

## Contents

<b>1</b>	<b>Building strings</b>	<b>2</b>
<b>2</b>	<b>Characters given by their position</b>	<b>3</b>
<b>3</b>	<b>String conditionals</b>	<b>5</b>
<b>4</b>	<b>Encoding functions</b>	<b>5</b>
<b>5</b>	<b>Internal string functions</b>	<b>7</b>
<b>6</b>	<b>Possibilities, and things to do</b>	<b>8</b>
<b>7</b>	<b>l3str implementation</b>	<b>9</b>
7.1	Helpers . . . . .	10
7.1.1	A function unrelated to strings . . . . .	10
7.1.2	Assigning strings . . . . .	10
7.1.3	Variables and constants . . . . .	10
7.1.4	Escaping spaces . . . . .	12
7.2	Characters given by their position . . . . .	13
7.3	String conditionals . . . . .	19
7.4	Conversions . . . . .	21
7.4.1	Producing one byte or character . . . . .	21
7.4.2	Mapping functions for conversions . . . . .	23
7.4.3	Error-reporting during conversion . . . . .	23
7.4.4	Framework for conversions . . . . .	24
7.4.5	Byte unescape and escape . . . . .	29
7.4.6	Native strings . . . . .	30
7.4.7	8-bit encodings . . . . .	32
7.5	Messages . . . . .	35
7.6	Deprecated string functions . . . . .	36
7.7	Escaping definition files . . . . .	36
7.7.1	Unescape methods . . . . .	36
7.7.2	Escape methods . . . . .	41
7.8	Encoding definition files . . . . .	43
7.8.1	UTF-8 support . . . . .	44
7.8.2	UTF-16 support . . . . .	49
7.8.3	UTF-32 support . . . . .	54
7.8.4	ISO 8859 support . . . . .	58

## Index

75

# The `I3str` package: manipulating strings of characters\*

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

Released 2012/02/07

L<sup>A</sup>T<sub>E</sub>X3 provides a set of functions to manipulate token lists as strings of characters, ignoring the category codes of those characters.

String variables are simply specialised token lists, but by convention should be named with the suffix `...str`. Such variables should contain characters with category code 12 (other), except spaces, which have category code 10 (blank space). All the functions in this module first convert their argument to a string for internal processing, and will not treat a token list or the corresponding string representation differently.

Most functions in this module come in three flavours:

- `\str_...:N...`, which expect a token list or string variable as their argument;
- `\str_...:n...`, taking any token list (or string) as an argument;
- `\str_..._ignore_spaces:n...`, which ignores any space encountered during the operation: these functions are faster than those which take care of escaping spaces appropriately;

## 1 Building strings

---

`\c_backslash_str`  
`\c_lbrace_str`  
`\c_rbrace_str`  
`\c_hash_str`  
`\c_tilde_str`  
`\c_percent_str`

---

Constant strings, containing a single character, with category code 12. Any character can be accessed as `\iow_char:N \langle character\rangle`.

---

\*This file describes v3330, last revised 2012/02/07.

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

---

\tl\_to\_str:N  $\star$   
\tl\_to\_str:n  $\star$

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , leaving the resulting tokens in the input stream.

**TeXhackers note:** The string representation of a token list may depend on the category codes in effect when it is evaluated, and the value of the \escapechar: for instance \tl\_to\_str:n {\a} normally produces the three character “backslash”, “lower-case a”, “space”, but it may also produce 1 or 2 characters depending on the escape character, and the category code of a. This impacts almost all functions in the module, which use \tl\_to\_str:n internally.

---

\str\_new:N  
\str\_new:c

Creates a new  $\langle str\ var \rangle$  or raises an error if the name is already taken. The declaration is global. The  $\langle str\ var \rangle$  will initially be empty.

---

\str\_const:Nn  
\str\_const:(Nx|cn|cx)

\str\_const:Nn  $\langle str\ var \rangle$   $\{\langle token\ list \rangle\}$

Creates a new constant  $\langle str\ var \rangle$  or raises an error if the name is already taken. The value of the  $\langle str\ var \rangle$  will be set globally to the  $\langle token\ list \rangle$ , converted to a string.

---

\str\_set:Nn  
\str\_set:(Nx|cn|cx)  
\str\_gset:Nn  
\str\_gset:(Nx|cn|cx)

\str\_set:Nn  $\langle str\ var \rangle$   $\{\langle token\ list \rangle\}$

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , and stores the result in  $\langle str\ var \rangle$ .

---

\str\_put\_left:Nn  
\str\_put\_left:(Nx|cn|cx)  
\str\_gput\_left:Nn  
\str\_gput\_left:(Nx|cn|cx)

\str\_put\_left:Nn  $\langle str\ var \rangle$   $\{\langle token\ list \rangle\}$

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , and prepends the result to  $\langle str\ var \rangle$ . The current contents of the  $\langle str\ var \rangle$  are not automatically converted to a string.

---

\str\_put\_right:Nn  
\str\_put\_right:(Nx|cn|cx)  
\str\_gput\_right:Nn  
\str\_gput\_right:(Nx|cn|cx)

\str\_put\_right:Nn  $\langle str\ var \rangle$   $\{\langle token\ list \rangle\}$

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , and appends the result to  $\langle str\ var \rangle$ . The current contents of the  $\langle str\ var \rangle$  are not automatically converted to a string.

## 2 Characters given by their position

---

\str\_length:N  $\star$   
\str\_length:n  $\star$   
\str\_length\_ignore\_spaces:n  $\star$

Leaves the length of the string representation of  $\langle token\ list \rangle$  in the input stream as an integer denotation. The functions differ in their treatment of spaces. In the case of \str\_length:N and \str\_length:n, all characters including spaces are counted. The \str\_length\_ignore\_spaces:n leaves the number of non-space characters in the input stream.

---

`\str_count_spaces:N` \*      `\str_count_spaces:n {<token list>}`  
`\str_count_spaces:n` \*      Leaves in the input stream the number of space characters in the string representation of  $\langle token\ list \rangle$ , as an integer denotation. Of course, this function has no `_ignore_spaces` variant.

---

`\str_head:N`                  \*      `\str_head:n {<token list>}`  
`\str_head:n`                  \*  
`\str_head_ignore_spaces:n` \*

Converts the  $\langle token\ list \rangle$  into a  $\langle string \rangle$ . The first character in the  $\langle string \rangle$  is then left in the input stream, with category code “other”. The functions differ in their treatment of spaces. In the case of `\str_head:N` and `\str_head:n`, a leading space is returned with category code 10 (blank space). The `\str_head_ignore_spaces:n` function leaves the first non-space character in the input stream. If the  $\langle token\ list \rangle$  is empty (or blank in the case of the `_ignore_spaces` variant), then nothing is left on the input stream.

---

`\str_tail:N`                  \*      `\str_tail:n {<token list>}`  
`\str_tail:n`                  \*  
`\str_tail_ignore_spaces:n` \*

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , removes the first character, and leaves the remaining characters (if any) in the input stream, with category codes 12 and 10 (for spaces). The functions differ in the case where the first character is a space: `\str_tail:N` and `\str_tail:n` will trim only that space, while `\str_tail_ignore_spaces:n` removes the first non-space character and any space before it. If the  $\langle token\ list \rangle$  is empty (or blank in the case of the `_ignore_spaces` variant), then nothing is left on the input stream.

---

`\str_item:Nn`                  \*      `\str_item:nn {<token list>} {<integer expression>}`  
`\str_item:nn`                  \*  
`\str_item_ignore_spaces:nn` \*

Converts the  $\langle token\ list \rangle$  to a  $\langle string \rangle$ , and leaves in the input stream the character in position  $\langle integer\ expression \rangle$  of the  $\langle string \rangle$ . In the case of `\str_item:Nn` and `\str_item:nn`, all characters including spaces are taken into account. The `\str_item_ignore_spaces:nn` function skips spaces in its argument. If the  $\langle integer\ expression \rangle$  is negative, characters are counted from the end of the  $\langle string \rangle$ . Hence,  $-1$  is the right-most character, etc., while  $0$  is the first (left-most) character.

---

<code>\str_substr:Nnn</code> <code>\str_substr:nnn</code> <code>\str_substr_ignore_spaces:nnn</code>	<code>* \str_substr:nnn {\langle token list\rangle} {\langle start index\rangle} {\langle end index\rangle}</code> <code>*</code> <code>*</code>
--	--

---

Converts the  $\langle token\ list\rangle$  to a  $\langle string\rangle$ , and leaves in the input stream the characters from the  $\langle start\ index\rangle$  inclusive to the  $\langle end\ index\rangle$  exclusive. Note that the length of the substring is equal to the difference between the two  $\langle indices\rangle$ . If either of  $\langle start\ index\rangle$  or  $\langle end\ index\rangle$  is negative, then it is incremented by the length of the list. If either of  $\langle start\ index\rangle$  or  $\langle end\ index\rangle$  is empty, it is replaced by the corresponding end-point of the string. Both  $\langle start\ index\rangle$  and  $\langle end\ index\rangle$  count from 0 for the first (left most) character. For instance,

```
\iow_term:x { \str_substr:nnn { abcdef } { 1 } { 4 } }
\iow_term:x { \str_substr:nnn { abcdef } { -4 } { } }
```

will print `bcd` and `cdef`.

### 3 String conditionals

---

<code>\str_if_eq_p:NN</code> <code>\str_if_eq_p:(nn Vn on no nV VV xx)</code> <code>\str_if_eq:NNTF</code> <code>\str_if_eq:(nn Vn on no nV VV xx)TF</code>	<code>* \str_if_eq_p:nn {\langle t1_1\rangle} {\langle t1_2\rangle}</code> <code>* \str_if_eq:nnTF {\langle t1_1\rangle} {\langle t1_2\rangle} {\langle true\ code\rangle} {\langle false\ code\rangle}</code> <code>*</code> <code>*</code>
--	---

---

Compares the two  $\langle token\ lists\rangle$  on a character by character basis, and is `true` if the two lists contain the same characters in the same order. Thus for example

```
\str_if_eq_p:xx { abc } { \tl_to_str:n { abc } }
```

is logically `true`. All versions of these functions are fully expandable (including those involving an `x`-type expansion).

### 4 Encoding functions

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 `TEX` 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>

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

<i>(Encoding)</i>	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
Empty	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.

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

---

```
\str_set_convert:Nnnn
\str_gset_convert:Nnnn
```

```
\str_set_convert:Nnnn <str var> {\<string>} {\<name1>} {\<name2>}
```

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  $\backslash l\_foo\_str$  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  $\langle escaping \rangle$  and  $\langle encoding \rangle$ , or if it cannot be reencoded in the  $\langle encoding \rangle$  and  $\langle escaping \rangle$  (for instance, if a character does not exist in the  $\langle encoding \rangle$ ). 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  $\langle encoding \rangle$ , or an encoding-specific replacement character, or the question mark character.

---

```
\str_set_convert:NnnnTF
\str_gset_convert:NnnnTF
```

```
\str_set_convert:NnnnTF <str var> {\<string>} {\<name1>} {\<name2>} {\<true code>} {\<false code>}
```

As  $\backslash str\_set\_convert:Nnnn$ , 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  $\backslash str\_set\_convert:Nnnn$ , 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.

## 5 Internal string functions

---

```
\c_max_char_int
```

The maximum valid character code, 255 for pdfTEX, and 1114111 for the two other engines.

---

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

---

```
\str_aux_hexadecimal_use:NTF \str_aux_hexadecimal_use:NTF <token> {\<true code>} {\<false code>}
```

If the *<token>* is a hexdecimal digit (upper case or lower case), its upper-case version is left in the input stream, *followed* by the *<true code>*. Otherwise, the *<false code>* is left in the input stream.

**TeXhackers note:** This function will fail if the escape character is a hexdecimal digit, or has a character code in the range [0, 5]. We are thus careful to set the escape character to a known value before using it.

---

```
\str_aux_octal_use:NTF \str_aux_octal_use:NTF <token> {\<true code>} {\<false code>}
```

If the *<token>* is an octal digit, it is left in the input stream, *followed* by the *<true code>*. Otherwise, the *<false code>* is left in the input stream.

**TeXhackers note:** This function will fail if the escape character is an octal digit. We are thus careful to set the escape character to a known value before using it.

---

```
\tl_to_other_str:n ★ \tl_to_other_str:n {<token list>}
```

Converts the *<token list>* to a *<other string>*, where spaces have category code “other”.

**TeXhackers note:** This function can be f-expanded without fear of losing a leading space, since spaces do not have category code 10 in its result. This function takes a time quadratic in the length of the string; `\str_gset_other:Nn` is faster but not expandable.

---

```
\str_gset_other:Nn \str_gset_other:Nn <tl var> {<token list>}
```

Converts the *<token list>* to an *<other string>*, where spaces have category code “other”, and assigns the result to the *<tl var>*, globally.

---

```
\str_if_contains_char:NNT ★ \str_if_contains_char:nNTF {<token list>} <char>  
\str_if_contains_char:(NNTF|nNTF) ★
```

Converts the *<token list>* to a *<string>* and tests whether the *<char>* is present in the *<string>*. Spaces are ignored.

---

```
\str_declare_eight_bit_encoding:nnn \str_declare_eight_bit_encoding:nnn {<name>} {<mapping>}  
 {<missing>}
```

This declares the encoding *<name>* to map bytes to Unicode characters according to the *<mapping>*, and map those bytes which are not mentionned in the *<mapping>* either to the replacement character (if they appear in *<missing>*), or to themselves.

## 6 Possibilities, and things to do

Encoding/escaping-related tasks.

- Describe the internal format in the code comments. Refuse code points in ["D800,"DFFF] in the internal representation?
- Add documentation about each encoding and escaping method, and add examples.
- The `hex` unescaping should raise an error for odd-length strings.
- Decide what bytes should be escaped in the `url` escaping. Perhaps !'()\*-./0123456789\_ are safe, and all other characters should be escaped?
- Automate generation of 8-bit mapping files.
- Change the framework for 8-bit encodings: for decoding from 8-bit to Unicode, use 256 integer registers; for encoding, use a tree-box.
- More encodings (see Heiko's `stringenc`). CESU?
- More escapings: shell escapes, lua escapes, etc?

Other string tasks.

- `\str_if_head_eq:nN`
- `\str_if_numeric/decimal/integer:n`, perhaps in `!3fp`?
- Should `\str_item:Nn` be `\str_char:Nn`?
- Should `\str_substr:Nnn` be `\str_range:Nnn`?
- Introduce `\str_slice:Nnnn` with a third “step” argument? Or should we simply have `\str_slice:Nn {string} {clist}`, where the `{clist}`'s items are either one integer expression, two integer expressions separated by `:`, or three integer expressions separated by `::`, cf. Python's extended slice syntax?

## 7 `!3str` implementation

```

1  (*initex | package)
2  \ProvidesExplPackage
3    {\ExplFileName}{\ExplFileVersion}{\ExplFileDescription}
4  \RequirePackage{l3tl-analysis,l3flag}

```

The following string-related functions are currently defined in `!3kernel`.

- `\str_if_eq:nn[ptF]` and variants,
- `\str_if_eq_return:on`, `\str_if_eq_return:xx`
- `\tl_to_str:n`, `\tl_to_str:N`, `\tl_to_str:c`,
- `\token_to_str:N`, `\cs_to_str:N`
- `\str_head:n`, `\str_head_aux:w`, (copied here)

- `\str_tail:n, \str_tail_aux:w`, (copied here)
- `\str_length_skip_spaces` (deprecated)
- `\str_length_loop:NNNNNNNNNN` (unchanged)

## 7.1 Helpers

### 7.1.1 A function unrelated to strings

`\use_ii_i:nn` A function used to swap its arguments.

```
5  \cs_if_exist:NF \use_ii_i:nn
6   { \cs_new:Npn \use_ii_i:nn #1#2 { #2 #1 } }
(End definition for \use_ii_i:nn. This function is documented on page ??.)
```

### 7.1.2 Assigning strings

`\str_set:Nn` Simply convert the token list inputs to *⟨strings⟩*.

```
7  \tl_map_inline:nn
8   {
9    { set }
10   { gset }
11   { const }
12   { put_left }
13   { gput_left }
14   { put_right }
15   { gput_right }
16 }
17 {
18  \cs_new_protected:cpx { str_ #1 :Nn } ##1##2
19   { \exp_not:c { tl_ #1 :Nx } ##1 { \exp_not:N \tl_to_str:n {##2} } }
20   \exp_args:Nc \cs_generate_variant:Nn { str_ #1 :Nn } { Nx , cn , cx }
21 }
```

(End definition for `\str_set:Nn` and others. These functions are documented on page ??.)

### 7.1.3 Variables and constants

Internal scratch space for some functions.

```
22 \cs_new_protected_nopar:Npn \str_tmp:w { }
23 \tl_new:N \l_str_internal_tl
24 \int_new:N \l_str_internal_int
(End definition for \str_tmp:w. This function is documented on page ??.)
```

The `\g_str_result_tl` variable is used to hold the result of various internal string operations (mostly conversions) which are typically performed in a group. The variable is global so that it remains defined outside the group, to be assigned to a user-provided variable.

```
25 \tl_new:N \g_str_result_tl
```

(End definition for `\g_str_result_t1`. This function is documented on page ??.)

`\c_forty_eight` We declare here some integer values which delimit ranges of ASCII characters of various types. This is mostly used in `l3regex`.

```
26 \int_const:Nn \c_forty_eight { 48 }
27 \int_const:Nn \c_fifty_eight { 58 }
28 \int_const:Nn \c_sixty_five { 65 }
29 \int_const:Nn \c_ninety_one { 91 }
30 \int_const:Nn \c_ninety_seven { 97 }
31 \int_const:Nn \c_one_hundred_twenty_three { 123 }
32 \int_const:Nn \c_one_hundred_twenty_seven { 127 }
```

(End definition for `\c_forty_eight` and others. These functions are documented on page ??.)

`\c_max_char_int` The maximum valid character code is 255 for pdfTEX, and 1114111 for other engines.

```
33 \int_const:Nn \c_max_char_int
34   { \pdftex_if_engine:TF { "FF" } { "10FFFF" } }
```

(End definition for `\c_max_char_int`. This function is documented on page ??.)

`\c_str_replacement_char_int` When converting, invalid bytes are replaced by the Unicode replacement character "FFFD.

```
35 \int_const:Nn \c_str_replacement_char_int { "FFFD" }
```

(End definition for `\c_str_replacement_char_int`. This function is documented on page ??.)

`\c_backslash_str` For all of those strings, `\cs_to_str:N` produces characters with the correct category code.

```
36 \tl_const:Nx \c_backslash_str { \cs_to_str:N \\ }
37 \tl_const:Nx \c_lbrace_str { \cs_to_str:N \{ }
38 \tl_const:Nx \c_rbrace_str { \cs_to_str:N \} }
39 \tl_const:Nx \c_hash_str { \cs_to_str:N \# }
40 \tl_const:Nx \c_tilde_str { \cs_to_str:N \~{} }
41 \tl_const:Nx \c_percent_str { \cs_to_str:N \% }
```

(End definition for `\c_backslash_str` and others. These functions are documented on page ??.)

`\g_str_file_prop` To avoid needing one file per encoding/escaping alias, we keep track of those in a property list.

```
42 \prop_new:N \g_str_file_prop
43 \prop_gput:Nnn \g_str_file_prop { latin1 } { iso88591 }
44 \prop_gput:Nnn \g_str_file_prop { latin2 } { iso88592 }
45 \prop_gput:Nnn \g_str_file_prop { latin3 } { iso88593 }
46 \prop_gput:Nnn \g_str_file_prop { latin4 } { iso88594 }
47 \prop_gput:Nnn \g_str_file_prop { latin5 } { iso88599 }
48 \prop_gput:Nnn \g_str_file_prop { latin6 } { iso885910 }
49 \prop_gput:Nnn \g_str_file_prop { latin7 } { iso885913 }
50 \prop_gput:Nnn \g_str_file_prop { latin8 } { iso885914 }
51 \prop_gput:Nnn \g_str_file_prop { latin9 } { iso885915 }
52 \prop_gput:Nnn \g_str_file_prop { latin10 } { iso885916 }
53 \prop_gput:Nnn \g_str_file_prop { utf16le } { utf16 }
54 \prop_gput:Nnn \g_str_file_prop { utf16be } { utf16 }
```

```

55 \prop_gput:Nnn \g_str_file_prop { utf32le } { utf32 }
56 \prop_gput:Nnn \g_str_file_prop { utf32be } { utf32 }
57 \prop_gput:Nnn \g_str_file_prop { hexadecimal } { hex }
(End definition for \g_str_file_prop. This function is documented on page ??.)
```

- \g\_str\_error\_bool** In conversion functions with a built-in conditional, errors are not reported directly to the user, but the information is collected in this boolean, used at the end to decide on which branch of the conditional to take.

```

58 \bool_new:N \g_str_error_bool
(End definition for \g_str_error_bool. This function is documented on page ??.)
```

- str\_byte** Conversions from one *<encoding>/<escaping>* pair to another are done within x-expanding assignments. Errors are signalled by raising the relevant flag.

```

59 \flag_new:n { str_byte }
60 \flag_new:n { str_error }
(End definition for str_byte and str_error. These functions are documented on page ??.)
```

#### 7.1.4 Escaping spaces

- \tl\_to\_other\_str:n** Replaces all spaces by “other” spaces, after converting the token list to a string via \tl\_to\_str:n. This function is f-expandable, storing the part of the string with escaped spaces between the \q\_mark and \q\_stop markers.

```

61 \group_begin:
62 \char_set_lccode:nn { '\* } { '\ }
63 \char_set_lccode:nn { '\A } { '\A }
64 \tl_to_lowercase:n
{
65   \group_end:
66   \cs_new:Npn \tl_to_other_str:n #1
   {
69     \exp_after:wN \tl_to_other_str_loop:w \tl_to_str:n {#1} ~ %
70     A ~ A ~ A ~ A ~ A ~ A ~ A ~ \q_mark \q_stop
71   }
72   \cs_new:Npn \tl_to_other_str_loop:w
73     #1 ~ #2 ~ #3 ~ #4 ~ #5 ~ #6 ~ #7 ~ #8 ~ #9 \q_stop
74   {
75     \if_meaning:w A #8
76       \tl_to_other_str_end:w
77     \fi:
78     \tl_to_other_str_loop:w
79     #9 #1 * #2 * #3 * #4 * #5 * #6 * #7 * #8 * \q_stop
80   }
81   \cs_new:Npn \tl_to_other_str_end:w \fi: #1 \q_mark #2 * A #3 \q_stop
82   { \fi: #2 }
83 }
```

(End definition for \tl\_to\_other\_str:n. This function is documented on page ??.)

`\str_gset_other:Nn` This function could be done by using `\tl_to_other_str:n` within an x-expansion, but  
`\str_gset_other_loop:w` that would take a time quadratic in the size of the string. Instead, we can “leave the  
`\str_gset_other_end:w` result behind us” in the input stream, to be captured into the expanding assignment.  
 This gives us a linear time.

```

84 \group_begin:
85 \char_set_lccode:nn { ‘* } { ‘\ }
86 \char_set_lccode:nn { ‘A } { ‘\A }
87 \tl_to_lowercase:n
88 {
89   \group_end:
90   \cs_new_protected:Npn \str_gset_other:Nn #1#2
91   {
92     \tl_gset:Nx #1
93     {
94       \exp_after:wN \str_gset_other_loop:w \tl_to_str:n {#2} ~ %
95       A ~ A ~ A ~ A ~ A ~ A ~ A ~ A ~ A ~ \q_stop
96     }
97   }
98   \cs_new:Npn \str_gset_other_loop:w
99   #1 ~ #2 ~ #3 ~ #4 ~ #5 ~ #6 ~ #7 ~ #8 ~ #9 ~
100  {
101    \if_meaning:w A #9
102      \str_gset_other_end:w
103    \fi:
104    #1 * #2 * #3 * #4 * #5 * #6 * #7 * #8 * #9
105    \str_gset_other_loop:w *
106  }
107  \cs_new:Npn \str_gset_other_end:w \fi: #1 * A #2 \q_stop
108  { \fi: #1 }
109 }
```

(End definition for `\str_gset_other:Nn`. This function is documented on page ??.)

## 7.2 Characters given by their position

`\str_count_spaces:N` To speed up this function, we grab 9 spaces in each step. The loop stops when the last  
`\str_count_spaces:n` argument is one of the trailing  $X\langle number\rangle$ , and that  $\langle number\rangle$  is added to the sum of 9  
`\str_count_spaces_loop:wwwwwww` that precedes, to adjust the result.

```

110 \cs_new:Npn \str_count_spaces:N
111   { \exp_args:No \str_count_spaces:n }
112 \cs_new:Npn \str_count_spaces:n #1
113   {
114     \int_eval:n
115     {
116       \exp_after:wN \str_count_spaces_loop:wwwwwwww
117       \tl_to_str:n {#1} ~
118       X 7 ~ X 6 ~ X 5 ~ X 4 ~ X 3 ~ X 2 ~ X 1 ~ X 0 ~ X -1 ~
119       \q_stop
120     }
```

```

121    }
122 \cs_new:Npn \str_count_spaces_loop:wwwwwwww #1~#2~#3~#4~#5~#6~#7~#8~#9~
123 {
124   \if_meaning:w X #9
125     \exp_after:wN \use_none_delimit_by_q_stop:w
126   \fi:
127   \c_nine + \str_count_spaces_loop:wwwwwwww
128 }

```

(End definition for `\str_count_spaces:N`. This function is documented on page ??.)

`\str_length:N`  
`\str_length:n`  
`\str_length_unsafe:n`  
`\str_length_ignore_spaces:n`  
`\str_length_aux:n`  
`\str_length_loop:NNNNNNNN`

To measure the length of a string we could first escape all spaces using `\tl_to_other_-str:N`, then measure the length of this token list. However, this would be quadratic in the length of the string, and we can do better. Namely, add the number of spaces (counted using the functions defined above) to the length ignoring spaces. To measure the length ignoring spaces we use the same technique as for counting spaces: loop, grabbing 9 characters at each step, and end as soon as we reach one of the 9 trailing items. The `_unsafe` variant expects a token list consisting entirely of category code 12 characters.

```

129 \cs_new_nopar:Npn \str_length:N { \exp_args:No \str_length:n }
130 \cs_new:Npn \str_length:n #1
131 {
132   \str_length_aux:n
133   {
134     \str_count_spaces:n {#1}
135     + \exp_after:wN \str_length_loop:NNNNNNNN \tl_to_str:n {#1}
136   }
137 }
138 \cs_new:Npn \str_length_unsafe:n #1
139 {
140   \str_length_aux:n
141   { \str_length_loop:NNNNNNNN #1 }
142 }
143 \cs_new:Npn \str_length_ignore_spaces:n #1
144 {
145   \str_length_aux:n
146   { \exp_after:wN \str_length_loop:NNNNNNNN \tl_to_str:n {#1} }
147 }
148 \cs_new:Npn \str_length_aux:n #1
149 {
150   \int_eval:n
151   {
152     #1
153     { X \c_eight } { X \c_seven } { X \c_six }
154     { X \c_five } { X \c_four } { X \c_three }
155     { X \c_two } { X \c_one } { X \c_zero }
156     \q_stop
157   }
158 }
159 \cs_set:Npn \str_length_loop:NNNNNNNN #1#2#3#4#5#6#7#8#9
160 {

```

```

161   \if_meaning:w X #9
162     \exp_after:wN \use_none_delimit_by_q_stop:w
163   \fi:
164   \c_nine + \str_length_loop:NNNNNNNNN
165 }

```

(End definition for `\str_length:N`. This function is documented on page ??.)

`\str_head:N` The `_ignore_spaces` variant is almost identical to `\tl_head:n`. As usual, `\str_head:N` expands its argument and hands it to `\str_head:n`. To circumvent the fact that TeX skips spaces when grabbing undelimited macro parameters, `\str_head_aux:w` takes an argument delimited by a space. If #1 starts with a non-space character, `\use_i_delimit_by_q_stop:nw` leaves that in the input stream. On the other hand, if #1 starts with a space, the `\str_head_aux:w` takes an empty argument, and the single (braced) space in the definition of `\str_head_aux:w` makes its way to the output. Finally, for an empty argument, the (braced) empty brace group in the definition of `\str_head:n` gives an empty result after passing through `\use_i_delimit_by_q_stop:nw`.

```

166 \cs_new_nopar:Npn \str_head:N { \exp_args:No \str_head:n }
167 \cs_set:Npn \str_head:n #1
168 {
169   \exp_after:wN \str_head_aux:w
170   \tl_to_str:n {#1}
171   { { } } ~ \q_stop
172 }
173 \cs_set:Npn \str_head_aux:w #1 ~ %
174 { \use_i_delimit_by_q_stop:nw #1 { ~ } }
175 \cs_new:Npn \str_head_ignore_spaces:n #1
176 {
177   \exp_after:wN \use_i_delimit_by_q_stop:nw
178   \tl_to_str:n {#1} { } \q_stop
179 }

```

(End definition for `\str_head:N`. This function is documented on page ??.)

`\str_tail:N` As when fetching the head of a string, the `_ignore_spaces` variant is similar to `\tl_tail:n`. The more commonly used `\str_tail:n` function is a little bit more convoluted: hitting the front of the string with `\reverse_if:N \if_charcode:w \scan_stop:` removes the first character (which necessarily makes the test true, since it cannot match `\scan_stop:`). The auxiliary function inserts the required `\fi:` to close the conditional, and leaves the tail of the string in the input string. The details are such that an empty string has an empty tail.

```

180 \cs_new_nopar:Npn \str_tail:N { \exp_args:No \str_tail:n }
181 \cs_set:Npn \str_tail:n #1
182 {
183   \exp_after:wN \str_tail_aux:w
184   \reverse_if:N \if_charcode:w
185   \scan_stop: \tl_to_str:n {#1} X X \q_stop
186 }
187 \cs_set:Npn \str_tail_aux:w #1 X #2 \q_stop { \fi: #1 }
188 \cs_new:Npn \str_tail_ignore_spaces:n #1

```

```

189  {
190    \exp_after:wN \str_tail_aux_ii:w
191    \tl_to_str:n {#1} X X \q_stop
192  }
193 \cs_new:Npn \str_tail_aux_ii:w #1 #2 X #3 \q_stop { #2 }
(End definition for \str_tail:N. This function is documented on page ??.)
```

`\str_skip_c_zero:w` Removes `max(#1,0)` characters from the input stream, and then leaves `\c_zero`. This should be expanded using `\tex_roman numeral:D`. We remove characters 8 at a time until there are at most 8 to remove. Then we do a dirty trick: the `\if_case:w` construction leaves between 0 and 8 times the `\or:` control sequence, and those `\or:` become arguments of `\str_skip_end_ii:NNNNNNNN`. If the number of characters to remove is 6, say, then there are two `\or:` left, and the 8 arguments of `\str_skip_end_ii:NNNNNNNN` are the two `\or:`, and 6 characters from the input stream, exactly what we wanted to remove. Then close the `\if_case:w` conditional with `\fi:`, and stop the initial expansion with `\c_zero` (see places where `\str_skip_c_zero:w` is called).

```

194 \cs_new:Npn \str_skip_c_zero:w #1;
195  {
196    \if_num:w \int_eval:w #1 > \c_eight
197      \exp_after:wN \str_skip_aux:NNNNNNNN
198    \else:
199      \exp_after:wN \str_skip_end:w
200      \int_use:N \int_eval:w
201    \fi:
202    #1 ;
203  }
204 \cs_new:Npn \str_skip_aux:NNNNNNNN #1; #2#3#4#5#6#7#8#9
205  { \exp_after:wN \str_skip_c_zero:w \int_use:N \int_eval:w #1 - \c_eight ; }
206 \cs_new:Npn \str_skip_end:w #1 ;
207  {
208    \exp_after:wN \str_skip_end_ii:NNNNNNNN
209    \if_case:w \if_num:w #1 > \c_zero #1 \else: 0 \fi: \exp_stop_f:
210    \or: \or: \or: \or: \or: \or: \or: \or:
211  }
212 \cs_new:Npn \str_skip_end_ii:NNNNNNNN #1#2#3#4#5#6#7#8 { \fi: \c_zero }
(End definition for \str_skip_c_zero:w. This function is documented on page ??.)
```

`\str_collect_delimit_by_q_stop:w` Collects `max(#1,0)` characters, and removes everything else until `\q_stop`. This is somewhat similar to `\str_skip_c_zero:w`, but this time we can only grab 7 characters at a time. At the end, we use an `\if_case:w` trick again, so that the 8 first arguments of `\str_collect_end_ii:nnnnnnnnw` are some `\or:`, followed by an `\if:`, followed by `#1` characters from the input stream. Simply leaving this in the input stream will close the conditional properly and the `\or:` disappear.

```

213 \cs_new:Npn \str_collect_delimit_by_q_stop:w #1;
214  {
215    \exp_after:wN \str_collect_aux:wn
216    \int_use:N \int_eval:w #1 ;
217  }
```

```

218    }
219 \cs_new:Npn \str_collect_aux:wn #1 ;
220 {
221   \if_num:w #1 > \c_seven
222     \exp_after:wN \str_collect_aux:wnNNNNNN
223   \else:
224     \exp_after:wN \str_collect_end:wn
225   \fi:
226   #1 ;
227 }
228 \cs_new:Npn \str_collect_aux:wnNNNNNN #1; #2 #3#4#5#6#7#8#9
229 {
230   \exp_after:wN \str_collect_aux:wn
231   \int_use:N \int_eval:w #1 - \c_seven ;
232   { #2 #3#4#5#6#7#8#9 }
233 }
234 \cs_new:Npn \str_collect_end:wn #1 ;
235 {
236   \exp_after:wN \str_collect_end_ii:nnnnnnnnw
237   \if_case:w \if_num:w #1 > \c_zero #1 \else: 0 \fi: \exp_stop_f:
238   \or: \or: \or: \or: \or: \or: \or: \or: \fi:
239 }
240 \cs_new:Npn \str_collect_end_ii:nnnnnnnnw #1#2#3#4#5#6#7#8 #9 \q_stop
241 { #1#2#3#4#5#6#7#8 }
(End definition for \str_collect_delimit_by_q_stop:w. This function is documented on page ??.)
```

\str\_item:Nn This is mostly shuffling arguments around to avoid measuring the length of the string more than once, and make sure that the parameters given to \str\_skip\_c\_zero:w are necessarily within the bounds of the length of the string. The \_ignore\_spaces function cheats a little bit in that it doesn't hand to \str\_item\_unsafe:nn an "other string".\str\_item:nn This is alright, as everything else is done with undelimited arguments.

```

242 \cs_new_nopar:Npn \str_item:Nn { \exp_args:No \str_item:nn }
243 \cs_new:Npn \str_item:nn #
244 {
245   \exp_last_unbraced:Nf \token_to_str:N
246   \exp_args:Nf \str_item_unsafe:nn { \tl_to_other_str:n {#1} }
247 }
248 \cs_new:Npn \str_item_ignore_spaces:nn #1
249 { \exp_args:No \str_item_unsafe:nn { \tl_to_str:n {#1} } }
250 \cs_new:Npn \str_item_unsafe:nn #1#2
251 {
252   \exp_after:wN \str_item_aux:ww
253   \int_use:N \int_eval:w #2 \exp_after:wN ;
254   \int_value:w \str_length_unsafe:n {#1} ;
255   #1
256   \q_stop
257 }
258 \cs_new:Npn \str_item_aux:ww #1; #2;
259 {
```

```
260 \int_compare:nNnTF {#1} < \c_zero
261 {
262     \int_compare:nNnTF {#1} < {-#2}
263     { \use_none_delimit_by_q_stop:w }
264     {
265         \exp_after:wN \use_i_delimit_by_q_stop:nw
266         \tex_roman numeral:D \str_skip_c_zero:w #1 + #2 ;
267     }
268 }
269 {
270     \int_compare:nNnTF {#1} < {#2}
271     {
272         \exp_after:wN \use_i_delimit_by_q_stop:nw
273         \tex_roman numeral:D \str_skip_c_zero:w #1 ;
274     }
275     { \use_none_delimit_by_q_stop:w }
276 }
```

(End definition for \str\_item:Nn. This function is documented on page ??.)

`\str_substr:Nnn` Sanitize the string. Then evaluate the arguments, replacing them by `\c_zero` or `\c_max_int` if they are empty. This is done by using the construction

```
\str_substr_ignore_spaces:nnn \int_use:N \int_eval:w #2 \c_zero \c_zero ;  
\str_substr_unsafe:nnn
```

`\str_substr_aux:nN` which expands to the value of #2, followed by `\c_zero \c_zero` ; if #2 is an expression, and expands to `0\c_zero` ; otherwise. The same is done to the end-point of the range. Then limit the range to be at most the length of the string (this avoids needing to check for the end of the string when grabbing characters). Afterwards, skip characters, then keep some more, and finally drop the end of the string.  
`\str_substr_aux:www`  
`\str_substr_aux:nw`  
`\str_substr_normalize_range:mm`

```

278 \cs_new_nopar:Npn \str_substr:Nnn { \exp_args:No \str_substr:nnn }
279 \cs_new:Npn \str_substr:nnn #1#2#3
280 {
281     \exp_args:Nf \tl_to_str:n
282     {
283         \exp_args:Nf \str_substr_unsafe:nnn
284         { \tl_to_other_str:n {#1} } {#2} {#3}
285     }
286 }
287 \cs_new:Npn \str_substr_ignore_spaces:nnn #1
288 { \exp_args:No \str_substr_unsafe:nnn { \tl_to_str:n {#1} } }
289 \cs_new:Npn \str_substr_unsafe:nnn #1#2#3
290 {
291     \exp_after:wN \str_substr_aux:www
292     \int_value:w \str_length_unsafe:n {#1} \exp_after:wN ;
293     \int_use:N \int_eval:w #2 + \c_zero \exp_after:wN ;
294     \int_use:N \int_eval:w
295     \exp_args:Nf \str_substr_aux:nN {#3} \c_max_int ;
296     #1
297     \q_stop

```

```

298    }
299 \cs_new:Npn \str_substr_aux:nN #1 #2
300   { \tl_if_empty:nTF {#1} {#2} {#1} }
301 \cs_new:Npn \str_substr_aux:www #1; #2; #3;
302   {
303     \exp_args:Nf \str_substr_aux:nnw
304       { \str_substr_normalize_range:nn {#2} {#1} }
305       { \str_substr_normalize_range:nn {#3} {#1} }
306   }
307 \cs_new:Npn \str_substr_aux:nnw #1#2
308   {
309     \exp_after:wN \str_collect_delimit_by_q_stop:w
310     \int_use:N \int_eval:w #2 - #1 \exp_after:wN ;
311     \tex_roman numeral:D \str_skip_c_zero:w #1 ;
312   }
313 \cs_new:Npn \str_substr_normalize_range:nn #1#2
314   {
315     \int_eval:n
316     {
317       \if_num:w #1 < \c_zero
318         \if_num:w #1 < - #2 \exp_stop_f:
319           \c_zero
320         \else:
321           #1 + #2
322         \fi:
323       \else:
324         \if_num:w #1 < #2 \exp_stop_f:
325           #1
326         \else:
327           #2
328         \fi:
329       \fi:
330     }
331   }

```

(End definition for `\str_substr:Nnn`. This function is documented on page ??.)

### 7.3 String conditionals

<code>\str_if_eq:NN</code> <code>\str_if_eq:nn</code> <code>\str_if_eq:xx</code>	The <code>nn</code> and <code>xx</code> variants are already defined in <code>l3basics</code> . Note that <code>\str_if_eq:NN</code> is different from <code>\tl_if_eq:NN</code> because it needs to ignore category codes.
--	---

```

332 \prg_new_conditional:Npnn \str_if_eq:NN #1#2 { p , TF , T , F }
333   {
334     \if_int_compare:w \pdfstrcmp:D { \tl_to_str:N #1 } { \tl_to_str:N #2 }
335     = \c_zero \prg_return_true: \else: \prg_return_false: \fi:
336   }

```

(End definition for `\str_if_eq:NN`. This function is documented on page ??.)

<code>\str_if_contains_char:NNT</code> <code>\str_if_contains_char:NNTF</code> <code>\str_if_contains_char:nNTF</code> <code>\str_if_contains_char_aux:NN</code> <code>\str_if_contains_char_true:</code>	Loop over the characters of the string, comparing character codes. The loop is broken if character codes match. Otherwise we return “false”.
---	--

```

337 \prg_new_if:NNN \str_if_contains_char:N #1#2 { T , TF }
338 {
339   \exp_after:wN \str_if_contains_char_aux:NN \exp_after:wN #2
340   \tl_to_str:n \exp_after:wN {#1} { \prg_map_break:n { ? \fi: } }
341   \prg_break_point:n { }
342   \prg_return_false:
343 }
344 \prg_new_if:NNN \str_if_contains_char:nN #1#2 { TF }
345 {
346   \exp_after:wN \str_if_contains_char_aux:NN \exp_after:wN #2
347   \tl_to_str:n {#1} { \prg_map_break:n { ? \fi: } }
348   \prg_break_point:n { }
349   \prg_return_false:
350 }
351 \cs_new:Npn \str_if_contains_char_aux:NN #1#2
352 {
353   \if_charcode:w #1 #2
354   \exp_after:wN \str_if_contains_char_true:
355   \fi:
356   \str_if_contains_char_aux:NN #1
357 }
358 \cs_new_nopar:Npn \str_if_contains_char_true:
359 { \prg_map_break:n { \prg_return_true: \use_none:n } }
(End definition for \str_if_contains_char:NNT and \str_if_contains_char:NNTF. These functions are documented on page ??.)
```

`\str_aux_octal_use:NTF` TeX dutifully detects octal digits for us: if #1 is an octal digit, then the right-hand side of the comparison is '1#1, greater than 1. Otherwise, the right-hand side stops as '1, and the conditional takes the `false` branch.

```

360 \prg_new_if:NNN \str_aux_octal_use:N #1 { TF }
361 {
362   \if_num:w \c_one < '1 \token_to_str:N #1 \exp_stop_f:
363   #1 \prg_return_true:
364   \else:
365   \prg_return_false:
366   \fi:
367 }
(End definition for \str_aux_octal_use:NTF. This function is documented on page 8.)
```

`\str_aux_hexadecimal_use:NTF` TeX detects uppercase hexadecimal digits for us (see `\str_aux_octal_use:NTF`), but not the lowercase letters, which we need to detect and replace by their uppercase counterpart.

```

368 \prg_new_if:NNN \str_aux_hexadecimal_use:N #1 { TF }
369 {
370   \if_num:w \c_two < "1 \token_to_str:N #1 \exp_stop_f:
371   #1 \prg_return_true:
372   \else:
373   \if_case:w \int_eval:w
374     \exp_after:wN ' \token_to_str:N #1 - 'a
375     \int_eval_end:
```

```

376      A
377      \or: B
378      \or: C
379      \or: D
380      \or: E
381      \or: F
382      \else:
383          \prg_return_false:
384          \exp_after:wN \use_none:n
385      \fi:
386      \prg_return_true:
387  \fi:
388 }

```

(End definition for `\str_aux_hexadecimal_use:NTF`. This function is documented on page 8.)

## 7.4 Conversions

### 7.4.1 Producing one byte or character

`\c_str_byte_0_tl` For each integer  $N$  in the range  $[0, 255]$ , we create a constant token list which holds three character tokens with category code other: the character with character code  $N$ , followed by the representation of  $N$  as two hexadecimal digits. The value  $-1$  is given a default token list which ensures that later functions give an empty result for the input  $-1$ . Simultaneously, we build a list of all bytes (from which we remove the null byte) in `\c_str_positive_bytes_tl`.

```

389 \group_begin:
390   \tl_clear:N \l_str_internal_tl
391   \char_set_catcode_other:n { \c_zero }
392   \tl_gset:Nx \g_str_result_tl { \tl_to_str:n { 0123456789ABCDEF } }
393   \exp_args:No \tl_map_inline:nn { \g_str_result_tl " }
394   { \char_set_lccode:nn {'#1} { \c_zero } }
395   \tl_map_inline:Nn \g_str_result_tl
396   {
397     \tl_map_inline:Nn \g_str_result_tl
398     {
399       \char_set_lccode:nn { \c_zero } {"#1##1}
400       \tl_to_lowercase:n
401       {
402         \tl_const:cx
403         { c_str_byte_ \int_eval:n {"#1##1} _tl }
404         { ^@ #1 ##1 }
405         \tl_put_right:Nn \l_str_internal_tl { ^@ }
406       }
407     }
408   }
409   \tl_const:Nx \c_str_positive_bytes_tl
410   { \exp_after:wN \use_none:n \l_str_internal_tl }
411 \group_end:
412 \tl_const:cn { c_str_byte_-1_tl } { { } \use_none:n { } }

```

*(End definition for \c\_str\_byte\_0\_t1, \c\_str\_byte\_1\_t1, and \c\_str\_byte\_255\_t1. These functions are documented on page ??.)*

\str\_output\_byte:n    Those functions must be used carefully: feeding them a value outside the range  $[-1, 255]$  will attempt to use the undefined token list variable \c\_str\_byte\_{number}\_t1. Assuming that the argument is in the right range, we expand the corresponding token list, and pick either the byte (first token) or the hexadecimal representations (second and third tokens). The value  $-1$  produces an empty result in both cases.

```

413 \cs_new:Npn \str_output_byte:n #1
414   { \str_output_byte:w #1 \str_output_end: }
415 \cs_new_nopar:Npn \str_output_byte:w
416   {
417     \exp_after:wN \exp_after:wN
418     \exp_after:wN \use_i:nnn
419     \cs:w c_str_byte_ \int_use:N \int_eval:w
420   }
421 \cs_new:Npn \str_output_hexadecimal:n #1
422   { \str_output_hexadecimal:w #1 \str_output_end: }
423 \cs_new_nopar:Npn \str_output_hexadecimal:w
424   {
425     \exp_after:wN \exp_after:wN
426     \exp_after:wN \use_none:n
427     \cs:w c_str_byte_ \int_use:N \int_eval:w
428   }
429 \cs_new_nopar:Npn \str_output_end:
430   { \int_eval_end: _tl \cs_end: }

```

*(End definition for \str\_output\_byte:n. This function is documented on page ??.)*

\str\_output\_byte\_pair\_be:n    Convert a number in the range  $[0, 65535]$  to a pair of bytes, either big-endian or little-endian.

```

431 \cs_new:Npn \str_output_byte_pair_be:n #1
432   {
433     \exp_args:Nf \str_output_byte_pair_aux:nnN
434     { \int_div_truncate:nn { #1 } { "100" } {#1} \use:nn
435   }
436 \cs_new:Npn \str_output_byte_pair_le:n #1
437   {
438     \exp_args:Nf \str_output_byte_pair_aux:nnN
439     { \int_div_truncate:nn { #1 } { "100" } {#1} \use_i:i:nn
440   }
441 \cs_new:Npn \str_output_byte_pair_aux:nnN #1#2#3
442   {
443     #3
444     { \str_output_byte:n { #1 } }
445     { \str_output_byte:n { #2 - #1 * "100" } }
446   }

```

*(End definition for \str\_output\_byte\_pair\_be:n. This function is documented on page ??.)*

### 7.4.2 Mapping functions for conversions

\str\_convert\_gmap:N This maps the function #1 over all characters in \g\_str\_result\_tl, which should be a byte string in most cases, sometimes a native string.

```

447 \cs_new_protected:Npn \str_convert_gmap:N #1
448 {
449     \tl_gset:Nx \g_str_result_tl
450     {
451         \exp_after:wN \str_convert_gmap_loop:NN
452         \exp_after:wN #1
453         \g_str_result_tl { ? \prg_map_break: }
454         \prg_break_point:n { }
455     }
456 }
457 \cs_new:Npn \str_convert_gmap_loop:NN #1#2
458 {
459     \use_none:n #2
460     #1#2
461     \str_convert_gmap_loop:NN #1
462 }
```

(End definition for \str\_convert\_gmap:N. This function is documented on page ??.)

\str\_convert\_gmap\_internal:N This maps the function #1 over all character codes in \g\_str\_result\_tl, which must be in the internal representation.

```

463 \cs_new_protected:Npn \str_convert_gmap_internal:N #1
464 {
465     \tl_gset:Nx \g_str_result_tl
466     {
467         \exp_after:wN \str_convert_gmap_internal_loop:Nww
468         \exp_after:wN #1
469         \g_str_result_tl \s_tl \q_stop \prg_map_break: \s_tl
470         \prg_break_point:n { }
471     }
472 }
473 \cs_new:Npn \str_convert_gmap_internal_loop:Nww #1 #2 \s_tl #3 \s_tl
474 {
475     \use_none_delimit_by_q_stop:w #3 \q_stop
476     #1 {#3}
477     \str_convert_gmap_internal_loop:Nww #1
478 }
```

(End definition for \str\_convert\_gmap\_internal:N. This function is documented on page ??.)

### 7.4.3 Error-reporting during conversion

\str\_if\_flag\_error:nnx \str\_if\_flag\_no\_error:nnx When converting using the function \str\_set\_convert:Nnnn, errors should be reported to the user after each step in the conversion. Errors are signalled by raising some flag (typically `str_error`), so here we test that flag: if it is raised, give the user an error,

otherwise remove the arguments. On the other hand, in the conditional functions `\str_set_convert:NnnnTF`, errors should be suppressed. This is done by changing `\str_if_flag_error:nnx` into `\str_if_flag_no_error:nnx` locally.

```

479 \cs_new_protected:Npn \str_if_flag_error:nnx #1
480   {
481     \flag_if_raised:nTF {#1}
482     { \msg_kernel_error:nnx { str } }
483     { \use_none:nn }
484   }
485 \cs_new_protected:Npn \str_if_flag_no_error:nnx #1#2#3
486   { \flag_if_raised:nT {#1} { \bool_gset_true:N \g_str_error_bool } }
(End definition for \str_if_flag_error:nnx. This function is documented on page ??.)
```

`\str_if_flag_times:nT` At the end of each conversion step, we raise all relevant errors as one error message, built on the fly. The height of each flag indicates how many times a given error was encountered. This function prints #2 followed by the number of occurrences of an error if it occurred, nothing otherwise.

```

487 \cs_new_protected:Npn \str_if_flag_times:nT #1#2
488   { \flag_if_raised:nT {#1} { #2~(x \flag_height:n {#1} ) } }
(End definition for \str_if_flag_times:nT. This function is documented on page ??.)
```

#### 7.4.4 Framework for conversions

Most functions in this module expect to be working with “native” strings. Strings can also be stored as bytes, in one of many encodings, for instance UTF8. The bytes themselves can be expressed in various ways in terms of TeX tokens, for instance as pairs of hexadecimal digits. The questions of going from arbitrary Unicode code points to bytes, and from bytes to tokens are mostly independent.

Conversions are done in four steps:

- “unescape” produces a string of bytes;
- “decode” takes in a string of bytes, and converts it to a list of Unicode characters in an internal representation, with items of the form

*⟨bytes⟩ \s\_tl ⟨Unicode code point⟩ \s\_tl*

where we have collected the *⟨bytes⟩* which combined to form this particular Unicode character, and the *⟨Unicode code point⟩* is in the range [0, "10FFFF].

- “encode” encodes the internal list of code points as a byte string in the new encoding;
- “escape” escapes bytes as requested.

The process is modified in case one of the encoding is empty (or the conversion function has been set equal to the empty encoding because it was not found): then the unescape or escape step is ignored, and the decode or encode steps work on tokens instead of bytes. Otherwise, each step must ensure that it passes a correct byte string or internal string to the next step.

\str\_set\_convert:Nnnn  
\str\_gset\_convert:Nnnn  
\str\_set\_convert:Nnnn  
\str\_gset\_convert:Nnnn  
\str\_convert\_aux\_i:nNNnnn

The input string is stored in \g\_str\_result\_tl, then we: unescape and decode; encode and escape; exit the group and store the result in the user's variable. The various conversion functions all act on \g\_str\_result\_tl. Errors are silenced for the conditional functions by redefining \str\_if\_flag\_error:nmx locally.

```

489 \cs_new_protected_nopar:Npn \str_set_convert:Nnnn
490   { \str_convert_aux_i:nNNnnn { } \tl_set_eq:NN }
491 \cs_new_protected_nopar:Npn \str_gset_convert:Nnnn
492   { \str_convert_aux_i:nNNnnn { } \tl_gset_eq:NN }
493 \prg_new_protected_conditional:Npnn
494   \str_set_convert:Nnnn #1#2#3#4 { T , F , TF }
495   {
496     \bool_gset_false:N \g_str_error_bool
497     \str_convert_aux_i:nNNnnn
498     { \cs_set_eq:NN \str_if_flag_error:nmx \str_if_flag_no_error:nmx }
499     \tl_set_eq:NN #1 {#2} {#3} {#4}
500     \bool_if:NTF \g_str_error_bool \prg_return_false: \prg_return_true:
501   }
502 \prg_new_protected_conditional:Npnn
503   \str_gset_convert:Nnnn #1#2#3#4 { T , F , TF }
504   {
505     \bool_gset_false:N \g_str_error_bool
506     \str_convert_aux_i:nNNnnn
507     { \cs_set_eq:NN \str_if_flag_error:nmx \str_if_flag_no_error:nmx }
508     \tl_gset_eq:NN #1 {#2} {#3} {#4}
509     \bool_if:NTF \g_str_error_bool \prg_return_false: \prg_return_true:
510   }
511 \cs_new_protected:Npn \str_convert_aux_i:nNNnnn #1#2#3#4#5#6
512   {
513     \group_begin:
514       #1
515       \str_gset_other:Nn \g_str_result_tl {#4}
516       \exp_after:wN \str_convert_aux_i:wwwnn
517         \tl_to_str:n {#5} /// \q_stop
518         { decode } { unescape }
519         \prg_do_nothing:
520         \str_convert_decode_:
521       \exp_after:wN \str_convert_aux_i:wwwnn
522         \tl_to_str:n {#6} /// \q_stop
523         { encode } { escape }
524         \use_i:i:nn
525         \str_convert_encode_:
526       \group_end:
527       #2 #3 \g_str_result_tl
528   }

```

(End definition for \str\_set\_convert:Nnnn and \str\_gset\_convert:Nnnn. These functions are documented on page ??.)

\str\_convert\_aux\_i:wwwnn  
\str\_convert\_aux\_iii>NNnNN

The task of \str\_convert\_aux\_i:wwwnn is to split  $\langle encoding \rangle / \langle escaping \rangle$  pairs into their components, #1 and #2. Calls to \str\_convert\_aux\_iv:nnn ensure that the cor-

responding conversion functions are defined. The third auxiliary does the main work.

- #1 is the encoding conversion function;
- #2 is the escaping function;
- #3 is the escaping name for use in an error message;
- #4 is `\prg_do_nothing`: for unescaping/decoding, and `\use_ii_i:nn` for encoding/escaping;
- #5 is the default encoding function (either “decode” or “encode”), for which there should be no escaping.

Let us ignore the native encoding for a second. In the unescaping/decoding phase, we want to do #2#1 in this order, and in the encoding/escaping phase, the order should be reversed: #4#2#1 does exactly that. If one of the encodings is the default (native), then the escaping should be ignored, with an error if any was given, and only the encoding, #1, should be performed.

```

529 \cs_new_protected:Npn \str_convert_aux_ii:wwnn
530   #1 / #2 // #3 \q_stop #4#5
531   {
532     \str_convert_aux_iv:nnn {enc} {#4} {#1}
533     \str_convert_aux_iv:nnn {esc} {#5} {#2}
534     \exp_args:Ncc \str_convert_aux_iii:NNnNN
535     { str_convert_#4_#1: } { str_convert_#5_#2: } {#2}
536   }
537 \cs_new_protected:Npn \str_convert_aux_iii:NNnNN #1#2#3#4#5
538   {
539     \if_meaning:w #1 #5
540       \tl_if_empty:nF {#3}
541         { \msg_kernel_error:nnx { str } { native-escaping } {#3} }
542       #1
543     \else:
544       #4 #2 #1
545     \fi:
546   }

```

(End definition for `\str_convert_aux_ii:wwnn`. This function is documented on page ??.)

`\str_convert_aux_iv:nnn`  
`\str_convert_aux_v:nnnn`

The arguments of `\str_convert_aux_iv:nnn` are: `enc` or `esc`, used to build filenames, the type of the conversion (unescape, decode, encode, escape), and the encoding or escaping name. If the function is already defined, no need to do anything. Otherwise, filter out all non-alphanumerics in the name, and lowercase it. Feed that, and the same three arguments, to `\str_convert_aux_v:nnnn`. The task is then to make sure that the conversion function `#3_#1` corresponding to the type `#3` and filtered name `#1` is defined, then set our initial conversion function `#3_#4` equal to that.

How do we get the `#3_#1` conversion to be defined if it isn’t? Two main cases.

First, if `#1` is a key in `\g_str_file_prop`, then the value `\l_str_internal_t1` tells us what file to load. Loading is skipped if the file was already read, *i.e.*, if the conversion

command based on `\l_str_internal_t1` already exists. Otherwise, try to load the file; if that fails, there is an error, use the default empty name instead.

Second, `#1` may be absent from the property list. The `\cs_if_exist:cF` test is automatically false, and we search for a file defining the encoding or escaping `#1` (this should allow third-party `.def` files). If the file is not found, there is an error, use the default empty name instead.

In all cases, the conversion based on `\l_str_internal_t1` is defined, so we can set the `#3_#1` function equal to that. In some cases (*e.g.*, `utf16be`), the `#3_#1` function is actually defined within the file we just loaded, and it is different from the `\l_str_internal_t1`-based function: we mustn't clobber that different definition.

```

547 \cs_new_protected:Npn \str_convert_aux_iv:nnn #1#2#3
548 {
549     \cs_if_exist:cF { str_convert_#2_#3: }
550     {
551         \exp_args:Nx \str_convert_aux_v:nnnn
552             { \str_convert_lowercase_alphanum:n {#3} }
553             {#1} {#2} {#3}
554     }
555 }
556 \cs_new_protected:Npn \str_convert_aux_v:nnnn #1#2#3#4
557 {
558     \cs_if_exist:cF { str_convert_#3_#1: }
559     {
560         \prop_get:NnNF \g_str_file_prop {#1} \l_str_internal_t1
561             { \tl_set:Nn \l_str_internal_t1 {#1} }
562         \cs_if_exist:cF { str_convert_#3_ \l_str_internal_t1 : }
563         {
564             \file_if_exist:nTF { 13str-#2- \l_str_internal_t1 .def }
565             {
566                 \group_begin:
567                     \str_load_catcodes:
568                     \file_input:n { 13str-#2- \l_str_internal_t1 .def }
569                 \group_end:
570             }
571             {
572                 \tl_clear:N \l_str_internal_t1
573                 \msg_kernel_error:nxxx { str } { unknown-#2 } { #4 } {#1}
574             }
575         }
576         \cs_if_exist:cF { str_convert_#3_#1: }
577         {
578             \cs_gset_eq:cc { str_convert_#3_#1: }
579             { str_convert_#3_ \l_str_internal_t1 : }
580         }
581     }
582     \cs_gset_eq:cc { str_convert_#3_#4: } { str_convert_#3_#1: }
583 }
```

(End definition for `\str_convert_aux_iv:nnn`. This function is documented on page ??.)

\str\_convert\_lowercase\_alphanum:n This function keeps only letters and digits, with upper case letters converted to lower case.  
\str\_convert\_lowercase\_alphanum\_loop:N

```

584 \cs_new:Npn \str_convert_lowercase_alphanum:n #1
585   {
586     \exp_after:wN \str_convert_lowercase_alphanum_loop:N
587     \tl_to_str:n {#1} { ? \prg_map_break: }
588     \prg_break_point:n { }
589   }
590 \cs_new:Npn \str_convert_lowercase_alphanum_loop:N #1
591   {
592     \use_none:n #1
593     \if_num:w '#1 < \c_ninety_one
594       \if_num:w '#1 < \c_sixty_five
595         \if_num:w \c_one < 1#1 \exp_stop_f:
596           #1
597         \fi:
598       \else:
599         \str_output_byte:n { '#1 + \c_thirty_two }
600       \fi:
601     \else:
602       \if_num:w '#1 < \c_one_hundred_twenty_three
603         \if_num:w '#1 < \c_ninety_seven
604           \else:
605             #1
606             \fi:
607           \fi:
608         \fi:
609       \str_convert_lowercase_alphanum_loop:N
610   }
(End definition for \str_convert_lowercase_alphanum:n. This function is documented on page ??.)
```

\str\_load\_catcodes: Since encoding files may be loaded at arbitrary places in a TeX document, including within verbatim mode, we set the catcodes of all characters appearing in any encoding definition file.

```

611 \cs_new_protected:Npn \str_load_catcodes:
612   {
613     \char_set_catcode_escape:N \\%
614     \char_set_catcode_group_begin:N \{
615     \char_set_catcode_group_end:N \}
616     \char_set_catcode_math_toggle:N \$%
617     \char_set_catcode_alignment:N \&%
618     \char_set_catcode_parameter:N \#
619     \char_set_catcode_math_superscript:N \^%
620     \char_set_catcode_ignore:N \ %
621     \char_set_catcode_space:N \~%
622     \tl_map_function:nN { abcdefghijklmnopqrstuvwxyz_ : ABCDEFILNPSTUX }
623       \char_set_catcode_letter:N
624     \tl_map_function:nN { 0123456789"?'*+-.(,) '!/<>[] ;= }
```

```

625     \char_set_catcode_other:N
626     \char_set_catcode_comment:N \%
627     \int_set:Nn \tex_endlinechar:D {32}
628 }
(End definition for \str_load_catcodes:. This function is documented on page ??.)
```

#### 7.4.5 Byte unescape and escape

Strings of bytes may need to be stored in auxiliary files in safe “escaping” formats. Each such escaping is only loaded as needed. By default, on input any non-byte is filtered out, while the output simply consists in letting bytes through.

\str\_filter\_bytes:n In the case of pdftEX, every character is a byte. For Unicode-aware engines, test the character code; non-bytes cause us to raise the flag **str\_byte**. Spaces have already been given the correct category code when this function is called.

```

629 \pdftex_if_engine:TF
630 { \cs_new_eq:NN \str_filter_bytes:n \use:n }
631 {
632     \cs_new:Npn \str_filter_bytes:n #1
633     {
634         \str_filter_bytes_aux:N #1
635         { ? \prg_map_break: }
636         \prg_break_point:n { }
637     }
638     \cs_new:Npn \str_filter_bytes_aux:N #1
639     {
640         \use_none:n #1
641         \if_num:w '#1 < 256 \exp_stop_f:
642             #1
643         \else:
644             \flag_raise:n { str_byte }
645         \fi:
646         \str_filter_bytes_aux:N
647     }
648 }
```

(End definition for \str\_filter\_bytes:n. This function is documented on page ??.)

\str\_convert\_unescape\_ : The simplest unescaping method removes non-bytes from \g\_str\_result\_tl.

```

649 \pdftex_if_engine:TF
650 { \cs_new_protected_nopar:Npn \str_convert_unescape_: { } }
651 {
652     \cs_new_protected_nopar:Npn \str_convert_unescape_:
653     {
654         \flag_clear:n { str_byte }
655         \tl_gset:Nx \g_str_result_tl
656         { \exp_args:No \str_filter_bytes:n \g_str_result_tl }
657         \str_if_flag_error:nnx { str_byte } { non-byte } { bytes }
658     }
659 }
```

```

660 \cs_new_eq:NN \str_convert_unescape_bytes: \str_convert_unescape_:
(End definition for \str_convert_unescape_:. This function is documented on page ??.)

```

\str\_convert\_escape\_: The simplest form of escape leaves the bytes from the previous step of the conversion unchanged.

```

661 \cs_new_protected_nopar:Npn \str_convert_escape_:
662 \cs_new_eq:NN \str_convert_escape_bytes: \str_convert_escape_:
(End definition for \str_convert_escape_:. This function is documented on page ??.)

```

#### 7.4.6 Native strings

\str\_convert\_decode\_: Convert each character to its character code, one at a time.

```

\str_decode_native_char:N
663 \cs_new_protected_nopar:Npn \str_convert_decode_:
664 { \str_convert_gmap:N \str_decode_native_char:N }
665 \cs_new:Npn \str_decode_native_char:N #1
666 { #1 \s_tl \int_value:w '#1 \s_tl }
(End definition for \str_convert_decode_:. This function is documented on page ??.)

```

\str\_convert\_encode\_: The conversion from an internal string to native character tokens is very different in pdfTeX and in other engines. For Unicode-aware engines, we need the definitions to be read when the null byte has category code 12, so we set that inside a group.

```

667 \group_begin:
668   \char_set_catcode_other:n { 0 }
669   \pdftex_if_engine:TF

```

\str\_encode\_native\_char:n Since pdfTeX only supports 8-bit characters, and we have a table of all bytes, the conversion can be done in linear time within an x-expanding assignment. Look out for character codes larger than 255, those characters are replaced by ?, and raise a flag, which then triggers a pdfTeX-specific error.

```

670 {
671   \cs_new_protected_nopar:Npn \str_convert_encode_:
672   {
673     \flag_clear:n { str_error }
674     \str_convert_gmap_internal:N \str_encode_native_char:n
675     \str_if_flag_error:nnx { str_error }
676     { pdfTeX-native-overflow } { }
677   }
678   \cs_new:Npn \str_encode_native_char:n #1
679   {
680     \if_num:w #1 < \c_two_hundred_fifty_six
681       \str_output_byte:n {#1}
682     \else:
683       \flag_raise:n { str_error }
684       ?
685     \fi:
686   }
687   \msg_kernel_new:nnnn { str } { pdfTeX-native-overflow }
688   { Character-code-too-large-for-pdfTeX. }

```

```

689 {
690     The~pdfTeX-engine~only~supports~8-bit~characters:~
691     valid~character~codes~are~in~the~range~[0,255].~
692     To~manipulate~arbitrary~Unicode,~use~LuaTeX~or~XeTeX.
693 }
694 }
```

\str\_encode\_native\_loop:w    In Unicode-aware engines, since building particular characters cannot be done expandably in TeX, we cannot hope to get a linear-time function. However, we get quite close by building the result one block of 255 characters at a time. Set the lowercase code of the ASCII character 1 to the first character code in the block, *etc.*, then lowercase the contents of \c\_str\_positive\_bytes\_tl. Unfortunately, this is complicated by the fact that TeX won't lowercase a character to the null character, thus the need for filtering: whenever the lowercase code of a given character is zero, it is replaced "by hand" by the null character. Unicode-aware engines will never incur an overflow because the internal string is guaranteed to only contain code points in [0, "10FFFF].

```

695 {
696     \cs_new_protected_nopar:Npn \str_convert_encode_:
697     {
698         \tex_lccode:D \c_zero \c_zero
699         \int_zero:N \l_str_internal_int
700         \tl_set_eq:NN \l_str_internal_tl \g_str_result_tl
701         \tl_gclear:N \g_str_result_tl
702         \exp_after:wN \str_encode_native_loop:w
703         \l_str_internal_tl \s_tl { \q_stop \prg_map_break: } \s_tl
704         \prg_break_point:n { \str_encode_native_flush: }
705     }
706     \cs_new_protected:Npn \str_encode_native_loop:w #1 \s_tl #2 \s_tl
707     {
708         \use_none_delimit_by_q_stop:w #2 \q_stop
709         \if_num:w \l_str_internal_int = \c_two_hundred_fifty_five
710             \str_encode_native_flush:
711             \l_str_internal_int \c_zero
712         \fi:
713         \tex_advance:D \l_str_internal_int \c_one
714         \tex_lccode:D \l_str_internal_int \int_eval:w #2 \int_eval_end:
715         \str_encode_native_loop:w
716     }
717     \cs_new_protected_nopar:Npn \str_encode_native_flush:
718     {
719         \use:x
720         {
721             \tl_to_lowercase:n
722             {
723                 \tl_gput_right:Nx \exp_not:N \g_str_result_tl
724                 {
725                     \exp_after:wN \str_encode_native_filter:N
726                     \c_str_positive_bytes_tl
727                     { ? = \c_zero \fi: \prg_map_break: }
```

```

728          \prg_break_point:n { }
729      }
730  }
731 }
732 }
733 \cs_new:Npn \str_encode_native_filter:N #1
734 {
735     \if_num:w '#1 > \l_str_internal_int
736         \exp_after:wN \prg_map_break:
737     \fi:
738     \if_num:w \tex_lccode:D '#1 = \c_zero
739         ^^@
740     \else:
741         #1
742     \fi:
743     \str_encode_native_filter:N
744 }
745 }
```

End the group to restore the catcode of the null byte.

```

746 \group_end:
(End definition for \str_convert_encode_:. This function is documented on page ??.)
```

#### 7.4.7 8-bit encodings

This section will be entirely rewritten: it is not yet clear in what situations 8-bit encodings are used, hence I don't know what exactly should be optimized. The current approach is reasonably efficient to convert long strings, and it scales well when using many different encodings. An approach based on csnames would have a smaller constant load time for each individual conversion, but has a large hash table cost. Using a range of \count registers works for decoding, but not for encoding: one possibility there would be to use a binary tree for the mapping of Unicode characters to bytes, stored as a box, one per encoding.

Since the section is going to be rewritten, documentation lacks.

All the 8-bit encodings which l3str supports rely on the same internal functions.

\str\_declare\_eight\_bit\_encoding:nnn All the 8-bit encoding definition file start with \str\_declare\_eight\_bit\_encoding:nnn {\langle encoding name\rangle} {\langle mapping\rangle} {\langle missing bytes\rangle}. The {\langle mapping\rangle} argument is a token list of pairs {\langle byte\rangle} {\langle Unicode\rangle} expressed in uppercase hexadecimal notation. The {\langle missing\rangle} argument is a token list of {\langle byte\rangle}. Every {\langle byte\rangle} which does not appear in the {\langle mapping\rangle} nor the {\langle missing\rangle} lists maps to the same code point in Unicode.

```

747 \cs_new_protected:Npn \str_declare_eight_bit_encoding:nnn #1#2#3
748 {
749     \tl_set:Nn \l_str_internal_tl {#1}
750     \cs_new_protected_nopar:cpn { str_convert_decode_#1: }
751     { \str_convert_decode_eight_bit:n {#1} }
752     \cs_new_protected_nopar:cpn { str_convert_encode_#1: }
753     { \str_convert_encode_eight_bit:n {#1} }
754     \tl_const:cn { c_str_encoding_#1_t1 } {#2}
```

```

755   \tl_const:cn { c_str_encoding_#1_missing_tl } {#3}
756 }
(End definition for \str_declare_eight_bit_encoding:nnn. This function is documented on page
8.)
\str_convert_decode_eight_bit:n
\str_decode_eight_bit_load:nn
\str_decode_eight_bit_load_missing:n
\str_decode_eight_bit_char:N
757 \cs_new_protected:Npn \str_convert_decode_eight_bit:n #1
758 {
759   \group_begin:
760   \int_zero:N \l_str_internal_int
761   \exp_last_unbraced:Nx \str_decode_eight_bit_load:nn
762   { \tl_use:c { c_str_encoding.#1_tl } }
763   { \q_stop \prg_map_break: } { }
764   \prg_break_point:n { }
765   \exp_last_unbraced:Nx \str_decode_eight_bit_load_missing:n
766   { \tl_use:c { c_str_encoding.#1_missing_tl } }
767   { \q_stop \prg_map_break: }
768   \prg_break_point:n { }
769   \flag_clear:n { str_error }
770   \str_convert_gmap:N \str_decode_eight_bit_char:N
771   \str_if_flag_error:nnx { str_error } { decode-8-bit } {#1}
772   \group_end:
773 }
774 \cs_new_protected:Npn \str_decode_eight_bit_load:nn #1#2
775 {
776   \use_none_delimit_by_q_stop:w #1 \q_stop
777   \tex_dimen:D "#1 = \l_str_internal_int sp \scan_stop:
778   \tex_skip:D \l_str_internal_int = "#1 sp \scan_stop:
779   \tex_toks:D \l_str_internal_int \exp_after:wN { \int_value:w "#2 }
780   \tex_advance:D \l_str_internal_int \c_one
781   \str_decode_eight_bit_load:nn
782 }
783 \cs_new_protected:Npn \str_decode_eight_bit_load_missing:n #1
784 {
785   \use_none_delimit_by_q_stop:w #1 \q_stop
786   \tex_dimen:D "#1 = \l_str_internal_int sp \scan_stop:
787   \tex_skip:D \l_str_internal_int = "#1 sp \scan_stop:
788   \tex_toks:D \l_str_internal_int \exp_after:wN
789   { \int_use:N \c_str_replacement_char_int }
790   \tex_advance:D \l_str_internal_int \c_one
791   \str_decode_eight_bit_load_missing:n
792 }
793 \cs_new:Npn \str_decode_eight_bit_char:N #1
794 {
795   #1 \s_tl
796   \if_num:w \tex_dimen:D '#1 < \l_str_internal_int
797     \if_num:w \tex_skip:D \tex_dimen:D '#1 = '#1 \exp_stop_f:
798       \tex_the:D \tex_toks:D \tex_dimen:D
799     \fi:
800   \fi:

```

```

801     \int_value:w '#1 \s_tl
802   }
(End definition for \str_convert_decode_eight_bit:n. This function is documented on page ??.)

\str_convert_encode_eight_bit:n
\str_encode_eight_bit_load:nn
\str_encode_eight_bit_char:n
\str_encode_eight_bit_char_aux:n

803 \cs_new_protected:Npn \str_convert_encode_eight_bit:n #1
804   {
805     \group_begin:
806       \int_zero:N \l_str_internal_int
807       \exp_last_unbraced:Nx \str_encode_eight_bit_load:nn
808         { \tl_use:c { c_str_encoding_#1_tl } }
809         { \q_stop \prg_map_break: } { }
810       \prg_break_point:n { }
811       \flag_clear:n { str_error }
812       \str_convert_gmap_internal:N \str_encode_eight_bit_char:n
813         \str_if_flag_error:nnx { str_error } { encode-8-bit } {#1}
814       \group_end:
815   }
816 \cs_new_protected:Npn \str_encode_eight_bit_load:nn #1#2
817   {
818     \use_none_delimit_by_q_stop:w #1 \q_stop
819     \tex_dimen:D "#2 = \l_str_internal_int sp \scan_stop:
820     \tex_skip:D \l_str_internal_int = "#2 sp \scan_stop:
821     \exp_args:NNf \tex_toks:D \l_str_internal_int
822       { \str_output_byte:n { "#1 } }
823     \tex_advance:D \l_str_internal_int \c_one
824     \str_encode_eight_bit_load:nn
825   }
826 \cs_new:Npn \str_encode_eight_bit_char:n #1
827   {
828     \if_num:w #1 > \c_max_register_int
829       \flag_raise:n { str_error }
830     \else:
831       \if_num:w \tex_dimen:D #1 < \l_str_internal_int
832         \if_num:w \tex_skip:D \tex_dimen:D #1 = #1 \exp_stop_f:
833           \tex_the:D \tex_toks:D \tex_dimen:D #1 \exp_stop_f:
834           \exp_after:wN \exp_after:wN \exp_after:wN \use_none:nn
835         \fi:
836       \fi:
837       \str_encode_eight_bit_char_aux:n {#1}
838     \fi:
839   }
840 \cs_new:Npn \str_encode_eight_bit_char_aux:n #1
841   {
842     \if_num:w #1 < \c_two_hundred_fifty_six
843       \str_output_byte:n {#1}
844     \else:
845       \flag_raise:n { str_error }
846     \fi:
847   }

```

(End definition for `\str_convert_encode_eight_bit:n`. This function is documented on page ??.)

## 7.5 Messages

General messages, and messages for the encodings and escapings loaded by default (“native”, and “bytes”).

```

848 \msg_kernel_new:nnn { str } { unknown-esc }
849   { Escaping-scheme-'#1'~(filtered:~'#2')~unknown. }
850 \msg_kernel_new:nnn { str } { unknown-enc }
851   { Encoding-scheme-'#1'~(filtered:~'#2')~unknown. }
852 \msg_kernel_new:nnnn { str } { native-escaping }
853   { The-'native'-encoding-scheme-does-not-support-any-escaping. }
854   {
855     Since-native-strings-do-not-consist-in-bytes,~
856     none-of-the-escaping-methods-make-sense.~
857     The-specified-escaping,-'#1',-will-be-ignored.
858   }
859 \msg_kernel_new:nnn { str } { file-not-found }
860   { File-'l3str-#1.def'-not-found. }

```

Message used when the “bytes” unescaping fails because the string given to `\str_set_convert:Nnnn` contains a non-byte. This cannot happen for the pdfTEX engine, since that engine only supports 8-bit characters. Messages used for other escapings and encodings are defined in each definition file.

```

861 \pdftex_if_engine:F
862   {
863     \msg_kernel_new:nnnn { str } { non-byte }
864       { String-invalid-in-escaping-'#1':-it-may-only-contain-bytes. }
865       {
866         Some-characters-in-the-string-you-asked-to-convert-are-not-
867         8-bit-characters.-Perhaps-the-string-is-a-'native'-Unicode-string?-
868         If-it-is,-try-using\\
869         \\
870         \iow_indent:n
871           {
872             \iow_char:N\str_set_convert:Nnnn \\
873             \ \ <str-var>-\{<string>-\}~\{-native-\}~\{-target-encoding->-\}
874           }
875         }
876     }

```

Those messages are used when converting to and from 8-bit encodings.

```

877 \msg_kernel_new:nnnn { str } { decode-8-bit }
878   { Invalid-string-in-encoding-'#1'. }
879   {
880     LaTeX-came-across-a-byte-which-is-not-defined-to-represent-
881     any-character-in-the-encoding-'#1'.
882   }
883 \msg_kernel_new:nnnn { str } { encode-8-bit }
884   { Unicode-string-cannot-be-converted-to-encoding-'#1'. }

```

```

885  {
886      The~encoding~'#1'~only~contains~a~subset~of~all~Unicode~characters.~
887      LaTeX~was~asked~to~convert~a~string~to~that~encoding,~but~that~
888      string~contains~a~character~that~'#1'~does~not~support.
889 }

```

## 7.6 Deprecated string functions

`\str_length_skip_spaces:N` The naming scheme is a little bit more consistent with “`ignore_spaces`” instead of “`skip_spaces`”.

```

890 \cs_set:Npn \str_length_skip_spaces:N
891   { \exp_args:No \str_length_skip_spaces:n }
892 \cs_set_eq:NN \str_length_skip_spaces:n \str_length_ignore_spaces:n
(End definition for \str_length_skip_spaces:N and \str_length_skip_spaces:n. These functions
are documented on page ??.)
893 ⟨/initex | package⟩

```

## 7.7 Escaping definition files

Several of those encodings are defined by the pdf file format. The following byte storage methods are defined:

- `bytes` (default), non-bytes are filtered out, and bytes are left untouched (this is defined by default);
- `hex` or `hexadecimal`, as per the pdftEX primitive `\pdfescapehex`
- `name`, as per the pdftEX primitive `\pdfescapename`
- `string`, as per the pdftEX primitive `\pdfescapestring`
- `url`, as per the percent encoding of urls.

### 7.7.1 Unescape methods

`\str_convert_unescape_hex:` Take chars two by two, and interpret each pair as the hexadecimal code for a byte.  
`\str_unescape_hex_i:N` Anything else than hexadecimal digits is ignored, raising the flag. A string which contains  
`\str_unescape_hex_ii:N` an odd number of hexadecimal digits gets 0 appended to it: this is equivalent to appending  
a 0 in all cases, and dropping it if it is alone.

```

894 {*hex}
895 \cs_new_protected_nopar:Npn \str_convert_unescape_hex:
896   {
897     \group_begin:
898       \flag_clear:n { str_error }
899       \int_set:Nn \tex_escapechar:D { 92 }
900       \tl_gset:Nx \g_str_result_tl
901       {
902         \str_output_byte:w "
903         \exp_last_unbraced:Nf \str_unescape_hex_i:N

```

```

904         { \tl_to_str:N \g_str_result_tl }
905         0 { ? 0 - \c_one \prg_map_break: }
906         \prg_break_point:n { \str_output_end: }
907     }
908     \str_if_flag_error:nnx { str_error } { unescape-hex } { }
909     \group_end:
910   }
911 \cs_new:Npn \str_unescape_hex_i:N #1
912   {
913     \use_none:n #1
914     \str_aux_hexadecimal_use:NTF #1
915     { \str_unescape_hex_ii:N }
916     {
917       \flag_raise:n { str_error }
918       \str_unescape_hex_i:N
919     }
920   }
921 \cs_new:Npn \str_unescape_hex_ii:N #1
922   {
923     \use_none:n #1
924     \str_aux_hexadecimal_use:NTF #1
925     {
926       \str_output_end:
927       \str_output_byte:w " \str_unescape_hex_i:N
928     }
929     {
930       \flag_raise:n { str_error }
931       \str_unescape_hex_ii:N
932     }
933   }
934 \msg_kernel_new:nnnn { str } { unescape-hex }
935 { String-invalid-in-escaping-'hex':-only-hexadecimal-digits-allowed. }
936 {
937   Some-characters-in-the-string-you-asked-to-convert-are-not-
938   hexadecimal-digits-(0-9,-A-F,-a-f)-nor-spaces.
939 }
940 (/hex)
(End definition for \str_convert_unescape_hex:. This function is documented on page ??.)

```

**\str\_convert\_unescape\_name:** The `\str_convert_unescape_name:` function replaces each occurrence of # followed by two hexadecimal digits in `\g_str_result_tl` by the corresponding byte. The `url` function is identical, with escape character % instead of #. Thus we define the two together. The arguments of `\str_tmp:w` are the character code of # or % in hexadecimal, the name of the main function to define, and the name of the auxiliary which performs the loop.

The looping auxiliary #3 finds the next escape character, reads the following two characters, and tests them. The test `\str_aux_hexadecimal_use:NTF` leaves the upper-case digit in the input stream, hence we surround the test with `\str_output_byte:w "`

and `\str_output_end`:. If both characters are hexadecimal digits, they should be removed before looping: this is done by `\use_i:nnn`. If one of the characters is not a hexadecimal digit, then feed "#1 to `\str_output_byte:w` to produce the escape character, raise the flag, and call the looping function followed by the two characters (remove `\use_i:nnn`).

```

941  {*name | url}
942  \cs_set_protected:Npn \str_tmp:w #1#2#3
943  {
944    \cs_new_protected:cpx { str_convert_unescape_#2: }
945    {
946      \group_begin:
947        \flag_clear:n { str_byte }
948        \flag_clear:n { str_error }
949        \int_set:Nn \tex_escapechar:D { 92 }
950        \tl_gset:Nx \g_str_result_tl
951        {
952          \exp_after:wn #3 \g_str_result_tl
953          #1 ? { ? \prg_map_break: }
954          \prg_break_point:n { }
955        }
956        \str_if_flag_error:nnx { str_byte } { non-byte } { #2 }
957        \str_if_flag_error:nnx { str_error } { unescape-#2 } { }
958      \group_end:
959    }
960    \cs_new:Npn #3 ##1#1##2##3
961    {
962      \str_filter_bytes:n {##1}
963      \use_none:n ##3
964      \str_output_byte:w "
965        \str_aux_hexadecimal_use:NTF ##2
966        {
967          \str_aux_hexadecimal_use:NTF ##3
968          { }
969          {
970            \flag_raise:n { str_error }
971            * \c_zero + '#1 \use_i:nn
972          }
973        }
974        {
975          \flag_raise:n { str_error }
976          0 + '#1 \use_i:nn
977        }
978        \str_output_end:
979        \use_i:nnn #3 ##2##3
980    }
981    \msg_kernel_new:nnnn { str } { unescape-#2 }
982    { String~invalid~in~escaping~'#2'. }
983    {
984      LaTeX-came-across-the-escape-character-'#1'-not-followed-by-

```

```

985         two~hexadecimal~digits.~This~is~invalid~in~the~escaping~'#2'.
986     }
987 }
988 </name | url>
989 <name>\exp_after:wN \str_tmp:w \c_hash_str { name }
990 <name> \str_unescape_name_loop:wNN
991 <url>\exp_after:wN \str_tmp:w \c_percent_str { url }
992 <url> \str_unescape_url_loop:wNN
(End definition for \str_convert_unescape_name:. This function is documented on page ??.)
```

\str\_convert\_unescape\_string:  
\str\_unescape\_string\_newlines:wN  
\str\_unescape\_string\_loop:wNNN  
\str\_unescape\_string\_repeat:NNNNNN

The **string** escaping is somewhat similar to the **name** and **url** escapings, with escape character **\**. The first step is to convert all three line endings, **^J**, **^M**, and **^M^J** to the common **^J**, as per the PDF specification. This step cannot raise the flag.

Then the following escape sequences are decoded.

- \n Line feed (10)
- \r Carriage return (13)
- \t Horizontal tab (9)
- \b Backspace (8)
- \f Form feed (12)
- \( Left parenthesis
- \) Right parenthesis
- \\\ Backslash

\ddd (backslash followed by 1 to 3 octal digits) Byte **ddd** (octal), subtracting 256 in case of overflow.

If followed by an end-of-line character, the backslash and the end-of-line are ignored. If followed by anything else, the backslash is ignored, raising the error flag.

```

993 /*string)
994 \group_begin:
995   \char_set_lccode:nn {'\*} {'\\}
996   \char_set_catcode_other:N \^J
997   \char_set_catcode_other:N \^M
998   \tl_to_lowercase:n
999   {
1000     \cs_new_protected_nopar:Npn \str_convert_unescape_string:
1001     {
1002       \group_begin:
1003         \flag_clear:n { str_byte }
1004         \flag_clear:n { str_error }
1005         \int_set:Nn \tex_escapechar:D { 92 }
1006         \tl_gset:Nx \g_str_result_tl
1007         {
```

```

1008     \exp_after:wN \str_unescape_string_newlines:wN
1009         \g_str_result_tl \prg_map_break: ^^M ?
1010         \prg_break_point:n { }
1011     }
1012     \tl_gset:Nx \g_str_result_tl
1013     {
1014         \exp_after:wN \str_unescape_string_loop:wNNN
1015             \g_str_result_tl * ?? { ? \prg_map_break: }
1016             \prg_break_point:n { }
1017         }
1018         \str_if_flag_error:nnx { str_byte } { non-byte } { string }
1019         \str_if_flag_error:nnx { str_error } { unescape-string } { }
1020         \group_end:
1021     }
1022     \cs_new:Npn \str_unescape_string_loop:wNNN #1 **#2#3#4
1023 }
1024 {
1025     \str_filter_bytes:n {#1}
1026     \use_none:n #4
1027     \str_output_byte:w '
1028     \str_aux_octal_use:NTF #2
1029     {
1030         \str_aux_octal_use:NTF #3
1031         {
1032             \str_aux_octal_use:NTF #4
1033             {
1034                 \if_int_compare:w #2 > \c_three
1035                     - 256
1036                 \fi:
1037                 \str_unescape_string_repeat:NNNNNN
1038             }
1039             { \str_unescape_string_repeat:NNNNNN ? }
1040         }
1041         { \str_unescape_string_repeat:NNNNNN ?? }
1042     }
1043     {
1044         \prg_case_str:xxn {#2}
1045         {
1046             { \c_backslash_str } { 134 }
1047             { ( } { 50 }
1048             { ) } { 51 }
1049             { r } { 15 }
1050             { f } { 14 }
1051             { n } { 12 }
1052             { t } { 11 }
1053             { b } { 10 }
1054             { ^^J } { 0 - \c_one }
1055         }
1056         {
1057             \flag_raise:n { str_error }

```

```

1058          0 - \c_one \use_i:nn
1059      }
1060  }
1061  \str_output_end:
1062  \use_i:nn \str_unescape_string_loop:wNNN #2#3#4
1063  }
1064 \cs_new:Npn \str_unescape_string_repeat:NNNNNN #1#2#3#4#5#6
1065  { \str_output_end: \str_unescape_string_loop:wNNN }
1066 \cs_new:Npn \str_unescape_string_newlines:wN #1 ^^M #2
1067  {
1068  #1
1069  \if_charcode:w ^^J #2 \else: ^^J \fi:
1070  \str_unescape_string_newlines:wN #2
1071  }
1072 \msg_kernel_new:nnnn { str } { unescape-string }
1073  { String~invalid-in~escaping~'string'. }
1074  {
1075  LaTeX~came~across~an~escape~character~'\c_backslash_str'~
1076  not~followed~by~any~of:~'n',~'r',~'t',~'b',~'f',~'('~,~')',~
1077  '\c_backslash_str',~one~to~three~octal~digits,~or~the~end~
1078  of~a~line.
1079  }
1080 \group_end:
1081 </string>
(End definition for \str_convert_unescape_string:. This function is documented on page ??.)
```

### 7.7.2 Escape methods

Currently, none of the escape methods can lead to errors, assuming that their input is made out of bytes.

\str\_convert\_escape\_hex: Loop and convert each byte to hexadecimal.

```

\str_escape_hex_char:N
1082  {*hex}
1083  \cs_new_protected_nopar:Npn \str_convert_escape_hex:
1084  { \str_convert_gmap:N \str_escape_hex_char:N }
1085  \cs_new:Npn \str_escape_hex_char:N #1
1086  { \str_output_hexadecimal:n { '#1' } }
1087 </hex>
(End definition for \str_convert_escape_hex:. This function is documented on page ??.)
```

\str\_convert\_escape\_name: \str\_escape\_name\_char:N \str\_if\_escape\_name:NTF \c\_str\_escape\_name\_str \c\_str\_escape\_name\_not\_str For each byte, test whether it should be output as is, or be “hash-encoded”. Roughly, bytes outside the range [“2A”, “7E] are hash-encoded. We keep two lists of exceptions: characters in \c\_str\_escape\_name\_not\_str are not hash-encoded, and characters in the \c\_str\_escape\_name\_str are encoded.

```

1088  {*name}
1089  \str_const:Nn \c_str_escape_name_not_str { ! " $ & ' } %$%
1090  \str_const:Nn \c_str_escape_name_str { {} / <> [] }
1091  \cs_new_protected_nopar:Npn \str_convert_escape_name:
1092  { \str_convert_gmap:N \str_escape_name_char:N }
```

```

1093 \cs_new:Npn \str_escape_name_char:N #1
1094 {
1095     \str_if_escape_name:NTF #1 {#1}
1096     { \c_hash_str \str_output_hexadecimal:n {'#1} }
1097 }
1098 \prg_new_conditional:Npnn \str_if_escape_name:N #1 { TF }
1099 {
1100     \if_num:w '#1 < "2A \exp_stop_f:
1101         \str_if_contains_char:NNTF \c_str_escape_name_not_str #1
1102         \prg_return_true: \prg_return_false:
1103     \else:
1104         \if_num:w '#1 > "7E \exp_stop_f:
1105             \prg_return_false:
1106         \else:
1107             \str_if_contains_char:NNTF \c_str_escape_name_str #1
1108             \prg_return_false: \prg_return_true:
1109         \fi:
1110     \fi:
1111 }
1112 /name
(End definition for \str_convert_escape_name:. This function is documented on page ??.)

```

**\str\_convert\_escape\_string:** Any character below (and including) space, and any character above (and including) **del**, are converted to octal. One backslash is added before each parenthesis and backslash.

```

1113 (*string)
1114 \str_const:Nx \c_str_escape_string_str
1115 { \c_backslash_str ( ) }
1116 \cs_new_protected_nopar:Npn \str_convert_escape_string:
1117 { \str_convert_gmap:N \str_escape_string_char:N }
1118 \cs_new:Npn \str_escape_string_char:N #1
1119 {
1120     \str_if_escape_string:NTF #1
1121     {
1122         \str_if_contains_char:NNT
1123         \c_str_escape_string_str #1
1124         { \c_backslash_str }
1125         #1
1126     }
1127     {
1128         \c_backslash_str
1129         \int_div_truncate:nn {'#1} {64}
1130         \int_mod:nn { \int_div_truncate:nn {'#1} \c_eight } \c_eight
1131         \int_mod:nn {'#1} \c_eight
1132     }
1133 }
1134 \prg_new_conditional:Npnn \str_if_escape_string:N #1 { TF }
1135 {
1136     \if_num:w '#1 < "21 \exp_stop_f:
1137         \prg_return_false:

```

```

1138     \else:
1139         \if_num:w '#1 > "7E \exp_stop_f:
1140             \prg_return_false:
1141         \else:
1142             \prg_return_true:
1143         \fi:
1144     \fi:
1145 }
1146 (/string)
(End definition for \str_convert_escape_string:. This function is documented on page ??.)
```

\str\_convert\_escape\_url: This function is similar to \str\_convert\_escape\_name:, escaping different characters.

```

\str_escape_url_char:N
\str_if_escape_url:NTF
1147 (*url)
1148 \cs_new_protected_nopar:Npn \str_convert_escape_url:
1149   { \str_convert_gmap:N \str_escape_url_char:N }
1150 \cs_new:Npn \str_escape_url_char:N #1
1151   {
1152     \str_if_escape_url:NTF #1 {#1}
1153     { \c_percent_str \str_output_hexadecimal:n { '#1 } }
1154   }
1155 \prg_new_conditional:Npnn \str_if_escape_url:N #1 { TF }
1156   {
1157     \if_num:w '#1 < "41 \exp_stop_f:
1158         \str_if_contains_char:nNTF { "-.<> } #1
1159         \prg_return_true: \prg_return_false:
1160     \else:
1161         \if_num:w '#1 > "7E \exp_stop_f:
1162             \prg_return_false:
1163         \else:
1164             \str_if_contains_char:nNTF { [ ] } #1
1165             \prg_return_false: \prg_return_true:
1166             \fi:
1167         \fi:
1168   }
1169 (/url)
(End definition for \str_convert_escape_url:. This function is documented on page ??.)
```

## 7.8 Encoding definition files

The `native` encoding is automatically defined. Other encodings are loaded as needed. The following encodings are supported:

- UTF-8;
- UTF-16, big-, little-endian, or with byte order mark;
- UTF-32, big-, little-endian, or with byte order mark;
- the ISO 8859 code pages, numbered from 1 to 16, skipping the nonexistent ISO 8859-12.

### 7.8.1 utf-8 support

1170 `(*utf8)`

`\str_convert_encode_utf8:` Loop through the internal string, and convert each character to its UTF-8 representation. The representation is built from the right-most (least significant) byte to the left-most (most significant) byte. Continuation bytes are in the range [128, 191], taking 64 different values, hence we roughly want to express the character code in base 64, shifting the first digit in the representation by some number depending on how many continuation bytes there are. In the range [0, 127], output the corresponding byte directly. In the range [128, 2047], output the remainder modulo 64, plus 128 as a continuation byte, then output the quotient (which is in the range [0, 31]), shifted by 192. In the next range, [2048, 65535], split the character code into residue and quotient modulo 64, output the residue as a first continuation byte, then repeat; this leaves us with a quotient in the range [0, 15], which we output shifted by 224. The last range, [65536, 1114111], follows the same pattern: once we realize that dividing twice by 64 leaves us with a number larger than 15, we repeat, producing a last continuation byte, and offset the quotient by 240 for the leading byte.

How is that implemented? `\str_encode_utf_viii_loop:wwnnw` takes successive quotients as its first argument, the quotient from the previous step as its second argument (except in step 1), the bound for quotients that trigger one more step or not, and finally the offset used if this step should produce the leading byte. Leading bytes can be in the ranges [0, 127], [192, 223], [224, 239], and [240, 247] (really, that last limit should be 244 because Unicode stops at the code point 1114111). At each step, if the quotient #1 is less than the limit #3 for that range, output the leading byte (#1 shifted by #4) and stop. Otherwise, we need one more step: use the quotient of #1 by 64, and #1 as arguments for the looping auxiliary, and output the continuation byte corresponding to the remainder #2 – 64#1 + 128. The bizarre construction `\c_minus_one + \c_zero *` removes the spurious initial continuation byte (better methods welcome).

```

1171 \cs_new_protected_nopar:cpn { str_convert_encode_utf8: }
1172   { \str_convert_gmap_internal:N \str_encode_utf_viii_char:n }
1173 \cs_new:Npn \str_encode_utf_viii_char:n #1
1174   {
1175     \str_encode_utf_viii_loop:wwnnw #1 ; \c_minus_one + \c_zero * ;
1176     { 128 } { \c_zero }
1177     { 32 } { 192 }
1178     { 16 } { 224 }
1179     { 8 } { 240 }
1180     \q_stop
1181   }
1182 \cs_new:Npn \str_encode_utf_viii_loop:wwnnw #1; #2; #3#4 #5 \q_stop
1183   {
1184     \if_num:w #1 < #3 \exp_stop_f:
1185       \str_output_byte:n { #1 + #4 }
1186       \exp_after:wN \use_none_delimit_by_q_stop:w
1187     \fi:
1188     \exp_after:wN \str_encode_utf_viii_loop:wwnnw
1189       \int_value:w \int_div_truncate:nn {#1} {64} ; #1 ;

```

```

1190      #5 \q_stop
1191      \str_output_byte:n { #2 - 64 * ( #1 - \c_two ) }
1192  }
(End definition for \str_convert_encode_utf8:. This function is documented on page ??.)
```

\l\_str\_missing\_flag    When decoding a string that is purportedly in the UTF-8 encoding, four different errors  
 \l\_str\_extra\_flag    can occur, signalled by a specific flag for each (we define those flags using \flag\_clear\_-  
 \l\_str\_overlong\_flag new:n rather than \flag\_new:n, because they are shared with other encoding definition  
 \l\_str\_overflow\_flag files).

- “Missing continuation byte”: a leading byte is not followed by the right number of continuation bytes.
- “Extra continuation byte”: a continuation byte appears where it was not expected, *i.e.*, not after an appropriate leading byte.
- “Overlong”: a Unicode character is expressed using more bytes than necessary, for instance, "C0"80 for the code point 0, instead of a single null byte.
- “Overflow”: this occurs when decoding produces Unicode code points greater than 1114111.

We only raise one L<sup>A</sup>T<sub>E</sub>X3 error message, combining all the errors which occurred. In the short message, the leading comma must be removed to get a grammatically correct sentence. In the long text, first remind the user what a correct UTF-8 string should look like, then add error-specific information.

```

1193 \flag_clear_new:n { str_missing }
1194 \flag_clear_new:n { str_extra }
1195 \flag_clear_new:n { str_overlong }
1196 \flag_clear_new:n { str_overflow }
1197 \msg_kernel_new:nnnn { str } { utf8-decode }
1198 {
1199   Invalid~UTF-8~string: \exp_last_unbraced:Nf \use_none:n
1200   \str_if_flag_times:nT { str_missing } { ,~missing~continuation~byte }
1201   \str_if_flag_times:nT { str_extra } { ,~extra~continuation~byte }
1202   \str_if_flag_times:nT { str_overlong } { ,~overlong~form }
1203   \str_if_flag_times:nT { str_overflow } { ,~code~point~too~large }
1204 .
1205 }
1206 {
1207   In~the~UTF-8~encoding,~each~Unicode~character~consists~in~
1208   1~to~4~bytes,~with~the~following~bit~pattern: \\%
1209   \iow_indent:
1210   {
1211     Code~point~\ \\ \\ <~128:~0xxxxxx \\%
1212     Code~point~\ \\ \\ <~2048:~110xxxx~10xxxxxx \\%
1213     Code~point~\ \\ <~65536:~1110xxxx~10xxxxxx~10xxxxxx \\%
1214     Code~point~ <~1114112:~11110xxx~10xxxxxx~10xxxxxx~10xxxxxx \\%
1215   }
1216 Bytes~of~the~form~10xxxxxx~are~called~continuation~bytes.
```

```

1217   \flag_if_raised:nT { str_missing }
1218   {
1219     \\
1220     A-leading-byte~(in~the~range~[192,255])~was~not~followed~by~
1221     the~appropriate~number~of~continuation~bytes.
1222   }
1223   \flag_if_raised:nT { str_extra }
1224   {
1225     \\
1226     LaTeX~came~across~a~continuation~byte~when~it~was~not~expected.
1227   }
1228   \flag_if_raised:nT { str_overlong }
1229   {
1230     \\
1231     Every~Unicode~code~point~must~be~expressed~in~the~shortest~
1232     possible~form.~For~instance,~'0xC0'~'0x83'~is~not~a~valid~
1233     representation~for~the~code~point~3.
1234   }
1235   \flag_if_raised:nT { str_overflow }
1236   {
1237     \\
1238     Unicode~limits~code~points~to~the~range~[0,1114111].
1239   }
1240 }

(End definition for \l_str_missing_flag and others. These functions are documented on page ??.)
```

```

\str_convert_decode_utf8:
\str_decode_utf_viii_start:N
  \str_decode_utf_viii_continuation:wwN
    \str_decode_utf_viii_aux:wNnwN
    \str_decode_utf_viii_overflow:w
\str_decode_utf_viii_end:
```

Decoding is significantly harder than encoding. As before, lower some flags, which are tested at the end (in bulk, to trigger at most one L<sup>A</sup>T<sub>E</sub>X3 error, as explained above). We expect successive multi-byte sequences of the form *<start byte> <continuation bytes>*. The *\_start* auxiliary tests the first byte:

- [0, "7F]: the byte stands alone, and is converted to its own character code;
- ["80, "BF]: unexpected continuation byte, raise the appropriate flag, and convert that byte to the replacement character "FFFD";
- ["C0, "FF]: this byte should be followed by some continuation byte(s).

In the first two cases, \use\_none\_delimit\_by\_q\_stop:w removes data that only the third case requires, namely the limits of ranges of Unicode characters which can be expressed with 1, 2, 3, or 4 bytes.

We can now concentrate on the multi-byte case and the *\_continuation* auxiliary. We expect #3 to be in the range ["80, "BF]. The test for this goes as follows: if the character code is less than "80, we compare it to ="C0, yielding **false**; otherwise to "C0, yielding **true** in the range ["80, "BF] and **false** otherwise. If we find that the byte is not a continuation range, stop the current slew of bytes, output the replacement character, and continue parsing with the *\_start* auxiliary, starting at the byte we just tested. Once we know that the byte is a continuation byte, leave it behind us in the

input stream, compute what code point the bytes read so far would produce, and feed that number to the `_aux` function.

The `_aux` function tests whether we should look for more continuation bytes or not. If the number it receives as `#1` is less than the maximum `#4` for the current range, then we are done: check for an overlong representation by comparing `#1` with the maximum `#3` for the previous range. Otherwise, we call the `_continuation` auxiliary again, after shifting the “current code point” by `#4` (maximum from the range we just checked).

Two additional tests are needed: if we reach the end of the list of range maxima and we are still not done, then we are faced with an overflow. Clean up, and again insert the code point “`FFFD`” for the replacement character. Also, every time we read a byte, we need to check whether we reached the end of the string. In a correct UTF-8 string, this happens automatically when the `_start` auxiliary leaves its first argument in the input stream: the end-marker begins with `\prg_map_break:`, which ends the loop. On the other hand, if the end is reached when looking for a continuation byte, the `\use_none:n #3` construction removes the first token from the end-marker, and leaves the `_end` auxiliary, which raises the appropriate error flag before ending the mapping.

```

1241 \cs_new_protected_nopar:cpn { str_convert_decode_utf8: }
1242 {
1243     \flag_clear:n { str_error }
1244     \flag_clear:n { str_missing }
1245     \flag_clear:n { str_extra }
1246     \flag_clear:n { str_overlong }
1247     \flag_clear:n { str_overflow }
1248     \tl_gset:Nx \g_str_result_tl
1249     {
1250         \exp_after:wN \str_decode_utf_viii_start:N \g_str_result_tl
1251         { \prg_map_break: \str_decode_utf_viii_end: }
1252         \prg_break_point:n { }
1253     }
1254     \str_if_flag_error:nnx { str_error } { utf8-decode } { }
1255 }
1256 \cs_new:Npn \str_decode_utf_viii_start:N #1
1257 {
1258     #1
1259     \if_num:w '#1 < "C0 \exp_stop_f:
1260         \s_tl
1261         \if_num:w '#1 < "80 \exp_stop_f:
1262             \int_value:w '#1
1263         \else:
1264             \flag_raise:n { str_extra }
1265             \flag_raise:n { str_error }
1266             \int_use:N \c_str_replacement_char_int
1267             \fi:
1268         \else:
1269             \exp_after:wN \str_decode_utf_viii_continuation:wwN
1270             \int_use:N \int_eval:w '#1 - "C0 \exp_after:wN \int_eval_end:
1271         \fi:
1272         \s_tl

```

```

1273   \use_none_delimit_by_q_stop:w {"80} {"800} {"10000} {"110000} \q_stop
1274   \str_decode_utf_viii_start:N
1275   }
1276 \cs_new:Npn \str_decode_utf_viii_continuation:wwN
1277   #1 \s_tl #2 \str_decode_utf_viii_start:N #3
1278   {
1279     \use_none:n #3
1280     \if_num:w '#3 < \if_num:w '#3 < "80 \exp_stop_f: - \fi: "C0 \exp_stop_f:
1281       #3
1282       \exp_after:wn \str_decode_utf_viii_aux:wnnnwN
1283       \int_use:N \int_eval:w
1284       #1 * "40 + '#3 - "80
1285       \exp_after:wn \int_eval_end:
1286     \else:
1287       \s_tl
1288       \flag_raise:n { str_missing }
1289       \flag_raise:n { str_error }
1290       \int_use:N \c_str_replacement_char_int
1291     \fi:
1292     \s_tl
1293     #2
1294     \str_decode_utf_viii_start:N #3
1295   }
1296 \cs_new:Npn \str_decode_utf_viii_aux:wnnnwN
1297   #1 \s_tl #2#3#4 #5 \str_decode_utf_viii_start:N #6
1298   {
1299     \if_num:w #1 < #4 \exp_stop_f:
1300       \s_tl
1301       \if_num:w #1 < #3 \exp_stop_f:
1302         \flag_raise:n { str_overlong }
1303         \flag_raise:n { str_error }
1304         \int_use:N \c_str_replacement_char_int
1305       \else:
1306         #1
1307       \fi:
1308     \else:
1309       \if_meaning:w \q_stop #5
1310         \str_decode_utf_viii_overflow:w #1
1311       \fi:
1312       \exp_after:wn \str_decode_utf_viii_continuation:wwN
1313       \int_use:N \int_eval:w #1 - #4 \exp_after:wn \int_eval_end:
1314     \fi:
1315     \s_tl
1316     #2 {#4} #5
1317     \str_decode_utf_viii_start:N
1318   }
1319 \cs_new:Npn \str_decode_utf_viii_overflow:w #1 \fi: #2 \fi:
1320   {
1321     \fi: \fi:
1322     \flag_raise:n { str_overflow }

```

```

1323     \flag_raise:n { str_error }
1324     \int_use:N \c_str_replacement_char_int
1325   }
1326 \cs_new_nopar:Npn \str_decode_utf_viii_end:
1327   {
1328     \s_tl
1329     \flag_raise:n { str_missing }
1330     \flag_raise:n { str_error }
1331     \int_use:N \c_str_replacement_char_int \s_tl
1332     \prg_map_break:
1333   }
(End definition for \str_convert_decode_utf8:. This function is documented on page ??.)
1334 
```

### 7.8.2 utf-16 support

The definitions are done in a category code regime where the bytes 254 and 255 used by the byte order mark have catcode 12.

```

1335 {*utf16}
1336 \group_begin:
1337   \char_set_catcode_other:N \^fe
1338   \char_set_catcode_other:N \^ff

```

\str\_convert\_encode\_utf16: When the endianness is not specified, it is big-endian by default, and we add a byte-order mark. Convert characters one by one in a loop, with different behaviours depending on the character code.

- [0, "D7FF]: converted to two bytes;
- ["D800, "DFFF] are used as surrogates: they cannot be converted and are replaced by the replacement character;
- ["E000, "FFFF]: converted to two bytes;
- ["10000, "10FFFF]: converted to a pair of surrogates, each two bytes.

For the duration of this operation, \str\_tmp:w is defined as a function to convert a number in the range [0, "FFFF] to a pair of bytes (either big endian or little endian), by feeding the quotient of the division of #1 by "100, followed by #1 to \str\_encode\_utf\_xvi\_be:nn or its le analog: those compute the remainder, and output two bytes for the quotient and remainder.

```

1339   \cs_new_protected_nopar:cpn { str_convert_encode_utf16: }
1340   {
1341     \str_encode_utf_xvi_aux:N \str_output_byte_pair_be:n
1342     \tl_gput_left:Nx \g_str_result_tl { ^fe ^ff }
1343   }
1344   \cs_new_protected_nopar:cpn { str_convert_encode_utf16be: }
1345   { \str_encode_utf_xvi_aux:N \str_output_byte_pair_be:n }
1346   \cs_new_protected_nopar:cpn { str_convert_encode_utf16le: }

```

```

1347 { \str_encode_utf_xvi_aux:N \str_output_byte_pair_le:n }
1348 \cs_new_protected:Npn \str_encode_utf_xvi_aux:N #1
1349 {
1350     \flag_clear:n { str_error }
1351     \cs_set_eq:NN \str_tmp:w #1
1352     \str_convert_gmap_internal:N \str_encode_utf_xvi_char:n
1353     \str_if_flag_error:nnx { str_error } { utf16-encode } { }
1354 }
1355 \cs_new:Npn \str_encode_utf_xvi_char:n #1
1356 {
1357     \if_int_compare:w #1 < "D800 \exp_stop_f:
1358         \str_tmp:w {#1}
1359     \else:
1360         \if_int_compare:w #1 < "10000 \exp_stop_f:
1361             \if_int_compare:w #1 < "E000 \exp_stop_f:
1362                 \flag_raise:n { str_error }
1363                 \str_tmp:w { \c_str_replacement_char_int }
1364             \else:
1365                 \str_tmp:w {#1}
1366             \fi:
1367         \else:
1368             \exp_args:Nf \str_tmp:w { \int_div_truncate:nn {#1} {"400} + "D800 }
1369             \exp_args:Nf \str_tmp:w { \int_mod:nn {#1} {"400} + "DC00 }
1370         \fi:
1371     \fi:
1372 }
(End definition for \str_convert_encode_utf16:, \str_convert_encode_utf16be:, and \str_convert_encode_utf16le:
These functions are documented on page ??.)

```

\l\_str\_missing\_flag When encoding a Unicode string to UTF-16, only one error can occur: code points in the range ["D800, "DFFF], corresponding to surrogates, cannot be encoded. We use the all-purpose flag **str\_error** to signal that error.

\l\_str\_extra\_flag When decoding a Unicode string which is purportedly in UTF-16, three errors can occur: a missing trail surrogate, an unexpected trail surrogate, and a string containing an odd number of bytes.

```

1373 \flag_clear_new:n { str_missing }
1374 \flag_clear_new:n { str_extra }
1375 \flag_clear_new:n { str_end }
1376 \msg_kernel_new:nnnn { str } { utf16-encode }
1377     { Unicode~string~cannot~be~expressed~in~UTF-16:~surrogate. }
1378     {
1379         Surrogate~code~points~(in~the~range~[U+D800,~U+DFFF])~
1380         can~be~expressed~in~the~UTF-8~and~UTF-32~encodings,~
1381         but~not~in~the~UTF-16~encoding.
1382     }
1383 \msg_kernel_new:nnnn { str } { utf16-decode }
1384     {
1385         Invalid~UTF-16~string: \exp_last_unbraced:Nf \use_none:n
1386         \str_if_flag_times:nT { str_missing } { ,~missing-trail-surrogate }

```

```

1387     \str_if_flag_times:nT { str_extra }      { ,~extra~trail~surrogate }
1388     \str_if_flag_times:nT { str_end }        { ,~odd~number~of~bytes }
1389     .
1390   }
1391   {
1392     In~the~UTF-16~encoding,~each~Unicode~character~is~encoded~as~
1393     2~or~4~bytes: \\
1394     \iow_indent:n
1395     {
1396       Code~point~in~[U+0000,~U+D7FF]:~two~bytes \\
1397       Code~point~in~[U+D800,~U+DBFF]:~illegal \\
1398       Code~point~in~[U+E000,~U+FFFF]:~two~bytes \\
1399       Code~point~in~[U+10000,~U+10FFFF]:~a~lead~surrogate~and~a~trail~surrogate \\
1400     }
1401   Lead~surrogates~are~pairs~of~bytes~in~the~range~[0xD800,~0xDBFF],~and~trail~surrogates~are~in~the~range~[0xDC00,~0xFFFF].
1402   \flag_if_raised:nT { str_missing }
1403   {
1404     \\
1405     \\
1406     A~lead~surrogate~was~not~followed~by~a~trail~surrogate.
1407   }
1408   \flag_if_raised:nT { str_extra }
1409   {
1410     \\
1411     \\
1412     LaTeX~came~across~a~trail~surrogate~when~it~was~not~expected.
1413   }
1414   \flag_if_raised:nT { str_end }
1415   {
1416     \\
1417     The~string~contained~an~odd~number~of~bytes.~This~is~invalid:~the~basic~code~unit~for~UTF-16~is~16~bits~(2~bytes).
1418   }
1419 }
1420
(End definition for \l_str_missing_flag, \l_str_extra_flag, and \l_str_end_flag. These functions are documented on page ??.)
```

`\str_convert_decode_utf16:` As for UTF-8, decoding UTF-16 is harder than encoding it. If the endianness is unknown, check the first two bytes: if those are "FE and "FF in either order, remove them and use the corresponding endianness, otherwise assume big-endianness. The three endianness cases are based on a common auxiliary whose first argument is 1 for big-endian and 2 for little-endian, and whose second argument, delimited by the scan mark `\s_stop`, is expanded once (the string may be long; passing `\g_str_result_tl` as an argument before expansion is cheaper).

The `\strr_decode_utf_xvi:Nw` function defines `\str_tmp:w` to take two arguments and return the character code of the first one if the string is big-endian, and the second one if the string is little-endian, then loops over the string using `\str_decode_utf_xvi_pair>NN` described below.

```
1421   \cs_new_protected_nopar:cpn { str_convert_decode_utf16be: }
```

```

1422 { \str_decode_utf_xvi:Nw 1 \g_str_result_tl \s_stop }
1423 \cs_new_protected_nopar:cpn { str_convert_decode_utf16le: }
1424 { \str_decode_utf_xvi:Nw 2 \g_str_result_tl \s_stop }
1425 \cs_new_protected_nopar:cpn { str_convert_decode_utf16: }
1426 {
1427   \exp_after:wN \str_decode_utf_xvi_bom>NN
1428     \g_str_result_tl \s_stop \s_stop \s_stop
1429 }
1430 \cs_new_protected:Npn \str_decode_utf_xvi_bom>NN #1#2
1431 {
1432   \str_if_eq:xxTF { #1#2 } { ^~ff ^~fe }
1433   { \str_decode_utf_xvi:Nw 2 }
1434   {
1435     \str_if_eq:xxTF { #1#2 } { ^~fe ^~ff }
1436     { \str_decode_utf_xvi:Nw 1 }
1437     { \str_decode_utf_xvi:Nw 1 #1#2 }
1438   }
1439 }
1440 \cs_new_protected:Npn \str_decode_utf_xvi:Nw #1#2 \s_stop
1441 {
1442   \flag_clear:n { str_error }
1443   \flag_clear:n { str_missing }
1444   \flag_clear:n { str_extra }
1445   \flag_clear:n { str_end }
1446   \cs_set:Npn \str_tmp:w ##1 ##2 { ' ## #1 }
1447   \tl_gset:Nx \g_str_result_tl
1448   {
1449     \exp_after:wN \str_decode_utf_xvi_pair>NN
1450       #2 \q_nil \q_nil
1451     \prg_break_point:n { }
1452   }
1453   \str_if_flag_error:n { str_error } { utf16-decode } { }
1454 }

```

(End definition for \str\_convert\_decode\_utf16:, \str\_convert\_decode\_utf16be:, and \str\_convert\_decode\_utf16le  
These functions are documented on page ??.)

```
\str_decode_utf_xvi_pair>NN
\str_decode_utf_xvi_quad>NNNN
\str_decode_utf_xvi_pair_end:Nw
\str_decode_utf_xvi_error:nNN
\str_decode_utf_xvi_extra:NNw
```

Bytes are read two at a time. At this stage, \str\_tmp:w #1#2 expands to the character code of the most significant byte, and we distinguish cases depending on which range it lies in:

- ["D8, "DB] signals a lead surrogate, and the integer expression yields 1 ( $\varepsilon$ -TeX rounds ties away from zero);
- ["DC, "DF] signals a trail surrogate, unexpected here, and the integer expression yields 2;
- any other value signals a code point in the Basic Multilingual Plane, which stands for itself, and the \if\_case:w construction expands to nothing (cases other than 1 or 2), leaving the relevant material in the input stream, followed by another call to the \_pair auxiliary.

The case of a lead surrogate is treated by the \_quad auxiliary, whose arguments #1, #2, #4 and #5 are the four bytes. We expect the most significant byte of #4#5 to be in the range ["DC, "DF] (trail surrogate). The test is similar to the test used for continuation bytes in the UTF-8 decoding functions. In the case where #4#5 is indeed a trail surrogate, leave #1#2#4#5 \s\_t1 ⟨code point⟩ \s\_t1, and remove the pair #4#5 before looping with \str\_decode\_utf\_xvi\_pair:NN. Otherwise, of course, complain about the missing surrogate.

Every time we read a pair of bytes, we test for the end-marker \q\_nil. When reaching the end, we additionally check that the string had an even length. Also, if the end is reached when expecting a trail surrogate, we treat that as a missing surrogate.

```

1455   \cs_new:Npn \str_decode_utf_xvi_pair:NN #1#2
1456   {
1457     \if_meaning:w \q_nil #2
1458       \str_decode_utf_xvi_pair_end:Nw #1
1459     \fi:
1460     \if_case:w
1461       \int_eval:w ( \str_tmp:w #1#2 - "D6 ) / \c_four \int_eval_end:
1462       \or: \exp_after:wN \str_decode_utf_xvi_quad>NNwNN
1463       \or: \exp_after:wN \str_decode_utf_xvi_extra>NNw
1464     \fi:
1465     #1#2 \s_t1
1466     \int_eval:n { "100 * \str_tmp:w #1#2 + \str_tmp:w #2#1 } \s_t1
1467     \str_decode_utf_xvi_pair:NN
1468   }
1469   \cs_new:Npn \str_decode_utf_xvi_quad>NNwNN
1470   #1#2 #3 \str_decode_utf_xvi_pair:NN #4#5
1471   {
1472     \if_meaning:w \q_nil #5
1473       \str_decode_utf_xvi_error:nNN { missing } #1#2
1474       \str_decode_utf_xvi_pair_end:Nw #4
1475     \fi:
1476     \if_num:w
1477       \if_num:w \str_tmp:w #4#5 < "DC \exp_stop_f:
1478         \c_zero = \c_one
1479       \else:
1480         \str_tmp:w #4#5 < "E0 \exp_stop_f:
1481       \fi:
1482       #1 #2 #4 #5 \s_t1
1483       \int_eval:n
1484       {
1485         ( "100 * \str_tmp:w #1#2 + \str_tmp:w #2#1 - "D800 ) * "400
1486         + "100 * \str_tmp:w #4#5 + \str_tmp:w #5#4 - "DC00
1487       }
1488       \s_t1
1489       \exp_after:wN \use_i:nnn
1490     \else:
1491       \str_decode_utf_xvi_error:nNN { missing } #1#2
1492     \fi:
1493     \str_decode_utf_xvi_pair:NN #4#5

```

```

1494     }
1495   \cs_new:Npn \str_decode_utf_xvi_pair_end:Nw #1 \fi:
1496   {
1497     \fi:
1498     \if_meaning:w \q_nil #1
1499     \else:
1500       \str_decode_utf_xvi_error:nNN { end } #1 \prg_do_nothing:
1501     \fi:
1502     \prg_map_break:
1503   }
1504   \cs_new:Npn \str_decode_utf_xvi_extra>NNw #1#2 \s_tl #3 \s_tl
1505   { \str_decode_utf_xvi_error:nNN { extra } #1#2 }
1506   \cs_new:Npn \str_decode_utf_xvi_error:nNN #1#2#3
1507   {
1508     \flag_raise:n { str_error }
1509     \flag_raise:n { str_#1 }
1510     #2 #3 \s_tl
1511     \int_use:N \c_str_replacement_char_int \s_tl
1512   }
(End definition for \str_decode_utf_xvi_pair:NN, \str_decode_utf_xvi_quad:NNwNN, and \str_decode_utf_xvi_pair:
These functions are documented on page ??.)

```

Restore the original catcodes of bytes 254 and 255.

```

1513 \group_end:
1514 
```

### 7.8.3 utf-32 support

The definitions are done in a category code regime where the bytes 0, 254 and 255 used by the byte order mark have catcode “other”.

```

1515 /*utf32*/
1516 \group_begin:
1517   \char_set_catcode_other:N \^00
1518   \char_set_catcode_other:N \^fe
1519   \char_set_catcode_other:N \^ff

```

\str\_convert\_encode\_utf32: Convert each integer in the comma-list \g\_str\_result\_tl to a sequence of four bytes.  
\str\_convert\_encode\_utf32be: The functions for big-endian and little-endian encodings are very similar, but the \str\_output\_byte:n instructions are reversed.

```

\str_encode_utf_xxxii_be:n
  \str_encode_utf_xxxii_be_aux:nn
\str_encode_utf_xxxii_le:n
  \str_encode_utf_xxxii_le_aux:nn
1520   \cs_new_protected_nopar:cpn { str_convert_encode_utf32: }
1521   {
1522     \str_convert_gmap_internal:N \str_encode_utf_xxxii_be:n
1523     \tl_gput_left:Nx \g_str_result_tl { ^00 ^00 ^fe ^ff }
1524   }
1525   \cs_new_protected_nopar:cpn { str_convert_encode_utf32be: }
1526   { \str_convert_gmap_internal:N \str_encode_utf_xxxii_be:n }
1527   \cs_new_protected_nopar:cpn { str_convert_encode_utf32le: }
1528   { \str_convert_gmap_internal:N \str_encode_utf_xxxii_le:n }
1529   \cs_new:Npn \str_encode_utf_xxxii_be:n #1
1530   {

```

```

1531     \exp_args:Nf \str_encode_utf_XXXii_be_aux:nn
1532         { \int_div_truncate:nn {#1} { "100" } {#1}
1533     }
1534 \cs_new:Npn \str_encode_utf_XXXii_be_aux:nn #1#2
1535     {
1536         ^^00
1537         \str_output_byte_pair_be:n {#1}
1538         \str_output_byte:n { #2 - #1 * "100" }
1539     }
1540 \cs_new:Npn \str_encode_utf_XXXii_le:n #1
1541     {
1542         \exp_args:Nf \str_encode_utf_XXXii_le_aux:nn
1543             { \int_div_truncate:nn {#1} { "100" } {#1}
1544     }
1545 \cs_new:Npn \str_encode_utf_XXXii_le_aux:nn #1#2
1546     {
1547         \str_output_byte:n { #2 - #1 * "100" }
1548         \str_output_byte_pair_le:n {#1}
1549         ^^00
1550     }

```

(End definition for `\str_convert_encode_utf32:`, `\str_convert_encode_utf32be:`, and `\str_convert_encode_utf32le`.  
These functions are documented on page ??.)

- str\_overflow** There can be no error when encoding in UTF-32. When decoding, the string may not have length  $4n$ , or it may contain code points larger than "10FFFF". The latter case often happens if the encoding was in fact not UTF-32, because most arbitrary strings are not valid in UTF-32.

```

1551 \flag_clear_new:n { str_overflow }
1552 \flag_clear_new:n { str_end }
1553 \msg_kernel_new:nnnn { str } { utf32-decode }
1554 {
    Invalid~UTF-32~string: \exp_last_unbraced:Nf \use_none:n
    \str_if_flag_times:nT { str_overflow } { ,~code~point~too~large }
    \str_if_flag_times:nT { str_end } { ,~truncated~string }
1555 .
1556 }
1557 {
    In~the~UTF-32~encoding,~every~Unicode~character~
    (in~the~range~[U+0000,~U+10FFFF])~is~encoded~as~4~bytes.
    \flag_if_raised:nT { str_overflow }
    {
        \\\\
        LaTeX~came~across~a~code~point~larger~than~1114111,~
        the~maximum~code~point~defined~by~Unicode.~
        Perhaps~the~string~was~not~encoded~in~the~UTF-32~encoding?
    }
    \flag_if_raised:nT { str_end }
    {
        \\\\

```

```

1573     The~length~of~the~string~is~not~a~multiple~of~4.~  

1574     Perhaps~the~string~was~truncated?  

1575   }  

1576 }

```

(End definition for `str_overflow` and `str_end`. These functions are documented on page ??.)

`\str_convert_decode_utf32:` The structure is similar to UTF-16 decoding functions. If the endianness is not given, test the first 4 bytes of the string (possibly `\s_stop` if the string is too short) for the presence of a byte-order mark. If there is a byte-order mark, use that endianness, and remove the 4 bytes, otherwise default to big-endian, and leave the 4 bytes in place. The `\str_decode_utf_xxxii:Nw` auxiliary receives 1 or 2 as its first argument indicating endianness, and the string to convert as its second argument (expanded or not). It sets `\str_tmp:w` to expand to the character code of either of its two arguments depending on endianness, then triggers the `_loop` auxiliary inside an x-expanding assignment to `\g_str_result_t1`.

`\str_convert_decode_utf32be:`

`\str_convert_decode_utf32le:`

`\str_decode_utf_xxxii_bom:NNNN`

`\str_decode_utf_xxxii:Nw`

`\str_decode_utf_xxxii_loop:NNNN`

`\str_decode_utf_xxxii_end:w`

The `_loop` auxiliary first checks for the end-of-string marker `\s_stop`, calling the `_end` auxiliary if appropriate. Otherwise, leave the *<4 bytes>* `\s_t1` behind, then check that the code point is not overflowing: the leading byte must be 0, and the following byte at most 16.

In the ending code, we check that there remains no byte: there should be nothing left until the first `\s_stop`. Break the map.

```

1577 \cs_new_protected_nopar:cpn { str_convert_decode_utf32be: }  

1578   { \str_decode_utf_xxxii:Nw 1 \g_str_result_t1 \s_stop }  

1579 \cs_new_protected_nopar:cpn { str_convert_decode_utf32le: }  

1580   { \str_decode_utf_xxxii:Nw 2 \g_str_result_t1 \s_stop }  

1581 \cs_new_protected_nopar:cpn { str_convert_decode_utf32: }  

1582   {  

1583     \exp_after:wN \str_decode_utf_xxxii_bom:NNNN \g_str_result_t1  

1584     \s_stop \s_stop \s_stop \s_stop  

1585   }  

1586 \cs_new_protected:Npn \str_decode_utf_xxxii_bom:NNNN #1#2#3#4  

1587   {  

1588     \str_if_eq:xxTF { #1#2#3#4 } { ^~ff ^~fe ^~00 ^~00 }  

1589       { \str_decode_utf_xxxii:Nw 2 }  

1590       {  

1591         \str_if_eq:xxTF { #1#2#3#4 } { ^~00 ^~00 ^~fe ^~ff }  

1592           { \str_decode_utf_xxxii:Nw 1 }  

1593           { \str_decode_utf_xxxii:Nw 1 #1#2#3#4 }  

1594       }  

1595   }  

1596 \cs_new_protected:Npn \str_decode_utf_xxxii:Nw #1#2 \s_stop  

1597   {  

1598     \flag_clear:n { str_overflow }  

1599     \flag_clear:n { str_end }  

1600     \flag_clear:n { str_error }  

1601     \cs_set:Npn \str_tmp:w ##1 ##2 { ' ## #1 }  

1602     \t1_gset:Nx \g_str_result_t1  

1603     {

```

```

1604         \exp_after:wN \str_decode_utf_xxxii_loop:NNNN
1605             #2 \s_stop \s_stop \s_stop \s_stop
1606             \prg_break_point:n { }
1607         }
1608     \str_if_flag_error:nx { str_error } { utf32-decode } { }
1609   }
1610 \cs_new:Npn \str_decode_utf_xxxii_loop:NNNN #1#2#3#4
1611   {
1612     \if_meaning:w \s_stop #4
1613       \exp_after:wN \str_decode_utf_xxxii_end:w
1614     \fi:
1615     #1#2#3#4 \s_tl
1616     \if_num:w \str_tmp:w #1#4 > \c_zero
1617       \flag_raise:n { str_overflow }
1618       \flag_raise:n { str_error }
1619       \int_use:N \c_str_replacement_char_int
1620     \else:
1621       \if_num:w \str_tmp:w #2#3 > \c_sixteen
1622         \flag_raise:n { str_overflow }
1623         \flag_raise:n { str_error }
1624         \int_use:N \c_str_replacement_char_int
1625     \else:
1626       \int_eval:n
1627         { \str_tmp:w #2#3*"10000 + \str_tmp:w #3#2*"100 + \str_tmp:w #4#1 }
1628     \fi:
1629   \fi:
1630   \s_tl
1631   \str_decode_utf_xxxii_loop:NNNN
1632 }
1633 \cs_new:Npn \str_decode_utf_xxxii_end:w #1 \s_stop
1634   {
1635     \tl_if_empty:nF {#1}
1636     {
1637       \flag_raise:n { str_end }
1638       \flag_raise:n { str_error }
1639       #1 \s_tl
1640       \int_use:N \c_str_replacement_char_int \s_tl
1641     }
1642     \prg_map_break:
1643   }
(End definition for \str_convert_decode_utf32:, \str_convert_decode_utf32be:, and \str_convert_decode_utf32le
These functions are documented on page ??.)
Restore the original catcodes of bytes 0, 254 and 255.

1644 \group_end:
1645 
```

#### 7.8.4 iso 8859 support

The ISO-8859-1 encoding exactly matches with the 256 first Unicode characters. For other 8-bit encodings of the ISO-8859 family, we keep track only of differences, and of unassigned bytes.

```
1646 /*iso88591)
1647 \str_declare_eight_bit_encoding:nnn { iso88591 }
1648 {
1649 }
1650 {
1651 }
1652 /iso88591)

1653 /*iso88592)
1654 \str_declare_eight_bit_encoding:nnn { iso88592 }
1655 {
1656   { A1 } { 0104 }
1657   { A2 } { 02D8 }
1658   { A3 } { 0141 }
1659   { A5 } { 013D }
1660   { A6 } { 015A }
1661   { A9 } { 0160 }
1662   { AA } { 015E }
1663   { AB } { 0164 }
1664   { AC } { 0179 }
1665   { AE } { 017D }
1666   { AF } { 017B }
1667   { B1 } { 0105 }
1668   { B2 } { 02DB }
1669   { B3 } { 0142 }
1670   { B5 } { 013E }
1671   { B6 } { 015B }
1672   { B7 } { 02C7 }
1673   { B9 } { 0161 }
1674   { BA } { 015F }
1675   { BB } { 0165 }
1676   { BC } { 017A }
1677   { BD } { 02DD }
1678   { BE } { 017E }
1679   { BF } { 017C }
1680   { CO } { 0154 }
1681   { C3 } { 0102 }
1682   { C5 } { 0139 }
1683   { C6 } { 0106 }
1684   { C8 } { 010C }
1685   { CA } { 0118 }
1686   { CC } { 011A }
1687   { CF } { 010E }
1688   { DO } { 0110 }
1689   { D1 } { 0143 }
```

```

1690 { D2 } { 0147 }
1691 { D5 } { 0150 }
1692 { D8 } { 0158 }
1693 { D9 } { 016E }
1694 { DB } { 0170 }
1695 { DE } { 0162 }
1696 { EO } { 0155 }
1697 { E3 } { 0103 }
1698 { E5 } { 013A }
1699 { E6 } { 0107 }
1700 { E8 } { 010D }
1701 { EA } { 0119 }
1702 { EC } { 011B }
1703 { EF } { 010F }
1704 { FO } { 0111 }
1705 { F1 } { 0144 }
1706 { F2 } { 0148 }
1707 { F5 } { 0151 }
1708 { F8 } { 0159 }
1709 { F9 } { 016F }
1710 { FB } { 0171 }
1711 { FE } { 0163 }
1712 { FF } { 02D9 }

1713 }
1714 {
1715 }
1716 </iso88592>

1717 /*iso88593)
1718 \str_declare_eight_bit_encoding:nnn { iso88593 }
1719 {
1720 { A1 } { 0126 }
1721 { A2 } { 02D8 }
1722 { A6 } { 0124 }
1723 { A9 } { 0130 }
1724 { AA } { 015E }
1725 { AB } { 011E }
1726 { AC } { 0134 }
1727 { AF } { 017B }
1728 { B1 } { 0127 }
1729 { B6 } { 0125 }
1730 { B9 } { 0131 }
1731 { BA } { 015F }
1732 { BB } { 011F }
1733 { BC } { 0135 }
1734 { BF } { 017C }
1735 { C5 } { 010A }
1736 { C6 } { 0108 }
1737 { D5 } { 0120 }
1738 { D8 } { 011C }

```

```

1739 { DD } { 016C }
1740 { DE } { 015C }
1741 { E5 } { 010B }
1742 { E6 } { 0109 }
1743 { F5 } { 0121 }
1744 { F8 } { 011D }
1745 { FD } { 016D }
1746 { FE } { 015D }
1747 { FF } { 02D9 }
1748 }
1749 {
1750 { A5 }
1751 { AE }
1752 { BE }
1753 { C3 }
1754 { DO }
1755 { E3 }
1756 { FO }
1757 }
1758 </iso88593>
1759 (*iso88594)
1760 \str_declare_eight_bit_encoding:nnn { iso88594 }
1761 {
1762 { A1 } { 0104 }
1763 { A2 } { 0138 }
1764 { A3 } { 0156 }
1765 { A5 } { 0128 }
1766 { A6 } { 013B }
1767 { A9 } { 0160 }
1768 { AA } { 0112 }
1769 { AB } { 0122 }
1770 { AC } { 0166 }
1771 { AE } { 017D }
1772 { B1 } { 0105 }
1773 { B2 } { 02DB }
1774 { B3 } { 0157 }
1775 { B5 } { 0129 }
1776 { B6 } { 013C }
1777 { B7 } { 02C7 }
1778 { B9 } { 0161 }
1779 { BA } { 0113 }
1780 { BB } { 0123 }
1781 { BC } { 0167 }
1782 { BD } { 014A }
1783 { BE } { 017E }
1784 { BF } { 014B }
1785 { CO } { 0100 }
1786 { C7 } { 012E }
1787 { C8 } { 010C }

```

```

1788 { CA } { 0118 }
1789 { CC } { 0116 }
1790 { CF } { 012A }
1791 { DO } { 0110 }
1792 { D1 } { 0145 }
1793 { D2 } { 014C }
1794 { D3 } { 0136 }
1795 { D9 } { 0172 }
1796 { DD } { 0168 }
1797 { DE } { 016A }
1798 { EO } { 0101 }
1799 { E7 } { 012F }
1800 { E8 } { 010D }
1801 { EA } { 0119 }
1802 { EC } { 0117 }
1803 { EF } { 012B }
1804 { FO } { 0111 }
1805 { F1 } { 0146 }
1806 { F2 } { 014D }
1807 { F3 } { 0137 }
1808 { F9 } { 0173 }
1809 { FD } { 0169 }
1810 { FE } { 016B }
1811 { FF } { 02D9 }
1812 }
1813 {
1814 }
1815 </iso88594>
1816 /*iso88595
1817 \str_declare_eight_bit_encoding:nnn { iso88595 }
1818 {
1819 { A1 } { 0401 }
1820 { A2 } { 0402 }
1821 { A3 } { 0403 }
1822 { A4 } { 0404 }
1823 { A5 } { 0405 }
1824 { A6 } { 0406 }
1825 { A7 } { 0407 }
1826 { A8 } { 0408 }
1827 { A9 } { 0409 }
1828 { AA } { 040A }
1829 { AB } { 040B }
1830 { AC } { 040C }
1831 { AE } { 040E }
1832 { AF } { 040F }
1833 { B0 } { 0410 }
1834 { B1 } { 0411 }
1835 { B2 } { 0412 }
1836 { B3 } { 0413 }

```

1837	{ B4 } { 0414 }
1838	{ B5 } { 0415 }
1839	{ B6 } { 0416 }
1840	{ B7 } { 0417 }
1841	{ B8 } { 0418 }
1842	{ B9 } { 0419 }
1843	{ BA } { 041A }
1844	{ BB } { 041B }
1845	{ BC } { 041C }
1846	{ BD } { 041D }
1847	{ BE } { 041E }
1848	{ BF } { 041F }
1849	{ CO } { 0420 }
1850	{ C1 } { 0421 }
1851	{ C2 } { 0422 }
1852	{ C3 } { 0423 }
1853	{ C4 } { 0424 }
1854	{ C5 } { 0425 }
1855	{ C6 } { 0426 }
1856	{ C7 } { 0427 }
1857	{ C8 } { 0428 }
1858	{ C9 } { 0429 }
1859	{ CA } { 042A }
1860	{ CB } { 042B }
1861	{ CC } { 042C }
1862	{ CD } { 042D }
1863	{ CE } { 042E }
1864	{ CF } { 042F }
1865	{ DO } { 0430 }
1866	{ D1 } { 0431 }
1867	{ D2 } { 0432 }
1868	{ D3 } { 0433 }
1869	{ D4 } { 0434 }
1870	{ D5 } { 0435 }
1871	{ D6 } { 0436 }
1872	{ D7 } { 0437 }
1873	{ D8 } { 0438 }
1874	{ D9 } { 0439 }
1875	{ DA } { 043A }
1876	{ DB } { 043B }
1877	{ DC } { 043C }
1878	{ DD } { 043D }
1879	{ DE } { 043E }
1880	{ DF } { 043F }
1881	{ EO } { 0440 }
1882	{ E1 } { 0441 }
1883	{ E2 } { 0442 }
1884	{ E3 } { 0443 }
1885	{ E4 } { 0444 }
1886	{ E5 } { 0445 }

```

1887      { E6 } { 0446 }
1888      { E7 } { 0447 }
1889      { E8 } { 0448 }
1890      { E9 } { 0449 }
1891      { EA } { 044A }
1892      { EB } { 044B }
1893      { EC } { 044C }
1894      { ED } { 044D }
1895      { EE } { 044E }
1896      { EF } { 044F }
1897      { F0 } { 2116 }
1898      { F1 } { 0451 }
1899      { F2 } { 0452 }
1900      { F3 } { 0453 }
1901      { F4 } { 0454 }
1902      { F5 } { 0455 }
1903      { F6 } { 0456 }
1904      { F7 } { 0457 }
1905      { F8 } { 0458 }
1906      { F9 } { 0459 }
1907      { FA } { 045A }
1908      { FB } { 045B }
1909      { FC } { 045C }
1910      { FD } { 00A7 }
1911      { FE } { 045E }
1912      { FF } { 045F }
1913  }
1914  {
1915  }
1916 </iso88595>
1917 /*iso88596)
1918 \str_declare_eight_bit_encoding:nnn { iso88596 }
1919  {
1920      { AC } { 060C }
1921      { BB } { 061B }
1922      { BF } { 061F }
1923      { C1 } { 0621 }
1924      { C2 } { 0622 }
1925      { C3 } { 0623 }
1926      { C4 } { 0624 }
1927      { C5 } { 0625 }
1928      { C6 } { 0626 }
1929      { C7 } { 0627 }
1930      { C8 } { 0628 }
1931      { C9 } { 0629 }
1932      { CA } { 062A }
1933      { CB } { 062B }
1934      { CC } { 062C }
1935      { CD } { 062D }

```

```

1936    { CE } { 062E }
1937    { CF } { 062F }
1938    { DO } { 0630 }
1939    { D1 } { 0631 }
1940    { D2 } { 0632 }
1941    { D3 } { 0633 }
1942    { D4 } { 0634 }
1943    { D5 } { 0635 }
1944    { D6 } { 0636 }
1945    { D7 } { 0637 }
1946    { D8 } { 0638 }
1947    { D9 } { 0639 }
1948    { DA } { 063A }
1949    { EO } { 0640 }
1950    { E1 } { 0641 }
1951    { E2 } { 0642 }
1952    { E3 } { 0643 }
1953    { E4 } { 0644 }
1954    { E5 } { 0645 }
1955    { E6 } { 0646 }
1956    { E7 } { 0647 }
1957    { E8 } { 0648 }
1958    { E9 } { 0649 }
1959    { EA } { 064A }
1960    { EB } { 064B }
1961    { EC } { 064C }
1962    { ED } { 064D }
1963    { EE } { 064E }
1964    { EF } { 064F }
1965    { FO } { 0650 }
1966    { F1 } { 0651 }
1967    { F2 } { 0652 }
1968    }
1969    {
1970    { A1 }
1971    { A2 }
1972    { A3 }
1973    { A5 }
1974    { A6 }
1975    { A7 }
1976    { A8 }
1977    { A9 }
1978    { AA }
1979    { AB }
1980    { AE }
1981    { AF }
1982    { B0 }
1983    { B1 }
1984    { B2 }
1985    { B3 }

```

```

1986 { B4 }
1987 { B5 }
1988 { B6 }
1989 { B7 }
1990 { B8 }
1991 { B9 }
1992 { BA }
1993 { BC }
1994 { BD }
1995 { BE }
1996 { CO }
1997 { DB }
1998 { DC }
1999 { DD }
2000 { DE }
2001 { DF }
2002 }
2003 </iso88596>
2004 <*iso88597>
2005 \str_declare_eight_bit_encoding:nnn { iso88597 }
2006 {
2007 { A1 } { 2018 }
2008 { A2 } { 2019 }
2009 { A4 } { 20AC }
2010 { A5 } { 20AF }
2011 { AA } { 037A }
2012 { AF } { 2015 }
2013 { B4 } { 0384 }
2014 { B5 } { 0385 }
2015 { B6 } { 0386 }
2016 { B8 } { 0388 }
2017 { B9 } { 0389 }
2018 { BA } { 038A }
2019 { BC } { 038C }
2020 { BE } { 038E }
2021 { BF } { 038F }
2022 { CO } { 0390 }
2023 { C1 } { 0391 }
2024 { C2 } { 0392 }
2025 { C3 } { 0393 }
2026 { C4 } { 0394 }
2027 { C5 } { 0395 }
2028 { C6 } { 0396 }
2029 { C7 } { 0397 }
2030 { C8 } { 0398 }
2031 { C9 } { 0399 }
2032 { CA } { 039A }
2033 { CB } { 039B }
2034 { CC } { 039C }

```

```
2035 { CD } { 039D }
2036 { CE } { 039E }
2037 { CF } { 039F }
2038 { DO } { 03A0 }
2039 { D1 } { 03A1 }
2040 { D3 } { 03A3 }
2041 { D4 } { 03A4 }
2042 { D5 } { 03A5 }
2043 { D6 } { 03A6 }
2044 { D7 } { 03A7 }
2045 { D8 } { 03A8 }
2046 { D9 } { 03A9 }
2047 { DA } { 03AA }
2048 { DB } { 03AB }
2049 { DC } { 03AC }
2050 { DD } { 03AD }
2051 { DE } { 03AE }
2052 { DF } { 03AF }
2053 { EO } { 03B0 }
2054 { E1 } { 03B1 }
2055 { E2 } { 03B2 }
2056 { E3 } { 03B3 }
2057 { E4 } { 03B4 }
2058 { E5 } { 03B5 }
2059 { E6 } { 03B6 }
2060 { E7 } { 03B7 }
2061 { E8 } { 03B8 }
2062 { E9 } { 03B9 }
2063 { EA } { 03BA }
2064 { EB } { 03BB }
2065 { EC } { 03BC }
2066 { ED } { 03BD }
2067 { EE } { 03BE }
2068 { EF } { 03BF }
2069 { FO } { 03C0 }
2070 { F1 } { 03C1 }
2071 { F2 } { 03C2 }
2072 { F3 } { 03C3 }
2073 { F4 } { 03C4 }
2074 { F5 } { 03C5 }
2075 { F6 } { 03C6 }
2076 { F7 } { 03C7 }
2077 { F8 } { 03C8 }
2078 { F9 } { 03C9 }
2079 { FA } { 03CA }
2080 { FB } { 03CB }
2081 { FC } { 03CC }
2082 { FD } { 03CD }
2083 { FE } { 03CE }
2084 }
```

```

2085      {
2086          { AE }
2087          { D2 }
2088      }
2089  </iso88597>
2090 {*iso88598}
2091 \str_declare_eight_bit_encoding:nnn { iso88598 }
2092  {
2093      { AA } { 00D7 }
2094      { BA } { 00F7 }
2095      { DF } { 2017 }
2096      { E0 } { 05D0 }
2097      { E1 } { 05D1 }
2098      { E2 } { 05D2 }
2099      { E3 } { 05D3 }
2100      { E4 } { 05D4 }
2101      { E5 } { 05D5 }
2102      { E6 } { 05D6 }
2103      { E7 } { 05D7 }
2104      { E8 } { 05D8 }
2105      { E9 } { 05D9 }
2106      { EA } { 05DA }
2107      { EB } { 05DB }
2108      { EC } { 05DC }
2109      { ED } { 05DD }
2110      { EE } { 05DE }
2111      { EF } { 05DF }
2112      { F0 } { 05E0 }
2113      { F1 } { 05E1 }
2114      { F2 } { 05E2 }
2115      { F3 } { 05E3 }
2116      { F4 } { 05E4 }
2117      { F5 } { 05E5 }
2118      { F6 } { 05E6 }
2119      { F7 } { 05E7 }
2120      { F8 } { 05E8 }
2121      { F9 } { 05E9 }
2122      { FA } { 05EA }
2123      { FD } { 200E }
2124      { FE } { 200F }
2125  }
2126  {
2127      { A1 }
2128      { BF }
2129      { C0 }
2130      { C1 }
2131      { C2 }
2132      { C3 }
2133      { C4 }

```

```

2134      { C5 }
2135      { C6 }
2136      { C7 }
2137      { C8 }
2138      { C9 }
2139      { CA }
2140      { CB }
2141      { CC }
2142      { CD }
2143      { CE }
2144      { CF }
2145      { D0 }
2146      { D1 }
2147      { D2 }
2148      { D3 }
2149      { D4 }
2150      { D5 }
2151      { D6 }
2152      { D7 }
2153      { D8 }
2154      { D9 }
2155      { DA }
2156      { DB }
2157      { DC }
2158      { DD }
2159      { DE }
2160      { FB }
2161      { FC }
2162    }
2163 </iso88598>
2164 /*iso88599)
2165 \str_declare_eight_bit_encoding:nnn { iso88599 }
2166  {
2167    { DO } { 011E }
2168    { DD } { 0130 }
2169    { DE } { 015E }
2170    { FO } { 011F }
2171    { FD } { 0131 }
2172    { FE } { 015F }
2173  }
2174  {
2175  }
2176 </iso88599>
2177 /*iso885910)
2178 \str_declare_eight_bit_encoding:nnn { iso885910 }
2179  {
2180    { A1 } { 0104 }
2181    { A2 } { 0112 }
2182    { A3 } { 0122 }

```

```

2183 { A4 } { 012A }
2184 { A5 } { 0128 }
2185 { A6 } { 0136 }
2186 { A8 } { 013B }
2187 { A9 } { 0110 }
2188 { AA } { 0160 }
2189 { AB } { 0166 }
2190 { AC } { 017D }
2191 { AE } { 016A }
2192 { AF } { 014A }
2193 { B1 } { 0105 }
2194 { B2 } { 0113 }
2195 { B3 } { 0123 }
2196 { B4 } { 012B }
2197 { B5 } { 0129 }
2198 { B6 } { 0137 }
2199 { B8 } { 013C }
2200 { B9 } { 0111 }
2201 { BA } { 0161 }
2202 { BB } { 0167 }
2203 { BC } { 017E }
2204 { BD } { 2015 }
2205 { BE } { 016B }
2206 { BF } { 014B }
2207 { C0 } { 0100 }
2208 { C7 } { 012E }
2209 { C8 } { 010C }
2210 { CA } { 0118 }
2211 { CC } { 0116 }
2212 { D1 } { 0145 }
2213 { D2 } { 014C }
2214 { D7 } { 0168 }
2215 { D9 } { 0172 }
2216 { E0 } { 0101 }
2217 { E7 } { 012F }
2218 { E8 } { 010D }
2219 { EA } { 0119 }
2220 { EC } { 0117 }
2221 { F1 } { 0146 }
2222 { F2 } { 014D }
2223 { F7 } { 0169 }
2224 { F9 } { 0173 }
2225 { FF } { 0138 }

2226 }
2227 {
2228 }
2229 </iso885910>
2230 {*iso885911}
2231 \str_declare_eight_bit_encoding:nnn { iso885911 }

```

```

2232    {
2233    { A1 } { OE01 }
2234    { A2 } { OE02 }
2235    { A3 } { OE03 }
2236    { A4 } { OE04 }
2237    { A5 } { OE05 }
2238    { A6 } { OE06 }
2239    { A7 } { OE07 }
2240    { A8 } { OE08 }
2241    { A9 } { OE09 }
2242    { AA } { OEOA }
2243    { AB } { OEOB }
2244    { AC } { OEOC }
2245    { AD } { OEOD }
2246    { AE } { OEOF }
2247    { AF } { OEOF }
2248    { BO } { OE10 }
2249    { B1 } { OE11 }
2250    { B2 } { OE12 }
2251    { B3 } { OE13 }
2252    { B4 } { OE14 }
2253    { B5 } { OE15 }
2254    { B6 } { OE16 }
2255    { B7 } { OE17 }
2256    { B8 } { OE18 }
2257    { B9 } { OE19 }
2258    { BA } { OE1A }
2259    { BB } { OE1B }
2260    { BC } { OE1C }
2261    { BD } { OE1D }
2262    { BE } { OE1E }
2263    { BF } { OE1F }
2264    { CO } { OE20 }
2265    { C1 } { OE21 }
2266    { C2 } { OE22 }
2267    { C3 } { OE23 }
2268    { C4 } { OE24 }
2269    { C5 } { OE25 }
2270    { C6 } { OE26 }
2271    { C7 } { OE27 }
2272    { C8 } { OE28 }
2273    { C9 } { OE29 }
2274    { CA } { OE2A }
2275    { CB } { OE2B }
2276    { CC } { OE2C }
2277    { CD } { OE2D }
2278    { CE } { OE2E }
2279    { CF } { OE2F }
2280    { DO } { OE30 }
2281    { D1 } { OE31 }

```

```

2282 { D2 } { OE32 }
2283 { D3 } { OE33 }
2284 { D4 } { OE34 }
2285 { D5 } { OE35 }
2286 { D6 } { OE36 }
2287 { D7 } { OE37 }
2288 { D8 } { OE38 }
2289 { D9 } { OE39 }
2290 { DA } { OE3A }
2291 { DF } { OE3F }
2292 { EO } { OE40 }
2293 { E1 } { OE41 }
2294 { E2 } { OE42 }
2295 { E3 } { OE43 }
2296 { E4 } { OE44 }
2297 { E5 } { OE45 }
2298 { E6 } { OE46 }
2299 { E7 } { OE47 }
2300 { E8 } { OE48 }
2301 { E9 } { OE49 }
2302 { EA } { OE4A }
2303 { EB } { OE4B }
2304 { EC } { OE4C }
2305 { ED } { OE4D }
2306 { EE } { OE4E }
2307 { EF } { OE4F }
2308 { F0 } { OE50 }
2309 { F1 } { OE51 }
2310 { F2 } { OE52 }
2311 { F3 } { OE53 }
2312 { F4 } { OE54 }
2313 { F5 } { OE55 }
2314 { F6 } { OE56 }
2315 { F7 } { OE57 }
2316 { F8 } { OE58 }
2317 { F9 } { OE59 }
2318 { FA } { OE5A }
2319 { FB } { OE5B }
2320 }
2321 {
2322 { DB }
2323 { DC }
2324 { DD }
2325 { DE }
2326 }
2327 </iso885911>
2328 {*iso885913}
2329 \str_declare_eight_bit_encoding:nnn { iso885913 }
2330 {

```

2331 { A1 } { 201D }  
2332 { A5 } { 201E }  
2333 { A8 } { 00D8 }  
2334 { AA } { 0156 }  
2335 { AF } { 00C6 }  
2336 { B4 } { 201C }  
2337 { B8 } { 00F8 }  
2338 { BA } { 0157 }  
2339 { BF } { 00E6 }  
2340 { CO } { 0104 }  
2341 { C1 } { 012E }  
2342 { C2 } { 0100 }  
2343 { C3 } { 0106 }  
2344 { C6 } { 0118 }  
2345 { C7 } { 0112 }  
2346 { C8 } { 010C }  
2347 { CA } { 0179 }  
2348 { CB } { 0116 }  
2349 { CC } { 0122 }  
2350 { CD } { 0136 }  
2351 { CE } { 012A }  
2352 { CF } { 013B }  
2353 { DO } { 0160 }  
2354 { D1 } { 0143 }  
2355 { D2 } { 0145 }  
2356 { D4 } { 014C }  
2357 { D8 } { 0172 }  
2358 { D9 } { 0141 }  
2359 { DA } { 015A }  
2360 { DB } { 016A }  
2361 { DD } { 017B }  
2362 { DE } { 017D }  
2363 { EO } { 0105 }  
2364 { E1 } { 012F }  
2365 { E2 } { 0101 }  
2366 { E3 } { 0107 }  
2367 { E6 } { 0119 }  
2368 { E7 } { 0113 }  
2369 { E8 } { 010D }  
2370 { EA } { 017A }  
2371 { EB } { 0117 }  
2372 { EC } { 0123 }  
2373 { ED } { 0137 }  
2374 { EE } { 012B }  
2375 { EF } { 013C }  
2376 { FO } { 0161 }  
2377 { F1 } { 0144 }  
2378 { F2 } { 0146 }  
2379 { F4 } { 014D }  
2380 { F8 } { 0173 }

```

2381 { F9 } { 0142 }
2382 { FA } { 015B }
2383 { FB } { 016B }
2384 { FD } { 017C }
2385 { FE } { 017E }
2386 { FF } { 2019 }
2387 }
2388 {
2389 }
2390 (/iso885913)
2391 (*iso885914)
2392 \str_declare_eight_bit_encoding:nnn { iso885914 }
2393 {
2394 { A1 } { 1E02 }
2395 { A2 } { 1E03 }
2396 { A4 } { 010A }
2397 { A5 } { 010B }
2398 { A6 } { 1EOA }
2399 { A8 } { 1E80 }
2400 { AA } { 1E82 }
2401 { AB } { 1EOB }
2402 { AC } { 1EF2 }
2403 { AF } { 0178 }
2404 { B0 } { 1E1E }
2405 { B1 } { 1E1F }
2406 { B2 } { 0120 }
2407 { B3 } { 0121 }
2408 { B4 } { 1E40 }
2409 { B5 } { 1E41 }
2410 { B7 } { 1E56 }
2411 { B8 } { 1E81 }
2412 { B9 } { 1E57 }
2413 { BA } { 1E83 }
2414 { BB } { 1E60 }
2415 { BC } { 1EF3 }
2416 { BD } { 1E84 }
2417 { BE } { 1E85 }
2418 { BF } { 1E61 }
2419 { D0 } { 0174 }
2420 { D7 } { 1E6A }
2421 { DE } { 0176 }
2422 { F0 } { 0175 }
2423 { F7 } { 1E6B }
2424 { FE } { 0177 }
2425 }
2426 {
2427 }
2428 (/iso885914)
2429 (*iso885915)

```

```

2430 \str_declare_eight_bit_encoding:nnn { iso885915 }
2431 {
2432   { A4 } { 20AC }
2433   { A6 } { 0160 }
2434   { A8 } { 0161 }
2435   { B4 } { 017D }
2436   { B8 } { 017E }
2437   { BC } { 0152 }
2438   { BD } { 0153 }
2439   { BE } { 0178 }
2440 }
2441 {
2442 }
2443 
```

2443

```

2444 (*iso885916)
2445 \str_declare_eight_bit_encoding:nnn { iso885916 }
2446 {
2447   { A1 } { 0104 }
2448   { A2 } { 0105 }
2449   { A3 } { 0141 }
2450   { A4 } { 20AC }
2451   { A5 } { 201E }
2452   { A6 } { 0160 }
2453   { A8 } { 0161 }
2454   { AA } { 0218 }
2455   { AC } { 0179 }
2456   { AE } { 017A }
2457   { AF } { 017B }
2458   { B2 } { 010C }
2459   { B3 } { 0142 }
2460   { B4 } { 017D }
2461   { B5 } { 201D }
2462   { B8 } { 017E }
2463   { B9 } { 010D }
2464   { BA } { 0219 }
2465   { BC } { 0152 }
2466   { BD } { 0153 }
2467   { BE } { 0178 }
2468   { BF } { 017C }
2469   { C3 } { 0102 }
2470   { C5 } { 0106 }
2471   { D0 } { 0110 }
2472   { D1 } { 0143 }
2473   { D5 } { 0150 }
2474   { D7 } { 015A }
2475   { D8 } { 0170 }
2476   { DD } { 0118 }
2477   { DE } { 021A }
2478   { E3 } { 0103 }

```

```

2479 { E5 } { 0107 }
2480 { F0 } { 0111 }
2481 { F1 } { 0144 }
2482 { F5 } { 0151 }
2483 { F7 } { 015B }
2484 { F8 } { 0171 }
2485 { FD } { 0119 }
2486 { FE } { 021B }
2487 }
2488 {
2489 }
2490 //iso885916)

```

## 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.

Symbols	
\# . . . . .	<i>39</i> , <u>618</u>
\\$ . . . . .	<u>616</u>
\% . . . . .	<i>41</i> , <u>626</u>
\ . . . . .	<i>617</i>
\* . . . . .	<i>62</i> , <u>85</u> , <u>995</u>
\\" . . . . .	<i>36</i> , <u>613</u> , <u>868</u> , <u>869</u> , <u>872</u> , <u>995</u> , <u>1208</u> , <u>1211</u> – <u>1214</u> , <u>1219</u> , <u>1225</u> , <u>1230</u> , <u>1237</u> , <u>1393</u> , <u>1396</u> – <u>1398</u> , <u>1400</u> , <u>1406</u> , <u>1411</u> , <u>1416</u> , <u>1565</u> , <u>1572</u>
\{ . . . . .	<i>37</i> , <u>614</u> , <u>873</u>
\} . . . . .	<i>38</i> , <u>615</u> , <u>873</u>
\^ . . . . .	<i>619</i> , <u>996</u> , <u>997</u> , <u>1337</u> , <u>1338</u> , <u>1517</u> – <u>1519</u>
\~ . . . . .	<i>40</i> , <u>621</u>
\_u . . . . .	<i>62</i> , <u>85</u> , <u>620</u> , <u>873</u> , <u>1211</u> – <u>1213</u>
<b>A</b>	
\A . . . . .	<i>63</i> , <u>86</u>
<b>B</b>	
\bool_gset_false:N . . . . .	<i>496</i> , <u>505</u>
\bool_gset_true:N . . . . .	<u>486</u>
\bool_if:NTF . . . . .	<i>500</i> , <u>509</u>
\bool_new:N . . . . .	<u>58</u>
<b>C</b>	
\c_backslash_str . . . . .	<i>2</i> , <u>36</u> , <u>36</u> , <u>1046</u> , <u>1075</u> , <u>1077</u> , <u>1115</u> , <u>1124</u> , <u>1128</u>
\c_eight . . . . .	<i>153</i> , <u>196</u> , <u>205</u> , <u>1130</u> , <u>1131</u>
\c_fifty_eight . . . . .	<u>26</u> , <u>27</u>
\c_five . . . . .	<u>154</u>
\c_forty_eight . . . . .	<u>26</u> , <u>26</u>
\c_four . . . . .	<u>154</u> , <u>1461</u>
\c_hash_str . . . . .	<i>2</i> , <u>36</u> , <u>39</u> , <u>989</u> , <u>1096</u>
\c_lbrace_str . . . . .	<i>2</i> , <u>36</u> , <u>37</u>
\c_max_char_int . . . . .	<u>7</u> , <u>33</u> , <u>33</u>
\c_max_int . . . . .	<u>295</u>
\c_max_register_int . . . . .	<u>828</u>
\c_minus_one . . . . .	<u>1175</u>
\c_nine . . . . .	<u>127</u> , <u>164</u>
\c_ninety_one . . . . .	<u>26</u> , <u>29</u> , <u>593</u>
\c_ninety_seven . . . . .	<u>26</u> , <u>30</u> , <u>603</u>
\c_one . . . . .	<u>155</u> , <u>362</u> , <u>595</u> , <u>713</u> , <u>780</u> , <u>790</u> , <u>823</u> , <u>905</u> , <u>1054</u> , <u>1058</u> , <u>1478</u>
\c_one_hundred_twenty_seven . . . . .	<u>26</u> , <u>32</u>
\c_one_hundred_twenty_three . . . . .	<u>26</u> , <u>31</u> , <u>602</u>
\c_percent_str . . . . .	<i>2</i> , <u>36</u> , <u>41</u> , <u>991</u> , <u>1153</u>
\c_rbrace_str . . . . .	<i>2</i> , <u>36</u> , <u>38</u>
\c_seven . . . . .	<i>153</i> , <u>221</u> , <u>231</u>
\c_six . . . . .	<u>153</u>
\c_sixteen . . . . .	<u>1621</u>
\c_sixty_five . . . . .	<u>26</u> , <u>28</u> , <u>594</u>
\c_str_byte_-1_t1 . . . . .	<u>389</u>
\c_str_byte_0_t1 . . . . .	<u>389</u>
\c_str_byte_1_t1 . . . . .	<u>389</u>
\c_str_byte_255_t1 . . . . .	<u>389</u>

\c_str_escape_name_not_str . . . . .	1173, 1182, 1256, 1276, 1296, 1319,
. . . . .	1355, 1455, 1469, 1495, 1504, 1506,
\c_str_escape_name_str .	1529, 1534, 1540, 1545, 1610, 1633
\c_str_escape_string_str	\cs_new_eq>NN . . . . . 630, 660, 662
\c_str_positive_bytes_t1 ..	\cs_new_nopar:Npn . . . . . 129, 166,
\c_str_replacement_char_int . . .	180, 242, 278, 358, 415, 423, 429, 1326
35, 389, 409, 726	\cs_new_protected:cpn . . . . . 944
135, 789, 1266, 1290, 1304, 1324,	\cs_new_protected:cpx . . . . . 18
1331, 1363, 1511, 1619, 1624, 1640	\cs_new_protected:Npn . . . 90, 447, 463,
\c_thirty_two . . . . .	479, 485, 487, 511, 529, 537, 547,
\c_three . . . . .	556, 611, 706, 747, 757, 774, 783,
\c_tilde_str . . . . .	803, 816, 1348, 1430, 1440, 1586, 1596
\c_two . . . . .	\cs_new_protected_nopar:cpn . . . . .
155, 370, 1191	750, 752, 1171, 1241,
\c_two_hundred_fifty_five . . . . .	1339, 1344, 1346, 1421, 1423, 1425,
709	1520, 1525, 1527, 1577, 1579, 1581
\c_two_hundred_fifty_six . . . . .	\cs_new_protected_nopar:Npn . . . 22, 489,
680, 842	491, 650, 652, 661, 663, 671, 696,
\c_zero 155, 209, 212, 237, 260, 293, 317,	717, 895, 1000, 1083, 1091, 1116, 1148
319, 335, 391, 394, 399, 698, 711,	\cs_set:Npn . . . . . 159,
727, 738, 971, 1175, 1176, 1478, 1616	167, 173, 181, 187, 890, 1446, 1601
\char_set_catcode_alignment:N . . . . .	\cs_set_eq>NN . . . . . 498, 507, 892, 1351
\char_set_catcode_comment:N . . . . .	\cs_set_protected:Npn . . . . . 942
\char_set_catcode_escape:N . . . . .	\cs_to_str:N . . . . . 36–41
\char_set_catcode_group_begin:N . . . . .	
\char_set_catcode_group_end:N . . . . .	
\char_set_catcode_ignore:N . . . . .	
\char_set_catcode_letter:N . . . . .	
\char_set_catcode_math_superscript:N	
. . . . .	E
\char_set_catcode_math_toggle:N . . . . .	\else: . . . . . 198,
\char_set_catcode_other:N . . . . .	209, 223, 237, 320, 323, 326, 335,
. . . . .	364, 372, 382, 543, 598, 601, 604,
\char_set_catcode_other:n . . . . .	643, 682, 740, 830, 844, 1069, 1103,
\char_set_catcode_parameter:N . . . . .	1106, 1138, 1141, 1160, 1163, 1263,
\char_set_catcode_space:N . . . . .	1268, 1286, 1305, 1308, 1359, 1364,
\char_set_lccode:nn . . . . .	1367, 1479, 1490, 1499, 1620, 1625
. . . . .	\exp_after:wN . . . . .
\cs:w . . . . .	69, 94, 116, 125, 135, 146, 162,
\cs_end: . . . . .	169, 177, 183, 190, 197, 199, 205,
\cs_generate_variant:Nn . . . . .	208, 215, 222, 224, 230, 236, 252,
\cs_gset_eq:cc . . . . .	253, 265, 272, 291–293, 309, 310,
\cs_if_exist:cF . . . . .	339, 340, 346, 354, 374, 384, 410,
\cs_if_exist:NF . . . . .	417, 418, 425, 426, 451, 452, 467,
\cs_new:Npn . . . . .	468, 516, 521, 586, 702, 725, 736,
6, 67, 72, 81, 98, 107, 110, 112, 122,	779, 788, 834, 952, 989, 991, 1008,
130, 138, 143, 148, 175, 188, 193,	1014, 1186, 1188, 1250, 1269, 1270,
194, 204, 206, 212, 213, 219, 228,	1282, 1285, 1312, 1313, 1427, 1449,
234, 240, 243, 248, 250, 258, 279,	1462, 1463, 1489, 1583, 1604, 1613
287, 289, 299, 301, 307, 313, 351,	\exp_args:Nc . . . . . 20
413, 421, 431, 436, 441, 457, 473,	\exp_args:Ncc . . . . . 534
584, 590, 632, 638, 665, 678, 733,	\exp_args:Nf . . . . . 246, 281, 283, 295,
793, 826, 840, 911, 921, 960, 1022,	303, 433, 438, 1368, 1369, 1531, 1542
1064, 1066, 1085, 1093, 1118, 1150,	\exp_args>NNf . . . . . 821

\exp_args:N	111, 129, 166, 180, 242, 249, 278, 288, 393, 656, 891	G
\exp_args:Nx	551	
\exp_last_unbraced:Nf	245, 903, 1199, 1385, 1555	
\exp_last_unbraced:Nx	761, 765, 807	
\exp_not:c	19	
\exp_not:N	19, 723	
\exp_stop:f	209, 237, 318, 324, 362, 370, 595, 641, 797, 832, 833, 1100, 1104, 1136, 1139, 1157, 1161, 1184, 1259, 1261, 1280, 1299, 1301, 1357, 1360, 1361, 1477, 1480	
\ExplFileVersion	3	
\ExplFileDescription	3	
\ExplFileName	3	
\ExplFileDate	3	
F		I
\fi:	77, 81, 82, 103, 107, 108, 126, 163, 187, 201, 209, 212, 225, 237, 238, 322, 328, 329, 335, 340, 347, 355, 366, 385, 387, 545, 597, 600, 606–608, 645, 685, 712, 727, 737, 742, 799, 800, 835, 836, 838, 846, 1036, 1069, 1109, 1110, 1143, 1144, 1166, 1167, 1187, 1267, 1271, 1280, 1291, 1307, 1311, 1314, 1319, 1321, 1366, 1370, 1371, 1459, 1464, 1475, 1481, 1492, 1495, 1497, 1501, 1614, 1628, 1629	
\file_if_exist:nTF	564	
\file_input:n	568	
\flag_clear:n	654, 673, 769, 811, 898, 947, 948, 1003, 1004, 1243– 1247, 1350, 1442–1445, 1598–1600	
\flag_clear_new:n	... 1193–1196, 1373–1375, 1551, 1552	
\flag_height:n	488	
\flag_if_raised:nT	... 486, 488, 1217, 1223, 1228, 1235, 1404, 1409, 1414, 1563, 1570	
\flag_if_raised:nTF	481	
\flag_new:n	59, 60	
\flag_raise:n	644, 683, 829, 845, 917, 930, 970, 975, 1057, 1264, 1265, 1288, 1289, 1302, 1303, 1322, 1323, 1329, 1330, 1362, 1508, 1509, 1617, 1618, 1622, 1623, 1637, 1638	
\g_str_error_bool	58, 58, 486, 496, 500, 505, 509	
\g_str_file_prop	42, 42–57, 560	
\g_str_result_tl	25, 25, 392, 393, 395, 397, 449, 453, 465, 469, 515, 527, 655, 656, 700, 701, 723, 900, 904, 950, 952, 1006, 1009, 1012, 1015, 1248, 1250, 1342, 1422, 1424, 1428, 1447, 1523, 1578, 1580, 1583, 1602	
\group_begin:	61, 84, 389, 513, 566, 667, 759, 805, 897, 946, 994, 1002, 1336, 1516	
\group_end:	66, 89, 411, 526, 569, 746, 772, 814, 909, 958, 1020, 1080, 1513, 1644	
I		
\if_case:w	209, 237, 373, 1460	
\if_charcode:w	184, 353, 1069	
\if_int_compare:w	334, 1034, 1357, 1360, 1361	
\if_meaning:w	75, 101, 124, 161, 539, 1309, 1457, 1472, 1498, 1612	
\if_num:w	196, 209, 221, 237, 317, 318, 324, 362, 370, 593–595, 602, 603, 641, 680, 709, 735, 738, 796, 797, 828, 831, 832, 842, 1100, 1104, 1136, 1139, 1157, 1161, 1184, 1259, 1261, 1280, 1299, 1301, 1476, 1477, 1616, 1621	
\int_compare:nNnTF	260, 262, 270	
\int_const:Nn	26–33, 35	
\int_div_truncate:nn	434, 439, 1129, 1130, 1189, 1368, 1532, 1543	
\int_eval:n	114, 150, 315, 403, 1466, 1483, 1626	
\int_eval:w	196, 200, 205, 216, 231, 253, 293, 294, 310, 373, 419, 427, 714, 1270, 1283, 1313, 1461	
\int_eval_end:	375, 430, 714, 1270, 1285, 1313, 1461	
\int_mod:nn	1130, 1131, 1369	
\int_new:N	24	
\int_set:Nn	627, 899, 949, 1005	
\int_use:N	200, 205, 216, 231, 253, 293, 294, 310, 419, 427, 789, 1266, 1270, 1283, 1290, 1304, 1313, 1324, 1331, 1511, 1619, 1624, 1640	

<pre>\int_value:w ..... 254, 292, 666, 779, 801, 1189, 1262 \int_zero:N ..... 699, 760, 806 \iow_char:N ..... 872 \iow_indent:n ..... 870, 1209, 1394</pre> <p style="text-align: center;"><b>L</b></p> <pre>\l_doc_pTF_name_tl ..... 5 \l_str_end_flag ..... 1373 \l_str_extra_flag ..... 1193, 1373 \l_str_internal_int ..... 22, 24, 699, 709, 711, 713, 714, 735, 760, 777–780, 786–788, 790, 796, 806, 819–821, 823, 831 \l_str_internal_tl ..... 22, 23, 390, 405, 410, 560–562, 564, 568, 572, 579, 700, 703, 749 \l_str_missing_flag ..... 1193, 1373 \l_str_overflow_flag ..... 1193 \l_str_overlong_flag ..... 1193</pre> <p style="text-align: center;"><b>M</b></p> <pre>\msg_kernel_error:nnx ..... 482, 541 \msg_kernel_error:nnxx ..... 573 \msg_kernel_new:nnn ..... 848, 850, 859 \msg_kernel_new:nnnn ..... 687, 852, 863, 877, 883, 934, 981, 1072, 1197, 1376, 1383, 1553</pre> <p style="text-align: center;"><b>O</b></p> <pre>\or: ..... 210, 238, 377–381, 1462, 1463</pre> <p style="text-align: center;"><b>P</b></p> <pre>\pdftex_if_engine:F ..... 861 \pdftex_if_engine:TF ..... 34, 629, 649, 669 \pdftex_strcmp:D ..... 334 \prg_break_point:n ..... 341, 348, 454, 470, 588, 636, 704, 728, 764, 768, 810, 906, 954, 1010, 1016, 1252, 1451, 1606 \prg_case_str:xxn ..... 1044 \prg_do_nothing: ..... 519, 1500 \prg_map_break: ..... 453, 469, 587, 635, 703, 727, 736, 763, 767, 809, 905, 953, 1009, 1015, 1251, 1332, 1502, 1642 \prg_map_break:n ..... 340, 347, 359 \prg_new_conditional:Npnn ..... 332, 337, 344, 360, 368, 1098, 1134, 1155 \prg_new_protected_conditional:Npnn ..... 493, 502</pre>	<pre>\prg_return_false: ..... 335, 342, 349, 365, 383, 500, 509, 1102, 1105, 1108, 1137, 1140, 1159, 1162, 1165 \prg_return_true: ..... 335, 359, 363, 371, 386, 500, 509, 1102, 1108, 1142, 1159, 1165 \prop_get:NnNF ..... 560 \prop_gput:Nnn ..... 43–57 \prop_new:N ..... 42 \ProvidesExplPackage ..... 2</pre> <p style="text-align: center;"><b>Q</b></p> <pre>\q_mark ..... 70, 81 \q_nil ..... 1450, 1457, 1472, 1498 \q_stop ..... 70, 73, 79, 81, 95, 107, 119, 156, 171, 178, 185, 187, 191, 193, 240, 256, 297, 469, 475, 517, 522, 530, 703, 708, 763, 767, 776, 785, 809, 818, 1180, 1182, 1190, 1273, 1309</pre> <p style="text-align: center;"><b>R</b></p> <pre>\RequirePackage ..... 4 \reverse_if:N ..... 184</pre> <p style="text-align: center;"><b>S</b></p> <pre>\s_stop .. 1422, 1424, 1428, 1440, 1578, 1580, 1584, 1596, 1605, 1612, 1633 \s_tl ..... 469, 473, 666, 703, 706, 795, 801, 1260, 1272, 1277, 1287, 1292, 1297, 1300, 1315, 1328, 1331, 1465, 1466, 1482, 1488, 1504, 1510, 1511, 1615, 1630, 1639, 1640 \scan_stop: 185, 777, 778, 786, 787, 819, 820 \str_aux_hexadecimal_use:N ..... 368 \str_aux_hexadecimal_use:NTF ..... 8, 368, 914, 924, 965, 967 \str_aux_octal_use:N ..... 360 \str_aux_octal_use:NTF ..... 8, 360, 1028, 1030, 1032 \str_byte ..... 59 \str_collect_aux:wn ... 213, 215, 219, 230 \str_collect_aux:wnNNNNNNN 213, 222, 228 \str_collect_delimit_by_q_stop:w ... 213, 213, 309 \str_collect_end:wn ..... 213, 224, 234 \str_collect_end_i:i:nnnnnnnw ..... 213, 236, 240 \str_const:cn ..... 7 \str_const:cx ..... 7 \str_const:Nn ..... 3, 7, 1089, 1090</pre>
--	--

```

\str_const:Nx ..... 7, 1114
\str_convert_aux_i:nNnnn ..... 489, 490, 492, 497, 506, 511
\str_convert_aux_ii:wwnn ..... 516, 521, 529, 529
\str_convert_aux_iii>NNnNN 529, 534, 537
\str_convert_aux_iv:nnn 532, 533, 547, 547
\str_convert_aux_v:nnnn .. 547, 551, 556
\str_convert_decode:_ ... 520, 663, 663
\str_convert_decode_eight_bit:n .... 751, 757, 757
\str_convert_decode_utf16: .... 1421
\str_convert_decode_utf16be: ... 1421
\str_convert_decode_utf16le: ... 1421
\str_convert_decode_utf32: .... 1577
\str_convert_decode_utf32be: ... 1577
\str_convert_decode_utf32le: ... 1577
\str_convert_decode_utf8: .... 1241
\str_convert_encode:_ . 525, 667, 671, 696
\str_convert_encode_eight_bit:n .... 753, 803, 803
\str_convert_encode_utf16: .... 1339
\str_convert_encode_utf16be: ... 1339
\str_convert_encode_utf16le: ... 1339
\str_convert_encode_utf32: .... 1520
\str_convert_encode_utf32be: ... 1520
\str_convert_encode_utf32le: ... 1520
\str_convert_encode_utf8: .... 1171
\str_convert_escape:_ ... 661, 661, 662
\str_convert_escape_bytes: ... 661, 662
\str_convert_escape_hex: ... 1082, 1083
\str_convert_escape_name: ... 1088, 1091
\str_convert_escape_string: . 1113, 1116
\str_convert_escape_url: ... 1147, 1148
\str_convert_gmap:N .... 447,
447, 664, 770, 1084, 1092, 1117, 1149
\str_convert_gmap_internal:N 463, 463,
674, 812, 1172, 1352, 1522, 1526, 1528
\str_convert_gmap_internal_loop:Nw 463
\str_convert_gmap_internal_loop:Nww
..... 467, 473, 477
\str_convert_gmap_loop:NN ....
..... 447, 451, 457, 461
\str_convert_lowercase_alphanum:n ..
..... 552, 584, 584
\str_convert_lowercase_alphanum_loop:N
..... 584, 586, 590, 609
\str_convert_unescape:_ 649, 650, 652, 660
\str_convert_unescape_bytes: .. 649, 660
\str_convert_unescape_hex: ... 894, 895
\str_convert_unescape_name: ..... 941
\str_convert_unescape_string: 993, 1000
\str_convert_unescape_url: ..... 941
\str_count_spaces:N ..... 4, 110, 110
\str_count_spaces:n ... 110, 111, 112, 134
\str_count_spaces_loop:wwwwwww ...
..... 110, 116, 122, 127
\str_declare_eight_bit_encoding:nnn
... 8, 747, 747, 1647, 1654, 1718,
1760, 1817, 1918, 2005, 2091, 2165,
2178, 2231, 2329, 2392, 2430, 2445
\str_decode_eight_bit_char:N ....
..... 757, 770, 793
\str_decode_eight_bit_load:nn ....
..... 757, 761, 774, 781
\str_decode_eight_bit_load_missing:n
..... 757, 765, 783, 791
\str_decode_native_char:N . 663, 664, 665
\str_decode_utf_viii_aux:wNnwN ....
..... 1241, 1282, 1296
\str_decode_utf_viii_continuation:wwN
..... 1241, 1269, 1276, 1312
\str_decode_utf_viii_end: 1241, 1251, 1326
\str_decode_utf_viii_overflow:w ....
..... 1241, 1310, 1319
\str_decode_utf_viii_start:N ....
..... 1241, 1250,
1256, 1274, 1277, 1294, 1297, 1317
\str_decode_utf_xvi:Nw .... 1421,
1422, 1424, 1433, 1436, 1437, 1440
\str_decode_utf_xvi_bom:NN ....
..... 1421, 1427, 1430
\str_decode_utf_xvi_error:nNN ....
... 1455, 1473, 1491, 1500, 1505, 1506
\str_decode_utf_xvi_extra:NNw ....
..... 1455, 1463, 1504
\str_decode_utf_xvi_pair:NN ....
... 1449, 1455, 1455, 1467, 1470, 1493
\str_decode_utf_xvi_pair_end:Nw ....
..... 1455, 1458, 1474, 1495
\str_decode_utf_xvi_quad:NNwNN ....
..... 1455, 1462, 1469
\str_decode_utf_xxxii:Nw .... 1577,
1578, 1580, 1589, 1592, 1593, 1596
\str_decode_utf_xxxii_bom:NNNN ....
..... 1577, 1583, 1586
\str_decode_utf_xxxii_end:w ....
..... 1577, 1613, 1633
\str_decode_utf_xxxii_loop:NNNN ....
..... 1577, 1604, 1610, 1631

```

```

\str_encode_eight_bit_char:n ..... 803, 812, 826
\str_encode_eight_bit_char_aux:n ... 803, 837, 840
\str_encode_eight_bit_load:nn .... 803, 807, 816, 824
\str_encode_native_char:n . 670, 674, 678
\str_encode_native_filter:N ..... 695, 725, 733, 743
\str_encode_native_flush: ..... 695, 704, 710, 717
\str_encode_native_loop:w ..... 695, 702, 706, 715
\str_encode_utf_viii_char:n ..... 1171, 1172, 1173
\str_encode_utf_viii_loop:wwnnw .... 1171, 1175, 1182, 1188
\str_encode_utf_xvi_aux:N ..... 1339, 1341, 1345, 1347, 1348
\str_encode_utf_xvi_char:n ..... 1339, 1352, 1355
\str_encode_utf_xxxii_be:n ..... 1520, 1522, 1526, 1529
\str_encode_utf_xxxii_be_aux:nn .... 1520, 1531, 1534
\str_encode_utf_xxxii_le:n ..... 1520, 1528, 1540
\str_encode_utf_xxxii_le_aux:nn .... 1520, 1542, 1545
\str_end ..... 1551
\str_error ..... 59
\str_escape_hex_char:N . 1082, 1084, 1085
\str_escape_name_char:N 1088, 1092, 1093
\str_escape_string_char:N 1113, 1117, 1118
\str_escape_url_char:N . 1147, 1149, 1150
\str_filter_bytes:n ..... 629, 630, 632, 656, 962, 1025
\str_filter_bytes_aux:N 629, 634, 638, 646
\str_gput_left:cn ..... 7
\str_gput_left:cx ..... 7
\str_gput_left:Nn ..... 3, 7
\str_gput_left:Nx ..... 7
\str_gput_right:cn ..... 7
\str_gput_right:cx ..... 7
\str_gput_right:Nn ..... 3, 7
\str_gput_right:Nx ..... 7
\str_gset:cn ..... 7
\str_gset:cx ..... 7
\str_gset:Nn ..... 3, 7
\str_gset:Nx ..... 7
\str_gset_convert:Nnnn . 7, 489, 491, 503
\str_gset_convert:NnnnTF ..... 7
\str_gset_other:Nn ..... 8, 84, 90, 515
\str_gset_other_end:w ..... 84, 102, 107
\str_gset_other_loop:w .. 84, 94, 98, 105
\str_head:N ..... 4, 166, 166
\str_head:n ..... 166, 166, 167
\str_head_aux:w ..... 166, 169, 173
\str_head_ignore_spaces:n .. 4, 166, 175
\str_if_contains_char>NN ..... 337
\str_if_contains_char:nN ..... 344
\str_if_contains_char:NNT .. 8, 337, 1122
\str_if_contains_char:NNTF 337, 1101, 1107
\str_if_contains_char:nNTF 337, 1158, 1164
\str_if_contains_char_aux>NN ..... 337, 339, 346, 351, 356
\str_if_contains_char_true: 337, 354, 358
\str_if_eq:NN ..... 332, 332
\str_if_eq:nn ..... 332
\str_if_eq:NNTF ..... 5
\str_if_eq:xx ..... 332
\str_if_eq:xxTF .. 1432, 1435, 1588, 1591
\str_if_escape_name:N ..... 1098
\str_if_escape_name:NTF .... 1088, 1095
\str_if_escape_string:N ..... 1134
\str_if_escape_string:NTF .. 1113, 1120
\str_if_escape_url:N ..... 1155
\str_if_escape_url:NTF .... 1147, 1152
\str_if_flag_error:nnx 479, 479, 498, 507, 657, 675, 771, 813, 908, 956, 957, 1018, 1019, 1254, 1353, 1453, 1608
\str_if_flag_no_error:nnx ..... 479, 485, 498, 507
\str_if_flag_times:nT ..... 487, 1200–1203, 1386–1388, 1556, 1557
\str_item:Nn ..... 4, 242, 242
\str_item:nn ..... 242, 242, 243
\str_item_aux:ww ..... 242, 252, 258
\str_item_ignore_spaces:nn .. 4, 242, 248
\str_item_unsafe:nn ... 242, 246, 249, 250
\str_length:N ..... 3, 129, 129
\str_length:n ..... 129, 129, 130
\str_length_aux:n 129, 132, 140, 145, 148
\str_length_ignore_spaces:n ..... 3, 129, 143, 892
\str_length_loop:NNNNNNNN ..... 129, 135, 141, 146, 159, 164
\str_length_skip_spaces:N ..... 890, 890
\str_length_skip_spaces:n .. 890, 891, 892
\str_length_unsafe:n .. 129, 138, 254, 292

```

\str_load_catcodes: .....	567, 611, 611	\str_tail_aux_i:w .....	180, 190, 193
\str_new:N .....	3	\str_tail_ignore_spaces:n .....	4, 180, 188
\str_output_byte:n .....	413, 413, 444, 445, 599, 681, 822, 843, 1185, 1191, 1538, 1547	\str_tmp:w .....	22, 22, 942, 989, 991, 1351, 1358, 1363, 1365, 1368, 1369, 1446, 1461, 1466, 1477, 1480, 1485, 1486, 1601, 1616, 1621, 1627
\str_output_byte:w .....	413, 414, 415, 902, 927, 964, 1027	\str_unescape_hex_i:N .....	894, 903, 911, 918, 927
\str_output_byte_pair_aux:nnN .....	431, 433, 438, 441	\str_unescape_hex_ii:N .....	894, 915, 921, 931
\str_output_byte_pair_be:n .....	431, 431, 1341, 1345, 1537	\str_unescape_name_loop:wNN .....	941, 990
\str_output_byte_pair_le:n .....	431, 436, 1347, 1548	\str_unescape_string_loop:wNNN .....	993, 1014, 1022, 1062, 1065
\str_output_end: .....	413, 414, 422, 429, 906, 926, 978, 1061, 1065	\str_unescape_string_newlines:wN .....	993, 1008, 1066, 1070
\str_output_hexadecimal:n .....	413, 421, 1086, 1096, 1153	\str_unescape_string_repeat:NNNNNN .....	993, 1037, 1039, 1041, 1064
\str_output_hexadecimal:w .....	413, 422, 423	\str_unescape_url_loop:wNN .....	941, 992
\str_overflow .....	1551		
\str_put_left:cn .....	7		
\str_put_left:cx .....	7		
\str_put_left:Nn .....	3, 7		
\str_put_left:Nx .....	7		
\str_put_right:cn .....	7		
\str_put_right:cx .....	7		
\str_put_right:Nn .....	3, 7	\tex_advance:D .....	713, 780, 790, 823
\str_put_right:Nx .....	7	\tex_dimen:D .....	777, 786, 796–798, 819, 831–833
\str_set:cn .....	7	\tex_endlinechar:D .....	627
\str_set:cx .....	7	\tex_escapechar:D .....	899, 949, 1005
\str_set:Nn .....	3, 7	\tex_lccode:D .....	698, 714, 738
\str_set:Nx .....	7	\tex_roman numeral:D .....	266, 273, 311
\str_set_convert:Nnn .....	7, 489, 489, 494	\tex_skip:D .....	778, 787, 797, 820, 832
\str_set_convert:NnnnTF .....	7	\tex_the:D .....	798, 833
\str_skip_aux:wNNNNNNNN .....	194, 197, 204	\tex_toks:D .....	779, 788, 798, 821, 833
\str_skip_c_zero:w .....	194, 194, 205, 266, 273, 311	\tl_clear:N .....	390, 572
\str_skip_end:w .....	194, 199, 206	\tl_const:cn .....	412, 754, 755
\str_skip_end_ii:NNNNNNNN .....	194, 208, 212	\tl_const:cx .....	402
\str_substr:Nnn .....	5, 278, 278	\tl_const:Nx .....	36–41, 409
\str_substr:nnn .....	278, 278, 279	\tl_gclear:N .....	701
\str_substr_aux:nN .....	278, 295, 299	\tl_gput_left:Nx .....	1342, 1523
\str_substr_aux:nnw .....	278, 303, 307	\tl_gput_right:Nx .....	723
\str_substr_aux:www .....	278, 291, 301	\tl_gset:Nx .....	92, 392, 449, 465, 655, 900, 950, 1006, 1012, 1248, 1447, 1602
\str_substr_ignore_spaces:nnn .....	5, 278, 287	\tl_gset_eq>NN .....	492, 508
\str_substr_normalize_range:nn .....	278, 304, 305, 313	\tl_if_empty:nF .....	540, 1635
\str_substr_unsafe:nnn .....	278, 283, 288, 289	\tl_if_empty:nTF .....	300
\str_tail:N .....	4, 180, 180	\tl_map_function:nN .....	622, 624
\str_tail:n .....	180, 180, 181	\tl_map_inline:Nn .....	395, 397
\str_tail_aux:w .....	180, 183, 187	\tl_map_inline:nn .....	7, 393

## T

TeX and L<sup>A</sup>T<sub>E</sub>X 2 $\epsilon$  commands:

\count .....	32
\escapechar .....	3
\pdfescapehex .....	36
\pdfescapename .....	6, 36
\pdfescapestring .....	6, 36
\tex_advance:D .....	713, 780, 790, 823
\tex_dimen:D .....	777, 786, 796–798, 819, 831–833
\tex_endlinechar:D .....	627
\tex_escapechar:D .....	899, 949, 1005
\tex_lccode:D .....	698, 714, 738
\tex_roman numeral:D .....	266, 273, 311
\tex_skip:D .....	778, 787, 797, 820, 832
\tex_the:D .....	798, 833
\tex_toks:D .....	779, 788, 798, 821, 833
\tl_clear:N .....	390, 572
\tl_const:cn .....	412, 754, 755
\tl_const:cx .....	402
\tl_const:Nx .....	36–41, 409
\tl_gclear:N .....	701
\tl_gput_left:Nx .....	1342, 1523
\tl_gput_right:Nx .....	723
\tl_gset:Nx .....	92, 392, 449, 465, 655, 900, 950, 1006, 1012, 1248, 1447, 1602
\tl_gset_eq>NN .....	492, 508
\tl_if_empty:nF .....	540, 1635
\tl_if_empty:nTF .....	300
\tl_map_function:nN .....	622, 624
\tl_map_inline:Nn .....	395, 397
\tl_map_inline:nn .....	7, 393

	U
\tl_new:N . . . . .	23, 25
\tl_put_right:Nn . . . . .	405
\tl_set:Nn . . . . .	561, 749
\tl_set_eq:NN . . . . .	490, 499, 700
\tl_to_lowercase:n . . .	64, 87, 400, 721, 998
\tl_to_other_str:n . . .	8, 61, 67, 246, 284
\tl_to_other_str_end:w . . .	61, 76, 81
\tl_to_other_str_loop:w . .	61, 69, 72, 78
\tl_to_str:N . . . . .	3, 334, 904
\tl_to_str:n . . . . .	19, 69, 94, 117, 135, 146, 170, 178, 185, 191, 249, 281, 288, 340, 347, 392, 517, 522, 587
\tl_use:c . . . . .	762, 766, 808
\token_to_str:N . . . . .	245, 362, 370, 374
\use:n . . . . .	630
\use:nn . . . . .	434
\use:x . . . . .	719
\use_i:nn . . . . .	971, 976, 1058, 1062
\use_i:nnn . . . . .	418, 979, 1489
\use_i_delimit_by_q_stop:nw . . . . .	174, 177, 265, 272
\use_i_i:mn . . . . .	5, 5, 6, 439, 524
\use_none:n . . . . .	359, 384, 410, 412, 426, 459, 592, 640, 913, 923, 963, 1026, 1199, 1279, 1385, 1555
\use_none:nn . . . . .	483, 834
\use_none_delimit_by_q_stop:w . . . . .	125, 162, 263, 275, 475, 708, 776, 785, 818, 1186, 1273