

# Development of LuaT<sub>E</sub>X-ja package

Hironori Kitagawa 北川 弘典

LuaT<sub>E</sub>X-ja project team [h\\_kitagawa2001@yahoo.co.jp](mailto:h_kitagawa2001@yahoo.co.jp)

KEYWORDS T<sub>E</sub>X, pT<sub>E</sub>X, LuaT<sub>E</sub>X, LuaT<sub>E</sub>X-ja, Japanese

ABSTRACT LuaT<sub>E</sub>X-ja package is a macro package for typesetting Japanese documents under LuaT<sub>E</sub>X. The package has more flexibility of typesetting than pT<sub>E</sub>X, which is widely used Japanese extension of T<sub>E</sub>X, and has corrected some unwanted features of pT<sub>E</sub>X. In this paper, we describe specifications, the current status and some internal processing methods of LuaT<sub>E</sub>X-ja.

## 1 Introduction

### 1.1 History

To typeset Japanese documents with T<sub>E</sub>X, ASCII pT<sub>E</sub>X [2] has been widely used in Japan. There are other methods—for example, using Omega and OTP [4], or with the CJK package—to do so, however, these alternative methods did not become majority. The author thinks that this is because pT<sub>E</sub>X enables us to produce high-quality documents (e.g., supporting vertical typesetting), and the appearance of pT<sub>E</sub>X is earlier than that of alternatives described above.

However, pT<sub>E</sub>X has been left behind from the extensions of T<sub>E</sub>X such as  $\epsilon$ -T<sub>E</sub>X and pdfT<sub>E</sub>X, and the diffusion of UTF-8 encoding. In recent years, the situation has become better, by development of ptexenc [19] by Nobuyuki Tsuchimura (土村展之),  $\epsilon$ -pT<sub>E</sub>X [7] by the author, and upT<sub>E</sub>X [18] by Takuji Tanaka (田中琢爾). However, continuing this approach, namely, to develop an engine extension localized for Japanese, is not wise. This approach needs lots of work for *each* engine. In addition, if we use LuaT<sub>E</sub>X, the necessity of an engine extension is getting smaller because LuaT<sub>E</sub>X has an ability to hook T<sub>E</sub>X's internal process by using Lua callbacks.

Before our LuaT<sub>E</sub>X-ja project, there were several experimental attempts to typeset Japanese documents with LuaT<sub>E</sub>X. Here we cite three examples:

- `luaums.sty` [8] developed by the author. This experimental package is for creating a certain Japanese-based presentation with LuaT<sub>E</sub>X.
- the `luajalayout` package [12], formerly known as the `jafontspec` package, by Kazuki Maeda (前田一貴). This package is based on L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub> and `fontspec` package.
- the `luajp-test` package [11], a test package made by Atsuhiro Kohda (香田温人), based on articles on the web page [5].

However, these packages are based on L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>, and do not have much ability to control the typesetting rule. And it is inefficient that more than one person separately develop similar packages. Development of the LuaT<sub>E</sub>X-ja package is started initially by the author and Kazuki Maeda, because of these situations.

## 1.2 Development policy of LuaT<sub>E</sub>X-ja

The first aim of LuaT<sub>E</sub>X-ja project was to implement features (from the ‘primitive’ level) of pT<sub>E</sub>X as macros under LuaT<sub>E</sub>X, therefore LuaT<sub>E</sub>X-ja is much affected by pT<sub>E</sub>X. However, as development proceeded, some technical/conceptual difficulties arose. Hence we changed the aim of the project as follows:

- *LuaT<sub>E</sub>X-ja offers at least the same flexibility of typesetting that pT<sub>E</sub>X has.*

We are not satisfied with the ability of producing (PDF) outputs conformed to JIS X 4051 [6], the Japanese Industrial Standard for typesetting, or to a technical note [20] by W3C; if one wants to produce very incoherent outputs for some reason, it should be possible. In this point, previous attempts of Japanese typesetting with LuaT<sub>E</sub>X which we cited in the previous subsection are inadequate. pT<sub>E</sub>X has some flexibility of typesetting, by changing internal parameters such as `\kanjiskip` or `\prebreakpenalty`, and by using custom JFM (Japanese TFM). Therefore we decided to include these functionality to LuaT<sub>E</sub>X-ja.

- *LuaT<sub>E</sub>X-ja isn't mere re-implementation or porting of pT<sub>E</sub>X; some (technically and/or conceptually) inconvenient features of pT<sub>E</sub>X are modified.*

We describe this point in more detail at the next section.

## 1.3 Overview of the processes

We describe an outline of LuaT<sub>E</sub>X-ja's process in order.

- In the `process_input_buffer` callback: treatment of line-break after a Japanese character (in Subsection 2.2).

- In the `hyphenate` callback: font replacement.

LuaT<sub>E</sub>X-ja looks into for each *glyph\_node*  $p$  in the horizontal list. If the character represented by  $p$  is considered as a Japanese character, the font used at  $p$  is replaced by the value of `\ltj@curjfont`, an attribute for ‘the current Japanese font’ at  $p$ .

Furthermore, the subtype of  $p$  is subtracted by 1 to suppress hyphenation around  $p$  by LuaT<sub>E</sub>X, because later processes of LuaT<sub>E</sub>X-ja take care of all things about Japanese characters.

- In `pre_linebreak_filter` and `hpack_filter` callbacks:

1. LuaT<sub>E</sub>X-ja has its own stack system, and the current horizontal list is traversed in this stage to determine what the level of LuaT<sub>E</sub>X-ja's internal stack at the end of the list is. We will discuss it in Subsection 5.2.
2. In this stage, LuaT<sub>E</sub>X-ja inserts glues/kerns for Japanese typesetting in the list. This is the core routine of LuaT<sub>E</sub>X-ja. We will discuss it in Subsections 2.4 and 2.5.

3. To make a match between a metric and a real font, sometimes adjustment of the position of (Japanese) glyphs are performed. We will discuss it in Subsection 5.3.
- In the `mlist_to_hlist` callback: treatment of Japanese characters in math formulas. This stage is similar to adjustment of the position of glyphs (see above), so we omit to describe this stage from this paper.

In this paper, a *alphabetic character* means a non-Japanese character. Similarly, we use the word an *alphabetic font* as the counterpart of a Japanese font.

## 1.4 Contents of this paper

Here we describe the contents of the rest of this paper briefly. In Section 2, we describe major differences between pTeX and LuaTeX-ja. The next section, Section 3, is concentrated on a problem how we distinguish between Japanese characters and alphabetic characters. In Section 4, we show current development status of the package. Finally, in Section 5, we describe some internal routines of LuaTeX-ja.

## 1.5 General information of the project

This LuaTeX-ja project is hosted by SourceForge.jp. The official wiki is located on <http://sourceforge.jp/projects/luatex-ja/wiki/>. There is no stable version on October 22, 2011, however a set of developer sources can be obtained from the git repository. Members of the project team are as follows (in random order): Hironori Kitagawa, Kazuki Maeda, Takayuki Yato, Yusuke Kuroki, Noriyuki Abe, Munehiro Yamamoto, Tomoaki Honda, and Shuzaburo Saito.

## 2 Major differences with pTeX

In this section, we explain several major differences between pTeX and our LuaTeX-ja. For general information of Japanese typesetting and the overview of pTeX, please see Okumura [14].

### 2.1 Names of control sequences

Because pTeX is an engine modification of Knuth's original TeX82 engine, some of the additional primitives take a form that is very difficult to be simulated by a macro. For example, an additional primitive `\prebreakpenalty<char_code> [=] <penalty>` in pTeX sets the amount of penalty inserted before a character whose code is `<char_code>` to `<penalty>`, and this form `\prebreakpenalty<char_code>` can be also used for retrieving the value.

Moreover, there are some internal parameters of pTeX which values of them at the end of a horizontal box or that of a paragraph are valid in whole box or paragraph. However, the implementation of these parameters in LuaTeX-ja is not so easy; we will discuss it in Subsection 5.2.

From the two problems discussed above, the assignment and retrieval of most parameters in LuaTeX-ja are summarized into the following three control sequences:

xあy    xあ y	<code>1 \font\x=IPAMincho \x</code> <code>2 \ltjsetparameter{jacharrange={-6}}x あ</code> <code>3 y\quad x あ</code> <code>4 y</code>
-------------	--

FIGURE 1. A notable sample showing the treatment of a line-break after a Japanese character.

- `\ltjsetparameter{<name>=<value>, ...}`: for local assignment.
- `\ltjglobalsetparameter`: for global assignment. Note that these two control sequences obey the value of `\globaldefs` primitive.
- `\ltjgetparameter{<name>}[<optional argument>]`: for retrieval. The returned value is always a string.

## 2.2 Line-break after a Japanese character

Japanese texts can break lines almost everywhere, in contrast with alphabetic texts can break lines only between words (or use hyphenation). Hence, pT<sub>E</sub>X's input processor is modified so that a line-break after a Japanese character doesn't emit a space. However, there is no way to customize the input processor of LuaT<sub>E</sub>X, other than to hack its CWEB-source. All a macro package can do is to modify an input line before when LuaT<sub>E</sub>X begin to process it, inside the `process_input_buffer` callback.

Hence, in LuaT<sub>E</sub>X-ja, a comment letter (we reserve U+FFFFF for this purpose) will be appended to an input line, if this line ends with a Japanese character.<sup>1</sup> One might jump to a conclusion that the treatment of a line-break by pT<sub>E</sub>X and that of LuaT<sub>E</sub>X-ja are totally same, however they are different in the respect that LuaT<sub>E</sub>X-ja's judgment whether a comment letter will be appended the line is done *before* the line is actually processed by LuaT<sub>E</sub>X.

Figure 1 shows an example of this situation; the command at the first line marks most of Japanese characters as 'non-Japanese characters'. In other words, from that command onward, the letter 'あ' will be treated as an alphabetic character by LuaT<sub>E</sub>X-ja. Then, it is natural to have a space between 'あ' and 'y' in the output (as the second example in the figure), where the actual output of the first example in the figure does not so. This is because 'あ' at the input line 2 is considered a Japanese character by LuaT<sub>E</sub>X-ja, when LuaT<sub>E</sub>X-ja does the decision whether U+FFFFF will be added to the input line 2.

## 2.3 Separation between 'real' fonts and metrics

Traditionally, most Japanese fonts used in typesetting are not proportional, that is, most glyphs have same size (in most cases, square-shaped). Hence, it is not rare that the contents of different JFMs are essentially same, and only differ in their names. For example, `min10.tfm` and `goth10.tfm`, which are JFMs shipped with pT<sub>E</sub>X for seriffed

1. Strictly speaking, it also requires that the catcode of the end-line character is 5 (*end-of-line*). This condition is useful under the verbatim environment.

---

```
\jfont\foo=file:ipam.ttf:jfm=ujis;script=latn;-kern;+jp04 at 12pt
\jfont\bar=psft:Ryumin-Light:jfm=ujis at 10pt
```

FIGURE 2. Typical declarations of Japanese fonts.

*mincho* family and sans-serifred *gothic* family, differ their FAMILY and FACE only. Moreover, *jis.tfm* and *jisg.tfm*, which is included in the *jis* font metric, which is used in *jsclasses* [13] by Haruhiko Okumura (奥村晴彦), are totally same as binary files. Considering this situation, we decided to separate ‘real’ fonts and metrics used for them in Lua $\TeX$ -ja. Typical declarations of Japanese fonts in the style of plain  $\TeX$  are shown in Figure 2. We would like to add several remarks:

- A control sequence `\jfont` must be used for Japanese fonts, instead of `\font`.
- Lua $\TeX$ -ja automatically loads the *luaotfload* package, so `file:` and `name:` prefixes, and various font features can be used as the first line in Figure 2.
- The `jfm` key specifies the metric for the font. In Figure 2, `\foo` and `\bar` will use a metric stored in a Lua script named `jfm-ujis.lua`. This metric is the standard metric in Lua $\TeX$ -ja, and is based on JFMs used in the *otf* package [16] (hence almost all characters are square-shaped).
- The `psft:` prefix can be used to specify name-only, non-embedded fonts. When one displays a pdf with these fonts, actual fonts which will be used for them depend on a pdf reader.

The specification of a metric for Lua $\TeX$ -ja is similar to that of a JFM (see [14]); characters are grouped into several classes, the size information of characters are specified for each class, and glue/kern insertions are specified for each pair of classes. Although the author have not tried, it may be possible to develop a program that ‘converts’ a JFM to a metric for Lua $\TeX$ -ja. Lua $\TeX$ -ja offers three metrics by default; `jfm-ujis.lua`, `jfm-jis.lua` based on the *jis* font metric, and `jfm-min.lua` based on old *min10.tfm*.

Note that `-kern` in features is important, because kerning information from a real font itself will clash with glue/kern information from the metric.

## 2.4 Insertion of glues/kerns for Japanese typesetting: timing

As described in [9], Lua $\TeX$ ’s kerning and ligaturing processes are totally different from those of  $\TeX$ 82.  $\TeX$ 82’s process is done just when a (sequence of) character is appended to the current list. Thus we can interrupt this process by writing as `f{}irm`. However, Lua $\TeX$ ’s process is *node-based*, that is, the process will be done when a horizontal box or a paragraph is ended, so `f{}irm` and `firm` yield same outputs under Lua $\TeX$ .

The situation for Japanese characters is more complicated. Glues (and kerns) which are needed for Japanese typesetting are divided into the following three categories:

- Glue (or kern) from the metric of Japanese fonts (*JFM glue*, for short).
- Default glue between a Japanese character and an alphabetic character (we say *xkanjiskip*, for short), usually 1/4 of full-width (*shibuaki*) with some stretch and shrink for justifying each line.

TABLE 1. Examples of differences between pT<sub>E</sub>X and LuaT<sub>E</sub>X-ja.

	(1)	(2)	(3)	(4)
Input	あ] {} [ ] \ / [	い』 \ / a	う) \ hbox{} (	え] \ special{} [
pT <sub>E</sub> X	あ] [ ] [	い』 a	う) (	え] [
LuaT <sub>E</sub> X-ja	あ] [ ] [	い』 a	う) (	え] [

FIGURE 3. Detail of the output of pT<sub>E</sub>X in the input (1) in Table 1.

- Default glue between two consecutive Japanese characters (*kanjiskip*, for short). The main reason of this glue is to enable breaking lines almost everywhere in Japanese texts. In most cases, its natural width is zero, and some stretch/shrink for justifying each line.

In pT<sub>E</sub>X, these three kinds of glues are treated differently. A JFM glue is inserted when a (sequence of) Japanese character is appended to the current list, same as the case of alphabetic characters in T<sub>E</sub>X82. This means that one can interrupt the insertion process by saying `{}`. A *xkanjiskip* is inserted just before ‘hpack’ or line-breaking of a paragraph; this timing is somewhat similar to that of LuaT<sub>E</sub>X’s kerning process. Finally, A *kanjiskip* is not appeared as a node anywhere; only appears implicitly in calculation of the width of a horizontal box, that of breaking lines, and the actual output process to a DVI file. These specifications have made pT<sub>E</sub>X’s behavior very hard to understand.

LuaT<sub>E</sub>X-ja inserts glues in all three categories simultaneously inside `hpack_filter` and `pre_linebreak_filter` callbacks. The reasons of this specification are to behave like alphabetic characters in LuaT<sub>E</sub>X (as described in the first paragraph in this subsection), and to clarify the specification for LuaT<sub>E</sub>X-ja’s process.

## 2.5 Insertion of glues/kerns for Japanese typesetting: specification

Now we will take a look at the insertion process itself through four points.

*Ignored nodes* As noted in the previous subsection, the insertion process in pT<sub>E</sub>X can be interrupted by saying `{}` or anything else.<sup>2</sup> This leads the second row in Table 1, or Figure 3. Here ‘the process is interrupted’ means that pT<sub>E</sub>X does not think the letter ‘]’ is followed by ‘[’, hence two half-width glues are inserted between ‘]’ and ‘[’, where the left one is from ‘]’ and the right one is from ‘[’.

On the other hand, in LuaT<sub>E</sub>X-ja, the process is done inside `hpack_filter` and `pre_linebreak_filter` callbacks. Hence, *anything that does not make any node will be ignored* in LuaT<sub>E</sub>X-ja, as shown in (1) in Table 1. LuaT<sub>E</sub>X-ja also ignores any nodes which does not make any contribution to current horizontal list—*ins\_node*, *adjust\_node*, *mark\_node*, *whatsit\_node* and *penalty\_node*—, as shown in (4).

2. This is why some tricks like `ちよ{}つと` for `min10.tfm` and other ‘old’ JFMs work.

By the way, around a *glyph\_node*  $p$  there may be some nodes attached to  $p$ . These are an accent and kerns for moving it to the right place, and a kern from the italic correction<sup>3</sup> for  $p$ . It is natural that these attachments should be ignored inside the process. Hence LuaTeX-ja takes this approach, as the latest version of pTeX (version p3.2). This explains (2) in the Table 1.

Summarizing the above, one should put an empty horizontal box `\hbox{}` to where he/she wants to interrupt the insertion process in LuaTeX-ja as (3) in the Table 1.

*Fonts with the same metric* Recall that LuaTeX-ja separates ‘real’ fonts and metrics, as in Subsection 2.3. Consider the following input, where all Japanese fonts use same metric (in LuaTeX-ja), and `\gt` selects *gothic* family for the current Japanese font family:

```
明朝) \gt (ゴシック
```

If the above input is processed by pTeX, because the insertion process is interrupt by `\gt`, the result looks like

```
明朝) (ゴシック
```

However this seems to be unnatural, since two Japanese fonts in the output use the same metric, i.e., the same typesetting rule. Hence, we decided that Japanese fonts with the same metric are treated as one font in the insertion process of LuaTeX-ja. Thus, the output from the above input in LuaTeX-ja looks like:

```
明朝) (ゴシック
```

One might have the situation that this default behavior is not suitable. LuaTeX-ja offers a way to handle this situation, but we leave it to the manual [10].

*Fonts with different metrics* The case where two adjacent Japanese characters use different metrics and/or different size is similar. Consider the following input where the *mincho* family and the *gothic* family use different metrics:

```
漢) \gt (漢) \large (大
```

As the previous paragraph, this input yields the following, by pTeX:

```
漢) (漢) (大
```

We had thought that amounts of spaces between parentheses in the above output are too much. Hence we have changed the default behavior of LuaTeX-ja, so that the amount of a glue between two Japanese characters with different metrics is the *average* of a glue from the left character and that from the right character. For example, Figure 4 shows the output from the above input. The width of glue indicated ‘(1)’ is  $(a/2 + a/2)/2 = 0.5a$ , and the width of glue indicated ‘(2)’ is  $(a/2 + 1.2a/2)/2 = 0.55a$ . This default behavior can be changed by `differentjfm` parameter of LuaTeX-ja.

*kanjiskip and xkanjiskip* In pTeX, the value of *xkanjiskip* is controlled by a skip named `\xkanjiskip`. A well-known defect of this implementation is that the value of

3. TeX82 (and LuaTeX) does not distinguish between explicit kern and a kern for italic correction. To distinguish them, an additional subtype for a kern is introduced in pTeX. On the other hand, LuaTeX-ja uses an additional attribute and redefines `\` to set this attribute.



FIGURE 4. Fonts with different metrics.

*xkanjiskip* is not connected with the size of the current Japanese font. It seems that EXTRASPACE, EXTRASTRETCH, EXTRASHRINK parameters in a JFM are reserved for specifying the default value of *xkanjiskip* in a unit of the design size, but p<sub>T</sub>E<sub>X</sub> did not use these parameters, actually.

Considering this situation of p<sub>T</sub>E<sub>X</sub>, Lua<sub>T</sub>E<sub>X</sub>-ja can use the value of *xkanjiskip* that specified in a metric. If the value of *xkanjiskip* on user side (this is the value of `xkanjiskip` parameter of `\ltjsetparameter`) is `\maxdimen`, then Lua<sub>T</sub>E<sub>X</sub>-ja uses the specification from the current used metric as the actual value of *xkanjiskip*. This description also applies for *kanjiskip*.

### 3 Distinction of characters

Since Lua<sub>T</sub>E<sub>X</sub> can handle Unicode characters natively, it is a major problem that how we distinguish Japanese characters and alphabetic characters. For example, the multiplication sign (U+00D7) exists both in ISO-8859-1 (hence in Latin-1 Supplement in Unicode) and in the basic Japanese character set JIS X 0208. It is not desirable that this character is always treated as an alphabetic character, because this symbol is often used in the sense of ‘negative’ in Japan.

#### 3.1 Character ranges

Before we describe the approach taken in Lua<sub>T</sub>E<sub>X</sub>-ja, we review the approach taken by up<sub>T</sub>E<sub>X</sub>. up<sub>T</sub>E<sub>X</sub> extends the `\kcatcode` primitive in p<sub>T</sub>E<sub>X</sub>, to use this primitive for setting how a character is treated among alphabetic characters (15), *kanji* (16), *kana* (17), *Hangul* (17), or *other CJK characters* (18). The assignment to `\kcatcode` can be done by a Unicode block.<sup>4</sup>

Lua<sub>T</sub>E<sub>X</sub>-ja adopted a different approach. There are many Unicode blocks in Basic Multilingual Plane which are not included in Japanese fonts, therefore it is inconvenient if we process by a Unicode block. Furthermore, JIS X 0208 are not just union of Unicode blocks; for example, the intersection of JIS X 0208 and Latin-1 Supplement is shown in Table 2. Considering these two points, to customize the range of Japanese characters in Lua<sub>T</sub>E<sub>X</sub>-ja, one has to define ranges of character codes in his/her source in advance.

We note that Lua<sub>T</sub>E<sub>X</sub>-ja offers two additional control sequences, `\ltjjachar` and `\ltjalchar`. They are similar to `\char` primitive, however `\ltjjachar` always yields a Japanese character, provided that the argument is more than or equal to 128, and `\ltjalchar` always yields an alphabetic character, regardless of the argument.

4. There are some exceptions. For example, U+FF00–FFEF (Halfwidth and Fullwidth Forms) are divided into three blocks in recent up<sub>T</sub>E<sub>X</sub>.



TABLE 2. Intersection of JIS X 0208 and Latin-1 Supplement.

§ (U+00A7), ¨ (U+00A8), ° (U+00B0), ± (U+00B1),  
 ´ (U+00B4), ¶ (U+00B6), × (U+00D7), ÷ (U+00F7)

TABLE 3. Predefined ranges in LuaTeX-ja.

- 1 (Additional) Latin characters which are not belonged in the range 8.
- 2 Greek and Cyrillic letters.
- 3 Punctuations and miscellaneous symbols.
- 4 Unicode blocks which does not intersect with Adobe-Japan1-6.
- 5 Surrogates and supplementary private use Areas.
- 6 Characters used in Japanese typesetting.
- 7 Characters possibly used in CJK typesetting, but not in Japanese.
- 8 Characters in Table 2.

### 3.2 Default setting of ranges

Patches for plain TeX and L<sup>A</sup>T<sub>E</sub>X<sub>2 $\epsilon$</sub>  of LuaTeX-ja predefine eight character ranges, as shown in Table 3. Almost of these ranges are just the union of Unicode blocks, and determined from the Adobe-Japan1-6 character collection [1], and JIS X 0208. Among these eight ranges, the ranges 2, 3, 6, 7, and 8 are considered ranges of Japanese characters, and others are considered ranges of alphabetic characters.<sup>5</sup> We remark on ranges 2 and 8:

*The range 2* JIS X 0208 includes Greek letters and Cyrillic letters, however, these letters cannot be used for typesetting Greek or Russian, of course. Hence it is reasonable that Greek letters and Cyrillic consist another character range.

*The range 8* If one wants to use 8-bit TFMs, such as T1 or TS1 encodings, he should mark this range 8 as a range of alphabetic characters by

```
\ltjsetParameter{jacharrange={-8}}
```

This is because some 8-bit TFMs have a glyph in this range; for example, the character ‘ $\mathcal{E}$ ’ is located at "D7 in the T1 encoding.

### 3.3 Control sequences producing Unicode characters

The *fontspec* package<sup>6</sup> offers various control sequences that produce Unicode characters. However, these control sequences as it stands cannot work correctly with the default range setting of LuaTeX-ja. For example, `\textquotedblleft` is just an abbreviation of `\char"201C\relax`, and the character U+201C (LEFT DOUBLE QUOTATION MARK) is treated as an Japanese character, because it belongs to the range 3. This problem is resolved by using `\tjchar` instead of the `\char` primitive. It is included

5. Note that ranges 3 and 8 are considered ranges of alphabetic characters in this paper.

6. Precisely saying, it is the *xunicode* package, originally a package for X<sub>Y</sub>L<sub>A</sub>T<sub>E</sub>X and automatically loaded by the *fontspec* package.

x, x, x, X, x	<code>1 x, \char`x, % depend on range setting</code> <code>2 \ltjalchar`x, % alphabetic char</code> <code>3 \ltjjachar`x, % Japanese char</code> <code>4 \texttimes % alph. char (by fontspec)</code>
---------------	--

FIGURE 5. Control sequences producing a Unicode character.

in an optional package named `luatexja-fontspec.sty`. Figure 5 shows several ways to typeset a character, both as a Japanese character and as an alphabetic characters.

The situation looks similar in math formulas, but in fact it differs. Each control sequence that represents an ordinary symbol defined by the *unicode-math* package is just synonym of a character. For example, the meaning of `\otimes` is just the character U+2297 (CIRCLED TIMES), which is included in the range 3. However, it is difficult to define a control sequence like `\ltjalUmathchar` as a counterpart of `\Umathchar`, since an input like `'\sum^{\ltjalUmathchar} \dots'` has to be permitted.

However, we couldn't develop a satisfactory solution to this problem in time for this paper, due to a lack of time. We are just testing a solution below:

- LuaT<sub>E</sub>X-ja has a list of character codes which will be always treated as alphabetic characters in math mode. Considering 8-bit TFMs for math symbols, this list includes natural numbers between "80 and "FF by default.
- Redefine internal commands defined in the *unicode-math* package so that codes of characters which are mentioned in the *unicode-math* package will be included in the list.

We would like to extend treatments described in this subsection to 8-bit font encodings, but we leave it to further development too.

## 4 Current status of development

At the moment, LuaT<sub>E</sub>X-ja can be used under plain T<sub>E</sub>X, and under L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>. Generally speaking, one only has to read `luatexja.sty`, by `\input` command or `\usepackage` (in L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub>), if you merely want to typeset Japanese characters. We look more details by parts.

### 4.1 'Engine extension'

The lowest part of LuaT<sub>E</sub>X-ja corresponds to the pT<sub>E</sub>X extension as *an engine extension of T<sub>E</sub>X*. We, the project members, think that this part is almost done. There is one more feature of LuaT<sub>E</sub>X-ja which we are going to explain:

*Shifting baseline* In order to make a match between Japanese fonts and alphabetic fonts, sometimes shifting the baseline of alphabetic characters may be needed. pT<sub>E</sub>X has a dimension `\ybaselineshift`, which corresponds to the amount of shifting down the baseline of alphabetic characters. This is useful for Japanese-based documents, but not for documents mainly in languages with alphabetic characters.

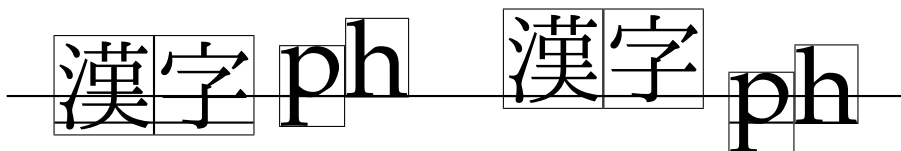


FIGURE 6. First example of shifting baseline.



FIGURE 7. Second example of shifting baseline.

Hence, Lua $\TeX$ -ja extends p $\TeX$ 's `\ybaselineshift` to Japanese characters. Namely, Lua $\TeX$ -ja offers two parameters, `yjabaselineshift` and `yalbaselineshift`, for the amount of shifting the baseline of Japanese characters and that of alphabetic characters, respectively.

An example output is shown in Figure 6. The left half is the output when `yjabaselineshift` is positive, hence the baseline of Japanese characters is shifted down. On the other hand, the right half is the output when `yalbaselineshift` is positive, hence the baseline of alphabetic characters is shifted down. Figure 7 shows an interesting use of these parameters.

Note that Lua $\TeX$ -ja doesn't support vertical typesetting, *tategaki*, for now.

## 4.2 Patches for plain $\TeX$ and $\LaTeX 2_{\epsilon}$

p $\TeX$  has a patch for plain  $\TeX$ , namely `ptex.tex`, that for  $\LaTeX 2_{\epsilon}$  macro (this patch and  $\LaTeX 2_{\epsilon}$  consist `p\LaTeX 2_{\epsilon}`), and `kinsoku.tex` which includes the default setting of *kinsoku shori*, the Japanese hyphenation. We ported them to Lua $\TeX$ -ja, except the codes related to vertical typesetting, because Lua $\TeX$ -ja doesn't support vertical typesetting yet. We remark one point related to the porting:

*Behavior of `\fontfamily`* The control sequence `\fontfamily` in p $\LaTeX 2_{\epsilon}$  changes the current alphabetic font family and/or the current Japanese font family, depending the argument. More concretely, `\fontfamily{<arg>}` changes the current alphabetic font family to `<arg>`, if and only if one of the following conditions are satisfied:

- An alphabetic font family named `<arg>` in *some* alphabetic encoding is already defined in the document.
- There exists an alphabetic encoding `<enc>` already defined in the document such that a font definition file `<enc><arg>.fd` (all lowercase) exists.

The same criterion is used for changing Japanese font family.

To work this behavior well, it is required that a list of all (alphabetic) encodings defined already in the document. However, since Lua $\TeX$ -ja is loaded as a

package, LuaT<sub>E</sub>X-ja cannot have this list. Hence LuaT<sub>E</sub>X-ja adopted a different approach, namely `\fontfamily{<arg>}` changes the current alphabetic font family to `<arg>`, if and only if:

- An alphabetic font family named `<arg>` in the current alphabetic encoding `<enc>` is already defined in the document.
- A font definition file `<enc><arg>.fd` (all lowercase) exists.

### 4.3 Classes for Japanese documents

To produce ‘high-quality’ Japanese documents, we need not only that Japanese characters are correctly placed, but also class files for Japanese documents. Two major families of classes are widely used in Japan: *jclasses* which is distributed with the official pT<sub>E</sub>X 2<sub>ε</sub> macros, and *jsclasses*. At the present, LuaT<sub>E</sub>X-ja simply contains their counterparts: *ltjclasses* and *ltjsclasses*. However, the policy on classes is not determined now, and we hope to have another family of classes which are useful for commercial printing. In the author’s opinion, *ltjclasses* is better to stay as an example of porting of class files for pT<sub>E</sub>X to LuaT<sub>E</sub>X-ja.

### 4.4 Patches for packages

Apart from patches for the L<sup>A</sup>T<sub>E</sub>X 2<sub>ε</sub> kernel and classes for Japanese documents, we need to make patches for several packages. At the present, we considered the following packages, and made patches or porting for the former two packages.

*The fontspec package* The *fontspec* package is built on NFSS2, hence control sequences offered by the *fontspec* package, such as `\setmainfont`, are only effective for alphabetic fonts if LuaT<sub>E</sub>X-ja is loaded. `luatexja-fontspec.sty` (not automatically loaded) offers these counterparts for Japanese fonts, with additional ‘j’ in the name of control sequences, such as `\setmainjfont`. As described in Subsection 3.3, it also includes a patch for control sequences producing Unicode characters.

*The otf package* This package is widely used in pT<sub>E</sub>X for typesetting characters which is not in JIS X 0208, and for using more than one weight in *mincho* and *gothic* font families. Therefore LuaT<sub>E</sub>X-ja supports features in the *otf* package, by loading `luatexja-otf.sty` manually. Note that characters by `\UTF{}` and `\CID{}` are not appended to the current list as a *glyph\_node*, to avoid from callbacks by the *luaotfload* package. We have another remark; `\CID` does not work with TrueType fonts, since `\CID` uses the conversion table between CID and the glyph order of the current Japanese font.

*The listings package* It is known for users of pT<sub>E</sub>X that there is a patch `jlisting.sty` for the *listings* package, to use Japanese characters in the `lstlisting` environment. Generally speaking, it also can be used in LuaT<sub>E</sub>X-ja. However, it seems to be that a Japanese character after a space does not receive any process of the *listings* package; this is inconvenient when we use the *showexpl* package.

There is another way to use characters whose code are above 256 with the *listings* package (described in [3]). However, this method is not suitable for Japanese, since the number of Japanese characters is very large. We hope that the *listings*

package will be able to handle all characters above 256 without any patch, in the future.

## 5 Implementation

### 5.1 Handling of Japanese fonts

In pTeX, there are three slots for maintaining current fonts, namely `\font` for alphabetic fonts, `\jfont` for Japanese fonts (in horizontal direction) and `\tfont` for Japanese fonts (in vertical direction). With these slots, we can manage the current font for alphabetic characters and that for Japanese characters separately in pTeX. However, LuaTeX has only one slot for maintaining the current font, as TeX82. This situation leads a problem: how can we maintain the ‘current Japanese font’?

There are three approaches for this problem. One approach is to make a mapping table from alphabetic fonts to corresponding Japanese fonts (here we don’t assume that NFSS2 is available). Another approach is that we always use composite fonts with alphabetic fonts and Japanese fonts. The third approach is that the information of the current Japanese font is stored in an attribute. We adopted the third approach, since LuaTeX-ja is much affected by pTeX as we noted in Subsection 1.2.

As in Figure 2, LuaTeX-ja uses `\jfont` for defining Japanese fonts, as pTeX. However, because the information of the current Japanese font is stored into an attribute, control sequences defined by `\jfont` (e.g., `\foo` and `\bar` in Figure 2) is not representing a font by the means of TeX82. In other words, each of these control sequences is just an assignment to an attribute, therefore they cannot be an argument of `\the`, `\fontname`, nor `\textfont`.

Callbacks by the *luaotfload* package, e.g., replacement of glyphs according to OpenType font features, are performed just after ‘Examination of stack level’ (see Subsections 1.3 and 5.2). Also note that calculation of character classes for each Japanese character is done *after* the these callbacks for now.

### 5.2 Stack management

As we noted in Subsection 2.1, parameters that the values at the end of a horizontal box or that of a paragraph are valid in whole box or paragraph, such as *kanjiskip*, cannot be implemented by internal integers or registers of other types in TeX. We explain it in this subsection.

Figure 8 is an extract of a CWEB-source `tex/packaging.w` of LuaTeX (SVN revision 4358). This function is called just when an explicit `\hbox{...}` or `\vbox{...}` is ended, and the function `filtered_hpack()` is where the `hpack_filter` and then the actual ‘hpack’ process are performed. Notice that the `unsave()` function is called before `filtered_hpack()`. This is the problem; because of `unsave()`, we can retrieve only the values of registers *outside* the box, even in the `hpack_filter` callback.

To cope with this problem, LuaTeX-ja has its own stack system, based on Lua codes in [17]. Furthermore, *whatsit* nodes whose `user_id` is 30112 (`stack_node`, for short) will be appended to the current horizontal list each time the current stack level is incremented, and their values are the values of `\currentgrouplevel` at that time. In the beginning

```

void package(int c)
{
    ...
    d = box_max_depth;
    unsave();
    save_ptr -= 4;
    if (cur_list.mode_field == -hmode) {
        cur_box = filtered_hpack(cur_list.head_field,
                                cur_list.tail_field, saved_value(1),
                                saved_level(1), grp, saved_level(2));
        subtype(cur_box) = HLIST_SUBTYPE_HBOX;
    } else {

```

FIGURE 8. An extract of a CWEB-source `tex/packaging.w` of LuaT<sub>E</sub>X.

of the `hpack_filter` callback, the list in question is traversed to determine whether the stack level at the end of the list and that outside the box coincides.

Let  $x$  be the value of `\currentgrouplevel`, and  $y$  be the current stack level, both inside the `hpack_filter` callback, i.e., outside a horizontal box. Consider a list which represents the content of the box, then we have:

- A *stack\_node* whose value is  $x + 1$  (because all materials in the box are included in a group `\hbox{...}`), the value of `\currentgrouplevel` inside the box is at least  $x + 1$  in the list corresponds to an assignment related to the stack system in just top-level of the list, like

```
\hbox{...(assignment)...}
```

In this case, the current stack level is incremented to  $y + 1$  after the assignment.

- A *stack\_node* whose value is more than  $x + 1$  in the list corresponds to an assignment inside another group contained in the box. For example, the following input creates a *stack\_node* whose value is  $x + 3 = (x + 1) + 2$ :

```
\hbox{...{...{...(assignment)}...}...}
```

Thus, we can conclude that the stack level at the end of the list is  $y + 1$ , if and only if there is a *stack\_node* whose value is  $x + 1$ . Otherwise, the stack level is just  $y$ .

### 5.3 Adjustment of the position of Japanese characters

The size of a glyph specified in a metric and that of a real font usually differ. For example, the letter ‘l’ is half-width in `jfm-ujis.lua` or `jis.tfm`, while this letter is full-width like ‘l’ in most TrueType fonts used in Japanese typesetting, such as IPA Mincho. Hence the adjustment of position of such glyphs is needed. In the context of pT<sub>E</sub>X, this process was performed using virtual fonts.

On the other hand, LuaT<sub>E</sub>X-ja does the adjustment by encapsulating a glyph into a horizontal box. There are two main reasons why we adopted this method; one is that we feared Lua codes for coexisting with callbacks by the `luaotfload` package would be large if we use virtual fonts, and the other is to cope with shifting of the baseline of characters at the same time.

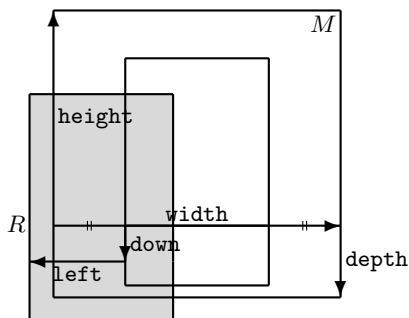


FIGURE 9. The position of the ‘real’ glyph.

Figure 9 shows the adjustment process. A large square  $M$  is the imaginary body specified in the metric, and a vertical rectangle is the imaginary body of a real glyph. First, the real glyph is aligned with respect to the width of  $M$ . In the figure, the real glyph is aligned ‘middle’; this setting is useful for the full-width middle dot ‘ $\cdot$ ’. We have other settings, ‘left’ and ‘right’. Furthermore, it is shifted according to the value of `left` and `down`, which are specified in the metric, for fine adjustment. The final position of the real glyph is shown by the gray rectangle  $R$ . If the amount of shifting the baseline is not zero,  $M$  (and hence the real glyph) is shifted by that amount.

We would like to remark briefly on the vertical position of a real glyph. A JFM (or a metric used in Lua $\TeX$ -ja) and a real font used for it may have different height or depth. In that case, it may look better if the real glyph is shifted vertically to match the height-depth ratio specified in the metric, while any vertical adjustment except the adjustment by the `down` value does not performed in the present implementation of Lua $\TeX$ -ja. This situation is carefully studied by Otake [15]. Here the policy on this problem is not determined now, however we would like to offer several solutions in future development.

#### 5.4 Further notes on metrics for Lua $\TeX$ -ja

*Proportional typesetting* Some fonts are proportional, that is, each glyphs in those fonts have its own width. An example of proportional fonts is IPA P Mincho. Using these fonts in p $\TeX$  is very hard, since one needs to make a dedicated JFM for a real font.

Lua $\TeX$ -ja supports these proportional fonts; specifying the width of a character class in a metric to “prop” makes the width of each character in this class that of a glyph in a real font. If no JFM glue is needed, one simply has to use `jfm-prop.lua`. The following is an example:

```

あいうえお
あいうえお

```

```

1 \jfont\pr=file:ipamp.ttf:jfm=prop at 3.25mm
2 あいうえお\\pr{あいうえお}

```

*Scaling by metrics* Because of virtual fonts, even if one specifies to use `min10.tfm` or `jis.tfm` at 10 pt in p $\TeX$ , the actual size of real fonts used in dviwares for these JFMs are 9.62216 pt. Hence, for example, if one wants to use 3.25 mm Japanese

fonts and 10 pt alphabetic fonts in pT<sub>E</sub>X, he/she needs to scale a Japanese font by

$$\frac{3.25 \text{ mm}}{10 \text{ pt} \cdot 0.962216} \simeq 0.961$$

in declarations of Japanese fonts.

LuaT<sub>E</sub>X-ja didn't support such scaling of glyphs by metrics, so one has to adjust the size argument for `\jfont` manually. Continuing the previous example, for using 3.25 mm Japanese fonts and 10 pt alphabetic fonts in LuaT<sub>E</sub>X-ja, he/she needs to scale a Japanese font by  $3.25 \text{ mm}/10 \text{ pt} \simeq 0.92487$ .

## 6 Conclusion

We have discussed about our LuaT<sub>E</sub>X-ja package, which is much affected by pT<sub>E</sub>X. For now, it can be used for experimental use, however there are much refinements which are needed for regular use. The author hopes that this paper and LuaT<sub>E</sub>X-ja project contribute the typesetting Japanese, and possibly other Asian languages, under LuaT<sub>E</sub>X.

## Acknowledgements

The author would like to thank Ken Nakano and Hideaki Togashi for their development and management of ASCII pT<sub>E</sub>X. The author is very grateful to Haruhiko Okumura for his leadership in the Japanese T<sub>E</sub>X community. The author is also very grateful to members of LuaT<sub>E</sub>X-ja project team for their valuable cooperation in development.

## References

1. Adobe Systems Incorporated, *Adobe-Japan1-6 Character Collection for CID-Keyed Fonts*, Technical Note #5078, 2004. <http://partners.adobe.com/public/developer/en/font/5078.Adobe-Japan1-6.pdf>
2. ASCII MEDIA WORKS, アスキー日本語 T<sub>E</sub>X (pT<sub>E</sub>X). <http://ascii.asciimw.jp/pb/ptex/>
3. John Baker, *Typesetting UTF8 APL code with the T<sub>E</sub>X l<sub>s</sub>tlisting package*. <http://bakerjd99.wordpress.com/2011/08/15/>
4. Jin-Hwan Cho and Haruhiko Okumura, *Typesetting CJK Languages with Omega, T<sub>E</sub>X, XML, and Digital Typography*, Lecture Notes in Computer Science, vol. 3130, Springer, 2004, 139–148.
5. Yannis Haralambous. *The Joy of LuaT<sub>E</sub>X*. <http://luatex.bluwiki.com/>
6. Japanese Industrial Standards Committee. *JIS X 4051: Formatting rules for Japanese documents*, 1993, 1995, 2004.
7. 北川弘典,  $\epsilon$ -pT<sub>E</sub>X についての wiki. <http://sourceforge.jp/projects/eptex/wiki/FrontPage>
8. 北川弘典, LuaT<sub>E</sub>X で日本語. <http://oku.edu.mie-u.ac.jp/tex/mod/forum/discuss.php?d=378>
9. LuaT<sub>E</sub>X development team, *The LuaT<sub>E</sub>X reference*. <http://www.luatex.org/svn/trunk/manual/luatexref-t.pdf> (snapshot of SVN trunk)



- 
10. LuaTeX-ja project team, *The LuaTeX-ja package*. Not completed for now. Available at `doc/man-en.pdf` (in English) or `doc/man-ja.pdf` (in Japanese) in the Git repository.
  11. 香田温人, LuaTeX と日本語. <http://www1.pm.tokushima-u.ac.jp/~kohda/tex/luatex-old.html>
  12. 前田一貴, luajalayout パッケージ—LuaTeXによる日本語組版—. <http://www-is.amp.i.kyoto-u.ac.jp/lab/kmaeda/lualatex/luajalayout/>
  13. 奥村晴彦, pTeX2 $\epsilon$ 新ドキュメントクラス. <http://oku.edu.mie-u.ac.jp/~okumura/jsclasses/>
  14. Haruhiko Okumura, *pTeX and Japanese Typesetting*, *The Asian Journal of T<sub>E</sub>X 2* (2008), 43–51.
  15. 乙部巖己, min10 フォントについて. <http://argent.shinshu-u.ac.jp/~otobe/tex/files/min10.pdf>
  16. 齋藤修三郎, Open Type Font 用 VF. <http://psitau.kitunebi.com/otf.html>
  17. Jonathan Sauer, *[Dev-luatex] tex.currentgrouplevel*. <http://www.ntg.nl/pipermail/dev-luatex/2008-August/001765.html>
  18. Takuji Tanaka, *upTeX, upL<sup>A</sup>T<sub>E</sub>X—unicode version of pTeX, pL<sup>A</sup>T<sub>E</sub>X*. [http://homepage3.nifty.com/ttk/comp/tex/uptex\\_en.html](http://homepage3.nifty.com/ttk/comp/tex/uptex_en.html)
  19. Nobuyuki Tsuchimura and Yusuke Kuroki, *Development of Japanese T<sub>E</sub>X Environment*, *The Asian Journal of T<sub>E</sub>X 2* (2008), 53–62.
  20. W3C Working Group, *Requirements for Japanese Text Layout*. <http://www.w3.org/TR/jlreq/>