

# lua-list-hyphen — Per-language listing of hyphenated words for Lua $\text{\LaTeX}$ \*

Alan J. Cain<sup>†</sup>

Released 2026-05-27

## Abstract

This Lua $\text{\LaTeX}$  package writes each word that has been hyphenated across lines to a file, using a different file for each language, for subsequent external checking.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Requirements</b>	<b>3</b>
<b>3</b>	<b>Installation</b>	<b>3</b>
<b>4</b>	<b>Getting started</b>	<b>3</b>
<b>5</b>	<b>Package options</b>	<b>4</b>
<b>6</b>	<b>Output format</b>	<b>5</b>
<b>7</b>	<b>Usage notes</b>	<b>5</b>
7.1	Languages . . . . .	5
7.2	Limitations . . . . .	5
<b>8</b>	<b>Implementation (<math>\text{\LaTeX}</math> package)</b>	<b>6</b>
8.1	Initial set-up . . . . .	6
8.2	Options . . . . .	6
8.3	Processing package options . . . . .	7
8.4	Lua backend . . . . .	8
8.5	Saving <code>babel</code> language names . . . . .	8
8.6	Processing and writing hyphenation lists . . . . .	8

---

\*This document describes v0.3.20, last revised 2026-05-27.

<sup>†</sup>`a.j.cain (AT) gmail.com`

<b>9</b>	<b>Implementation (Lua backend)</b>	<b>9</b>
9.1	Debugging function . . . . .	9
9.2	Table key constants . . . . .	9
9.3	Segment type . . . . .	10
9.4	Node ID and subtype constants . . . . .	10
9.5	Output constants . . . . .	10
9.6	Utility functions . . . . .	11
9.7	Getting text from nodes . . . . .	11
9.8	String manipulation . . . . .	14
9.9	Pre-linebreak processing . . . . .	15
9.10	Post-linebreak processing . . . . .	18
9.11	Callbacks . . . . .	23
9.12	Language settings . . . . .	24
9.13	Processing hyphenation lists . . . . .	24
	9.13.1 Comparisons and equality checks . . . . .	25
	9.13.2 Sorting . . . . .	26
	9.13.3 Deduplication . . . . .	27
	9.13.4 Combined processing . . . . .	28
9.14	Writing . . . . .	28
9.15	Export public functions . . . . .	32

## Index 33

## 1 Introduction

TeX’s algorithm for finding points where a word can be hyphenated is good, but not perfect.<sup>1</sup> The present author writes in British English, where the valid division points can depend on both the pronunciation of a word and its internal structure (and hence its etymology). Currently, TeX’s pattern-based approach produces *bio-lo-gic*, *bio-logy*, *bio-lo-gist*, rather than the standard *bio-logic*, *biol-ogy*, *biolo-gist*.<sup>2</sup> To deal with such cases, at least a substantially larger number of patterns would be required than are available at present. There are also various words where the valid division points in British English cannot be deduced from their spelling alone: for instance, the verbs *at-trib-ute*, *pre-sent*, *pro-duce*, *re-cord* have different division points from the orthographically identical nouns *at-tri-bute*, *pres-ent*, *prod-uce*, *rec-ord*. For another example, compare *cur-ric-ulum vitae* and *school cur-ricu-lum*.

Easy checking of the chosen hyphenations is desirable. With LuaTeX, it is possible to extract the hyphenated words. The LuaLaTeX package `lua-check-hyphen` offers this facility. It checks hyphenated words against a whitelist, visually flags unknown hyphenations, and writes unknown hyphenations to a file. But it was first written in 2012, when LuaTeX was at an earlier stage of development, and so it has certain problems, such as with words containing ligatures. It also lacks multi-language support.

This LuaLaTeX package, `lua-list-hyphen`, uses some ideas from `lua-check-hyphen` but was written from scratch to work with a modern LuaTeX. It simply writes hyphenated

<sup>1</sup>For a description of the algorithm and its limitations, see Knuth’s account in Appendix H of *The TeXbook* (Addison-Wesley, 2021. ISBN: 978-0-201-13447-6)

<sup>2</sup>See the *New Oxford Spelling Dictionary*, which is the authority for word divisions in British English (Oxford University Press, 2005. ISBN: 978-0-19-860881-3).

words from each language to a separate file, so that they can be checked (manually or by an external program).

[The author has written a simple Python application `hyphenassist`<sup>3</sup> that checks the listed hyphenations against a dictionary of valid divisions and allows the user to quickly choose to add entries to the division dictionary, add hyphenation exceptions, or ignore particular hyphenations. He has used this program in conjunction with code incorporated into this package to check hyphenations in his own books.<sup>4</sup>]

**Licence.** `lua-list-hyphen` is released under the L<sup>A</sup>T<sub>E</sub>X Project Public Licence v1.3c or later.<sup>5</sup>

**Acknowledgements.** The author thanks Keno Wehr for corrections and comments on the documentation.

**Feature requests and bug reports** The development code and issue tracker are hosted at Codeberg.<sup>6</sup>

## 2 Requirements

`lua-list-hyphen` requires

- (1) LuaL<sup>A</sup>T<sub>E</sub>X,
- (2) a recent L<sup>A</sup>T<sub>E</sub>X kernel with `expl3` support (any kernel version since 2020-02-02 should suffice).

It does not depend on any other packages, but will interface with `babel` or `polyglossia` (if one of them is loaded) to determine language names.

## 3 Installation

To install `lua-list-hyphen` manually, run `luatex lua-list-hyphen.ins` and copy `lua-list-hyphen.sty` and `lua-list-hyphen.lua` to somewhere LuaL<sup>A</sup>T<sub>E</sub>X can find them.

## 4 Getting started

Simply load the package; the hyphenated words are by default written to the file `\jobname-⟨lang-id⟩.hyph`, without being sorted or having duplicates removed. The `⟨lang-id⟩` is either a LuaL<sup>A</sup>T<sub>E</sub>X numerical language ID, or a `babel` or `polyglossia` name of the language, if one of these packages is in use. The prefix `\jobname-` and the extension `.hyph` can be customized; see [Section 5](#).

---

<sup>3</sup>URL: <https://codeberg.org/ajcain/hyphenassist>.

<sup>4</sup>In particular, *Form & Number: A History of Mathematical Beauty*. URL: [https://archive.org/details/cain\\_formandnumber\\_ebook\\_large](https://archive.org/details/cain_formandnumber_ebook_large).

<sup>5</sup>URL: <https://www.latex-project.org/lppl.txt>

<sup>6</sup>URL: <https://codeberg.org/ajcain/lua-list-hyphen>

## 5 Package options

**verbose** The boolean option **verbose** controls how much information is written to the file about each hyphenated word. When **true**, for each hyphenated word, both the undivided original and the divided word are written out, as well as the page number on which the hyphenated word appears (or, more precisely, begins) and the undivided word in context (as specified by the **context** keys; see below). When **false**, only the hyphenated word is written. (*Default: false*)

**context** Integer options controlling how many words before (**context-before**) and after (**context-after**) the hyphenated word are written as context when **verbose=true**. The key **context** is simply a shortcut for setting **context-before** and **context-after** to the same value. (*Default: 2*)

**unique** The option **unique** controls removal of duplicates from the list of hyphenated words written out. It can be set to one of the following three values:

**none:** Duplicate hyphenations are not removed.

**case:** Hyphenations that are duplicate (case-sensitively) are removed. In this case, the hyphenations **geo-metry** and **Geo-metry** are considered to be distinct.

**nocase:** Hyphenations that are duplicate (case-insensitively) are removed. In this case, the hyphenations **geo-metry** and **Geo-metry** are considered to be duplicates. The case of each listed hyphenation will be that of the first appearance of that hyphenation.

Note that removal of duplicates is unaffected by the page number or context that is written out when **verbose=true**. (*Default: none*)

**sort** The option **sort** controls sorting of the list of hyphenated words. It can be set to one of the following three values:

**none:** Hyphenations appear in the same order as they occur in the document, or, if duplicates are removed, in the order of first appearance in the document.

**case:** Hyphenations are sorted case-sensitively. In this case, **Geo-metry** precedes **geo-meter**.

**nocase:** Hyphenations are sorted case-insensitively. In this case, **geo-meter** precedes **Geo-metry**.

(*Default: none*)

**include-non-output** Boolean option determining whether hyphenated words that are never written to the page are listed. (For instance, a hyphenated word might occur in text that a package temporarily typesets into a box, measures, and then discards.) (*Default: false*)

The two options **prefix** and **extension** specify the files to which hyphenations are written. Between the prefix and the extension is either a LuaTeX numerical language ID, or a **babel** or **polyglossia** name of the language, if one of these packages is in use.

**prefix** The **prefix** is the part of the file name to which the list of hyphenated words is written, before the language ID. (*Default: \jobname-* (note the hyphen).)

**extension** The extension of the file (including the **.**) to which the list of hyphenated words for each language is written. (*Default: .hyph*)

**debug** The boolean option **debug** controls whether debugging information is written to the terminal. (*Default: false*)

## 6 Output format

Each output file begins with a header (each line of which begins with a ‘comment’ symbol %) that includes information about the language and the package options that were used. Each line of the remainder of the file describes one hyphenation.

When `verbose=false`, the line contains only the hyphenated word.

When `verbose=true`, the line contains the original undivided word, the hyphenated word, the page number where the hyphenated word appears (or, to be precise, begins), and the context in which the hyphenated word appears. Each part of the output is padded so that the various lines align. The original and undivided words are separated by the ASCII ‘arrow’ `->`; the page number is prefixed by `p.`; and the context is surrounded by (straight) quotation marks `" "`. If the hyphenation was never written to the page, `p.<page>` is replaced by `<none>`. (This can only happen with `include-non-output=true`.)

## 7 Usage notes

### 7.1 Languages

To determine the language of a word, `lua-list-hyphen` looks at what language is applied at the first possible hyphenation point, first considering the part of the word before it, then the part after it. In the (presumably rare) case of a ‘mixed-language’ word like ‘near-Zugzwang’ being specified (using, for example, `babel`) with `near-\foreignlanguage{german}{Zugzwang}`, it would be assigned to the language in which ‘near-’ is set.

Duplicates are removed within each language. If the same hyphenation occurs in two different languages, it will appear in both files, regardless of the value of the `unique` package option.

### 7.2 Limitations

`lua-list-hyphen` uses LuaTeX’s built-in Unicode functions for pattern matching and converting between upper and lower case. These functions are based on the `slnunicode` library. This library has not been updated for some time and is based on an out-of-date version of the Unicode standard. Thus there may be problems with languages added to Unicode more recently. Hyphenated words from such languages should still be listed, but may contain extraneous characters (such as adjacent punctuation) and may not be sorted correctly. Users may prefer to leave sorting and removal of duplicates to an external program that adheres to the current Unicode standard.

## 8 Implementation (L<sup>A</sup>T<sub>E</sub>X package)

```

1 <{*package>
2 <@@=lualisthyphen>

```

### 8.1 Initial set-up

Package identification/version information.

```

3 \NeedsTeXFormat{LaTeX2e}[2020-02-02]
4 \ProvidesExplPackage{lua-list-hyphen}{2026-05-27}{0.3.20}{
5   {Listing hyphenated words for LuaLaTeX}

```

Check that Lua<sub>T</sub><sub>E</sub>X is in use.

```

6 \sys_if_engine luatex:F
7   {
8     \msg_new:nnn{ lua-list-hyphen }{ luatex_required }
9     { LuaLaTeX~required.~Package~loading~will~abort. }
10    \msg_critical:nn{ lua-list-hyphen }{ luatex_required }
11  }

```

### 8.2 Options

`\l_lualisthyphen_verbose_bool` Boolean option to indicate whether lists of hyphenations should be written verbosely.

```

12 \keys_define:nn { lua-list-hyphen }{
13   verbose .bool_set:N = \l_lualisthyphen_verbose_bool,
14 }

```

*(End of definition for \l\_lualisthyphen\_verbose\_bool.)*

`\l_lualisthyphen_context_before_int` Integer options to determine the number of words before and after a hyphenation shown as context in verbose output.

```

15 \keys_define:nn { lua-list-hyphen }{
16   context-before .int_set:N = \l_lualisthyphen_context_before_int,
17   context-before .initial:n = { 2 },
18   context-after .int_set:N = \l_lualisthyphen_context_after_int,
19   context-after .initial:n = { 2 },
20   context .code:n = {
21     \keys_set:nn{ lua-list-hyphen }{
22       context-before=#1,
23       context-after=#1,
24     }
25   },
26
27 }

```

*(End of definition for \l\_lualisthyphen\_context\_before\_int and \l\_lualisthyphen\_context\_after\_int.)*

`\l_lualisthyphen_unique_int` Choice option to indicate whether lists of hyphenations should have duplicates removed, case-sensitively or case-insensitively.

```

28 \int_new:N\l_lualisthyphen_unique_int
29 \keys_define:nn { lua-list-hyphen }{
30   unique .choices:nn = { none, case, nocase }{
31     \int_set:Nn\l_lualisthyphen_unique_int{ \l_keys_choice_int - 1 }
32   },
33 }

```

(End of definition for \l\_\_lualisthyphen\_unique\_int.)

`\l__lualisthyphen_sort_int` Choice option to indicate whether lists of hyphenations should be sorted, case-sensitively or case-insensitively.

```

34 \int_new:N\l__lualisthyphen_sort_int
35 \keys_define:nn { lua-list-hyphen }{
36   sort .choices:nn = { none, case, nocase }{
37     \int_set:Nn\l__lualisthyphen_sort_int{ \l_keys_choice_int - 1 }
38   },
39 }

```

(End of definition for \l\_\_lualisthyphen\_sort\_int.)

`\l__lualisthyphen_include_non_output_bool` Boolean option to indicate whether lists of hyphenations should include those that are never output to the page.

```

40 \keys_define:nn { lua-list-hyphen }{
41   include-non-output .bool_set:N = \l__lualisthyphen_include_non_output_bool,
42 }

```

(End of definition for \l\_\_lualisthyphen\_include\_non\_output\_bool.)

`\l__lualisthyphen_file_prefix_str` String option for the prefix of files to which hyphenations are wrtitten.

```

43 \keys_define:nn { lua-list-hyphen }{
44   prefix .str_set:N = \l__lualisthyphen_file_prefix_str,
45   prefix .initial:e = { \c_sys_jobname_str- },
46 }

```

(End of definition for \l\_\_lualisthyphen\_file\_prefix\_str.)

`\l__lualisthyphen_file_extension_str` String option for the extension of files to which hyphenations are wrtitten.

```

47 \keys_define:nn { lua-list-hyphen }{
48   extension .str_set:N = \l__lualisthyphen_file_extension_str,
49   extension .initial:n = { .hyph },
50 }

```

(End of definition for \l\_\_lualisthyphen\_file\_extension\_str.)

`\l__lualisthyphen_debug_int` Option to specify whether debug information is written to the terminal. Not intended for end users.

```

51 \int_new:N\l__lualisthyphen_debug_int
52 \keys_define:nn { lua-list-hyphen }{
53   debug .code:n = {\int_set_eq:NN\l__lualisthyphen_debug_int\c_one_int}
54 }

```

(End of definition for \l\_\_lualisthyphen\_debug\_int.)

### 8.3 Processing package options

Process package options.

```

55 \ProcessKeyOptions [ lua-list-hyphen ]
    Convert boolean options to integers (which can be accessed from Lua).
56 \int_new:N\l__lualisthyphen_verbose_int
57 \bool_if:NT\l__lualisthyphen_verbose_bool
58   { \int_set_eq:NN\l__lualisthyphen_verbose_int\c_one_int }
59 \int_new:N\l__lualisthyphen_include_non_output_int
60 \bool_if:NT\l__lualisthyphen_include_non_output_bool
61   { \int_set_eq:NN\l__lualisthyphen_include_non_output_int\c_one_int }

```

## 8.4 Lua backend

Load the Lua backend.

```
62 \lua_now:n{
63   lualisthyphen = require('lua-list-hyphen')
64 }
```

## 8.5 Saving babel language names

At `enddocument/afterlastpage`, if possible save babel's language names. (polyglossia's names can be found directly from Lua.)

```
65 \hook_gput_code:nnn{ enddocument/afterlastpage }{ lua-list-hyphen } {
66   \__lualisthyphen_babel_save_language_names:
67 }
```

`\__lualisthyphen_babel_save_language_names:` If babel is in use, get language names from `\bbl@languages`.

```
68 \cs_new:Npn \__lualisthyphen_babel_save_language_names:
69 {
70   \cs_if_exist:NT\bbl@languages
71   {
```

Iterate through `\bbl@languages` to get language names. Items stored in this macro are quadruples prefixed with `\bbl@elt`, so locally redefine this latter macro to an auxiliary function that passes language ID/name pairs to the Lua backend.

```
72     \group_begin:
73     \cs_set_eq:NN
74       \bbl@elt
75       \__lualisthyphen_babel_save_language_names_elt:nnnn
76     \bbl@languages
77     \group_end:
78   }
79 }
```

*(End of definition for `\__lualisthyphen_babel_save_language_names:.`)*

`\__lualisthyphen_babel_save_language_names_elt:nnnn`

Auxiliary function that takes a quadruple stored in `\bbl@languages` and passes language ID/name pairs to the Lua backend.

```
80 \cs_new:Npn \__lualisthyphen_babel_save_language_names_elt:nnnn #1#2#3#4
81 {
82   \lua_now:n{
83     lualisthyphen.babel_save_language_name(#2,'#1')
84   }
85 }
```

*(End of definition for `\__lualisthyphen_babel_save_language_names_elt:nnnn.`)*

## 8.6 Processing and writing hyphenation lists

At `enddocument/info`, process and output the hyphenations that have been found.

```
86 \hook_gput_code:nnn{ enddocument/info }{ lua-list-hyphen } {
87   \__lualisthyphen_process_write_hyphenation_lists:ee
88   {\str_use:N\l__lualisthyphen_file_prefix_str}
89   {\str_use:N\l__lualisthyphen_file_extension_str}
90 }
```



sthyphen\_process\_write\_hyphenation\_lists:nn

Sort the list of hyphenations into separate lists for each language, sort and deduplicate them as required, and write them to files with prefix given in the first parameter and suffix in the second.

```
91 \cs_new:Npn \__lualisthyphen_process_write_hyphenation_lists:nn #1#2
92 {
93   \lua_now:e{
94     lualisthyphen.process_write_hyphenation_lists(
95       '\luaescapestring{#1}',
96       '\luaescapestring{#2}'
97     )
98   }
99 }
100 \cs_generate_variant:Nn
101   \__lualisthyphen_process_write_hyphenation_lists:nn
102   { ee }

(End of definition for \__lualisthyphen_process_write_hyphenation_lists:nn.)
103 </package>
```

## 9 Implementation (Lua backend)

104 <lua>

### 9.1 Debugging function

debug Debugging function. Defined according to the package option `debug` to either do nothing or write debugging information.

```
105 local debug
106
107 if tex.count['l__lualisthyphen_debug_int'] == 0 then
108   debug = function(s)
109     end
110 else
111   debug = function(s)
112     print('lua-list-hyphen DEBUG: ' .. s)
113   end
114 end
```

(End of definition for debug.)

### 9.2 Table key constants

Keys for tables containing hyphenatable/hyphenated word data.

```
115 local KEY_TYPE = 'type'
116 local KEY_WORD = 'word'
117 local KEY_LANG = 'lang'
118 local KEY_DIVISION = 'division'
119 local KEY_INDEX = 'index'
120 local KEY_CONTEXT = 'context'
121 local KEY_PAGE = 'page'
```

### 9.3 Segment type

Constants for types of segments found while scanning hlist before linebreaking.

```
122 local SEGMENT_WORD = 0
123 local SEGMENT_SPACE = 1
124 local SEGMENT_MATH = 2
```

### 9.4 Node ID and subtype constants

Define constants for the node IDs that need to be recognized.

```
125 local NODE_ID_HLIST = node.id('hlist')
126 local NODE_ID_DISC = node.id('disc')
127 local NODE_ID_GLUE = node.id('glue')
128 local NODE_ID_KERN = node.id('kern')
129 local NODE_ID_MARGIN_KERN = node.id('margin_kern')
130 local NODE_ID_GLYPH = node.id('glyph')
131 local NODE_ID_MATH = node.id('math')
```

Define constants for the kern node subtypes that have to be recognized. (There seems to be no automatic way to get the numerical value from the subtype text other than searching the `node.subtype(<node type>)` tables.)

```
132 local NODE_KERN_SUBTYPE_FONTKERN
133 local NODE_KERN_SUBTYPE_USERKERN
134 for k,v in pairs(node.subtypes('kern')) do
135   if v == 'fontkern' then
136     NODE_KERN_SUBTYPE_FONTKERN = k
137   elseif v == 'userkern' then
138     NODE_KERN_SUBTYPE_USERKERN = k
139   end
140 end
```

Define constants for the math node subtypes.

```
141 local NODE_MATH_SUBTYPE_BEGIN
142 local NODE_MATH_SUBTYPE_END
143 for k,v in pairs(node.subtypes('math')) do
144   if v == 'beginmath' then
145     NODE_MATH_SUBTYPE_BEGIN = k
146   elseif v == 'endmath' then
147     NODE_MATH_SUBTYPE_END = k
148   end
149 end
```

### 9.5 Output constants

Constants for output.

```
150 local STR_MATH = '[MATH] '
151 local STR_SPACE = ' '
152 local STR_SPACE_TWO = '  '
153 local STR_ARROW = '-> '
154 local STR_PAGE_PREFIX = 'p. '
155 local STR_PAGE_NONE = '<none>'
156 local STR_QUOTE_OPEN = '"'
157 local STR_QUOTE_CLOSE = '"'
```

## 9.6 Utility functions

`list_filter` Take a list `t` and remove from it any elements for which the function `f` does not return true. (The index `j` is always the destination index to which a ‘keep’ element is moved.)<sup>7</sup>

```
158 local function list_filter(t, f)
159   local j = 1
160   local n = #t
161
162   for i=1,n do
163     if (f(t[i])) then
164       if (i ~= j) then
165         t[j] = t[i]
166         t[i] = nil
167       end
168       j = j + 1
169     else
170       t[i] = nil
171     end
172   end
173
174 end
```

*(End of definition for list\_filter.)*

`list_uniq` Take a list `t` and remove from it adjacent elements for which the function `f` returns true. (The index `j` is always the last ‘kept’ element.)

```
175 local function list_uniq(t, f)
176   local j = 1
177   local n = #t
178
179   for i=2,n do
180     if (f(t[i],t[j])) then
181       t[i] = nil
182     else
183       j = i
184     end
185   end
186
187   list_filter(
188     t,
189     function(a) return a end
190   )
191 end
```

*(End of definition for list\_uniq.)*

## 9.7 Getting text from nodes

Getting the components of the ligatures that have Unicode code points can be problematic, at least for some fonts, so define a lookup table for these cases.

```
192 local LIGATURE_TEXT = {
193   [0xfb00] = 'ff',
```

---

<sup>7</sup>Code adapted from <https://stackoverflow.com/a/53038524>.

```

194     [0xfb01] = 'fi',
195     [0xfb02] = 'fl',
196     [0xfb03] = 'ffi',
197     [0xfb04] = 'ffl',
198 }

```

Cache to save table lookups when extracting text.

```

199 local font_characters = {}

```

Extracting text from nodes uses two functions that call each other, so the names have to be defined ahead of time.

```

200 local get_node_text
201 local get_nodelist_text

```

`get_node_text` Return the text content of a glyph node (which might be a normal glyph, a ligature, etc.).

```

202 get_node_text = function(n)
203
204     if n.id == NODE_ID_GLYPH then
205
206         local ligature_text = LIGATURE_TEXT[n.char]
207         if ligature_text ~= nil then
208             return ligature_text
209         elseif n.components then
210             return get_nodelist_text(n.components)
211         else
212             -- See [https://tug.org/pipermail/luatex/2018-March/006786.html]
213             local characters = font_characters[n.font]
214             if not characters then
215                 characters = fonts.hashes.identifiers[n.font].characters
216                 font_characters[n.font] = characters
217             end

```

In looking up the text, things can go wrong (e.g. the glyph might not be part of the version of unicode than is supported by the underlying `slnunicode` library). So use `pcall` and return an empty string if there is an error.

```

218         local ok,retval = pcall(
219             function ()
220                 return utf8.char(tonumber(characters[n.char].tounicode,16))
221             end
222         )
223         if ok then
224             return retval
225         else
226             return ''
227         end
228     end
229
230     elseif n.id == NODE_ID_DISC then
231
232         if n.replace then
233             return get_nodelist_text(n.replace)
234         else
235             return ''
236         end

```

```

237
238     else
239         return ''
240     end
241
242 end

```

*(End of definition for get\_node\_text.)*

**get\_nodelist\_text** Return the text content of the glyph nodes in the list starting at **head** up to and including the node **last**, or up to the end of the list if **last** is not specified.

```

243 get_nodelist_text = function (head,last)
244
245     local text = ''
246
247     for item in node.traverse(head) do
248
249         text = text .. get_node_text(item)
250
251         if item == last then
252             break
253         end
254     end
255
256     return text
257
258 end

```

*(End of definition for get\_nodelist\_text.)*

**is\_possible\_word\_node** Return boolean indicating if node **n** could be part of a word. Assume that **glyph**, **disc**, and **margin\_kern** nodes could be part of a word, as could a **kern** node with subtype **fontkern**.

```

259 local function is_possible_word_node(n)
260
261     return (
262         n.id == NODE_ID_GLYPH
263         or
264         n.id == NODE_ID_DISC
265         or
266         (n.id == NODE_ID_KERN and n.subtype == NODE_KERN_SUBTYPE_FONTKERN)
267         or
268         n.id == NODE_ID_MARGIN_KERN
269     )
270
271 end

```

*(End of definition for is\_possible\_word\_node.)*

**is\_possible\_space\_node** Return boolean indicating if node **n** could be part of a space. Assume that **glue** nodes could be part of a space, as could a **kern** node with subtype **userkern**.

```

272 local function is_possible_space_node(n)
273
274     return (

```

```

275     n.id == NODE_ID_GLUE
276     or
277     (n.id == NODE_ID_KERN and n.subtype == NODE_KERN_SUBTYPE_USERKERN)
278 )
279
280 end

```

*(End of definition for is\_possible\_space\_node.)*

## 9.8 String manipulation

`trim_nonlettershyphens_both` Remove characters other than letters and hyphens from both the start and end of a string.

```

281 local function trim_nonlettershyphens_both(s)
282
283     return unicode.utf8.match(s, '^[^%a-]*([^-]*)[^%a-]*$')
284
285 end

```

*(End of definition for trim\_nonlettershyphens\_both.)*

`trim_nonlettershyphens_start` Remove characters other than letters and hyphens from the start of a string.

```

286 local function trim_nonlettershyphens_start(s)
287
288     return unicode.utf8.match(s, '^[^%a-]*([^-]*)$')
289
290 end

```

*(End of definition for trim\_nonlettershyphens\_start.)*

`trim_nonlettershyphens_end` Remove characters other than letters and hyphens from the end of a string.

```

291 local function trim_nonlettershyphens_end(s)
292
293     return unicode.utf8.match(s, '^(.*)[^%a-]*$')
294
295 end

```

*(End of definition for trim\_nonlettershyphens\_end.)*

`rpadd` Return string `s` padded on the right with spaces to length `n`.

```

296 local function rpadd(s,n)
297
298     return s .. unicode.utf8.rep(STR_SPACE,n - unicode.utf8.len(s))
299
300 end

```

*(End of definition for rpadd.)*

`lpadd` Return string `s` padded on the left with spaces to length `n`.

```

301 local function lpadd(s,n)
302
303     return unicode.utf8.rep(STR_SPACE,n - unicode.utf8.len(s)) .. s
304
305 end

```

*(End of definition for lpadd.)*

## 9.9 Pre-linebreak processing

Before each line has been broken, find all potential division points and store the words in which they occur, linking each potential break point to the corresponding word.

Declare a new attribute, which will be used to store in each disc node the index of the corresponding word in the table `hlist_segment_list`.

```
306 local hyphen_attr = luatexbase.new_attribute('hyphen_attr')
```

Table to hold segments (word/space/math) in the hlist that will be broken. This table will be cleared after the post-linebreak processing.

```
307 local hlist_segment_list = {}
```

`get_first_glyph_lang` Return the lang attribute of the first glyph in the the part of the list starting n that could be part of a word. (Currently unused; see the documentation of `get_disc_lang`.)

```
308 -- local function get_first_glyph_lang(n)
309
310 --     local item = n
311 --     while item and is_possible_word_node(item) do
312 --         if item.id == NODE_ID_GLYPH then
313 --             return item.lang
314 --         end
315 --         item = item.next
316 --     end
317
318 --     return nil
319
320 -- end
```

*(End of definition for `get_first_glyph_lang`.)*

`get_disc_lang` Try to find the language ID in force at a given disc node by looking at (1) the last glyph in the word before the disc node; (2) the first glyph in the word after the disc node. Default to language ID 0.

(Looking at `replace`, `pre`, `post` is possible, but is unreliable and so disabled for the present. The author has encountered the situation where an explicit hyphen results in the hyphen characters in `replace` and `pre` having different language IDs. He has not had time to investigate how this arises from the interaction of `babel/polyglossia` and `LuaATEX`.)

```
321 local function get_disc_lang(n)
322
323 -- lang = get_first_glyph_lang(n.replace)
324 -- if lang then
325 --     print(lang)
326 --     return lang
327 -- end
328
329 -- lang = get_first_glyph_lang(n.pre)
330 -- if lang then
331 --     print(lang)
332 --     return lang
333 -- end
334
335 -- lang = get_first_glyph_lang(n.post)
```

```

336 -- if lang then
337 --   return lang
338 -- end
339
340 local item

```

Before the disc node.

```

341 item = n
342 while item and is_possible_word_node(item) do
343   if item.id == NODE_ID_GLYPH then
344     return item.lang
345   end
346   item = item.prev
347 end

```

After the disc node.

```

348 item = n
349 while item and is_possible_word_node(item) do
350   if item.id == NODE_ID_GLYPH then
351     return item.lang
352   end
353   item = item.next
354 end
355
356 return 0
357
358 end

```

*(End of definition for get\_disc\_lang.)*

`pre_linebreak` Extract segments (word/space/math) from the hlist at `hlist_head` and store appropriate data in `hlist_segment_list`. For spaces and math, this is just the existence of a segment. For a word, store its text and its language ID (as determined by `get_disc_lang`). Also, for each disc node, assign the index of the word in `hlist_segment_list` to its `hyphen_attr` attribute (declared above).

```

359 local function pre_linebreak(hlist_head,groupcode)
360
361   local word_start_node = nil
362   local segment_count = 0
363   local lang = nil
364
365   debug('Pre-linebreak processing start')
366

```

Per § 8.2 of the LuaTeX manual, there can be cases where `.prev` is invalid, so run `node.slide()` to set them correctly.

```

367   node.slide(hlist_head)
368
369   local item = hlist_head
370   while item do

```

If `item` is a math node (which must have subtype `beginmath`, unless something has changed the node list), skip the math and add [MATH] to `hlist_segment_list`.

```

371     if item.id == NODE_ID_MATH then
372       assert(item.subtype == NODE_MATH_SUBTYPE_BEGIN)

```



```

373     while not (
374         item.id == NODE_ID_MATH and item.subtype == NODE_MATH_SUBTYPE_END
375     ) do
376         item = item.next
377     end
378     item = item.next
379
380     segment_count = segment_count + 1
381     hlist_segment_list[segment_count] = {
382         [KEY_TYPE] = SEGMENT_MATH
383     }
384
385     goto continue
386 end

```

If `item` is a possible word node, read the whole word, setting the `hyphen_attr` of any disc nodes to `segment_count`, and adding the word to `hlist_segment_list`.

```

387     if is_possible_word_node(item) then
388         word_start_node = item
389         segment_count = segment_count + 1
390         while item and is_possible_word_node(item) do
391

```

When the first disc node is found, find the language of the word.

```

392             if item.id == NODE_ID_DISC then
393                 if not lang then
394                     lang = get_disc_lang(item)
395                 end
396                 node.set_attribute(item, hyphen_attr, segment_count)
397             end
398
399             item = item.next
400         end

```

`item` should be a node, because even after the last word node, the `hlist` will contain something. But just in case, check and find the last node using `node.tail` if necessary. This latter case should be very rare, so it is more efficient to recalculate here if necessary rather than having an extra assignment to store the previous node in the while loop.

```

401         local word_end_node
402         if item then
403             word_end_node = item.prev
404         else
405             word_end_node = node.tail(word_start_node)
406         end
407
408         local word = get_nodelist_text(word_start_node, word_end_node)
409         hlist_segment_list[segment_count] = {
410             [KEY_TYPE] = SEGMENT_WORD,
411             [KEY_WORD] = word,
412             [KEY_LANG] = lang,
413         }
414
415         word_start_node = nil
416         lang = nil
417

```

```

418     goto continue
419 end

```

If `item` is a node that could be part of a space, add a space to the segment list.

```

420     if is_possible_space_node(item) then
421         segment_count = segment_count + 1
422
423         while item and is_possible_space_node(item) do
424             item = item.next
425         end
426
427         hlist_segment_list[segment_count] = {
428             [KEY_TYPE] = SEGMENT_SPACE
429         }
430
431         goto continue
432     end

```

If `item` is anything else, just move on.

```

433     item = item.next
434
435     ::continue::
436 end
437
438 debug('Pre-linebreak processing finish')
439
440 return true
441 end

```

*(End of definition for pre\_linebreak.)*

## 9.10 Post-linebreak processing

After linebreaking, look for a discretionary node at the end of each line, which indicates that a word has been divided between the end of that line and the start of the next. Extract the two word-pieces from the lines and store them, together with the undivided word and its context in the appropriate language table. Also insert a whatsit to that will set the page number when the hyphenation is written out.

`get_used_disc` If at the tail of the hlist at `hlist_head` (which will be a line) there is a disc node not followed by a glyph node, return that disc node. Otherwise return `nil`. Only search up to at most `USED_DISC_SEARCH_LIMIT` nodes from the tail.

```

442 local USED_DISC_SEARCH_LIMIT = 20
443
444 local function get_used_disc(hlist_head)
445
446     debug(' Looking for disc node at end of line')
447
448     local item = node.tail(hlist_head)
449
450     local count = 0
451     while item and item.id ~= NODE_ID_GLYPH and count < USED_DISC_SEARCH_LIMIT do
452         if item.id == NODE_ID_DISC then
453             return item

```

```

454     end
455     item = item.prev
456     count = count + 1
457 end
458
459 return nil
460
461 end

```

*(End of definition for get\_used\_disc.)*

`get_disc_word_start` Return the node starting the word that includes a given disc node `n`, or `nil` if there is no such node.

```

462 local function get_disc_word_start(hlist_head,n)
463
464     local item = n
465
466     while item do
467         local prev = item.prev
468
469         if not (prev and is_possible_word_node(prev)) then
470             return item
471         end
472
473         item = prev
474     end
475
476     return nil
477 end

```

*(End of definition for get\_disc\_word\_start.)*

`get_next_hlist` Return the next hlist in the list containing the given node `n`, or `nil` if there is no such hlist node.

```

478 local function get_next_hlist(n)
479
480     local item = n.next
481
482     while item do
483         if item.id == NODE_ID_HLIST then
484             return item
485         end
486         item = item.next
487     end
488
489     return nil
490
491 end

```

*(End of definition for get\_next\_hlist.)*

`get_line_first_word` Return the first word in the hlist at `hlist_head`, or `nil` if there is no such word.

```

492 local function get_line_first_word(hlist_head)

```

`word_start_node` is either nil or the (glyph) node that starts the word.

```

493   local word_start_node = nil
494
495   for item in node.traverse(hlist_head) do
496
497     if item.id == NODE_ID_GLYPH then
498       if not word_start_node then
499         word_start_node = item
500       end
501     end
502
503     if not is_possible_word_node(item) then
504       if word_start_node then
505         return get_nodelist_text(word_start_node,item.prev)
506       end
507     end
508
509   end

```

It is possible that the word ends at the end of the hlist, so check if a word has been started.

```

510   if word_start_node then
511     return get_nodelist_text(word_start_node,node.tail(hlist_head))
512   else
513     return nil
514   end
515 end

```

*(End of definition for `get_line_first_word`.)*

`get_context` Return a string assembled from the part of `hlist_segment_list` before or after `index` according to `incr` (which must be  $\pm 1$ ) up a maximum of `target_word_count` words.

```

516 local function get_context(index,incr,target_word_count)
517
518   local result = ''
519   local word_count = 0
520
521   local i = index + incr
522   while (
523     i > 0 and i <= #hlist_segment_list and word_count < target_word_count
524   ) do
525     local t = hlist_segment_list[i]
526
527     local item
528
529     if t[KEY_TYPE] == SEGMENT_WORD then
530       item = t[KEY_WORD]
531       word_count = word_count + 1
532     elseif t[KEY_TYPE] == SEGMENT_SPACE then
533       item = STR_SPACE
534     elseif t[KEY_TYPE] == SEGMENT_MATH then
535       item = STR_MATH
536     end
537

```

```

538     if incr > 0 then
539         result = result .. item
540     else
541         result = item .. result
542     end
543
544     i = i + incr
545 end
546
547 return result
548
549 end

```

(End of definition for `get_context`.)

Count and list for hyphenated words. Each entry in the list will be a table containing the original word, the hyphenation, the language, the index of the table in the list (which is needed later for stable sorting and sorting into the original order), and the context.

```

550 local hyphenation_list = {}
551 local hyphenation_count = 0

```

`check_line_hyphenation` Check whether there is a hyphenated word at the end of the given `hlist`; if so, save the word to `hyphenation_list`.

```

552 local function check_line_hyphenation(hlist)

```

First, is there a disc node not followed by a glyph node at the end of the list?

```

553     local last_disc = get_used_disc(hlist.head)
554     if not last_disc then
555         debug(' No disc node found at end of line')
556         return
557     end
558     debug(' Disc node found at end of line')

```

Get the undivided word and its language from `hlist_segment_list`.

```

559     local hyphenation_index = node.has_attribute(last_disc, hyphen_attr)
560     local t = hlist_segment_list[hyphenation_index]
561     assert(t)
562     assert(t[KEY_TYPE] == SEGMENT_WORD)
563     local word = t[KEY_WORD]
564     local lang = t[KEY_LANG]

```

`word` might be something other than a genuine word, such as an ISBN (with hyphen separators). So only proceed if it contains at least one letter.

```

565     if not unicode.utf8.match(word, '%a') then
566         debug(' Divided "word" contains no letters')
567         return
568     end

```

There should always be a next line, since there is a disc node at the end of `hlist`, but check anyway.

```

569     local next_line = get_next_hlist(hlist)
570
571     if not next_line then
572         debug(' No following line found (which should not happen)')
573         return
574     end

```

For the pre-linebreak part of the word, get the word that ends the line, and trim any leading non-letters. This could leave an empty word; for example, if *n*-dimensional is broken at the hyphen, the word ending the line is just the hyphen. If an empty word is left, just use the non-trimmed result.

```

575 local pre = get_nodelist_text(get_disc_word_start(hlist.head,last_disc))
576 local pre_temp = trim_nonlettershyphens_start(pre)
577 if pre_temp ~= '' then
578   pre = pre_temp
579 end

```

For the post-linebreak part, just get the word at the start of the next line, and trim and trailing non-letters.

```

580 local post = trim_nonlettershyphens_end(get_line_first_word(next_line.head))

```

Compute the context and then trim any unwanted symbols from the word itself.

```

581 local context =
582   get_context(
583     hyphenation_index,-1,tex.count['l__lualisthyphen_context_before_int']
584   )
585   .. word ..
586   get_context(
587     hyphenation_index,1,tex.count['l__lualisthyphen_context_after_int']
588   )
589
590 word = trim_nonlettershyphens_both(word)
591
592 debug(
593   ' Hyphenated word found: "' .. word .. '" -> "' .. pre .. '<>' .. post .. '"'
594 )

```

Store everything (except the page number on which the hyphenated word appears, which is not yet known) in the hyphenation list.

```

595 hyphenation_count = hyphenation_count + 1
596 hyphenation_list[hyphenation_count] = {
597   [KEY_LANG] = lang,
598   [KEY_WORD] = word,
599   [KEY_DIVISION] = pre .. post,
600   [KEY_INDEX] = hyphenation_count,
601   [KEY_CONTEXT] = context,
602 }

```

Add a whatsit to record the page number when the page with the hyphenation is shipped out. This information also serves to distinguish hyphenations that are written to the page from those that occur in (e.g.) boxes that are discarded without being written to the page.

```

603 late_lua_n = node.new('whatsit','late_lua')
604 late_lua_n.data =
605   'lualisthyphen.set_hyphenation_page(' .. hyphenation_count .. ',tex.count["c@page"])'
606
607 node.insert_after(hlist.head,last_disc,late_lua_n)
608
609 end

```

*(End of definition for check\_line\_hyphenation.)*

`set_hyphenation_page` Set the page on which the hyphenation with the given index appears.

```
610 local function set_hyphenation_page(index,page)
611
612     hyphenation_list[index][KEY_PAGE] = page
613
614 end
```

*(End of definition for set\_hyphenation\_page.)*

`post_linebreak` For every line in the vlist at `vlist_head`, check whether there is a hyphenated word at the end.

```
615 local function post_linebreak(vlist_head,groupcode)
616
617     debug('Post-linebreak processing start')
618
619     for item in node.traverse(vlist_head) do
620         if item.id == NODE_ID_HLIST then
621             node.slide(item.head)
622         end
623     end
```

Now actually look for hyphenations.

```
623     local line_no = 0
624
625     for item in node.traverse(vlist_head) do
626
627         if item.id == NODE_ID_HLIST then
628             line_no = line_no + 1
629             debug(' Line no.' .. line_no)
630             check_line_hyphenation(item)
631         end
632
633     end
634
635     hlist_segment_list = {}
636
637     debug('Post-linebreak processing end')
638
639     return true
640
641 end
```

*(End of definition for post\_linebreak.)*

## 9.11 Callbacks

Add `pre_linebreak` and `post_linebreak` to the relevant callbacks.

```
642 local LUA_LIST_HYPHEN_PRE_LINEBREAK = 'LUA_LIST_HYPHEN_PRE_LINEBREAK'
643 luatexbase.add_to_callback(
644     'pre_linebreak_filter',
645     pre_linebreak,
646     LUA_LIST_HYPHEN_PRE_LINEBREAK
```

```

647 )
648
649 local LUA_LIST_HYPHEN_POST_LINEBREAK = 'LUA_LIST_HYPHEN_POST_LINEBREAK'
650 luatexbase.add_to_callback(
651   'post_linebreak_filter',
652   post_linebreak,
653   LUA_LIST_HYPHEN_POST_LINEBREAK
654 )

```

## 9.12 Language settings

Table mapping language IDs to textual names.

```

655 local language_table = {}

```

Populating `language_table` is done differently for `babel` and `polyglossia`. If `babel` is in use, the  $\text{\LaTeX}$  frontend iterates through `\bbl@languages` and calls `babel_save_language_name`. If `polyglossia` is in use, `language_table` is populated by `polyglossia_get_language_names`, which is called just before the hyphenation lists are written.

`babel_save_language_name` Store the association of a language ID to `babel`'s textual name, if no name has been assigned to that ID already.

```

656 local function babel_save_language_name(lang_id,name)
657
658   if not language_table[lang_id] then
659     language_table[lang_id] = name
660   end
661
662 end

```

*(End of definition for babel\_save\_language\_name.)*

`polyglossia_get_language_names` If `polyglossia` has been loaded, use it to build the table mapping language IDs to textual names.

```

663 local function polyglossia_get_language_names()
664
665   if not polyglossia then
666     return
667   end
668
669   for name,language in pairs(polyglossia.newloader_loaded_languages) do
670     language_table[lang.id(language)] = name
671   end
672
673 end

```

*(End of definition for polyglossia\_get\_language\_names.)*

## 9.13 Processing hyphenation lists

Before writing out hyphenation lists, remove duplicates and/or perform sorting, in accordance with the set options.



### 9.13.1 Comparisons and equality checks

`equal_hyphenation_case_sensitive` Equality check for deduplicating the list of hyphenations case-sensitively.

```
674 local function equal_hyphenation_case_sensitive(a,b)
675   return (
676     a[KEY_WORD] == b[KEY_WORD]
677     and
678     a[KEY_DIVISION] == b[KEY_DIVISION]
679   )
680 end
```

*(End of definition for equal\_hyphenation\_case\_sensitive.)*

`equal_hyphenation_case_insensitive` Equality check for deduplicating the list of hyphenations case-insensitively.

```
681 local function equal_hyphenation_case_insensitive(a,b)
682   return (
683     unicode.utf8.lower(a[KEY_WORD]) == unicode.utf8.lower(b[KEY_WORD])
684     and
685     unicode.utf8.lower(a[KEY_DIVISION]) == unicode.utf8.lower(b[KEY_DIVISION])
686   )
687 end
```

*(End of definition for equal\_hyphenation\_case\_insensitive.)*

`lessthan_hyphenation_case_sensitive` Comparison for sorting the list of hyphenations case-sensitively.

The comparison of index keys ensures that the sorting is stable.

```
688 local function lessthan_hyphenation_case_sensitive(a,b)
689   return (
690     a[KEY_WORD] < b[KEY_WORD]
691     or
692     (
693       a[KEY_WORD] == b[KEY_WORD]
694       and
695       a[KEY_DIVISION] < b[KEY_DIVISION]
696     )
697     or
698     (
699       a[KEY_WORD] == b[KEY_WORD]
700       and
701       a[KEY_DIVISION] == b[KEY_DIVISION]
702       and
703       a[KEY_INDEX] < b[KEY_INDEX]
704     )
705   )
706 end
```

*(End of definition for lessthan\_hyphenation\_case\_sensitive.)*

`lessthan_hyphenation_case_insensitive` Comparison for sorting the list of hyphenations case-insensitively.

The comparison of index keys ensures that the sorting is stable.

```
707 local function lessthan_hyphenation_case_insensitive(a,b)
708   return (
709     unicode.utf8.lower(a[KEY_WORD]) < unicode.utf8.lower(b[KEY_WORD])
710     or
711     (
```

```

712     unicode.utf8.lower(a[KEY_WORD]) == unicode.utf8.lower(b[KEY_WORD])
713     and
714     unicode.utf8.lower(a[KEY_DIVISION]) < unicode.utf8.lower(b[KEY_DIVISION])
715 )
716 or
717 (
718     unicode.utf8.lower(a[KEY_WORD]) == unicode.utf8.lower(b[KEY_WORD])
719     and
720     unicode.utf8.lower(a[KEY_DIVISION]) < unicode.utf8.lower(b[KEY_DIVISION])
721     and
722     a[KEY_INDEX] < b[KEY_INDEX]
723 )
724 )
725 end

```

*(End of definition for lessthan\_hyphenation\_case\_insensitive.)*

### 9.13.2 Sorting

`sort_hyphenation_list_none` Sort `hyphenation_list` into its original order of appearance.

```

726 local function sort_hyphenation_list_none(hyphenation_list)
727     table.sort(
728         hyphenation_list,
729         function(a,b)
730             return a[KEY_INDEX] < b[KEY_INDEX]
731         end
732     )
733 end

```

*(End of definition for sort\_hyphenation\_list\_none.)*

`sort_hyphenation_list_case` Sort `hyphenation_list` case-sensitively.

```

734 local function sort_hyphenation_list_case(hyphenation_list)
735     table.sort(
736         hyphenation_list,
737         lessthan_hyphenation_case_sensitive
738     )
739 end

```

*(End of definition for sort\_hyphenation\_list\_case.)*

`sort_hyphenation_list_nocase` Sort `hyphenation_list` case-insensitively.

```

740 local function sort_hyphenation_list_nocase(hyphenation_list)
741     table.sort(
742         hyphenation_list,
743         lessthan_hyphenation_case_insensitive
744     )
745 end

```

*(End of definition for sort\_hyphenation\_list\_nocase.)*

`process_lang_hyphenation_list_sort` Select the appropriate function for sorting.

```

746 local sort_hyphenation_list
747 if tex.count['l__lualisthyphen_sort_int'] == 1 then
748     sort_hyphenation_list = sort_hyphenation_list_case

```

```

749 elseif tex.count['l__lualisthyphen_sort_int'] == 2 then
750   sort_hyphenation_list = sort_hyphenation_list_nocase
751 else
752   sort_hyphenation_list = sort_hyphenation_list_none
753 end

```

*(End of definition for process\_lang\_hyphenation\_list\_sort.)*

### 9.13.3 Deduplication

deduplicate\_hyphenation\_list\_none Dummy function; does not deduplicate hyphenation\_list.

```

754 local function deduplicate_hyphenation_list_none(hyphenation_list)
755 end

```

*(End of definition for deduplicate\_hyphenation\_list\_none.)*

deduplicate\_hyphenation\_list\_case Remove duplicates from hyphenation\_list case-sensitively.

```

756 local function deduplicate_hyphenation_list_case(hyphenation_list)
757   table.sort(
758     hyphenation_list,
759     lessthan_hyphenation_case_sensitive
760   )
761   list_uniq(
762     hyphenation_list,
763     equal_hyphenation_case_sensitive
764   )
765 end

```

*(End of definition for deduplicate\_hyphenation\_list\_case.)*

deduplicate\_hyphenation\_list\_nocase Remove duplicates from hyphenation\_list case-insensitively.

```

766 local function deduplicate_hyphenation_list_nocase(hyphenation_list)
767   table.sort(
768     hyphenation_list,
769     lessthan_hyphenation_case_insensitive
770   )
771   list_uniq(
772     hyphenation_list,
773     equal_hyphenation_case_insensitive
774   )
775 end

```

*(End of definition for deduplicate\_hyphenation\_list\_nocase.)*

deduplicate\_hyphenation\_list Select the appropriate function for whether duplicates would be removed.

```

776 local deduplicate_hyphenation_list
777 if tex.count['l__lualisthyphen_unique_int'] == 1 then
778   deduplicate_hyphenation_list = deduplicate_hyphenation_list_case
779 elseif tex.count['l__lualisthyphen_unique_int'] == 2 then
780   deduplicate_hyphenation_list = deduplicate_hyphenation_list_nocase
781 else
782   deduplicate_hyphenation_list = deduplicate_hyphenation_list_none
783 end

```

*(End of definition for deduplicate\_hyphenation\_list.)*

### 9.13.4 Combined processing

```
process_lang_hyphenation_list Remove duplicates and sort hyphenation_list.
784 local function process_lang_hyphenation_list(hyphenation_list)
785     deduplicate_hyphenation_list(hyphenation_list)
786     sort_hyphenation_list(hyphenation_list)
787 end
```

*(End of definition for process\_lang\_hyphenation\_list.)*

## 9.14 Writing

```
write_lang_hyphenation_list_standard Write out just the hyphenated words in hyphenation_list to file handle f.
788 local function write_lang_hyphenation_list_standard(f,hyphenation_list,widths)
789
790     for i,v in ipairs(hyphenation_list) do
791
792         if v then
793             f:write(v[KEY_DIVISION] .. '\n')
794         end
795
796     end
797
798 end
```

*(End of definition for write\_lang\_hyphenation\_list\_standard.)*

```
write_lang_hyphenation_list_verbose Write out all hyphenation information in hyphenation_list to file handle f, in columns
as specified in widths.
```

```
799 local function write_lang_hyphenation_list_verbose(f,hyphenation_list,widths)
800
801     local cols_word = widths[KEY_WORD]
802     local cols_division = widths[KEY_DIVISION]
803     local cols_page = widths[KEY_PAGE]
804
805     for i,v in ipairs(hyphenation_list) do
806
807         if v then
808             local page = v[KEY_PAGE]
809             if page then
810                 page = STR_PAGE_PREFIX .. page
811             else
812                 page = STR_PAGE_NONE
813             end
814
815             f:write(
816                 rpad(v[KEY_WORD],cols_word)
817                 .. STR_ARROW
818                 .. rpad(v[KEY_DIVISION],cols_division)
819                 .. STR_SPACE_TWO
820                 .. lpad(page,cols_page)
```

It is possible for KEY\_PAGE not to have been set, for instance if the hyphenation occurred in a box that was never output.

```

821         .. STR_SPACE
822         .. STR_QUOTE_OPEN
823         .. v[KEY_CONTEXT]
824         .. STR_QUOTE_CLOSE
825         .. '\n'
826     )
827 end
828
829 end
830
831 end

```

*(End of definition for write\_lang\_hyphenation\_list\_verbose.)*

`write_lang_hyphenation_list` Set `write_lang_hyphenation_list` to be either `write_lang_hyphenation_list_standard` or `write_lang_hyphenation_list_verbose`, depending on the package options.

```

832 local write_lang_hyphenation_list
833 if tex.count['l__lualisthyphen_verbose_int'] == 0 then
834     write_lang_hyphenation_list = write_lang_hyphenation_list_standard
835 else
836     write_lang_hyphenation_list = write_lang_hyphenation_list_verbose
837 end

```

*(End of definition for write\_lang\_hyphenation\_list.)*

Compute a settings description to insert into file headers.

```

838 local settings_desc
839 if tex.count['l__lualisthyphen_verbose_int'] == 0 then
840     settings_desc = 'verbose=false'
841 else
842     settings_desc = 'verbose=true'
843     .. ',context-before=' .. tex.count['l__lualisthyphen_context_before_int']
844     .. ',context-after=' .. tex.count['l__lualisthyphen_context_after_int']
845 end
846 if tex.count['l__lualisthyphen_include_non_output_int'] == 0 then
847     settings_desc = settings_desc .. ',include-non-output=false'
848 else
849     settings_desc = settings_desc .. ',include-non-output=true'
850 end
851 local NONE_CASE_NOCASE = {
852     [0] = 'none',
853     [1] = 'case',
854     [2] = 'nocase'
855 }
856 settings_desc = settings_desc
857 .. ',sort=' .. NONE_CASE_NOCASE[tex.count['l__lualisthyphen_sort_int']]
858 .. ',unique=' .. NONE_CASE_NOCASE[tex.count['l__lualisthyphen_unique_int']]

```

`get_hyphenation_file_path` Get the file to which the list of hyphenated words will be written, based on the given prefix, extension, lang\_name, and taking into account any specified output directory for LuaTeX, and with a file header.

```

859 local function get_hyphenation_file_path(prefix,extension,lang_name)
860
861     local hyphenation_file_path = prefix .. tostring(lang_name) .. extension

```

```

862
863     if not status.output_directory then
864         return hyphenation_file_path
865     end
866
867     if string.sub(status.output_directory,-1,-1) == '/' then
868         hyphenation_file_path = status.output_directory
869         .. hyphenation_file_path
870     else
871         hyphenation_file_path = status.output_directory
872         .. '/' .. hyphenation_file_path
873     end
874
875     return hyphenation_file_path
876
877 end

```

*(End of definition for get\_hyphenation\_file\_path.)*

process\_write\_lang\_hyphenation\_list

Process and write out the `hyphenation_list` (which will be for the language with the numerical `lang_id`) to a file with the given `prefix` and `extension`, using `widths` for the ‘columns’ in verbose mode.

```

878 local function process_write_lang_hyphenation_list(
879     prefix,extension,lang_id,hyphenation_list,widths
880 )
881
882     process_lang_hyphenation_list(hyphenation_list)
883
884     local lang_name = language_table[lang_id]
885     local lang_desc
886     if not lang_name then
887         lang_name = lang_id
888         lang_desc = 'language with ID ' .. lang_id
889     else
890         lang_desc = 'language "' .. lang_name .. '" (ID ' .. lang_id .. ')'
891     end
892
893     local f = io.open(get_hyphenation_file_path(prefix,extension,lang_name),'w')
894
895     f:write('% Chosen hyphenations for ' .. lang_desc .. '\n')
896     f:write('% Generated by lua-list-hyphen (' .. settings_desc .. ')\n')
897
898     write_lang_hyphenation_list(f,hyphenation_list,widths)
899     f:close()
900
901 end

```

*(End of definition for process\_write\_lang\_hyphenation\_list.)*

process\_write\_hyphenation\_lists

Sort `hyphenation_list` into per-language lists and write them out to separate files.

```

902 local function process_write_hyphenation_lists(prefix,extension)
903
904     local lang_hyphenation_table = {}
905     local lang_widths_table = {}

```

Iterate through all the stored hyphenations. Sort them into per-language lists (creating the list the first time each language is encountered) and also storing the maximum width of values, for output alignment.

```

906   for _,h in pairs(hyphenation_list) do
907
908       if h[KEY_PAGE] or tex.count['l__lualisthyphen_include_non_output_int'] == 1 then
909
910           local lang = h[KEY_LANG]
911
912           local t = lang_hyphenation_table[lang]
913           if not t then
914               lang_hyphenation_table[lang] = {}
915               t = lang_hyphenation_table[lang]
916           end
917
918           local widths = lang_widths_table[lang]
919           if not widths then
920               lang_widths_table[lang] = {
921                   [KEY_WORD] = 0,
922                   [KEY_DIVISION] = 0,
923                   [KEY_PAGE] = 0
924               }
925               widths = lang_widths_table[lang]
926           end
927
928           widths[KEY_WORD] = math.max(
929               widths[KEY_WORD],
930               unicode.utf8.len(h[KEY_WORD])
931           )
932           widths[KEY_DIVISION] = math.max(
933               widths[KEY_DIVISION],
934               unicode.utf8.len(h[KEY_DIVISION])
935           )
936           widths[KEY_PAGE] = math.max(
937               widths[KEY_PAGE],
938               unicode.utf8.len(tostring(h[KEY_PAGE]))
939           )
940
941           table.insert(t,h)
942
943       end
944
945   end

```

Adjust the maximum width for the page output, since there is a prefix and a ‘no page’ indicator to consider.

```

946   for _,widths in pairs(lang_widths_table) do
947       widths[KEY_PAGE] = math.max(
948           widths[KEY_PAGE] + unicode.utf8.len(STR_PAGE_PREFIX),
949           unicode.utf8.len(STR_PAGE_NONE)
950       )
951   end

```

If polyglossia is in use, populate language\_table.

```

952   polyglossia_get_language_names()

```

For each language, process and write out its hyphenations to a file.

```
953   for k,v in pairs(lang_hyphenation_table) do
954     process_write_lang_hyphenation_list(prefix,extension,k,v,lang_widths_table[k])
955   end
956
957 end
```

*(End of definition for process\_write\_hyphenation\_lists.)*

## 9.15 Export public functions

Finally, make available the functions that will be called from the L<sup>A</sup>T<sub>E</sub>X frontend using `\lua_now:n`.

```
958 return {
959   process_write_hyphenation_lists = process_write_hyphenation_lists,
960   set_hyphenation_page = set_hyphenation_page,
961   babel_save_language_name = babel_save_language_name,
962 }
963 \lua
```



# Index

The italic numbers denote the pages where the corresponding entry is described, numbers underlined point to the definition, all others indicate the places where it is used.

<b>B</b>	
babel commands:	
babel_save_language_name .....	<u>656</u>
bool commands:	
\bool_if:NTF .....	57, 60
<b>C</b>	
check commands:	
check_line_hyphenation .....	<u>552</u>
context (option) .....	4
context-after (option) .....	4
context-before (option) .....	4
cs commands:	
\cs_generate_variant:Nn .....	100
\cs_if_exist:NTF .....	70
\cs_new:Npn .....	68, 80, 91
\cs_set_eq:NN .....	73
<b>D</b>	
debug (option) .....	4
debug .....	<u>105</u>
deduplicate commands:	
deduplicate_hyphenation_list ...	<u>776</u>
deduplicate_hyphenation_list_-	
case .....	<u>756</u>
deduplicate_hyphenation_list_-	
nocase .....	<u>766</u>
deduplicate_hyphenation_list_-	
none .....	<u>754</u>
<b>E</b>	
equal commands:	
equal_hyphenation_case_insensitive	
.....	<u>681</u>
equal_hyphenation_case_sensitive	<u>674</u>
extension (option) .....	4
<b>F</b>	
\foreignlanguage .....	5
<b>G</b>	
get commands:	
get_context .....	<u>516</u>
get_disc_lang .....	<u>321</u>
get_disc_word_start .....	<u>462</u>
get_first_glyph_lang .....	<u>308</u>
get_hyphenation_file_path .....	<u>859</u>
get_line_first_word .....	<u>492</u>
get_next_hlist .....	<u>478</u>
get_node_text .....	<u>202</u>
get_nodelist_text .....	<u>243</u>
get_used_disc .....	<u>442</u>
group commands:	
\group_begin: .....	72
\group_end: .....	77
<b>H</b>	
hook commands:	
\hook_gput_code:nnn .....	65, 86
<b>I</b>	
include-non-output (option) .....	4
int commands:	
\int_new:N .....	28, 34, 51, 56, 59
\int_set:Nn .....	31, 37
\int_set_eq:NN .....	53, 58, 61
\c_one_int .....	53, 58, 61
is commands:	
is_possible_space_node .....	<u>272</u>
is_possible_word_node .....	<u>259</u>
<b>J</b>	
\jobname .....	3, 4
<b>K</b>	
keys commands:	
\l_keys_choice_int .....	31, 37
\keys_define:nn .....	
.....	12, 15, 29, 35, 40, 43, 47, 52
\keys_set:nn .....	21
<b>L</b>	
lessthan commands:	
lessthan_hyphenation_case_-	
insensitive .....	<u>707</u>
lessthan_hyphenation_case_-	
sensitive .....	<u>688</u>
list commands:	
list_filter .....	<u>158</u>
list_uniq .....	<u>175</u>
lpad .....	<u>301</u>
lua commands:	
\lua_now:n .....	32, 62, 82, 93
\luaescapestring .....	95, 96
lualisthyphen internal commands:	
__lualisthyphen_babel_save_-	
language_names: .....	66, <u>68</u> , 68
__lualisthyphen_babel_save_-	
language_names_elt:nnnn	75, <u>80</u> , 80

