `%+` means “a plus sign” (+ being magic, you have to escape it to refer to it), in *Interpreter's* eye the two patterns have the same length: two. Finally, one should be aware that patterns declared with a `nomagic` entry set to `true` are sorted after they've been transformed (so that their real length might not be obvious). So classes are needed when patterns need a proper ordering no matter their lengths. For instance, some patterns should always be declared first, as they protect input from *Interpreter* (see next section), while others might need to be declared last, as they rely on what previous patterns might have done. Besides, classes are metatables for the patterns they contain.

### interpreter.default\_class (Default: 1)

All patterns belong to a class, even though you may omit the `class` entry when declaring one. In this case, the pattern is assigned to the class denoted by this number.

### interpreter.set\_class(number, table)

Defines class `number` as `table`. Classes don't need to be defined beforehand for patterns to be added to them (rather, *Interpreter* defines them implicitly when needed). However, classes are also metatables for the patterns, so that if there lacks an entry in a pattern's table, the class's entry is used if it exists. The function returns a table.

### Protecting input

Sometimes you want *Interpreter* to refrain from interpreting; that is most useful for verbatim code, for instance. There are various ways to do that.

### interpreter.active (Default: true)

A boolean switching *Interpreter* on and off. Beware, the switching applies only starting at the next paragraph.

### interpreter.protect([line])

A function protecting all or part of the current paragraph. If `line` is given, it should be a number  $n$ , and line  $n$  in the current paragraph will be protected; without `line`, the entire paragraph is protected. Protecting means that the patterns not yet searched for will be ignored. For instance, if you want material to be read verbatim when surrounded with

`<code>` and `</code>`, you can declare a pattern as follows:

```
local function verbatim (buffer)
  buffer[1] = "\\verbatim"
  buffer[#buffer - 1] = "\\endverbatim"
  interpreter.protect()
end
interpreter.add_pattern{
  pattern = "^%s*<code>%*s$",
  call    = verbatim,
  class   = 1
}
```

This code is extremely simplified: it assumes that `<code>` and `</code>` starts and ends the paragraph and that `</code>` isn't the last line of the file (otherwise it'd also be the last line in the paragraph, whereas here the last one is the paragraph boundary). An important point is that the pattern belongs to the first class, so it is called before all other patterns (provided there is no shorter pattern in class 1) and prevents them from doing anything, since the entire paragraph is protected. (Typesetting the material as verbatim material obviously depends on the `\verbatim` macro, not on *Interpreter*.)

### `interpreter.escape`

A character which prevents patterns from being replaced if immediately preceded by it. As an example, if `interpreter.escape = '_'`, and `*text*` denotes italic, then `*text*` will produce *text* while `_*text*` will produce *\*text\**. Once a paragraph has been processed, *Interpreter* removes all escape characters. Only one character can be an escape character.

### `interpreter.protector(left[, right])` (*right defaults to left*)

Defines two characters to protect what they surround. In other words, *Interpreter* replaces patterns only if the match isn't found between `left` and `right`. Unlike the escape character, you can define as many protectors as you wish; and unlike the escape character again, *Interpreter* doesn't remove them once the paragraph has been processed, so you must take care of them. For instance:

```
interpreter.protector('')
```

```
interpreter.add_pattern{
  pattern = '"(.-)"',
  replace = '\\verb`%1`',
  class   = 0
}
```

Anything between double quotes will be left untouched; then, when the paragraph has been processed for all other classes, a pattern in class 0 calls the `\verb` command to take care of the argument. Note that the protectors should enclose what they protect without coinciding with it; this is not the case here, which is why the pattern is applied.

### `interpreter.direct`

(*Default: two percent signs then I and at least one space*)

A string, actually a pattern, signalling that the line which it begins should be processed as Lua code. The default is `%%I%S+`, i.e. `%I` followed by at least one space. The pattern shouldn't declare itself as attached to the beginning of the line (as in `^%%I%S+`) because they will be matched at the beginning of the line only anyway. The line is processed with the `loadstring` function, and then turned into an empty line. For instance:

```
%I interpreter.active = false
This won't be interpreted...
%i interpreter.active = true
```

As this example shows, lines flagged with `interpreter.direct` don't obey `interpreter.active` and are always processed as described above.

### Technical stuff

You don't have to bother with this section if you don't mind how *Interpreter* does its job; actually you won't learn much anyway.

### `interpreter.reset()`

A function which resets everything to default and deletes classes. It is used when calling `\interpretfile` so that new interpretations start from zero.

### `interpreter.register(function)`

A function called to put *Interpreter*'s main function into the `post_linebreak_filter` callback; you can redefine it at

will. If it is undefined, `callback.register()` is used, unless `luatexbase.add_to_callback()` is detected. (The detection takes place at the first call to `\interpretfile`, so there is no need to load *Interpreter* after `luatexbase`.)

### `interpreter.unregister(function)`

A function called to remove *Interpreter*'s main function from the `post_linebreak_filter` callback. It works similarly to the previous one.

### An example: `i-doc.lua`

Here's a description of `i-doc.lua`, the file containing the interpretation used for *Interpreter*'s documentation. Remember that none of the TeX macros used here is defined by *Interpreter*; instead, they are my own and should be adapted if necessary. Also several options taken here are far from optimal but are convenient examples.

Shorthands for often used functions.

```
local gsub, match = string.gsub, string.match
local add_pattern = interpreter.add_pattern
local nomagic     = interpreter.nomagic
```

Class 1 and 2 will be used for verbatim (thus protecting) and "normal" patterns go into class 3 or higher.

```
interpreter.default_class = 3
```

The reader might have observed that `interpreter-doc.txt` begins with a table of contents. This table is useful for the source file only, and isn't typeset by TeX, because the following pattern suppresses it: the entire paragraph containing `TABLE OF CONTENTS` on a line of its own is deleted. Protecting the paragraph is useless, but it makes things a little bit faster because the paragraph won't be pointlessly searched for other patterns.

```
local function contents (buffer)
  for n in ipairs(buffer) do
    buffer[n] = ""
  end
  interpreter.protect()
end
add_pattern{
  pattern = "^%s*TABLE OF CONTENTS%s*$",
```

```
  call    = contents,
  class   = 1
}
```

Sections headers are typeset as

```
===== section_tag
=== Section title =====
=====
```

The first and third line are decorations and they are removed. The `section_tag` is meant for the source only again (linking the section to the table of contents). I could have used it to create PDF destinations, but that seemed unnecessary in such a small file. The associated pattern is: at least four equals signs.

```
add_pattern{
  pattern = "^====+.*",
  replace = ""
}
```

The middle line is spotted with the tree equals sign at the beginning of the line (the previous pattern being longer, the decoration lines have been already removed and they won't be taken for section titles). The signs are removed and replaced with `\section{` and `}`.

```
local function section (buffer, num)
  local l = buffer[num]
  l = gsub(l, "^===%s*", "\\section{")
  l = gsub(l, "%s*=%s*", "}")
  buffer[num] = l
end
add_pattern{
  pattern = "^===",
  call    = section
}
```

The following pattern simply turns *Interpreter* into `\ital{Interpreter}`. The meaning of the `\ital` command is obvious, I suppose. Note the offset: starting at the backslash, this leads to the `n` in *Interpreter*, thus avoiding matching the pattern again. The Lua notation with double square brackets is used for strings with no escape character (hence `\ital`

and not `\ital` as would be necessary with a simple string).

```
add_pattern{
  pattern = "Interpreter",
  replace = [[\ital{Interpreter}]],
  offset  = 7
}
```

Turning `tex` into `TEX`. This illustrates the use of a function as `replace`; the point is that `\tex` should be suffixed with a space if initially followed by anything but a space or end of line (so as not to form a control sequence with the following letters), and it should be suffixed with a control space if initially followed by a space or end of line (so as to avoid gobbling the space). So the function checks the second capture. Note that simply replacing `tex` with `\tex{}` would be much simpler, but less instructive!

```
local function maketex (tex, next)
  if next == " " or next == "" then
    return [[\tex\ ]]
  else
    return [[\tex ]] .. next
  end
end
add_pattern{
  pattern = "(Tex)(.?)",
  replace = maketex,
  offset  = 2
}
```

The following turns `<text>` into `<text>` and `_text_` into `text`. Setting a class just so the patterns inherit the `nomagic` feature is of course an overkill, but that's an example.

```
interpreter.set_class(4, {nomagic = true})
add_pattern{
  pattern = "<...>",
  replace = [[\arg{%1}]],
  class   = 4
}
add_pattern{
  pattern = "_..._",
  replace = [[\ital{%1}]],
  class   = 4 }
}
```

I use double quotes as protectors; they are replaced with a `\verb` command at the very end of the processing (with class 0).

```
interpreter.protector('')
add_pattern{
  pattern = nomagic'"..."',
  replace = [[\verb`%1`]],
  class   = 0
}
```

The description of functions (in red in the PDF file) are handled with the `\describe` macro, which takes the function as its first argument and additional information as its second one (typeset in italics in the PDF file). In the source, it is simply marked as

```
> function (arguments) [Additional information]
```

with `[Additional information]` sometimes missing (i.e. there is no empty pairs of square brackets). Descriptions of entries in pattern tables follows the same syntax, except the line begins with `>>`. So the pattern first spots lines beginning with `>[>]` followed by at least one space, adds an empty pair of brackets at the end if there isn't any, and turn the whole into `\describe`. The number of `>` symbols sets `\describe`'s third argument, which specifies the level of the bookmark.

```
local function describe (buffer, num)
  local l = buffer[num]
  if not match (l, "%[.--%]s*$") then
    l = l .. " []"
  end
  local le = match(l, ">>") and 4 or 3
  buffer[num] = gsub(l, ">+%s+(-)%s+%[(-)%]",
    [[\describe{%1}{%2}{}] .. le .. ""])
end
add_pattern{
  pattern = "^>+%s+",
  call    = describe
}
```

Here's how multiline verbatim is handled; in the source it is simply marked by indenting the line with ten spaces; thus code is easily spotted when reading the source without

useless and annoying `</code>` or anything similar to mark it. To be properly processed by  $\TeX$ , the code should be surrounded by `\verbatim` and `\verbatim/` (my way of signalling blocks). Those must be on their own lines, so we insert a line at the beginning and at the end of the paragraph: for the closing `\verbatim/`, we can simply replace the last line of the paragraph, which is the boundary line, unless we're at the end of the file. But for the opening `\verbatim` a line must be added at the beginning of the paragraph; thus line numbers in the original source file and in its processed version don't match anymore, and this might be annoying when  $\TeX$  reports errors. Besides, blank `verbatim` lines aren't handled correctly and create a new `verbatim` block instead. So this way of marking `verbatim` material is good for small documents, but explicit marking is cleaner and more powerful (albeit not so good-looking in the source file).

Note that the `verbatim` pattern belongs to class 2 and the entire paragraph is protected, so *Interpreter* leaves it alone afterward (remember the default class is 3). Of course, the first ten space characters are removed.

```
local function verbatim (buffer)
  for n, l in ipairs(buffer) do
    buffer[n] = gsub(l, "%s%s%s%s%s%s%s%s", "")
  end
  table.insert(buffer, 1, [[\verbatim]])
  if gsub(buffer[#buffer],
    interpreter.paragraph, "") == "" then
    buffer[#buffer] = [[\verbatim/]]
  else
    table.insert(buffer, [[\verbatim/]])
  end
  interpreter.protect()
end
add_pattern{
  pattern = "^%s%s%s%s%s%s%s%s",
  call    = verbatim,
  class   = 2
}
```

And now comes the fun part. I wanted `i-doc.lua` to be self-describing. The source of what you're reading right now isn't `interpreter-doc.txt`, but `i-doc.lua` itself input in the latter file with

```
\interpreterfile{doc}{i-doc.lua}
```

How should code and comment be organized in `i-doc.lua`? Well, there is little choice, since the file is a normal Lua file: comment lines should be prefixed with `--` or surrounded with `--[[` and `--]]`. I chose the latter option, which is simpler. But normal code should also be typeset as `verbatim` material; I could have begun all lines with ten spaces, but that would have seemed strange. Instead, `--]]` is turned into `\source` and `\source/` is added at the end of the paragraph (`\source` is just `\verbatim` with a different layout). Which means all paragraphs have the same structure: comments between `--[[` and `--]]` and code immediately following (`--[[` is simply removed). The pattern is in class 1 and the paragraph is protected, so that lines indented with ten spaces or more aren't touched by the previous `verbatim` pattern (in class 2).

```
local function autoverbatim (buffer, line)
  buffer[line] = [[\source]]
  for n = line + 1, #buffer do
    interpreter.protect(n)
  end
  if gsub(buffer[#buffer],
    interpreter.paragraph, "") == "" then
    buffer[#buffer] = [[\source/]]
  else
    table.insert(buffer, [[\source/]])
  end
end
add_pattern{
  pattern = nomagic"%^--]",
  call    = autoverbatim,
  class   = 1
}
local function test ()
  return ""
end
add_pattern{
  pattern = nomagic"%^--[[",
  replace = test,
}
```

*Typeset with Lua $\TeX$  0.71 in Chaparral Pro and Lucida Console*