

## Contents

<b>1</b>	<b>Encoding and escaping schemes</b>	<b>2</b>
<b>2</b>	<b>Conversion functions</b>	<b>4</b>
<b>3</b>	<b>Internal string functions</b>	<b>4</b>
<b>4</b>	<b>Possibilities, and things to do</b>	<b>5</b>
	<b>Index</b>	<b>5</b>

# The `l3str-convert` package: string encoding conversions\*

The L<sup>A</sup>T<sub>E</sub>X3 Project<sup>†</sup>

Released 2017/02/20

## 1 Encoding and escaping schemes

Traditionally, string encodings only specify how strings of characters should be stored as bytes. However, the resulting lists of bytes are often to be used in contexts where only a restricted subset of bytes are permitted (*e.g.*, PDF string objects, URLs). Hence, storing a string of characters is done in two steps.

- The code points (“character codes”) are expressed as bytes following a given “encoding”. This can be UTF-16, ISO 8859-1, *etc.* See Table 1 for a list of supported encodings.<sup>1</sup>
- Bytes are translated to T<sub>E</sub>X tokens through a given “escaping”. Those are defined for the most part by the `pdf` file format. See Table 2 for a list of escaping methods supported.<sup>2</sup>

---

\*This file describes v6969, last revised 2017/02/20.

<sup>†</sup>E-mail: [latex-team@latex-project.org](mailto:latex-team@latex-project.org)

<sup>1</sup>Encodings and escapings will be added as they are requested.

Table 1: Supported encodings. Non-alphanumeric characters are ignored, and capital letters are lower-cased before searching for the encoding in this list.

$\langle Encoding \rangle$	description
<code>utf8</code>	UTF-8
<code>utf16</code>	UTF-16, with byte-order mark
<code>utf16be</code>	UTF-16, big-endian
<code>utf16le</code>	UTF-16, little-endian
<code>utf32</code>	UTF-32, with byte-order mark
<code>utf32be</code>	UTF-32, big-endian
<code>utf32le</code>	UTF-32, little-endian
<code>iso88591, latin1</code>	ISO 8859-1
<code>iso88592, latin2</code>	ISO 8859-2
<code>iso88593, latin3</code>	ISO 8859-3
<code>iso88594, latin4</code>	ISO 8859-4
<code>iso88595</code>	ISO 8859-5
<code>iso88596</code>	ISO 8859-6
<code>iso88597</code>	ISO 8859-7
<code>iso88598</code>	ISO 8859-8
<code>iso88599, latin5</code>	ISO 8859-9
<code>iso885910, latin6</code>	ISO 8859-10
<code>iso885911</code>	ISO 8859-11
<code>iso885913, latin7</code>	ISO 8859-13
<code>iso885914, latin8</code>	ISO 8859-14
<code>iso885915, latin9</code>	ISO 8859-15
<code>iso885916, latin10</code>	ISO 8859-16
<code>clist</code>	comma-list of integers
$\langle empty \rangle$	native (Unicode) string

Table 2: Supported escapings. Non-alphanumeric characters are ignored, and capital letters are lower-cased before searching for the escaping in this list.

$\langle Escaping \rangle$	description
<code>bytes</code> , or <code>empty</code>	arbitrary bytes
<code>hex</code> , <code>hexadecimal</code>	byte = two hexadecimal digits
<code>name</code>	see <code>\pdfescapename</code>
<code>string</code>	see <code>\pdfescapestring</code>
<code>url</code>	encoding used in URLs

## 2 Conversion functions

---

<code>\str_set_convert:Nnnn</code> <code>\str_gset_convert:Nnnn</code>	<code>\str_set_convert:Nnnn &lt;str var&gt; {&lt;string&gt;} {&lt;name 1&gt;} {&lt;name 2&gt;}</code>
---	---

---

This function converts the  $\langle string \rangle$  from the encoding given by  $\langle name 1 \rangle$  to the encoding given by  $\langle name 2 \rangle$ , and stores the result in the  $\langle str var \rangle$ . Each  $\langle name \rangle$  can have the form  $\langle encoding \rangle$  or  $\langle encoding \rangle / \langle escaping \rangle$ , where the possible values of  $\langle encoding \rangle$  and  $\langle escaping \rangle$  are given in Tables 1 and 2, respectively. The default escaping is to input and output bytes directly. The special case of an empty  $\langle name \rangle$  indicates the use of “native” strings, 8-bit for pdfTeX, and Unicode strings for the other two engines.

For example,

`\str_set_convert:Nnnn \l_foo_str { Hello! } { } { utf16/hex }`

results in the variable `\l_foo_str` holding the string `FEFF00480065006C006C006F0021`. This is obtained by converting each character in the (native) string `Hello!` to the UTF-16 encoding, and expressing each byte as a pair of hexadecimal digits. Note the presence of a (big-endian) byte order mark “FEFF”, which can be avoided by specifying the encoding `utf16be/hex`.

An error is raised if the  $\langle string \rangle$  is not valid according to the  $\langle escaping 1 \rangle$  and  $\langle encoding 1 \rangle$ , or if it cannot be reencoded in the  $\langle encoding 2 \rangle$  and  $\langle escaping 2 \rangle$  (for instance, if a character does not exist in the  $\langle encoding 2 \rangle$ ). Erroneous input is replaced by the Unicode replacement character “FFFD”, and characters which cannot be reencoded are replaced by either the replacement character “FFFD” if it exists in the  $\langle encoding 2 \rangle$ , or an encoding-specific replacement character, or the question mark character.

---

<code>\str_set_convert:NnnnTF</code> <code>\str_gset_convert:NnnnTF</code>	<code>\str_set_convert:NnnnTF &lt;str var&gt; {&lt;string&gt;} {&lt;name 1&gt;} {&lt;name 2&gt;} {&lt;true code&gt;} {&lt;false code&gt;}</code>
---	--

---

As `\str_set_convert:Nnnn`, converts the  $\langle string \rangle$  from the encoding given by  $\langle name 1 \rangle$  to the encoding given by  $\langle name 2 \rangle$ , and assigns the result to  $\langle str var \rangle$ . Contrarily to `\str_set_convert:Nnnn`, the conditional variant does not raise errors in case the  $\langle string \rangle$  is not valid according to the  $\langle name 1 \rangle$  encoding, or cannot be expressed in the  $\langle name 2 \rangle$  encoding. Instead, the  $\langle false code \rangle$  is performed.

## 3 Internal string functions

---

<code>\__str_gset_other:Nn</code>	<code>\__str_gset_other:Nn &lt;tl var&gt; {&lt;token list&gt;}</code>
-----------------------------------	---

---

Converts the  $\langle token list \rangle$  to an  $\langle other string \rangle$ , where spaces have category code “other”, and assigns the result to the  $\langle tl var \rangle$ , globally.

---

<code>\__str_hexadecimal_use:NTF</code>	<code>\__str_hexadecimal_use:NTF &lt;token&gt; {&lt;true code&gt;} {&lt;false code&gt;}</code>
---	--

---

If the  $\langle token \rangle$  is a hexadecimal digit (upper case or lower case), its upper-case version is left in the input stream, followed by the  $\langle true code \rangle$ . Otherwise, the  $\langle false code \rangle$  is left in the input stream.

**TeXhackers note:** This function fails on some inputs if the escape character is a hexadecimal digit. We are thus careful to set the escape character to a known (safe) value before using it.

