

Contents

1	Building strings	2
2	Accessing the contents of a string	3
3	String conditionals	5
4	Viewing strings	6
5	Scratch strings	6
6	Encoding functions	7
7	Internal string functions	9
8	Possibilities, and things to do	9
9	I3str implementation	10
9.1	Helpers	10
9.1.1	A function unrelated to strings	10
9.1.2	Assigning strings	11
9.1.3	Variables and constants	11
9.1.4	Escaping spaces	13
9.2	Characters given by their position	14
9.3	String conditionals	20
9.4	Viewing strings	22
9.5	Conversions	22
9.5.1	Producing one byte or character	22
9.5.2	Mapping functions for conversions	24
9.5.3	Error-reporting during conversion	24
9.5.4	Framework for conversions	25
9.5.5	Byte unescape and escape	30
9.5.6	Native strings	31
9.5.7	8-bit encodings	32
9.6	Messages	35
9.7	Deprecated string functions	36
9.8	Escaping definition files	36
9.8.1	Unescape methods	37
9.8.2	Escape methods	42
9.9	Encoding definition files	44
9.9.1	UTF-8 support	44
9.9.2	UTF-16 support	49
9.9.3	UTF-32 support	55
9.9.4	ISO 8859 support	58

The `I3str` package: manipulating strings of characters*

The L^AT_EX3 Project[†]

Released 2012/11/24

L^AT_EX3 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_max_char_int`

The maximum valid character code, 255 for pdfL^AT_EX, and 1114111 for X_GL^AT_EX and LuaL^AT_EX.

`\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 v4339, last revised 2012/11/24.

†E-mail: latex-team@latex-project.org

\tl_to_str:N ★
\tl_to_str:n ★

\tl_to_str:N {tl var}
\tl_to_str:n {{token list}}

Converts the *<token list>* to a *<string>*, 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

\str_new:N {str var}

Creates a new *<str var>* or raises an error if the name is already taken. The declaration is global. The *<str var>* will initially be empty.

\str_const:Nn
\str_const:(Nx|cn|cx)

\str_const:Nn {str var} {{token list}}

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

\str_set:Nn
\str_set:(Nx|cn|cx)
\str_gset:Nn
\str_gset:(Nx|cn|cx)

\str_set:Nn {str var} {{token list}}

Converts the *<token list>* to a *<string>*, and stores the result in *<str var>*.

\str_put_left:Nn
\str_put_left:(Nx|cn|cx)
\str_gput_left:Nn
\str_gput_left:(Nx|cn|cx)

\str_put_left:Nn {str var} {{token list}}

Converts the *<token list>* to a *<string>*, and prepends the result to *<str var>*. The current contents of the *<str var>* 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 {str var} {{token list}}

Converts the *<token list>* to a *<string>*, and appends the result to *<str var>*. The current contents of the *<str var>* are not automatically converted to a string.

2 Accessing the contents of a string

\str_count:N ★\str_count:n {{token list}}
\str_count:n ★
\str_count_ignore_spaces:n ★

Leaves the number of tokens in the string representation of *<token list>* in the input stream as an integer denotation. The functions differ in their treatment of spaces. In the case of \str_count:N and \str_count:n, all characters including spaces are counted. The \str_count_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 1 is the first (left-most) character.

```
\str_substr:Nnn          * \str_substr:nnn {\<token list>} {\<start index>} {\<end index>}
\str_substr:nnn          *
\str_substr_ignore_spaces:nnn *
```

Converts the *<token list>* to a *<string>*, and leaves in the input stream the characters from the *<start index>* inclusive to the *<end index>* exclusive. Note that the token count of the substring is equal to the difference between the two *<indices>*. If either of *<start index>* or *<end index>* is negative, then it is incremented by the token count of the list. If either of *<start index>* or *<end index>* is empty, it is replaced by the corresponding end-point of the string. Both *<start index>* and *<end index>* count from 1 for the first (left most) character. For instance,

```
\iow_term:x { \str_substr:nnn { abcdef } { 2 } { 5 } }
\iow_term:x { \str_substr:nnn { abcdef } { -4 } { } }
```

will print `bcd` and `cdef` to the terminal.

3 String conditionals

```
\str_if_eq_p:nn          * \str_if_eq_p:nnn {\<t1>} {\<t2>}
\str_if_eq_p:(Vn|on|no|nV|VV) * \str_if_eq:nnTF {\<t1>} {\<t2>} {\<true code>} {\<false code>}
\str_if_eq:nnTF          *
\str_if_eq:(Vn|on|no|nV|VV)TF *
```

Compares the two *<token lists>* 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:no { abc } { \tl_to_str:n { abc } }
```

is logically `true`.

```
\str_if_eq_x_p:nn * \str_if_eq_x_p:nnn {\<t1>} {\<t2>}
\str_if_eq_x:nnTF *
```

New: 2012-06-05

Compares the full expansion of two *<token lists>* 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_x_p:nn { abc } { \tl_to_str:n { abc } }
```

is logically `true`.

<code>\str_case:nnn *</code>	<code>\str_case:nnn {<test string>}</code>
<code>\str_case:onn *</code>	{
<hr/>	<code>{<string case1>} {<code case1>}</code>
<code>New: 2012-06-03</code>	<code>{<string case2>} {<code case2>}</code>
	<code>...</code>
	<code>{<string case_n>} {<code case_n>}</code>
	}
	{<else case>}

This function compares the *<test string>* in turn with each of the *<string cases>*. If the two are equal (as described for `\str_if_eq:nnTF` then the associated *<code>* is left in the input stream. If none of the tests are `true` then the `else` code will be left in the input stream.

<code>\str_case_x:nnn *</code>	<code>\str_case_x:nnn {<test string>}</code>
<hr/>	{
<code>New: 2012-06-05</code>	<code>{<string case1>} {<code case1>}</code>
	<code>{<string case2>} {<code case2>}</code>
	<code>...</code>
	<code>{<string case_n>} {<code case_n>}</code>
	}
	{<else case>}

This function compares the full expansion of the *<test string>* in turn with the full expansion of the *<string cases>*. If the two full expansions are equal (as described for `\str_if_eq:nnTF` then the associated *<code>* is left in the input stream. If none of the tests are `true` then the `else` code will be left in the input stream. The *<test string>* is expanded in each comparison, and must always yield the same result: for example, random numbers must not be used within this string.

4 Viewing strings

<code>\str_show:N</code>	<code>\str_show:N <tl var></code>
<code>\str_show:(c n)</code>	Displays the content of the <i><str var></i> on the terminal.

5 Scratch strings

<code>\l_tmpa_str</code>	Scratch strings for local assignment. These are never used by the kernel code, and so
<code>\l_tmpb_str</code>	are safe for use with any L ^A T _E X3-defined function. However, they may be overwritten by
	other non-kernel code and so should only be used for short-term storage.

<code>\g_tmpa_str</code>	Scratch strings for global assignment. These are never used by the kernel code, and so
<code>\g_tmpb_str</code>	are safe for use with any L ^A T _E X3-defined function. However, they may be overwritten by
	other non-kernel code and so should only be used for short-term storage.

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.

6 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.¹
- 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.²

¹Encodings and escapings will be added as they are requested.

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
bytes, or empty	arbitrary bytes
hex, hexadecimal	byte = two hexadecimal digits
name	see \pdfescapename
string	see \pdfescapestring
url	encoding used in URLs

\str_set_convert:Nnnn
\str_gset_convert:Nnnn

\str_set_convert:Nnnn *str var* {\i<string>} {\i<name 1>} {\i<name 2>}

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

For example,

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

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

An error is raised if the *<string>* is not valid according to the *<escaping 1>* and *<encoding 1>*, or if it cannot be reencoded in the *<encoding 2>* and *<escaping 2>* (for instance, if a character does not exist in the *<encoding 2>*). 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 *<encoding 2>*, 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* {\i<string>} {\i<name 1>} {\i<name 2>} {\i<true code>} {\i<false code>}

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

7 Internal string functions

`__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_hexadecimal_use:NTF`

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

If the *<token>* is a hexadecimal 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 fails on some inputs if the escape character is a hexadecimal digit. We are thus careful to set the escape character to a known (safe) value before using it.

`__str_output_byte:n *`

```
\__str_output_byte:n {<intexpr>}
```

Expands to a character token with category other and character code equal to the value of *<intexpr>*. The value of *<intexpr>* must be in the range $[-1, 255]$, and any value outside this range results in undefined behaviour. The special value -1 is used to produce an empty result.

8 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-token count 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.

- Expandable `\str_if_in:nTF?`
- `\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?
- Analog of `printf?`

9 |3str implementation

```

1  (*initex | package)
2  <@=str>
3  \ProvidesExplPackage
4    {\ExplFileName}{\ExplFileVersion}{\ExplFileDescription}
5  \RequirePackage{l3tl-analysis,l3tl-build,l3flag}

```

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

- `\str_if_eq:nn[ptF]` and variants,
- `\str_if_eq_x_return:on`, `\str_if_eq_x_return:nn`
- `\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:w`, (copied here)
- `\str_tail:n`, `_str_tail:w`, (copied here)
- `\str_count_ignore_spaces` (unchanged)
- `\str_count_loop:NNNNNNNN` (unchanged)

9.1 Helpers

9.1.1 A function unrelated to strings

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

```

6  \cs_if_exist:NF \use_i:i:nn
7    { \cs_new:Npn \use_i:i:nn #1#2 { #2 #1 } }
(End definition for \use_i:i:nn)

```

9.1.2 Assigning strings

\str_new:N A string is simply a token list.
\str_new:c

```

8 \cs_new_eq:NN \str_new:N \tl_new:N
9 \cs_generate_variant:Nn \str_new:N { c }

```

(End definition for `\str_new:N` and `\str_new:c`. These functions are documented on page ??.)

\str_set:Nn Simply convert the token list inputs to *{strings}*.

```

10 \tl_map_inline:nN
11   {
12     { set }
13     { gset }
14     { const }
15     { put_left }
16     { gput_left }
17     { put_right }
18     { gput_right }
19   }
20   {
21     \cs_new_protected:cp { str_ #1 :Nn } ##1##2
22     { \exp_not:c { tl_ #1 :Nx } ##1 { \exp_not:N \tl_to_str:n {##2} } }
23     \exp_args:Nc \cs_generate_variant:Nn { str_ #1 :Nn } { Nx , cn , cx }
24   }

```

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

9.1.3 Variables and constants

Internal scratch space for some functions.

```

25 \cs_new_protected_nopar:Npn \__str_tmp:w { }
26 \tl_new:N \l__str_internal_tl
27 \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.

```
28 \tl_new:N \g__str_result_tl
```

(End definition for `\g__str_result_tl`. This variable is documented on page ??.)

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

```

29 \int_const:Nn \c_forty_eight { 48 }
30 \int_const:Nn \c_fifty_eight { 58 }
31 \int_const:Nn \c_sixty_five { 65 }
32 \int_const:Nn \c_ninety_one { 91 }
33 \int_const:Nn \c_ninety_seven { 97 }
34 \int_const:Nn \c_one_hundred_twenty_three { 123 }
35 \int_const:Nn \c_one_hundred_twenty_seven { 127 }

```

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

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

```
36 \int_const:Nn \c_max_char_int  
37   { \pdfTeX_if_engine:TF { "FF" } { "10FFFF" } }
```

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

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

```
38 \int_const:Nn \c__str_replacement_char_int { "FFFD" }
```

(End definition for `\c__str_replacement_char_int` This variable is documented on page ??.)

\c_backslash_str **\c_lbrace_str** **\c_rbrace_str** **\c_hash_str** **\c_tilde_str** **\c_percent_str** For all of those strings, use `\cs_to_str:N` to get characters with the correct category code.

```
39 \tl_const:Nx \c_backslash_str { \cs_to_str:N \\ }  
40 \tl_const:Nx \c_lbrace_str { \cs_to_str:N \{ }  
41 \tl_const:Nx \c_rbrace_str { \cs_to_str:N \} }  
42 \tl_const:Nx \c_hash_str { \cs_to_str:N \# }  
43 \tl_const:Nx \c_tilde_str { \cs_to_str:N \~ }  
44 \tl_const:Nx \c_percent_str { \cs_to_str:N \% }
```

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

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

```
45 \prop_new:N \g__str_alias_prop  
46 \prop_gput:Nnn \g__str_alias_prop { latin1 } { iso88591 }  
47 \prop_gput:Nnn \g__str_alias_prop { latin2 } { iso88592 }  
48 \prop_gput:Nnn \g__str_alias_prop { latin3 } { iso88593 }  
49 \prop_gput:Nnn \g__str_alias_prop { latin4 } { iso88594 }  
50 \prop_gput:Nnn \g__str_alias_prop { latin5 } { iso88599 }  
51 \prop_gput:Nnn \g__str_alias_prop { latin6 } { iso885910 }  
52 \prop_gput:Nnn \g__str_alias_prop { latin7 } { iso885913 }  
53 \prop_gput:Nnn \g__str_alias_prop { latin8 } { iso885914 }  
54 \prop_gput:Nnn \g__str_alias_prop { latin9 } { iso885915 }  
55 \prop_gput:Nnn \g__str_alias_prop { latin10 } { iso885916 }  
56 \prop_gput:Nnn \g__str_alias_prop { utf16le } { utf16 }  
57 \prop_gput:Nnn \g__str_alias_prop { utf16be } { utf16 }  
58 \prop_gput:Nnn \g__str_alias_prop { utf32le } { utf32 }  
59 \prop_gput:Nnn \g__str_alias_prop { utf32be } { utf32 }  
60 \prop_gput:Nnn \g__str_alias_prop { hexadecimal } { hex }
```

(End definition for `\g__str_alias_prop` This variable 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.

```
61 \bool_new:N \g__str_error_bool
```

(End definition for `\g__str_error_bool` This variable 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.

```
62 \flag_new:n { str_byte }
63 \flag_new:n { str_error }
```

(End definition for *str_byte* and *str_error*. These variables are documented on page ??.)

9.1.4 Escaping spaces

__str_to_other:n Converts the *<token list>* to a *<other string>*, where spaces have category code “other”. First apply *\tl_to_str:n*, then replace all spaces by “other” spaces, 8 at a time, storing the converted part of the string between the *\q_mark* and *\q_stop* markers. 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 token count of the string; *__str_gset_other:Nn* is faster but not expandable.

```
64 \group_begin:
65 \char_set_lccode:nn { '*' } { '\_'}
66 \char_set_lccode:nn { 'A' } { '\A'}
67 \tl_to_lowercase:n
68 {
69   \group_end:
70   \cs_new:Npn \_\_str_to_other:n #1
71   {
72     \exp_after:wN \_\_str_to_other_loop:w \tl_to_str:n {#1} ~ %
73     A ~ A ~ A ~ A ~ A ~ A ~ A ~ \q_mark \q_stop
74   }
75   \cs_new:Npn \_\_str_to_other_loop:w
76   #1 ~ #2 ~ #3 ~ #4 ~ #5 ~ #6 ~ #7 ~ #8 ~ #9 \q_stop
77   {
78     \if_meaning:w A #8
79       \_\_str_to_other_end:w
80     \fi:
81     \_\_str_to_other_loop:w
82     #9 #1 * #2 * #3 * #4 * #5 * #6 * #7 * #8 * \q_stop
83   }
84   \cs_new:Npn \_\_str_to_other_end:w \fi: #1 \q_mark #2 * A #3 \q_stop
85   { \fi: #2 }
86 }
```

(End definition for *__str_to_other:n*. This function is documented on page ??.)

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

```
87 \group_begin:
88 \char_set_lccode:nn { '*' } { '\_'}
89 \char_set_lccode:nn { 'A' } { '\A'}
90 \tl_to_lowercase:n
91 {
```

```

92   \group_end:
93   \cs_new_protected:Npn \__str_gset_other:Nn #1#2
94   {
95     \tl_gset:Nx #1
96     {
97       \exp_after:wN \__str_gset_other_loop:w \tl_to_str:n {#2} ~ %
98       A ~ A ~ A ~ A ~ A ~ A ~ A ~ A ~ A ~ \q_stop
99     }
100   }
101   \cs_new:Npn \__str_gset_other_loop:w
102   #1 ~ #2 ~ #3 ~ #4 ~ #5 ~ #6 ~ #7 ~ #8 ~ #9 ~
103   {
104     \if_meaning:w A #9
105       \__str_gset_other_end:w
106     \fi:
107     #1 * #2 * #3 * #4 * #5 * #6 * #7 * #8 * #9
108     \__str_gset_other_loop:w *
109   }
110   \cs_new:Npn \__str_gset_other_end:w \fi: #1 * A #2 \q_stop
111   { \fi: #1 }
112 }
```

(End definition for `__str_gset_other:Nn`. This function is documented on page 9.)

9.2 Characters given by their position

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

```

113 \cs_new:Npn \str_count_spaces:N
114   { \exp_args:No \str_count_spaces:n }
115 \cs_new:Npn \str_count_spaces:n #1
116   {
117     \int_eval:n
118     {
119       \exp_after:wN \__str_count_spaces_loop:wwwwwwww
120       \tl_to_str:n {#1} ~
121       X 7 ~ X 6 ~ X 5 ~ X 4 ~ X 3 ~ X 2 ~ X 1 ~ X 0 ~ X -1 ~
122       \q_stop
123     }
124   }
125 \cs_new:Npn \__str_count_spaces_loop:wwwwwwww #1~#2~#3~#4~#5~#6~#7~#8~#9~
126   {
127     \if_meaning:w X #9
128       \exp_after:wN \use_none_delimit_by_q_stop:w
129     \fi:
130     \c_nine + \__str_count_spaces_loop:wwwwwwww
131   }
```

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

```

\str_count:N      To measure the token count of a string we could first escape all spaces using \__str_-
\str_count:n      to_other:o, then measure the count of this token list. However, this would be quadratic
\__str_count_unsafe:n   in the length of the string, and we can do better. Namely, add the number of spaces
\str_count_ignore_spaces:n (counted using the functions defined above) to the length ignoring spaces. To measure the
\__str_count_loop:NNNNNNNNN 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.

132 \cs_new_nopar:Npn \str_count:N { \exp_args:No \str_count:n }
133 \cs_new:Npn \str_count:n #1
134 {
135     \__str_count:n
136     {
137         \str_count_spaces:n {#1}
138         + \exp_after:wN \__str_count_loop:NNNNNNNNN \tl_to_str:n {#1}
139     }
140 }
141 \cs_new:Npn \__str_count_unsafe:n #1
142 {
143     \__str_count:n
144     { \__str_count_loop:NNNNNNNNN #1 }
145 }
146 \cs_new:Npn \str_count_ignore_spaces:n #1
147 {
148     \__str_count:n
149     { \exp_after:wN \__str_count_loop:NNNNNNNNN \tl_to_str:n {#1} }
150 }
151 \cs_new:Npn \__str_count:n #1
152 {
153     \int_eval:n
154     {
155         #1
156         { X \c_eight } { X \c_seven } { X \c_six }
157         { X \c_five } { X \c_four } { X \c_three }
158         { X \c_two } { X \c_one } { X \c_zero }
159         \q_stop
160     }
161 }
162 \cs_set:Npn \__str_count_loop:NNNNNNNNN #1#2#3#4#5#6#7#8#9
163 {
164     \if_meaning:w X #9
165     \exp_after:wN \use_none_delimit_by_q_stop:w
166     \fi:
167     \c_nine + \__str_count_loop:NNNNNNNNN
168 }

(End definition for \str_count:N This function is documented on page 3.)
```

```

\str_head:N      The _ignore_spaces variant is almost identical to \tl_head:n. As usual, \str_head:N
\str_head:n      expands its argument and hands it to \str_head:n. To circumvent the fact that TeX
\str_head_ignore_spaces:n
\__str_head:w
```

skips spaces when grabbing undelimited macro parameters, `__str_head: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:w` takes an empty argument, and the single (braced) space in the definition of `__str_head: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`.

```

169 \cs_new_nopar:Npn \str_head:N { \exp_args:No \str_head:n }
170 \cs_set:Npn \str_head:n #1
171 {
172     \exp_after:wN \__str_head:w
173     \tl_to_str:n {#1}
174     { { } } ~ \q_stop
175 }
176 \cs_set:Npn \__str_head:w #1 ~ %
177 { \use_i_delimit_by_q_stop:nw #1 { ~ } }
178 \cs_new:Npn \str_head_ignore_spaces:n #1
179 {
180     \exp_after:wN \use_i_delimit_by_q_stop:nw
181     \tl_to_str:n {#1} { } \q_stop
182 }
```

(End definition for `\str_head:N` This function is documented on page 4.)

`\str_tail:N`

`\str_tail:n`

`\str_tail_ignore_spaces:n`

`__str_tail_auxi:w`

`__str_tail_auxii:w`

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.

```

183 \cs_new_nopar:Npn \str_tail:N { \exp_args:No \str_tail:n }
184 \cs_set:Npn \str_tail:n #1
185 {
186     \exp_after:wN \__str_tail_auxi:w
187     \reverse_if:N \if_charcode:w
188     \scan_stop: \tl_to_str:n {#1} X X \q_stop
189 }
190 \cs_set:Npn \__str_tail_auxi:w #1 X #2 \q_stop { \fi: #1 }
191 \cs_new:Npn \str_tail_ignore_spaces:n #1
192 {
193     \exp_after:wN \__str_tail_auxii:w
194     \tl_to_str:n {#1} X X \q_stop
195 }
196 \cs_new:Npn \__str_tail_auxii:w #1 #2 X #3 \q_stop { #2 }
```

(End definition for `\str_tail:N` This function is documented on page 4.)

`__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

`__str_skip_loop:wNNNNNNNN`

`__str_skip_end:w`

`__str_skip_end>NNNNNNNN`

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: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: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).

```

197 \cs_new:Npn \__str_skip_c_zero:w #1;
198 {
199     \if_int_compare:w \__int_eval:w #1 > \c_eight
200         \exp_after:wN \__str_skip_loop:wNNNNNNNN
201     \else:
202         \exp_after:wN \__str_skip_end:w
203         \int_use:N \__int_eval:w
204     \fi:
205     #1 ;
206 }
207 \cs_new:Npn \__str_skip_loop:wNNNNNNNN #1; #2#3#4#5#6#7#8#9
208 { \exp_after:wN \__str_skip_c_zero:w \int_use:N \__int_eval:w #1 - \c_eight ; }
209 \cs_new:Npn \__str_skip_end:w #1 ;
210 {
211     \exp_after:wN \__str_skip_end:NNNNNNNN
212     \if_case:w \if_int_compare:w #1 > \c_zero #1 \else: 0 \fi: \exp_stop_f:
213     \or: \or: \or: \or: \or: \or: \or: \or:
214 }
215 \cs_new:Npn \__str_skip_end: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`
`__str_collect_loop:wn`
`__str_collect_loop:wnNNNNNNN`
`__str_collect_end:wn`
`__str_collect_end:nnnnnnnnw`

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:nnnnnnnnw` are some `\or:`, followed by an `\fi:`, followed by `#1` characters from the input stream. Simply leaving this in the input stream will close the conditional properly and the `\or:` disappear.

```

216 \cs_new:Npn \__str_collect_delimit_by_q_stop:w #1;
217 {
218     \exp_after:wN \__str_collect_loop:wn
219     \int_use:N \__int_eval:w #1 ;
220     {}
221 }
222 \cs_new:Npn \__str_collect_loop:wn #1 ;
223 {
224     \if_int_compare:w #1 > \c_seven
225         \exp_after:wN \__str_collect_loop:wnNNNNNNN
226     \else:
227         \exp_after:wN \__str_collect_end:wn
228     \fi:
229     #1 ;
```

```

230    }
231 \cs_new:Npn \__str_collect_loop:wnNNNNNNN #1; #2 #3#4#5#6#7#8#9
232 {
233     \exp_after:wN \__str_collect_loop:wn
234     \int_use:N \__int_eval:w #1 - \c_seven ;
235     { #2 #3#4#5#6#7#8#9 }
236 }
237 \cs_new:Npn \__str_collect_end:wn #1 ;
238 {
239     \exp_after:wN \__str_collect_end:nnnnnnnnw
240     \if_case:w \if_int_compare:w #1 > \c_zero #1 \else: 0 \fi: \exp_stop_f:
241     \or: \or: \or: \or: \or: \or: \or: \or: \fi:
242 }
243 \cs_new:Npn \__str_collect_end:nnnnnnnnw #1#2#3#4#5#6#7#8 #9 \q_stop
244 { #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". This is alright, as everything else is done with undelimited arguments.

```

245 \cs_new_nopar:Npn \str_item:Nn { \exp_args:No \str_item:nn }
246 \cs_new:Npn \str_item:nn #1#2
247 {
248     \exp_args:Nf \tl_to_str:n
249     {
250         \exp_args:Nf \__str_item_unsafe:nn
251         { \__str_to_other:n {#1} } {#2}
252     }
253 }
254 \cs_new:Npn \str_item_ignore_spaces:nn #1
255 { \exp_args:No \__str_item_unsafe:nn { \tl_to_str:n {#1} } }
256 \cs_new:Npn \__str_item_unsafe:nn #1#2
257 {
258     \exp_after:wN \__str_item:ww
259     \int_use:N \__int_eval:w #2 \exp_after:wN ;
260     \__int_value:w \__str_count_unsafe:n {#1} ;
261     { } #1
262     \q_stop
263 }
264 \cs_new:Npn \__str_item:ww #1; #2;
265 {
266     \int_compare:nNnTF {#1} < \c_zero
267     {
268         \int_compare:nNnTF {#1} < {-#2}
269         { \use_none_delimit_by_q_stop:w }
270     }
271     \exp_after:wN \use_i_delimit_by_q_stop:nw

```

```

272           \tex_roman numeral:D \__str_skip_c_zero:w #1 + #2 + \c_one ;
273       }
274   }
275   {
276     \int_compare:nNnTF {#1} > {#2}
277     { \use_none_delimit_by_q_stop:w }
278     {
279       \exp_after:wN \use_i_delimit_by_q_stop:nw
280       \tex_roman numeral:D \__str_skip_c_zero:w #1 ;
281     }
282   }
283 }
```

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

```

\str_substr:Nnn
\str_substr:nnn
\str_substr_ignore_spaces:nnn
\__str_substr_unsafe:nnn
\__str_substr:nN
\__str_substr:www
\__str_substr:nnw
\__str_substr_normalize_range:nn
```

Sanitize the string. Then evaluate the arguments, replacing them by `\c_zero` or `\c_max_int` if they are empty. 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.

```

284 \cs_new_nopar:Npn \str_substr:Nnn { \exp_args:No \str_substr:nnn }
285 \cs_new:Npn \str_substr:nnn #1#2#3
286   {
287     \exp_args:Nf \tl_to_str:n
288     {
289       \exp_args:Nf \__str_substr_unsafe:nnn
290       { \__str_to_other:n {#1} } {#2} {#3}
291     }
292   }
293 \cs_new:Npn \str_substr_ignore_spaces:nnn #1
294   { \exp_args:No \__str_substr_unsafe:nnn { \tl_to_str:n {#1} } }
295 \cs_new:Npn \__str_substr_unsafe:nnn #1#2#3
296   {
297     \exp_after:wN \__str_substr:www
298     \__int_value:w \__str_count_unsafe:n {#1} \exp_after:wN ;
299     \int_use:N \__int_eval:w #2 + \c_zero \exp_after:wN ;
300     \int_use:N \__int_eval:w
301       \exp_args:Nf \__str_substr:nN {#3} \c_max_int ;
302       { } #1
303     \q_stop
304   }
305 \cs_new:Npn \__str_substr:nN #1 #2
306   { \tl_if_empty:nTF {#1} {#2} {#1} }
307 \cs_new:Npn \__str_substr:www #1; #2; #3;
308   {
309     \exp_args:Nf \__str_substr:nnw
310     { \__str_substr_normalize_range:nn {#2} {#1} }
311     { \__str_substr_normalize_range:nn {#3} {#1} }
312   }
313 \cs_new:Npn \__str_substr:nnw #1#2
314   {
```

```

315   \exp_after:wN \__str_collect_delimit_by_q_stop:w
316   \int_use:N \__int_eval:w #2 + \c_one - #1 \exp_after:wN ;
317   \tex_roman numeral:D \__str_skip_c_zero:w #1 ;
318 }
319 \cs_new:Npn \__str_substr_normalize_range:nn #1#2
320 {
321   \int_eval:n
322   {
323     \if_int_compare:w #1 < \c_zero
324       \if_int_compare:w #1 < - #2 \exp_stop_f:
325         \c_zero
326       \else:
327         #1 + #2 + \c_one
328       \fi:
329     \else:
330       \if_int_compare:w #1 > #2 \exp_stop_f:
331         #2
332       \else:
333         #1
334       \fi:
335     \fi:
336   }
337 }

```

(End definition for `\str_substr:Nnn` This function is documented on page 5.)

9.3 String conditionals

`\str_if_eq:p:NN`
`\str_if_eq:NNTF`

`\str_if_eq_p:nn`

```

338 \prg_new_conditional:Npnn \str_if_eq:NN #1#2 { p , TF , T , F }
339 {
340   \if_int_compare:w \pdftex_strcmp:D { \tl_to_str:N #1 } { \tl_to_str:N #2 }
341     = \c_zero \prg_return_true: \else: \prg_return_false: \fi:
342 }

```

(End definition for `\str_if_eq:NN` These functions are documented on page 5.)

`\str_case:nnn`

`\str_case:onn`

`\str_case_x:nnn`

Defined in `l3basics` at present.
(End definition for `\str_case:nnn`, `\str_case:onn`, and `\str_case_x:nnn` These functions are documented on page 6.)

`__str_if_contains_char:NNT`

`__str_if_contains_char:NNTF`

`__str_if_contains_char:nNTF`

`__str_if_contains_char_aux>NN`

`__str_if_contains_char_true:`

Expects the `<token list>` to be an `<other string>`: the caller is responsible for ensuring that no (too)-special catcodes remain. Spaces with catcode 10 are ignored. Loop over the characters of the string, comparing character codes. The loop is broken if character codes match. Otherwise we return “false”.

```

343 \prg_new_conditional:Npnn \__str_if_contains_char:NN #1#2 { T , TF }
344 {
345   \exp_after:wN \__str_if_contains_char_aux:NN \exp_after:wN #2
346   #1 { \__prg_break:n { ? \fi: } }
347   \__prg_break_point:

```

```

348     \prg_return_false:
349 }
350 \prg_new_if:nPnn \__str_if_contains_char:nN #1#2 { TF }
351 {
352     \__str_if_contains_char_aux:NN #2 #1 { \__prg_break:n { ? \fi: } }
353     \__prg_break_point:
354     \prg_return_false:
355 }
356 \cs_new:Npn \__str_if_contains_char_aux:NN #1#2
357 {
358     \ifCharCode:w #1 #2
359     \exp_after:wN \__str_if_contains_char_true:
360     \fi:
361     \__str_if_contains_char_aux:NN #1
362 }
363 \cs_new_nopar:Npn \__str_if_contains_char_true:
364 { \__prg_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_octal_use:NTF` If the $\langle token \rangle$ is an octal digit, it is left in the input stream, *followed* by the $\langle true\ code \rangle$. Otherwise, the $\langle false\ code \rangle$ is left in the input stream.

TeXhackers note: This function 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. 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.

```

365 \prg_new_if:nPnn \__str_octal_use:N #1 { TF }
366 {
367     \ifIntCompare:w \c_one < '1 \tokenToStr:N #1 \exp_stop_f:
368     #1 \prg_return_true:
369     \else:
370     \prg_return_false:
371     \fi:
372 }
```

(End definition for __str_octal_use:NTF)

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

```

373 \prg_new_if:nPnn \__str_hexadecimal_use:N #1 { TF }
374 {
375     \ifIntCompare:w \c_two < "1 \tokenToStr:N #1 \exp_stop_f:
376     #1 \prg_return_true:
377     \else:
378     \ifCase:w \__int_eval:w
379         \exp_after:wN ` \tokenToStr:N #1 - `a
380         \__int_eval_end:
```

```

381      A
382      \or: B
383      \or: C
384      \or: D
385      \or: E
386      \or: F
387      \else:
388          \prg_return_false:
389          \exp_after:wN \use_none:n
390      \fi:
391      \prg_return_true:
392  \fi:
393 }
(End definition for \__str_hexadecimal_use:NTF)

```

9.4 Viewing strings

```

\str_show:n Displays a string on the terminal.
\str_show:N
\str_show:c
394 \cs_new_eq:NN \str_show:n \tl_show:n
395 \cs_new_eq:NN \str_show:N \tl_show:N
396 \cs_generate_variant:Nn \str_show:N { c }
(End definition for \str_show:n, \str_show:N, and \str_show:c These functions are documented on
page ??.)

```

9.5 Conversions

9.5.1 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 .

```

397 \group_begin:
398   \char_set_catcode_other:n { \c_zero }
399   \tl_set:Nx \l__str_internal_tl { \tl_to_str:n { 0123456789ABCDEF } }
400   \exp_args:No \tl_map_inline:nn { \l__str_internal_tl " }
401   { \char_set_lccode:nn {'#1} { \c_zero } }
402   \tl_map_inline:Nn \l__str_internal_tl
403   {
404     \tl_map_inline:Nn \l__str_internal_tl
405     {
406       \char_set_lccode:nn { \c_zero } {"#1##1}
407       \tl_to_lowercase:n
408       {
409         \tl_const:cx
410         { \c__str_byte_ \int_eval:n {"#1##1} _tl }
411         { ^~@ #1 ##1 }
412       }
413     }

```

```

414     }
415 \group_end:
416 \tl_const:cn { c__str_byte_-1_tl } { { } \use_none:n { } }
(End definition for \c__str_byte_0_tl, \c__str_byte_1_tl, and \c__str_byte_255_tl 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 $\backslash c_{\text{--}}\text{str_byte}_{\langle \text{number} \rangle}\text{tl}$. 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.

```

417 \cs_new:Npn \__str_output_byte:n #1
418   { \__str_output_byte:w #1 \__str_output_end: }
419 \cs_new_nopar:Npn \__str_output_byte:w
420   {
421     \exp_after:wN \exp_after:wN
422     \exp_after:wN \use_i:nnn
423     \cs:w c__str_byte_ \int_use:N \__int_eval:w
424   }
425 \cs_new:Npn \__str_output_hexadecimal:n #1
426   { \__str_output_hexadecimal:w #1 \__str_output_end: }
427 \cs_new_nopar:Npn \__str_output_hexadecimal:w
428   {
429     \exp_after:wN \exp_after:wN
430     \exp_after:wN \use_none:n
431     \cs:w c__str_byte_ \int_use:N \__int_eval:w
432   }
433 \cs_new_nopar:Npn \__str_output_end:
434   { \__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.

```

435 \cs_new:Npn \__str_output_byte_pair_be:n #1
436   {
437     \exp_args:Nf \__str_output_byte_pair:nnN
438       { \int_div_truncate:nn { #1 } { "100" } { #1 } \use:nn
439     }
440 \cs_new:Npn \__str_output_byte_pair_le:n #1
441   {
442     \exp_args:Nf \__str_output_byte_pair:nnN
443       { \int_div_truncate:nn { #1 } { "100" } { #1 } \use_i:i:nn
444     }
445 \cs_new:Npn \__str_output_byte_pair:nnN #1#2#3
446   {
447     #3
448     { \__str_output_byte:n { #1 } }
449     { \__str_output_byte:n { #2 - #1 * "100" } }
450   }
(End definition for \__str_output_byte_pair_be:n This function is documented on page ??.)
```

9.5.2 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.

```

451 \cs_new_protected:Npn \_\_str_convert_gmap:N #1
452 {
453     \tl_gset:Nx \g\_str_result_tl
454     {
455         \exp_after:wN \_\_str_convert_gmap_loop:NN
456         \exp_after:wN #1
457         \g\_str_result_tl { ? \_\_prg_break: }
458         \_\_prg_break_point:
459     }
460 }
461 \cs_new:Npn \_\_str_convert_gmap_loop:NN #1#2
462 {
463     \use_none:n #2
464     #1#2
465     \_\_str_convert_gmap_loop:NN #1
466 }
```

(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.

```

467 \cs_new_protected:Npn \_\_str_convert_gmap_internal:N #1
468 {
469     \tl_gset:Nx \g\_str_result_tl
470     {
471         \exp_after:wN \_\_str_convert_gmap_internal_loop:Nww
472         \exp_after:wN #1
473         \g\_str_result_tl \s\_tl \q_stop \_\_prg_break: \s\_tl
474         \_\_prg_break_point:
475     }
476 }
477 \cs_new:Npn \_\_str_convert_gmap_internal_loop:Nww #1 #2 \s\_tl #3 \s\_tl
478 {
479     \use_none_delimit_by_q_stop:w #3 \q_stop
480     #1 {#3}
481     \_\_str_convert_gmap_internal_loop:Nww #1
482 }
```

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

9.5.3 Error-reporting during conversion

`__str_if_flag_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 `@@_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.

```

483 \cs_new_protected:Npn \_str_if_flag_error:nnx #1
484   {
485     \flag_if_raised:nTF {#1}
486     { \msg_kernel_error:nnx { str } }
487     { \use_none:nn }
488   }
489 \cs_new_protected:Npn \_str_if_flag_no_error:nnx #1#2#3
490   { \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.

```

491 \cs_new_protected:Npn \_str_if_flag_times:nT #1#2
492   { \flag_if_raised:nT {#1} { #2~(x \flag_height:n {#1} ) } }
(End definition for \_str_if_flag_times:nT)
```

9.5.4 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:NnnnTF
\str_gset_convert:NnnnTF
\__str_convert: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:nnx` locally.

```

493 \cs_new_protected_nopar:Npn \str_set_convert:Nnnn
494   { \__str_convert:nNNnnn { } \tl_set_eq:NN }
495 \cs_new_protected_nopar:Npn \str_gset_convert:Nnnn
496   { \__str_convert:nNNnnn { } \tl_gset_eq:NN }
497 \prg_new_protected_conditional:Npnn
498   \str_set_convert:Nnnn #1#2#3#4 { T , F , TF }
499   {
500     \bool_gset_false:N \g__str_error_bool
501     \__str_convert:nNNnnn
502       { \cs_set_eq:NN \__str_if_flag_error:nnx \__str_if_flag_no_error:nnx }
503       \tl_set_eq:NN #1 {#2} {#3} {#4}
504     \bool_if:NTF \g__str_error_bool \prg_return_false: \prg_return_true:
505   }
506 \prg_new_protected_conditional:Npnn
507   \str_gset_convert:Nnnn #1#2#3#4 { T , F , TF }
508   {
509     \bool_gset_false:N \g__str_error_bool
510     \__str_convert:nNNnnn
511       { \cs_set_eq:NN \__str_if_flag_error:nnx \__str_if_flag_no_error:nnx }
512       \tl_gset_eq:NN #1 {#2} {#3} {#4}
513     \bool_if:NTF \g__str_error_bool \prg_return_false: \prg_return_true:
514   }
515 \cs_new_protected:Npn \__str_convert:nNNnnn #1#2#3#4#5#6
516   {
517     \group_begin:
518       #1
519       \__str_gset_other:Nn \g__str_result_tl {#4}
520       \exp_after:wN \__str_convert:wwnnn
521         \tl_to_str:n {#5} /// \q_stop
522         { decode } { unescape }
523         \prg_do_nothing:
524         \__str_convert_decode_:
525       \exp_after:wN \__str_convert:wwnnn
526         \tl_to_str:n {#6} /// \q_stop
527         { encode } { escape }
528         \use_i:i:nn
529         \__str_convert_encode_:
530       \group_end:
531       #2 #3 \g__str_result_tl
532   }

```

(End definition for `\str_set_convert:Nnnn` and `\str_gset_convert:Nnnn`. These functions are documented on page 8.)

```

\__str_convert:wwnnn
\__str_convert>NNnNN

```

The task of `__str_convert:wwnnn` is to split `<encoding>/<escaping>` pairs into their components, #1 and #2. Calls to `__str_convert:nnn` ensure that the corresponding

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_i: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.

```

533 \cs_new_protected:Npn \__str_convert:wwnn
534   #1 / #2 // #3 \q_stop #4#5
535   {
536     \__str_convert:nnn {enc} {#4} {#1}
537     \__str_convert:nnn {esc} {#5} {#2}
538     \exp_args:Ncc \__str_convert>NNnNN
539     { __str_convert_#4:#1: } { __str_convert_#5:#2: } {#2}
540   }
541 \cs_new_protected:Npn \__str_convert:NNnNN #1#2#3#4#5
542   {
543     \if_meaning:w #1 #5
544       \tl_if_empty:nF {#3}
545         { \__msg_kernel_error:nnx { str } { native-escaping } {#3} }
546       #1
547     \else:
548       #4 #2 #1
549     \fi:
550   }

```

(End definition for `__str_convert:wwnn`. This function is documented on page 8.)

`__str_convert:nnn`
`__str_convert:nnnn`

The arguments of `__str_convert: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: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_alias_prop`, then the value `\l__str_internal_tl` 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.

```

551 \cs_new_protected:Npn \__str_convert:n {#1} {#3}
552 {
553     \cs_if_exist:cF { __str_convert_#2_#3: }
554     {
555         \exp_args:Nx \__str_convert:nnn
556         { \__str_convert_lowercase_alphanum:n {#3} }
557         {#1} {#2} {#3}
558     }
559 }
560 \cs_new_protected:Npn \__str_convert:nnnn {#1} {#2} {#3} {#4}
561 {
562     \cs_if_exist:cF { __str_convert_#3_#1: }
563     {
564         \prop_get:NnNF \g__str_alias_prop {#1} \l__str_internal_t1
565         { \tl_set:Nn \l__str_internal_t1 {#1} }
566         \cs_if_exist:cF { __str_convert_#3_ \l__str_internal_t1 : }
567         {
568             \file_if_exist:nTF { l3str-#2- \l__str_internal_t1 .def }
569             {
570                 \group_begin:
571                     \__str_load_catcodes:
572                     \file_input:n { l3str-#2- \l__str_internal_t1 .def }
573                 \group_end:
574             }
575             {
576                 \tl_clear:N \l__str_internal_t1
577                 \__msg_kernel_error:nnxx { str } { unknown-#2 } { #4 } {#1}
578             }
579         }
580         \cs_if_exist:cF { __str_convert_#3_#1: }
581         {
582             \cs_gset_eq:cc { __str_convert_#3_#1: }
583             { __str_convert_#3_ \l__str_internal_t1 : }
584         }
585     }
586     \cs_gset_eq:cc { __str_convert_#3_#4: } { __str_convert_#3_#1: }
587 }
```

(End definition for `__str_convert:n` This function is documented on page 8.)

`_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`

```

588 \cs_new:Npn \_str_convert_lowercase_alphanum:n #1
589   {
590     \exp_after:wN \_str_convert_lowercase_alphanum_loop:N
591     \tl_to_str:n {#1} { ? \_prg_break: }
592     \_prg_break_point:
593   }
594 \cs_new:Npn \_str_convert_lowercase_alphanum_loop:N #1
595   {
596     \use_none:n #1
597     \if_int_compare:w '#1 < \c_ninety_one
598       \if_int_compare:w '#1 < \c_sixty_five
599         \if_int_compare:w \c_one < 1#1 \exp_stop_f:
600           #1
601         \fi:
602       \else:
603         \_str_output_byte:n { '#1 + \c_thirty_two }
604       \fi:
605     \else:
606       \if_int_compare:w '#1 < \c_one_hundred_twenty_three
607         \if_int_compare:w '#1 < \c_ninety_seven
608           \else:
609             #1
610             \fi:
611           \fi:
612         \fi:
613       \_str_convert_lowercase_alphanum_loop:N
614   }

```

(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.

```

615 \cs_new_protected:Npn \_str_load_catcodes:
616   {
617     \char_set_catcode_escape:N \\%
618     \char_set_catcode_group_begin:N \{
619     \char_set_catcode_group_end:N \}
620     \char_set_catcode_math_toggle:N \$%
621     \char_set_catcode_alignment:N \&%
622     \char_set_catcode_parameter:N \#%
623     \char_set_catcode_math_superscript:N ^%
624     \char_set_catcode_ignore:N \ %
625     \char_set_catcode_space:N \~%
626     \tl_map_function:nN { abcdefghijklmnopqrstuvwxyz_ : ABCDEFILNPSTUX }%
627       \char_set_catcode_letter:N %
628     \tl_map_function:nN { 0123456789"?'*+-()., '!/>[] ;= }%
629       \char_set_catcode_other:N %

```

```

630     \char_set_catcode_comment:N \%%
631     \int_set:Nn \tex_endlinechar:D {32}
632 }
(End definition for \__str_load_catcodes:)
```

9.5.5 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.

```

633 \pdftex_if_engine:TF
634   { \cs_new_eq:NN \__str_filter_bytes:n \use:n }
635   {
636     \cs_new:Npn \__str_filter_bytes:n #1
637     {
638       \__str_filter_bytes_aux:N #1
639       { ? \__prg_break: }
640       \__prg_break_point:
641     }
642     \cs_new:Npn \__str_filter_bytes_aux:N #1
643     {
644       \use_none:n #1
645       \if_int_compare:w '#1 < 256 \exp_stop_f:
646         #1
647       \else:
648         \flag_raise:n { str_byte }
649       \fi:
650       \__str_filter_bytes_aux:N
651     }
652   }
(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`.

```

\__str_convert_unescape_bytes:
653 \pdftex_if_engine:TF
654   { \cs_new_protected_nopar:Npn \__str_convert_unescape_: { } }
655   {
656     \cs_new_protected_nopar:Npn \__str_convert_unescape_:
657     {
658       \flag_clear:n { str_byte }
659       \tl_gset:Nx \g__str_result_tl
660       { \exp_args:No \__str_filter_bytes:n \g__str_result_tl }
661       \__str_if_flag_error:nnx { str_byte } { non-byte } { bytes }
662     }
663   }
664 \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.
`__str_convert_escape_bytes_`:

```
665 \cs_new_protected_nopar:Npn \__str_convert_escape_: { }
666 \cs_new_eq:NN \__str_convert_escape_bytes: \__str_convert_escape_:
(End definition for \__str_convert_escape_: This function is documented on page ??.)
```

9.5.6 Native strings

`__str_convert_decode_`: Convert each character to its character code, one at a time.

```
667 \cs_new_protected_nopar:Npn \__str_convert_decode_:
668 { \__str_convert_gmap:N \__str_decode_native_char:N }
669 \cs_new:Npn \__str_decode_native_char:N #1
670 { #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.

```
671 \group_begin:
672   \char_set_catcode_other:n { 0 }
673   \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.

```
674 {
675   \cs_new_protected_nopar:Npn \__str_convert_encode_:
676   {
677     \flag_clear:n { str_error }
678     \__str_convert_gmap_internal:N \__str_encode_native_char:n
679     \__str_if_flag_error:nnx { str_error }
680     { pdfTeX-native-overflow } { }
681   }
682   \cs_new:Npn \__str_encode_native_char:n #1
683   {
684     \if_int_compare:w #1 < \c_two_hundred_fifty_six
685       \__str_output_byte:n {#1}
686     \else:
687       \flag_raise:n { str_error }
688       ?
689     \fi:
690   }
691   \msg_kernel_new:nnnn { str } { pdfTeX-native-overflow }
692   { Character-code-too-large-for-pdfTeX. }
693 }
```

```

694     The-pdfTeX-engine-only-supports-8-bit-characters:-
695     valid-character-codes-are-in-the-range-[0,255].-
696     To-manipulate-arbitrary-Unicode,-use-LuaTeX-or-XeTeX.
697   }
698 }
```

```
\__str_encode_native_loop:w
\__str_encode_native_flush:
\__str_encode_native_filter:N
```

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 using the l3tl-build module, which abuses \toks to reach an almost linear time. Use the standard lowercase trick to produce an arbitrary character from the null character, and add that character to the end of the token list being built. At the end of the loop, put the token list together with __tl_build_end:. Note that we use an x-expanding assignment because it is slightly faster. Unicode-aware engines will never incur an overflow because the internal string is guaranteed to only contain code points in [0, "10FFFF].

```

699 {
700   \cs_new_protected_nopar:Npn \__str_convert_encode_:
701   {
702     \int_zero:N \l__tl_build_offset_int
703     \__tl_gbuild_x:Nw \g__str_result_tl
704     \exp_after:wN \__str_encode_native_loop:w
705     \g__str_result_tl \s__tl { \q_stop \__prg_break: } \s__tl
706     \__prg_break_point:
707     \__tl_build_end:
708   }
709   \cs_new_protected:Npn \__str_encode_native_loop:w #1 \s__tl #2 \s__tl
710   {
711     \use_none_delimit_by_q_stop:w #2 \q_stop
712     \tex_lccode:D \l__str_internal_int \__int_eval:w #2 \__int_eval_end:
713     \tl_to_lowercase:n { \__tl_build_one:n { ^@ } }
714     \__str_encode_native_loop:w
715   }
716 }
```

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

```
717 \group_end:
```

(End definition for __str_convert_encode_: This function is documented on page ??.)

9.5.7 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 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.

All the 8-bit encoding definition file start with _str_declare_eight_bit_encoding:nnn {{encoding name}} {{mapping}} {{missing bytes}}. The *<mapping>* argument is a token list of pairs {{byte}} {{Unicode}} expressed in uppercase hexadecimal notation. The *<missing>* argument is a token list of {{byte}}. Every *<byte>* which does not appear in the *<mapping>* nor the *<missing>* lists maps to the same code point in Unicode.

```

718 \cs_new_protected:Npn \_str_declare_eight_bit_encoding:nnn #1#2#3
719   {
720     \tl_set:Nn \l__str_internal_tl {#1}
721     \cs_new_protected_nopar:cpn { _str_convert_decode:#1: }
722       { \_str_convert_decode_eight_bit:n {#1} }
723     \cs_new_protected_nopar:cpn { _str_convert_encode:#1: }
724       { \_str_convert_encode_eight_bit:n {#1} }
725     \tl_const:cn { c__str_encoding_#1_tl } {#2}
726     \tl_const:cn { c__str_encoding_#1_missing_tl } {#3}
727   }
728 (End definition for \_str_declare_eight_bit_encoding:nnn)

```

```

\_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
    \_str_decode_eight_bit_load:nnn
      \_str_decode_eight_bit_load_missing:n
        \_str_decode_eight_bit_load:nn
          \group_begin:
            \int_zero:N \l__str_internal_int
            \exp_last_unbraced:Nx \_str_decode_eight_bit_load:nn
              { \tl_use:c { c__str_encoding_#1_tl } }
              { \q_stop \_prg_break: } { }
            \_prg_break_point:
            \exp_last_unbraced:Nx \_str_decode_eight_bit_load_missing:n
              { \tl_use:c { c__str_encoding_#1_missing_tl } }
              { \q_stop \_prg_break: }
            \_prg_break_point:
            \flag_clear:n { str_error }
            \_str_convert_gmap:N \_str_decode_eight_bit_char:N
            \_str_if_flag_error:nnx { str_error } { decode-8-bit } {#1}
          \group_end:
        }
      \cs_new_protected:Npn \_str_decode_eight_bit_load:nn #1#2
      {
        \use_none_delimit_by_q_stop:w #1 \q_stop
        \tex_dimen:D "#1 = \l__str_internal_int sp \scan_stop:
        \tex_skip:D \l__str_internal_int = "#1 sp \scan_stop:
        \tex_toks:D \l__str_internal_int \exp_after:wN { \int_value:w "#2 }
        \tex_advance:D \l__str_internal_int \c_one
        \_str_decode_eight_bit_load:nn
      }
    \cs_new_protected:Npn \_str_decode_eight_bit_load_missing:n #1
    {

```

```

756 \use_none_delimit_by_q_stop:w #1 \q_stop
757 \tex_dimen:D "#1 = \l__str_internal_int sp \scan_stop:
758 \tex_skip:D \l__str_internal_int = "#1 sp \scan_stop:
759 \tex_toks:D \l__str_internal_int \exp_after:wN
760   { \int_use:N \c__str_replacement_char_int }
761 \tex_advance:D \l__str_internal_int \c_one
762 \__str_decode_eight_bit_load_missing:n
763 }
764 \cs_new:Npn \__str_decode_eight_bit_char:N #1
765 {
766   #1 \s__tl
767   \if_int_compare:w \tex_dimen:D '#1 < \l__str_internal_int
768     \if_int_compare:w \tex_skip:D \tex_dimen:D '#1 = '#1 \exp_stop_f:
769       \tex_the:D \tex_toks:D \tex_dimen:D
770     \fi:
771   \fi:
772   \__int_value:w '#1 \s__tl
773 }

```

(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
774 \cs_new_protected:Npn \__str_convert_encode_eight_bit:n #1
775 {
776   \group_begin:
777     \int_zero:N \l__str_internal_int
778     \exp_last_unbraced:Nx \__str_encode_eight_bit_load:nn
779     { \tl_use:c { c__str_encoding_#1_tl } }
780     { \q_stop \__prg_break: } { }
781     \__prg_break_point:
782     \fclear:n { str_error }
783     \__str_convert_gmap_internal:N \__str_encode_eight_bit_char:n
784     \__str_if_flag_error:nnx { str_error } { encode-8-bit } {#1}
785   \group_end:
786 }
787 \cs_new_protected:Npn \__str_encode_eight_bit_load:nn #1#2
788 {
789   \use_none_delimit_by_q_stop:w #1 \q_stop
790   \tex_dimen:D "#2 = \l__str_internal_int sp \scan_stop:
791   \tex_skip:D \l__str_internal_int = "#2 sp \scan_stop:
792   \exp_args:NNf \tex_toks:D \l__str_internal_int
793   { \__str_output_byte:n { "#1 } }
794   \tex_advance:D \l__str_internal_int \c_one
795   \__str_encode_eight_bit_load:nn
796 }
797 \cs_new:Npn \__str_encode_eight_bit_char:n #1
798 {
799   \if_int_compare:w #1 > \c_max_register_int
800     \fraise:n { str_error }
801   \else:
802     \if_int_compare:w \tex_dimen:D #1 < \l__str_internal_int

```

```

803     \if_int_compare:w \tex_skip:D \tex_dimen:D #1 = #1 \exp_stop_f:
804         \tex_the:D \tex_toks:D \tex_dimen:D #1 \exp_stop_f:
805             \exp_after:wN \exp_after:wN \exp_after:wN \use_none:nn
806             \fi:
807             \fi:
808             \_str_encode_eight_bit_char_aux:n {#1}
809             \fi:
810         }
811 \cs_new:Npn \_str_encode_eight_bit_char_aux:n #1
812 {
813     \if_int_compare:w #1 < \c_two_hundred_fifty_six
814         \_str_output_byte:n {#1}
815     \else:
816         \flag_raise:n { str_error }
817     \fi:
818 }

```

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

9.6 Messages

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

```

819 \_msg_kernel_new:nnn { str } { unknown-esc }
820   { Escaping~scheme~'#1'~(filtered:'#2')~unknown. }
821 \_msg_kernel_new:nnn { str } { unknown-enc }
822   { Encoding~scheme~'#1'~(filtered:'#2')~unknown. }
823 \_msg_kernel_new:nnnn { str } { native-escaping }
824   { The~'native'~encoding~scheme~does~not~support~any~escaping. }
825   {
826       Since~native~strings~do~not~consist~in~bytes,~
827       none~of~the~escaping~methods~make~sense.~
828       The~specified~escaping,~'#1',~will~be~ignored.
829   }
830 \_msg_kernel_new:nnn { str } { file-not-found }
831   { File~'13str-#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.

```

832 \pdftex_if_engine:F
833 {
834   \_msg_kernel_new:nnnn { str } { non-byte }
835   { String~invalid~in~escaping~'#1':~it~may~only~contain~bytes. }
836   {
837       Some~characters~in~the~string~you~asked~to~convert~are~not~
838       8-bit~characters.~Perhaps~the~string~is~a~'native'~Unicode~string?~
839       If~it~is,~try~using\\
840       \\

```

```

841     \iow_indent:n
842     {
843         \iow_char:N\str_set_convert:Nnnn \\
844         \ \ <str-var>~\{~<string>~\}~\{~native~\}~\{~target-encoding>~\}
845     }
846 }
847 }
```

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

```

848 \__msg_kernel_new:nnnn { str } { decode-8-bit }
849 { Invalid-string-in-encoding-'#1'. }
850 {
851     LaTeX-came-across-a-byte-which-is-not-defined-to-represent-
852     any-character-in-the-encoding-'#1'.
853 }
854 \__msg_kernel_new:nnnn { str } { encode-8-bit }
855 { Unicode-string-cannot-be-converted-to-encoding-'#1'. }
856 {
857     The-encoding-'#1'-only-contains-a-subset-of-all-Unicode-characters.-
858     LaTeX-was-asked-to-convert-a-string-to-that-encoding,-but-that-
859     string-contains-a-character-that-'#1'-does-not-support.
860 }
```

9.7 Deprecated string functions

Deprecated 2012-05-13 for removal by 2012-08-31.

```

\str_length:N
\str_length:n
\str_length_ignore_spaces:n
861 \cs_new_eq:NN \str_length:N \str_count:N
862 \cs_new_eq:NN \str_length:n \str_count:n
863 \cs_new_eq:NN \str_length_ignore_spaces:n \str_count_ignore_spaces:n
(End definition for \str_length:N, \str_length:n, and \str_length_ignore_spaces:n These functions
are documented on page ??.)
864 
```

9.8 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.

9.8.1 Unescape methods

```
\__str_convert_unescape_hex:  
  \__str_unescape_hex_auxi:N  
  \__str_unescape_hex_auxii:N
```

Take chars two by two, and interpret each pair as the hexadecimal code for a byte. Anything else than hexadecimal digits is ignored, raising the flag. A string which contains 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.

```
865  {*hex}  
866  \cs_new_protected_nopar:Npn \__str_convert_unescape_hex:  
867  {  
868  \group_begin:  
869  \flag_clear:n { str_error }  
870  \int_set:Nn \tex_escapechar:D { 92 }  
871  \tl_gset:Nx \g__str_result_tl  
872  {  
873  \__str_output_byte:w "  
874  \exp_last_unbraced:Nf \__str_unescape_hex_auxi:N  
875  { \tl_to_str:N \g__str_result_tl }  
876  0 { ? 0 - \c_one \__prg_break: }  
877  \__prg_break_point:  
878  \__str_output_end:  
879  }  
880  \__str_if_flag_error:n { str_error } { unescape-hex } { }  
881  \group_end:  
882  }  
883  \cs_new:Npn \__str_unescape_hex_auxi:N #1  
884  {  
885  \use_none:n #1  
886  \__str_hexadecimal_use:NTF #1  
887  { \__str_unescape_hex_auxii:N }  
888  {  
889  \flag_raise:n { str_error }  
890  \__str_unescape_hex_auxi:N  
891  }  
892  }  
893  \cs_new:Npn \__str_unescape_hex_auxii:N #1  
894  {  
895  \use_none:n #1  
896  \__str_hexadecimal_use:NTF #1  
897  {  
898  \__str_output_end:  
899  \__str_output_byte:w " \__str_unescape_hex_auxi:N  
900  }  
901  {  
902  \flag_raise:n { str_error }  
903  \__str_unescape_hex_auxii:N  
904  }  
905  }  
906  \__msg_kernel_new:nnnn { str } { unescape-hex }  
907  { String-invalid-in-escaping-'hex':~only~hexadecimal~digits~allowed. }
```

```

908     {
909         Some~characters~in~the~string~you~asked~to~convert~are~not~
910         hexadecimal~digits~(0-9,~A-F,~a-f)~nor~spaces.
911     }
912 
```

(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_t1` 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_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`).

```

913  /*name | url*/
914  \cs_set_protected:Npn \__str_tmp:w #1#2#3
915  {
916      \cs_new_protected:cpx { \__str_convert_unescape_#2: }
917      {
918          \group_begin:
919              \flag_clear:n { str_byte }
920              \flag_clear:n { str_error }
921              \int_set:Nn \tex_escapechar:D { 92 }
922              \tl_gset:Nx \g__str_result_t1
923              {
924                  \exp_after:wN #3 \g__str_result_t1
925                  #1 ? { ? \__prg_break: }
926                  \__prg_break_point:
927              }
928              \__str_if_flag_error:nnx { str_byte } { non-byte } { #2 }
929              \__str_if_flag_error:nnx { str_error } { unescape-#2 } { }
930          \group_end:
931      }
932      \cs_new:Npn #3 ##1#1##2##3
933      {
934          \__str_filter_bytes:n {##1}
935          \use_none:n ##3
936          \__str_output_byte:w "
937          \__str_hexadecimal_use:NTF ##2
938          {
939              \__str_hexadecimal_use:NTF ##3

```

```

940           { }
941           {
942             \flag_raise:n { str_error }
943             * \c_zero + '#1 \use_i:nn
944           }
945         }
946         {
947           \flag_raise:n { str_error }
948           0 + '#1 \use_i:nn
949         }
950         \__str_output_end:
951         \use_i:nnn #3 ##2##3
952       }
953     \__msg_kernel_new:nnnn { str } { unescape-#2 }
954     { String~invalid~in~escaping~'#2'. }
955     {
956       LaTeX~came~across~the~escape~character~'#1'~not~followed~by~
957       two~hexadecimal~digits.~This~is~invalid~in~the~escaping~'#2'.
958     }
959   }
960   </name | url>
961   (*name)
962   \exp_after:wN \__str_tmp:w \c_hash_str { name }
963   \__str_unescape_name_loop:wNN
964   </name>
965   (*url)
966   \exp_after:wN \__str_tmp:w \c_percent_str { url }
967   \__str_unescape_url_loop:wNN
968   </url>
(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.

```

969 /*string)
970 \group_begin:
971   \char_set_lccode:nn {'\*} {'\\}
972   \char_set_catcode_other:N \^J
973   \char_set_catcode_other:N \^M
974   \tl_to_lowercase:n
975   {
976     \cs_new_protected_nopar:Npn \__str_convert_unescape_string:
977     {
978       \group_begin:
979         \flag_clear:n { str_byte }
980         \flag_clear:n { str_error }
981         \int_set:Nn \tex_escapechar:D { 92 }
982         \tl_gset:Nx \g__str_result_tl
983         {
984           \exp_after:wN \__str_unescape_string_newlines:wN
985             \g__str_result_tl \__prg_break: ^^M ?
986             \__prg_break_point:
987         }
988         \tl_gset:Nx \g__str_result_tl
989         {
990           \exp_after:wN \__str_unescape_string_loop:wNN
991             \g__str_result_tl * ?? { ? \__prg_break: }
992             \__prg_break_point:
993         }
994         \__str_if_flag_error:nnx { str_byte } { non-byte } { string }
995         \__str_if_flag_error:nnx { str_error } { unescape-string } { }
996       \group_end:
997     }
998     \cs_new:Npn \__str_unescape_string_loop:wNNN #1 **#2#3#4
999   }
1000   {
1001     \__str_filter_bytes:n {#1}
1002     \use_none:n #4
1003     \__str_output_byte:w '
1004     \__str_octal_use:NTF #2
1005     {
1006       \__str_octal_use:NTF #3
1007       {
1008         \__str_octal_use:NTF #4
1009         {
1010           \if_int_compare:w #2 > \c_three
1011             - 256
1012           \fi:
1013           \__str_unescape_string_repeat:NNNNNN

```

```

1014         }
1015         { \__str_unescape_string_repeat:NNNNNN ? }
1016     }
1017     { \__str_unescape_string_repeat:NNNNNN ?? }
1018 }
1019 {
1020     \str_case_x:n:nnn {#2}
1021     {
1022         { \c_backslash_str } { 134 }
1023         { ( } { 50 }
1024         { ) } { 51 }
1025         { r } { 15 }
1026         { f } { 14 }
1027         { n } { 12 }
1028         { t } { 11 }
1029         { b } { 10 }
1030         { ^^J } { 0 - \c_one }
1031     }
1032     {
1033         \flag_raise:n { str_error }
1034         0 - \c_one \use_i:nn
1035     }
1036 }
1037     \__str_output_end:
1038     \use_i:nn \__str_unescape_string_loop:wNNN #2#3#4
1039 }
1040 \cs_new:Npn \__str_unescape_string_repeat:NNNNNN #1#2#3#4#5#6
1041     { \__str_output_end: \__str_unescape_string_loop:wNNN }
1042 \cs_new:Npn \__str_unescape_string_newlines:wN #1 ^^M #2
1043     {
1044     #1
1045     \if_charcode:w ^^J #2 \else: ^^J \fi:
1046     \__str_unescape_string_newlines:wN #2
1047 }
1048 \__msg_kernel_new:nnnn { str } { unescape-string }
1049     { String~invalid-in~escaping~'string'. }
1050     {
1051     LaTeX-came-across-an-escape-character-'\\c_backslash_str'-
1052     not-followed-by-any-of:-'n',-'r',-'t',-'b',-'f',-'( ',' ')',-
1053     '\\c_backslash_str',-one-to-three-octal-digits,-or-the-end-
1054     of-a-line.
1055 }
1056 \group_end:
1057 </string>
(End definition for \__str_convert_unescape_string: This function is documented on page ??.)
```

9.8.2 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:  
  \_\_str\_escape\_hex\_char:N  
  1058  (*hex)  
  1059  \cs_new_protected_nopar:Npn \_\_str\_convert\_escape\_hex:  
  1060  { \_\_str\_convert\_gmap:N \_\_str\_escape\_hex\_char:N }  
  1061  \cs_new:Npn \_\_str\_escape\_hex\_char:N #1  
  1062  { \_\_str\_output_hexadecimal:n { '#1' } }  
  1063  
```

(End definition for `__str_convert_escape_hex`: This function is documented on page ??.)

`__str_convert_escape_name:` 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.

```
\c\_str\_escape\_name\_not\_str  
  1064  (*name)  
  1065  \str_const:Nn \c\_str\_escape\_name\_not\_str { ! " $ & ' } %$  
  1066  \str_const:Nn \c\_str\_escape\_name\_str { {} />[] }  
  1067  \cs_new_protected_nopar:Npn \_\_str\_convert\_escape\_name:  
  1068  { \_\_str\_convert\_gmap:N \_\_str\_escape\_name\_char:N }  
  1069  \cs_new:Npn \_\_str\_escape\_name\_char:N #1  
  1070  {  
  1071    \_\_str\_if\_escape\_name:NTF #1 {#1}  
  1072    { \c\_hash\_str \_\_str\_output_hexadecimal:n { '#1' } }  
  1073  }  
  1074  \prg_new_conditional:Npnn \_\_str\_if\_escape\_name:N #1 { TF }  
  1075  {  
  1076    \if_int_compare:w '#1 < "2A \exp_stop_f:  
  1077    \_\_str\_if\_contains\_char:NNTF \c\_str\_escape\_name\_not\_str #1  
  1078    \prg_return_true: \prg_return_false:  
  1079  \else:  
  1080    \if_int_compare:w '#1 > "7E \exp_stop_f:  
  1081    \prg_return_false:  
  1082  \else:  
  1083    \_\_str\_if\_contains\_char:NNTF \c\_str\_escape\_name\_str #1  
  1084    \prg_return_false: \prg_return_true:  
  1085    \fi:  
  1086  \fi:  
  1087  }  
  1088  
```

(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.

```
\c\_str\_escape\_string\_str  
  1089  (*string)  
  1090  \str_const:Nx \c\_str\_escape\_string\_str
```

```

1091 { \c_backslash_str ( ) }
1092 \cs_new_protected_nopar:Npn \__str_convert_escape_string:
1093 { \__str_convert_gmap:N \__str_escape_string_char:N }
1094 \cs_new:Npn \__str_escape_string_char:N #1
1095 {
1096     \__str_if_escape_string:NTF #1
1097     {
1098         \__str_if_contains_char:NNT
1099             \c__str_escape_string_str #1
1100             { \c_backslash_str }
1101             #1
1102     }
1103     {
1104         \c_backslash_str
1105         \int_div_truncate:nn {'#1} {64}
1106         \int_mod:nn { \int_div_truncate:nn {'#1} \c_eight } \c_eight
1107         \int_mod:nn {'#1} \c_eight
1108     }
1109 }
1110 \prg_new_conditional:Npnn \__str_if_escape_string:N #1 { TF }
1111 {
1112     \if_int_compare:w '#1 < "21 \exp_stop_f:
1113         \prg_return_false:
1114     \else:
1115         \if_int_compare:w '#1 > "7E \exp_stop_f:
1116             \prg_return_false:
1117         \else:
1118             \prg_return_true:
1119         \fi:
1120     \fi:
1121 }
1122 </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.

```

1123 (*url)
1124 \cs_new_protected_nopar:Npn \__str_convert_escape_url:
1125 { \__str_convert_gmap:N \__str_escape_url_char:N }
1126 \cs_new:Npn \__str_escape_url_char:N #1
1127 {
1128     \__str_if_escape_url:NTF #1 {'#1}
1129     { \c_percent_str \__str_output_hexadecimal:n {'#1} }
1130 }
1131 \prg_new_conditional:Npnn \__str_if_escape_url:N #1 { TF }
1132 {
1133     \if_int_compare:w '#1 < "41 \exp_stop_f:
1134         \__str_if_contains_char:nNTF { "-.<>" } #1
1135             \prg_return_true: \prg_return_false:
1136     \else:
1137         \if_int_compare:w '#1 > "7E \exp_stop_f:

```

```

1138          \prg_return_false:
1139      \else:
1140          \__str_if_contains_char:nNTF { [ ] } #1
1141          \prg_return_false: \prg_return_true:
1142          \fi:
1143      \fi:
1144  }
1145 
```

(End definition for `__str_convert_escape_url`: This function is documented on page ??.)

9.9 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.

9.9.1 utf-8 support

```
1146 (*utf8)
```

`__str_convert_encode_utf8:`

```

\__str_encode_utf_viii_char:n
\__str_encode_utf_viii_loop:wwnnw
```

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_vii_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 $\backslash c_minus_one + \backslash c_zero *$ removes the spurious initial continuation byte (better methods welcome).

```

1147 \cs_new_protected_nopar:cpn { __str_convert_encode_utf8: }
1148   { __str_convert_gmap_internal:N __str_encode_utf_viii_char:n }
1149 \cs_new:Npn __str_encode_utf_viii_char:n #1
1150   {
1151     __str_encode_utf_viii_loop:wwnw #1 ; \c_minus_one + \c_zero * ;
1152     { 128 } { \c_zero }
1153     { 32 } { 192 }
1154     { 16 } { 224 }
1155     { 8 } { 240 }
1156     \q_stop
1157   }
1158 \cs_new:Npn __str_encode_utf_viii_loop:wwnw #1; #2; #3#4 #5 \q_stop
1159   {
1160     \if_int_compare:w #1 < #3 \exp_stop_f:
1161       __str_output_byte:n { #1 + #4 }
1162       \exp_after:wN \use_none_delimit_by_q_stop:w
1163     \fi:
1164     \exp_after:wN __str_encode_utf_viii_loop:wwnw
1165       __int_value:w \int_div_truncate:nn {#1} {64} ; #1 ;
1166       #5 \q_stop
1167       __str_output_byte:n { #2 - 64 * ( #1 - \c_two ) }
1168   }

```

(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 can occur, signalled by a specific flag for each (we define those flags using `\flag_clear_new:n` rather than `\flag_new:n`, because they are shared with other encoding definition 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^AT_EX3 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.

```

1169 \flag_clear_new:n { str_missing }
1170 \flag_clear_new:n { str_extra }
1171 \flag_clear_new:n { str_overlong }
1172 \flag_clear_new:n { str_overflow }
1173 \__msg_kernel_new:nnn { str } { utf8-decode }
1174 {
1175   Invalid~UTF-8~string: \exp_last_unbraced:Nf \use_none:n
1176   \_\_str_if_flag_times:nT { str_missing } { ,~missing~continuation~byte }
1177   \_\_str_if_flag_times:nT { str_extra } { ,~extra~continuation~byte }
1178   \_\_str_if_flag_times:nT { str_overlong } { ,~overlong~form }
1179   \_\_str_if_flag_times:nT { str_overflow } { ,~code~point~too~large }
1180 .
1181 }
1182 {
1183   In~the~UTF-8~encoding,~each~Unicode~character~consists~in~
1184   1~to~4~bytes,~with~the~following~bit~pattern: \\
1185   \iow_indent:n
1186   {
1187     Code~point~\ \ \ \ <~128:~0xxxxxx \\ \
1188     Code~point~\ \ \ <~2048:~110xxxx~10xxxxxx \\ \
1189     Code~point~\ \ <~65536:~1110xxxx~10xxxxxx~10xxxxxx \\ \
1190     Code~point~ <~1114112:~11110xxx~10xxxxxx~10xxxxxx~10xxxxxx \\ \
1191   }
1192 Bytes~of~the~form~10xxxxxx~are~called~continuation~bytes.
1193 \flag_if_raised:nT { str_missing }
1194 {
1195   \\\\
1196   A~leading~byte~(in~the~range~[192,255])~was~not~followed~by~
1197   the~appropriate~number~of~continuation~bytes.
1198 }
1199 \flag_if_raised:nT { str_extra }
1200 {
1201   \\\\
1202   LaTeX~came~across~a~continuation~byte~when~it~was~not~expected.
1203 }
1204 \flag_if_raised:nT { str_overlong }
1205 {
1206   \\\\
1207   Every~Unicode~code~point~must~be~expressed~in~the~shortest~
1208   possible~form.~For~instance,~'0xC0'~'0x83'~is~not~a~valid~
1209   representation~for~the~code~point~3.
1210 }
1211 \flag_if_raised:nT { str_overflow }
1212 {
1213   \\\\
1214   Unicode~limits~code~points~to~the~range~[0,1114111].
1215 }
1216 }

(End definition for \l__str_missing_flag and others. These variables 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:wNmNw
        \_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^AT_EX3 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_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.

```

1217 \cs_new_protected_nopar:cpxn { __str_convert_decode_utf8: }
1218 {
1219     \flag_clear:n { str_error }
1220     \flag_clear:n { str_missing }
1221     \flag_clear:n { str_extra }
1222     \flag_clear:n { str_overlong }
1223     \flag_clear:n { str_overflow }
1224     \tl_gset:Nx \g__str_result_tl
1225     {
1226         \exp_after:wN \__str_decode_utf_viii_start:N \g__str_result_tl

```

```

1227          { \_prg_break: \_str_decode_utf_viii_end: }
1228          \_prg_break_point:
1229      }
1230      \_str_if_flag_error:n { str_error } { utf8-decode } { }
1231  }
1232 \cs_new:Npn \_str_decode_utf_viii_start:N #1
1233  {
1234      #
1235      \if_int_compare:w '#1 < "C0 \exp_stop_f:
1236          \s__tl
1237          \if_int_compare:w '#1 < "80 \exp_stop_f:
1238              \int_value:w '#1
1239          \else:
1240              \flag_raise:n { str_extra }
1241              \flag_raise:n { str_error }
1242              \int_use:N \c__str_replacement_char_int
1243          \fi:
1244      \else:
1245          \exp_after:wn \_str_decode_utf_viii_continuation:wwN
1246          \int_use:N \_int_eval:w '#1 - "C0 \exp_after:wn \_int_eval_end:
1247      \fi:
1248      \s__tl
1249      \use_none_delimit_by_q_stop:w {"80} {"800} {"10000} {"110000} \q_stop
1250      \_str_decode_utf_viii_start:N
1251  }
1252 \cs_new:Npn \_str_decode_utf_viii_continuation:wwN
1253     #1 \s__tl #2 \_str_decode_utf_viii_start:N #3
1254  {
1255      \use_none:n #3
1256      \if_int_compare:w '#3 <
1257          \if_int_compare:w '#3 < "80 \exp_stop_f: - \fi:
1258          "C0 \exp_stop_f:
1259          #
1260          \exp_after:wn \_str_decode_utf_viii_aux:wNnnwN
1261          \int_use:N \_int_eval:w
1262              #1 * "40 + '#3 - "80
1263          \exp_after:wn \_int_eval_end:
1264      \else:
1265          \s__tl
1266          \flag_raise:n { str_missing }
1267          \flag_raise:n { str_error }
1268          \int_use:N \c__str_replacement_char_int
1269      \fi:
1270      \s__tl
1271      #
1272      \_str_decode_utf_viii_start:N #3
1273  }
1274 \cs_new:Npn \_str_decode_utf_viii_aux:wNnnwN
1275     #1 \s__tl #2#3#4 #5 \_str_decode_utf_viii_start:N #6
1276  {

```

```

1277 \if_int_compare:w #1 < #4 \exp_stop_f:
1278   \s__tl
1279   \if_int_compare:w #1 < #3 \exp_stop_f:
1280     \flag_raise:n { str_overflow }
1281     \flag_raise:n { str_error }
1282     \int_use:N \c__str_replacement_char_int
1283   \else:
1284     #1
1285   \fi:
1286 \else:
1287   \if_meaning:w \q_stop #5
1288     \__str_decode_utf_viii_overflow:w #1
1289   \fi:
1290   \exp_after:wN \__str_decode_utf_viii_continuation:wwN
1291   \int_use:N \__int_eval:w #1 - #4 \exp_after:wN \__int_eval_end:
1292 \fi:
1293 \s__tl
1294 #2 {#4} #5
1295 \__str_decode_utf_viii_start:N
1296 }
1297 \cs_new:Npn \__str_decode_utf_viii_overflow:w #1 \fi: #2 \fi:
1298 {
1299   \fi: \fi:
1300   \flag_raise:n { str_overflow }
1301   \flag_raise:n { str_error }
1302   \int_use:N \c__str_replacement_char_int
1303 }
1304 \cs_new_nopar:Npn \__str_decode_utf_viii_end:
1305 {
1306   \s__tl
1307   \flag_raise:n { str_missing }
1308   \flag_raise:n { str_error }
1309   \int_use:N \c__str_replacement_char_int \s__tl
1310   \__prg_break:
1311 }
(End definition for \__str_convert_decode_utf8: This function is documented on page ??.)
1312 
```

9.9.2 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.

```

1313 {*utf16}
1314 \group_begin:
1315   \char_set_catcode_other:N \^^fe
1316   \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
`__str_convert_encode_utf16be:`
`__str_convert_encode_utf16le:`
`__str_encode_utf_xvi_aux:N`
`__str_encode_utf_xvi_char:n`

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. The magic "D7C0 is "D800 - "10000/"400.

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.

```

1317   \cs_new_protected_nopar:cpn { __str_convert_encode_utf16: }
1318   {
1319     \__str_encode_utf_xvi_aux:N \__str_output_byte_pair_be:n
1320     \tl_gput_left:Nx \g__str_result_tl { ^^fe ^^ff }
1321   }
1322   \cs_new_protected_nopar:cpn { __str_convert_encode_utf16be: }
1323   {
1324     \__str_encode_utf_xvi_aux:N \__str_output_byte_pair_be:n
1325   }
1326   \cs_new_protected_nopar:cpn { __str_convert_encode_utf16le: }
1327   {
1328     \__str_encode_utf_xvi_aux:N \__str_output_byte_pair_le:n
1329   }
1330   \cs_new_protected:Npn \__str_encode_utf_xvi_aux:N #1
1331   {
1332     \flag_clear:n { str_error }
1333     \cs_set_eq:NN \__str_tmp:w #1
1334     \__str_convert_gmap_internal:N \__str_encode_utf_xvi_char:n
1335     \__str_if_flag_error:nnx { str_error } { utf16-encode } { }
1336   }
1337   \cs_new:Npn \__str_encode_utf_xvi_char:n #1
1338   {
1339     \if_int_compare:w #1 < "D800 \exp_stop_f:
1340       \__str_tmp:w {#1}
1341     \else:
1342       \if_int_compare:w #1 < "10000 \exp_stop_f:
1343         \if_int_compare:w #1 < "E000 \exp_stop_f:
1344           \flag_raise:n { str_error }
1345           \__str_tmp:w { \c__str_replacement_char_int }
1346         \else:
1347           \__str_tmp:w {#1}
1348         \fi:
1349       \else:
1350         \exp_args:Nf \__str_tmp:w { \int_div_truncate:nn {#1} {"400} + "D7C0 }
1351         \exp_args:Nf \__str_tmp:w { \int_mod:nn {#1} {"400} + "DC00 }
1352       \fi:
1353     \fi:
1354   }

```

(End definition for `_str_convert_encode_utf16:`, `_str_convert_encode_utf16be:`, and `_str_convert_encode_utf16l`.
These functions are documented on page ??.)

<code>\l__str_missing_flag</code>	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 <code>@@_error</code> to signal that error.
<code>\l__str_extra_flag</code>	
<code>\l__str_end_flag</code>	

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.

```

1351   \flag_clear_new:n { str_missing }
1352   \flag_clear_new:n { str_extra }
1353   \flag_clear_new:n { str_end }
1354   \_msg_kernel_new:nnnn { str } { utf16-encode }
1355   { Unicode~string~cannot~be~expressed~in~UTF-16:~surrogate. }
1356   {
1357     Surrogate~code~points~(in~the~range~[U+D800,~U+DFFF])~
1358     can~be~expressed~in~the~UTF-8~and~UTF-32~encodings,~
1359     but~not~in~the~UTF-16~encoding.
1360   }
1361   \_msg_kernel_new:nnnn { str } { utf16-decode }
1362   {
1363     Invalid~UTF-16~string: \exp_last_unbraced:Nf \use_none:n
1364     \_str_if_flag_times:nT { str_missing } { ,~missing~trail~surrogate }
1365     \_str_if_flag_times:nT { str_extra } { ,~extra~trail~surrogate }
1366     \_str_if_flag_times:nT { str_end } { ,~odd~number~of~bytes }
1367   .
1368 }
1369 {
1370   In~the~UTF-16~encoding,~each~Unicode~character~is~encoded~as~
1371   2~or~4~bytes: \\
1372   \iow_indent:n
1373   {
1374     Code~point~in~[U+0000,~U+D7FF]:~two~bytes \\
1375     Code~point~in~[U+D800,~U+DFFF]:~illegal \\
1376     Code~point~in~[U+E000,~U+FFFF]:~two~bytes \\
1377     Code~point~in~[U+10000,~U+10FFFF]:~
1378       a~lead~surrogate~and~a~trail~surrogate \\
1379   }
1380   Lead~surrogates~are~pairs~of~bytes~in~the~range~[0xD800,~0xDBFF],~
1381   and~trail~surrogates~are~in~the~range~[0xDC00,~0xDFFF].
1382   \flag_if_raised:nT { str_missing }
1383   {
1384     \\\\
1385     A~lead~surrogate~was~not~followed~by~a~trail~surrogate.
1386   }
1387   \flag_if_raised:nT { str_extra }
1388   {
1389     \\\\
1390     LaTeX~came~across~a~trail~surrogate~when~it~was~not~expected.

```

```

1391     }
1392     \flag_if_raised:nT { str_end }
1393     {
1394         \\\\
1395         The~string~contained~an~odd~number~of~bytes.~This~is~invalid:~
1396         the~basic~code~unit~for~UTF-16~is~16~bits~(2~bytes).
1397     }
1398 }
(End definition for \l__str_missing_flag, \l__str_extra_flag, and \l__str_end_flag These variables 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 `__str_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.

```

1399 \cs_new_protected_nopar:cpn { __str_convert_decode_utf16be: }
1400   { __str_decode_utf_xvi:Nw 1 \g__str_result_tl \s__stop }
1401 \cs_new_protected_nopar:cpn { __str_convert_decode_utf16le: }
1402   { __str_decode_utf_xvi:Nw 2 \g__str_result_tl \s__stop }
1403 \cs_new_protected_nopar:cpn { __str_convert_decode_utf16: }
1404   {
1405     \exp_after:wN __str_decode_utf_xvi_bom>NN
1406     \g__str_result_tl \s__stop \s__stop \s__stop
1407   }
1408 \cs_new_protected:Npn __str_decode_utf_xvi_bom>NN #1#2
1409   {
1410     \str_if_eq_x:nnTF { #1#2 } { ^ff ^fe }
1411       { __str_decode_utf_xvi:Nw 2 }
1412       {
1413         \str_if_eq_x:nnTF { #1#2 } { ^fe ^ff }
1414           { __str_decode_utf_xvi:Nw 1 }
1415           { __str_decode_utf_xvi:Nw 1 #1#2 }
1416       }
1417   }
1418 \cs_new_protected:Npn __str_decode_utf_xvi:Nw #1#2 \s__stop
1419   {
1420     \flag_clear:n { str_error }
1421     \flag_clear:n { str_missing }
1422     \flag_clear:n { str_extra }
1423     \flag_clear:n { str_end }
1424     \cs_set:Npn __str_tmp:w ##1 ##2 { ' ## #1 }
1425     \tl_gset:Nx \g__str_result_tl
```

```

1426    {
1427        \exp_after:wN \__str_decode_utf_xvi_pair:NN
1428            #2 \q_nil \q_nil
1429            \__prg_break_point:
1430        }
1431    \__str_if_flag_error:n { str_error } { utf16-decode } { }
1432 }
(End definition for \__str_convert_decode_utf16:, \__str_convert_decode_utf16be:, and \__str_convert_decode_utf16l:
These functions are documented on page ??.)
```

Bytes are read two at a time. At this stage, `\@@_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 (ε -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.

The magic number "D7F7 is such that "D7F7*"400 = "D800*"400+"DC00-"10000.

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.

```

1433 \cs_new:Npn \__str_decode_utf_xvi_pair:NN #1#2
1434 {
1435     \if_meaning:w \q_nil #2
1436         \__str_decode_utf_xvi_pair_end:Nw #1
1437     \fi:
1438     \if_case:w
1439         \__int_eval:w ( \__str_tmp:w #1#2 - "D6 ) / \c_four \__int_eval_end:
1440     \or: \exp_after:wN \__str_decode_utf_xvi_quad>NNwNN
1441     \or: \exp_after:wN \__str_decode_utf_xvi_extra>NNw
1442     \fi:
1443     #1#2 \s_t1
1444     \int_eval:n { "100 * \__str_tmp:w #1#2 + \__str_tmp:w #2#1 } \s_t1
1445     \__str_decode_utf_xvi_pair:NN
```

```

1446    }
1447 \cs_new:Npn \__str_decode_utf_xvi_quad:NNwNN
1448     #1#2 #3 \__str_decode_utf_xvi_pair:NN #4#5
1449 {
1450   \if_meaning:w \q_nil #5
1451     \__str_decode_utf_xvi_error:nNN { missing } #1#2
1452     \__str_decode_utf_xvi_pair_end:Nw #4
1453   \fi:
1454   \if_int_compare:w
1455     \if_int_compare:w \__str_tmp:w #4#5 < "DC \exp_stop_f:
1456       \c_zero = \c_one
1457     \else:
1458       \__str_tmp:w #4#5 < "E0 \exp_stop_f:
1459     \fi:
1460     #1 #2 #4 #5 \s__tl
1461   \int_eval:n
1462   {
1463     ( "100 * \__str_tmp:w #1#2 + \__str_tmp:w #2#1 - "D7F7 ) * "400
1464     + "100 * \__str_tmp:w #4#5 + \__str_tmp:w #5#4
1465   }
1466   \s__tl
1467   \exp_after:wN \use_i:nnn
1468 \else:
1469   \__str_decode_utf_xvi_error:nNN { missing } #1#2
1470 \fi:
1471 \__str_decode_utf_xvi_pair:NN #4#5
1472 }
1473 \cs_new:Npn \__str_decode_utf_xvi_pair_end:Nw #1 \fi:
1474 {
1475   \fi:
1476   \if_meaning:w \q_nil #1
1477   \else:
1478     \__str_decode_utf_xvi_error:nNN { end } #1 \prg_do_nothing:
1479   \fi:
1480   \__prg_break:
1481 }
1482 \cs_new:Npn \__str_decode_utf_xvi_extra>NNw #1#2 \s__tl #3 \s__tl
1483   { \__str_decode_utf_xvi_error:nNN { extra } #1#2 }
1484 \cs_new:Npn \__str_decode_utf_xvi_error:nNN #1#2#3
1485 {
1486   \flag_raise:n { str_error }
1487   \flag_raise:n { str_#1 }
1488   #2 #3 \s__tl
1489   \int_use:N \c__str_replacement_char_int \s__tl
1490 }

```

(End definition for `__str_decode_utf_xvi_pair:NN`, `__str_decode_utf_xvi_quad:NNwNN`, and `__str_decode_utf_xvi_pai`
These functions are documented on page ??.)

Restore the original catcodes of bytes 254 and 255.

```
1491 \group_end:
```

```
1492 〈/utf16〉
```

9.9.3 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”.

```
1493  {*utf32}
1494  \group_begin:
1495  \char_set_catcode_other:N \^00
1496  \char_set_catcode_other:N \^fe
1497  \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. 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
1498  \cs_new_protected_nopar:cpn { __str_convert_encode_utf32: }
1499  {
1500    \__str_convert_gmap_internal:N \__str_encode_utf_xxxii_be:n
1501    \tl_gput_left:Nx \g__str_result_tl { ^00 ^00 ^fe ^ff }
1502  }
1503  \cs_new_protected_nopar:cpn { __str_convert_encode_utf32be: }
1504  { \__str_convert_gmap_internal:N \__str_encode_utf_xxxii_be:n }
1505  \cs_new_protected_nopar:cpn { __str_convert_encode_utf32le: }
1506  { \__str_convert_gmap_internal:N \__str_encode_utf_xxxii_le:n }
1507  \cs_new:Npn \__str_encode_utf_xxxii_be:n #1
1508  {
1509    \exp_args:Nf \__str_encode_utf_xxxii_be_aux:nn
1510    { \int_div_truncate:nn {#1} { "100 } } {#1}
1511  }
1512  \cs_new:Npn \__str_encode_utf_xxxii_be_aux:nn #1#2
1513  {
1514    ^00
1515    \__str_output_byte_pair_be:n {#1}
1516    \__str_output_byte:n { #2 - #1 * "100 }
1517  }
1518  \cs_new:Npn \__str_encode_utf_xxxii_le:n #1
1519  {
1520    \exp_args:Nf \__str_encode_utf_xxxii_le_aux:nn
1521    { \int_div_truncate:nn {#1} { "100 } } {#1}
1522  }
1523  \cs_new:Npn \__str_encode_utf_xxxii_le_aux:nn #1#2
1524  {
1525    \__str_output_byte:n { #2 - #1 * "100 }
1526    \__str_output_byte_pair_le:n {#1}
1527    ^00
1528  }
```

(End definition for `__str_convert_encode_utf32:`, `__str_convert_encode_utf32be:`, and `__str_convert_encode_utf32l`. 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.

```

1529   \flag_clear_new:n { str_overflow }
1530   \flag_clear_new:n { str_end }
1531   \__msg_kernel_new:nmm { str } { utf32-decode }
1532   {
1533     Invalid~UTF-32~string: \exp_last_unbraced:Nf \use_none:n
1534     \__str_if_flag_times:nT { str_overflow } { ,~code-point-too-large }
1535     \__str_if_flag_times:nT { str_end } { ,~truncated-string }
1536   .
1537 }
1538 {
1539   In~the~UTF-32~encoding, ~every~Unicode~character~
1540   (in~the~range~[U+0000,~U+10FFFF])~is~encoded~as~4~bytes.
1541   \flag_if_raised:nT { str_overflow }
1542   {
1543     \\\\
1544     LaTeX~came~across~a~code~point~larger~than~1114111,~
1545     the~maximum~code~point~defined~by~Unicode.~
1546     Perhaps~the~string~was~not~encoded~in~the~UTF-32~encoding?
1547   }
1548   \flag_if_raised:nT { str_end }
1549   {
1550     \\\\
1551     The~length~of~the~string~is~not~a~multiple~of~4.~
1552     Perhaps~the~string~was~truncated?
1553   }
1554 }
```

(End definition for `str_overflow` and `str_end`. These variables 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:bom:NNNN` 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_tl`.

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_tl` 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.

```
1555   \cs_new_protected_nopar:cpx { __str_convert_decode_utf32be: }
```

```

1556     { \__str_decode_utf_xxxii:Nw 1 \g__str_result_tl \s__stop }
1557 \cs_new_protected_nopar:cpn { __str_convert_decode_utf32le: }
1558     { \__str_decode_utf_xxxii:Nw 2 \g__str_result_tl \s__stop }
1559 \cs_new_protected_nopar:cpn { __str_convert_decode_utf32: }
1560     {
1561         \exp_after:wN \__str_decode_utf_xxxii_bom:NNNN \g__str_result_tl
1562             \s__stop \s__stop \s__stop \s__stop
1563     }
1564 \cs_new_protected:Npn \__str_decode_utf_xxxii_bom:NNNN #1#2#3#4
1565     {
1566         \str_if_eq_x:nnTF { #1#2#3#4 } { ^~ff ^~fe ^~00 ^~00 }
1567             { \__str_decode_utf_xxxii:Nw 2 }
1568             {
1569                 \str_if_eq_x:nnTF { #1#2#3#4 } { ^~00 ^~00 ^~fe ^~ff }
1570                     { \__str_decode_utf_xxxii:Nw 1 }
1571                     { \__str_decode_utf_xxxii:Nw 1 #1#2#3#4 }
1572             }
1573     }
1574 \cs_new_protected:Npn \__str_decode_utf_xxxii:Nw #1#2 \s__stop
1575     {
1576         \flag_clear:n { str_overflow }
1577         \flag_clear:n { str_end }
1578         \flag_clear:n { str_error }
1579         \cs_set:Npn \__str_tmp:w ##1 ##2 { ' ## #1 }
1580         \tl_gset:Nx \g__str_result_tl
1581             {
1582                 \exp_after:wN \__str_decode_utf_xxxii_loop:NNNN
1583                     #2 \s__stop \s__stop \s__stop \s__stop
1584                     \__prg_break_point:
1585             }
1586         \__str_if_flag_error:nnx { str_error } { utf32-decode } { }
1587     }
1588 \cs_new:Npn \__str_decode_utf_xxxii_loop:NNNN #1#2#3#4
1589     {
1590         \if_meaning:w \s__stop #4
1591             \exp_after:wN \__str_decode_utf_xxxii_end:w
1592         \fi:
1593         #1#2#3#4 \s__tl
1594         \if_int_compare:w \__str_tmp:w #1#4 > \c_zero
1595             \flag_raise:n { str_overflow }
1596             \flag_raise:n { str_error }
1597             \int_use:N \c__str_replacement_char_int
1598         \else:
1599             \if_int_compare:w \__str_tmp:w #2#3 > \c_sixteen
1600                 \flag_raise:n { str_overflow }
1601                 \flag_raise:n { str_error }
1602                 \int_use:N \c__str_replacement_char_int
1603             \else:
1604                 \int_eval:n
1605                     { \__str_tmp:w #2#3*"10000 + \__str_tmp:w #3#2*"100 + \__str_tmp:w #4#1 }
```

```

1606          \fi:
1607          \fi:
1608          \s__tl
1609          \_\_str\_decode\_utf\_xxxii\_loop:NNNN
1610      }
1611      \cs_new:Npn \_\_str\_decode\_utf\_xxxii\_end:w #1 \s__stop
1612      {
1613          \tl_if_empty:nF {#1}
1614          {
1615              \flag_raise:n { str_end }
1616              \flag_raise:n { str_error }
1617              #1 \s__tl
1618              \int_use:N \c__str_replacement_char_int \s__tl
1619          }
1620          \_\_prg_break:
1621      }
(End definition for \_\_str\_convert\_decode\_utf32:, \_\_str\_convert\_decode\_utf32be:, and \_\_str\_convert\_decode\_utf32l
These functions are documented on page ??.)
    Restore the original catcodes of bytes 0, 254 and 255.

1622 \group_end:
1623 </utf32>

```

9.9.4 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.

```

1624 (*iso88591)
1625 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso88591 }
1626 {
1627 }
1628 {
1629 }
1630 </iso88591>
1631 (*iso88592)
1632 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso88592 }
1633 {
1634     { A1 } { 0104 }
1635     { A2 } { 02D8 }
1636     { A3 } { 0141 }
1637     { A5 } { 013D }
1638     { A6 } { 015A }
1639     { A9 } { 0160 }
1640     { AA } { 015E }
1641     { AB } { 0164 }
1642     { AC } { 0179 }
1643     { AE } { 017D }
1644     { AF } { 017B }

```

```

1645 { B1 } { 0105 }
1646 { B2 } { 02DB }
1647 { B3 } { 0142 }
1648 { B5 } { 013E }
1649 { B6 } { 015B }
1650 { B7 } { 02C7 }
1651 { B9 } { 0161 }
1652 { BA } { 015F }
1653 { BB } { 0165 }
1654 { BC } { 017A }
1655 { BD } { 02DD }
1656 { BE } { 017E }
1657 { BF } { 017C }
1658 { CO } { 0154 }
1659 { C3 } { 0102 }
1660 { C5 } { 0139 }
1661 { C6 } { 0106 }
1662 { C8 } { 010C }
1663 { CA } { 0118 }
1664 { CC } { 011A }
1665 { CF } { 010E }
1666 { DO } { 0110 }
1667 { D1 } { 0143 }
1668 { D2 } { 0147 }
1669 { D5 } { 0150 }
1670 { D8 } { 0158 }
1671 { D9 } { 016E }
1672 { DB } { 0170 }
1673 { DE } { 0162 }
1674 { EO } { 0155 }
1675 { E3 } { 0103 }
1676 { E5 } { 013A }
1677 { E6 } { 0107 }
1678 { E8 } { 010D }
1679 { EA } { 0119 }
1680 { EC } { 011B }
1681 { EF } { 010F }
1682 { F0 } { 0111 }
1683 { F1 } { 0144 }
1684 { F2 } { 0148 }
1685 { F5 } { 0151 }
1686 { F8 } { 0159 }
1687 { F9 } { 016F }
1688 { FB } { 0171 }
1689 { FE } { 0163 }
1690 { FF } { 02D9 }

1691 }
1692 {
1693 }
1694 </iso88592>

```

```

1695  /*iso88593*/
1696  \__str_declare_eight_bit_encoding:nnn { iso88593 }
1697  {
1698      { A1 } { 0126 }
1699      { A2 } { 02D8 }
1700      { A6 } { 0124 }
1701      { A9 } { 0130 }
1702      { AA } { 015E }
1703      { AB } { 011E }
1704      { AC } { 0134 }
1705      { AF } { 017B }
1706      { B1 } { 0127 }
1707      { B6 } { 0125 }
1708      { B9 } { 0131 }
1709      { BA } { 015F }
1710      { BB } { 011F }
1711      { BC } { 0135 }
1712      { BF } { 017C }
1713      { C5 } { 010A }
1714      { C6 } { 0108 }
1715      { D5 } { 0120 }
1716      { D8 } { 011C }
1717      { DD } { 016C }
1718      { DE } { 015C }
1719      { E5 } { 010B }
1720      { E6 } { 0109 }
1721      { F5 } { 0121 }
1722      { F8 } { 011D }
1723      { FD } { 016D }
1724      { FE } { 015D }
1725      { FF } { 02D9 }
1726  }
1727  {
1728      { A5 }
1729      { AE }
1730      { BE }
1731      { C3 }
1732      { D0 }
1733      { E3 }
1734      { F0 }
1735  }
1736  
```

```

1737  /*iso88594*/
1738  \__str_declare_eight_bit_encoding:nnn { iso88594 }
1739  {
1740      { A1 } { 0104 }
1741      { A2 } { 0138 }
1742      { A3 } { 0156 }
1743      { A5 } { 0128 }

```

```

1744 { A6 } { 013B }
1745 { A9 } { 0160 }
1746 { AA } { 0112 }
1747 { AB } { 0122 }
1748 { AC } { 0166 }
1749 { AE } { 017D }
1750 { B1 } { 0105 }
1751 { B2 } { 02DB }
1752 { B3 } { 0157 }
1753 { B5 } { 0129 }
1754 { B6 } { 013C }
1755 { B7 } { 02C7 }
1756 { B9 } { 0161 }
1757 { BA } { 0113 }
1758 { BB } { 0123 }
1759 { BC } { 0167 }
1760 { BD } { 014A }
1761 { BE } { 017E }
1762 { BF } { 014B }
1763 { CO } { 0100 }
1764 { C7 } { 012E }
1765 { C8 } { 010C }
1766 { CA } { 0118 }
1767 { CC } { 0116 }
1768 { CF } { 012A }
1769 { DO } { 0110 }
1770 { D1 } { 0145 }
1771 { D2 } { 014C }
1772 { D3 } { 0136 }
1773 { D9 } { 0172 }
1774 { DD } { 0168 }
1775 { DE } { 016A }
1776 { EO } { 0101 }
1777 { E7 } { 012F }
1778 { E8 } { 010D }
1779 { EA } { 0119 }
1780 { EC } { 0117 }
1781 { EF } { 012B }
1782 { FO } { 0111 }
1783 { F1 } { 0146 }
1784 { F2 } { 014D }
1785 { F3 } { 0137 }
1786 { F9 } { 0173 }
1787 { FD } { 0169 }
1788 { FE } { 016B }
1789 { FF } { 02D9 }
1790 }
1791 {
1792 }
1793 </iso88594>

```

```

1794  /*iso88595
1795  \__str_declare_eight_bit_encoding:nnn { iso88595 }
1796  {
1797      { A1 } { 0401 }
1798      { A2 } { 0402 }
1799      { A3 } { 0403 }
1800      { A4 } { 0404 }
1801      { A5 } { 0405 }
1802      { A6 } { 0406 }
1803      { A7 } { 0407 }
1804      { A8 } { 0408 }
1805      { A9 } { 0409 }
1806      { AA } { 040A }
1807      { AB } { 040B }
1808      { AC } { 040C }
1809      { AE } { 040E }
1810      { AF } { 040F }
1811      { B0 } { 0410 }
1812      { B1 } { 0411 }
1813      { B2 } { 0412 }
1814      { B3 } { 0413 }
1815      { B4 } { 0414 }
1816      { B5 } { 0415 }
1817      { B6 } { 0416 }
1818      { B7 } { 0417 }
1819      { B8 } { 0418 }
1820      { B9 } { 0419 }
1821      { BA } { 041A }
1822      { BB } { 041B }
1823      { BC } { 041C }
1824      { BD } { 041D }
1825      { BE } { 041E }
1826      { BF } { 041F }
1827      { C0 } { 0420 }
1828      { C1 } { 0421 }
1829      { C2 } { 0422 }
1830      { C3 } { 0423 }
1831      { C4 } { 0424 }
1832      { C5 } { 0425 }
1833      { C6 } { 0426 }
1834      { C7 } { 0427 }
1835      { C8 } { 0428 }
1836      { C9 } { 0429 }
1837      { CA } { 042A }
1838      { CB } { 042B }
1839      { CC } { 042C }
1840      { CD } { 042D }
1841      { CE } { 042E }
1842      { CF } { 042F }
1843      { DO } { 0430 }

```

```

1844 { D1 } { 0431 }
1845 { D2 } { 0432 }
1846 { D3 } { 0433 }
1847 { D4 } { 0434 }
1848 { D5 } { 0435 }
1849 { D6 } { 0436 }
1850 { D7 } { 0437 }
1851 { D8 } { 0438 }
1852 { D9 } { 0439 }
1853 { DA } { 043A }
1854 { DB } { 043B }
1855 { DC } { 043C }
1856 { DD } { 043D }
1857 { DE } { 043E }
1858 { DF } { 043F }
1859 { EO } { 0440 }
1860 { E1 } { 0441 }
1861 { E2 } { 0442 }
1862 { E3 } { 0443 }
1863 { E4 } { 0444 }
1864 { E5 } { 0445 }
1865 { E6 } { 0446 }
1866 { E7 } { 0447 }
1867 { E8 } { 0448 }
1868 { E9 } { 0449 }
1869 { EA } { 044A }
1870 { EB } { 044B }
1871 { EC } { 044C }
1872 { ED } { 044D }
1873 { EE } { 044E }
1874 { EF } { 044F }
1875 { F0 } { 2116 }
1876 { F1 } { 0451 }
1877 { F2 } { 0452 }
1878 { F3 } { 0453 }
1879 { F4 } { 0454 }
1880 { F5 } { 0455 }
1881 { F6 } { 0456 }
1882 { F7 } { 0457 }
1883 { F8 } { 0458 }
1884 { F9 } { 0459 }
1885 { FA } { 045A }
1886 { FB } { 045B }
1887 { FC } { 045C }
1888 { FD } { 00A7 }
1889 { FE } { 045E }
1890 { FF } { 045F }
1891 }
1892 {
1893 }

```

```

1894 </iso88595>
1895 /*iso88596)
1896 \__str_declare_eight_bit_encoding:nnn { iso88596 }
1897 {
1898     { AC } { 060C }
1899     { BB } { 061B }
1900     { BF } { 061F }
1901     { C1 } { 0621 }
1902     { C2 } { 0622 }
1903     { C3 } { 0623 }
1904     { C4 } { 0624 }
1905     { C5 } { 0625 }
1906     { C6 } { 0626 }
1907     { C7 } { 0627 }
1908     { C8 } { 0628 }
1909     { C9 } { 0629 }
1910     { CA } { 062A }
1911     { CB } { 062B }
1912     { CC } { 062C }
1913     { CD } { 062D }
1914     { CE } { 062E }
1915     { CF } { 062F }
1916     { D0 } { 0630 }
1917     { D1 } { 0631 }
1918     { D2 } { 0632 }
1919     { D3 } { 0633 }
1920     { D4 } { 0634 }
1921     { D5 } { 0635 }
1922     { D6 } { 0636 }
1923     { D7 } { 0637 }
1924     { D8 } { 0638 }
1925     { D9 } { 0639 }
1926     { DA } { 063A }
1927     { EO } { 0640 }
1928     { E1 } { 0641 }
1929     { E2 } { 0642 }
1930     { E3 } { 0643 }
1931     { E4 } { 0644 }
1932     { E5 } { 0645 }
1933     { E6 } { 0646 }
1934     { E7 } { 0647 }
1935     { E8 } { 0648 }
1936     { E9 } { 0649 }
1937     { EA } { 064A }
1938     { EB } { 064B }
1939     { EC } { 064C }
1940     { ED } { 064D }
1941     { EE } { 064E }
1942     { EF } { 064F }

```

```

1943      { F0 } { 0650 }
1944      { F1 } { 0651 }
1945      { F2 } { 0652 }
1946    }
1947  {
1948      { A1 }
1949      { A2 }
1950      { A3 }
1951      { A5 }
1952      { A6 }
1953      { A7 }
1954      { A8 }
1955      { A9 }
1956      { AA }
1957      { AB }
1958      { AE }
1959      { AF }
1960      { BO }
1961      { B1 }
1962      { B2 }
1963      { B3 }
1964      { B4 }
1965      { B5 }
1966      { B6 }
1967      { B7 }
1968      { B8 }
1969      { B9 }
1970      { BA }
1971      { BC }
1972      { BD }
1973      { BE }
1974      { CO }
1975      { DB }
1976      { DC }
1977      { DD }
1978      { DE }
1979      { DF }
1980    }
1981  
```

```
</iso88596>
```

```
1982 (*iso88597)
```

```
1983 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso88597 }
```

```
1984  {

```

```
1985      { A1 } { 2018 }

```

```
1986      { A2 } { 2019 }

```

```
1987      { A4 } { 20AC }

```

```
1988      { A5 } { 20AF }

```

```
1989      { AA } { 037A }

```

```
1990      { AF } { 2015 }

```

```
1991      { B4 } { 0384 }

```

1992	{ B5 } { 0385 }
1993	{ B6 } { 0386 }
1994	{ B8 } { 0388 }
1995	{ B9 } { 0389 }
1996	{ BA } { 038A }
1997	{ BC } { 038C }
1998	{ BE } { 038E }
1999	{ BF } { 038F }
2000	{ CO } { 0390 }
2001	{ C1 } { 0391 }
2002	{ C2 } { 0392 }
2003	{ C3 } { 0393 }
2004	{ C4 } { 0394 }
2005	{ C5 } { 0395 }
2006	{ C6 } { 0396 }
2007	{ C7 } { 0397 }
2008	{ C8 } { 0398 }
2009	{ C9 } { 0399 }
2010	{ CA } { 039A }
2011	{ CB } { 039B }
2012	{ CC } { 039C }
2013	{ CD } { 039D }
2014	{ CE } { 039E }
2015	{ CF } { 039F }
2016	{ DO } { 03A0 }
2017	{ D1 } { 03A1 }
2018	{ D3 } { 03A3 }
2019	{ D4 } { 03A4 }
2020	{ D5 } { 03A5 }
2021	{ D6 } { 03A6 }
2022	{ D7 } { 03A7 }
2023	{ D8 } { 03A8 }
2024	{ D9 } { 03A9 }
2025	{ DA } { 03AA }
2026	{ DB } { 03AB }
2027	{ DC } { 03AC }
2028	{ DD } { 03AD }
2029	{ DE } { 03AE }
2030	{ DF } { 03AF }
2031	{ EO } { 03B0 }
2032	{ E1 } { 03B1 }
2033	{ E2 } { 03B2 }
2034	{ E3 } { 03B3 }
2035	{ E4 } { 03B4 }
2036	{ E5 } { 03B5 }
2037	{ E6 } { 03B6 }
2038	{ E7 } { 03B7 }
2039	{ E8 } { 03B8 }
2040	{ E9 } { 03B9 }
2041	{ EA } { 03BA }

```

2042 { EB } { 03BB }
2043 { EC } { 03BC }
2044 { ED } { 03BD }
2045 { EE } { 03BE }
2046 { EF } { 03BF }
2047 { FO } { 03C0 }
2048 { F1 } { 03C1 }
2049 { F2 } { 03C2 }
2050 { F3 } { 03C3 }
2051 { F4 } { 03C4 }
2052 { F5 } { 03C5 }
2053 { F6 } { 03C6 }
2054 { F7 } { 03C7 }
2055 { F8 } { 03C8 }
2056 { F9 } { 03C9 }
2057 { FA } { 03CA }
2058 { FB } { 03CB }
2059 { FC } { 03CC }
2060 { FD } { 03CD }
2061 { FE } { 03CE }
2062 }
2063 {
2064 { AE }
2065 { D2 }
2066 }
2067 </iso88597>
2068 /*iso88598*/
2069 \_str_declare_eight_bit_encoding:nnn { iso88598 }
2070 {
2071 { AA } { 00D7 }
2072 { BA } { 00F7 }
2073 { DF } { 2017 }
2074 { EO } { 05D0 }
2075 { E1 } { 05D1 }
2076 { E2 } { 05D2 }
2077 { E3 } { 05D3 }
2078 { E4 } { 05D4 }
2079 { E5 } { 05D5 }
2080 { E6 } { 05D6 }
2081 { E7 } { 05D7 }
2082 { E8 } { 05D8 }
2083 { E9 } { 05D9 }
2084 { EA } { 05DA }
2085 { EB } { 05DB }
2086 { EC } { 05DC }
2087 { ED } { 05DD }
2088 { EE } { 05DE }
2089 { EF } { 05DF }
2090 { FO } { 05E0 }

```

```

2091 { F1 } { 05E1 }
2092 { F2 } { 05E2 }
2093 { F3 } { 05E3 }
2094 { F4 } { 05E4 }
2095 { F5 } { 05E5 }
2096 { F6 } { 05E6 }
2097 { F7 } { 05E7 }
2098 { F8 } { 05E8 }
2099 { F9 } { 05E9 }
2100 { FA } { 05EA }
2101 { FD } { 200E }
2102 { FE } { 200F }
2103 }
2104 {
2105 { A1 }
2106 { BF }
2107 { CO }
2108 { C1 }
2109 { C2 }
2110 { C3 }
2111 { C4 }
2112 { C5 }
2113 { C6 }
2114 { C7 }
2115 { C8 }
2116 { C9 }
2117 { CA }
2118 { CB }
2119 { CC }
2120 { CD }
2121 { CE }
2122 { CF }
2123 { DO }
2124 { D1 }
2125 { D2 }
2126 { D3 }
2127 { D4 }
2128 { D5 }
2129 { D6 }
2130 { D7 }
2131 { D8 }
2132 { D9 }
2133 { DA }
2134 { DB }
2135 { DC }
2136 { DD }
2137 { DE }
2138 { FB }
2139 { FC }
2140 }

```

```

2141  </iso88598>
2142  (*iso88599)
2143  \_\_str\_declare\_eight\_bit\_encoding:nnn { iso88599 }
2144  {
2145      { DO } { 011E }
2146      { DD } { 0130 }
2147      { DE } { 015E }
2148      { FO } { 011F }
2149      { FD } { 0131 }
2150      { FE } { 015F }
2151  }
2152  {
2153  }
2154  </iso88599>
2155  (*iso885910)
2156  \_\_str\_declare\_eight\_bit\_encoding:nnn { iso885910 }
2157  {
2158      { A1 } { 0104 }
2159      { A2 } { 0112 }
2160      { A3 } { 0122 }
2161      { A4 } { 012A }
2162      { A5 } { 0128 }
2163      { A6 } { 0136 }
2164      { A8 } { 013B }
2165      { A9 } { 0110 }
2166      { AA } { 0160 }
2167      { AB } { 0166 }
2168      { AC } { 017D }
2169      { AE } { 016A }
2170      { AF } { 014A }
2171      { B1 } { 0105 }
2172      { B2 } { 0113 }
2173      { B3 } { 0123 }
2174      { B4 } { 012B }
2175      { B5 } { 0129 }
2176      { B6 } { 0137 }
2177      { B8 } { 013C }
2178      { B9 } { 0111 }
2179      { BA } { 0161 }
2180      { BB } { 0167 }
2181      { BC } { 017E }
2182      { BD } { 2015 }
2183      { BE } { 016B }
2184      { BF } { 014B }
2185      { CO } { 0100 }
2186      { C7 } { 012E }
2187      { C8 } { 010C }
2188      { CA } { 0118 }
2189      { CC } { 0116 }

```

```

2190 { D1 } { 0145 }
2191 { D2 } { 014C }
2192 { D7 } { 0168 }
2193 { D9 } { 0172 }
2194 { EO } { 0101 }
2195 { E7 } { 012F }
2196 { E8 } { 010D }
2197 { EA } { 0119 }
2198 { EC } { 0117 }
2199 { F1 } { 0146 }
2200 { F2 } { 014D }
2201 { F7 } { 0169 }
2202 { F9 } { 0173 }
2203 { FF } { 0138 }
2204 }
2205 {
2206 }
2207 </iso885910>
2208 {*iso885911}
2209 \_str_declare_eight_bit_encoding:nnn { iso885911 }
2210 {
2211 { A1 } { OE01 }
2212 { A2 } { OE02 }
2213 { A3 } { OE03 }
2214 { A4 } { OE04 }
2215 { A5 } { OE05 }
2216 { A6 } { OE06 }
2217 { A7 } { OE07 }
2218 { A8 } { OE08 }
2219 { A9 } { OE09 }
2220 { AA } { OEOA }
2221 { AB } { OEOB }
2222 { AC } { OEOC }
2223 { AD } { OEOD }
2224 { AE } { OEOE }
2225 { AF } { OEOF }
2226 { B0 } { OE10 }
2227 { B1 } { OE11 }
2228 { B2 } { OE12 }
2229 { B3 } { OE13 }
2230 { B4 } { OE14 }
2231 { B5 } { OE15 }
2232 { B6 } { OE16 }
2233 { B7 } { OE17 }
2234 { B8 } { OE18 }
2235 { B9 } { OE19 }
2236 { BA } { OE1A }
2237 { BB } { OE1B }
2238 { BC } { OE1C }

```

2239 { BD } { OE1D }
2240 { BE } { OE1E }
2241 { BF } { OE1F }
2242 { CO } { OE20 }
2243 { C1 } { OE21 }
2244 { C2 } { OE22 }
2245 { C3 } { OE23 }
2246 { C4 } { OE24 }
2247 { C5 } { OE25 }
2248 { C6 } { OE26 }
2249 { C7 } { OE27 }
2250 { C8 } { OE28 }
2251 { C9 } { OE29 }
2252 { CA } { OE2A }
2253 { CB } { OE2B }
2254 { CC } { OE2C }
2255 { CD } { OE2D }
2256 { CE } { OE2E }
2257 { CF } { OE2F }
2258 { DO } { OE30 }
2259 { D1 } { OE31 }
2260 { D2 } { OE32 }
2261 { D3 } { OE33 }
2262 { D4 } { OE34 }
2263 { D5 } { OE35 }
2264 { D6 } { OE36 }
2265 { D7 } { OE37 }
2266 { D8 } { OE38 }
2267 { D9 } { OE39 }
2268 { DA } { OE3A }
2269 { DF } { OE3F }
2270 { EO } { OE40 }
2271 { E1 } { OE41 }
2272 { E2 } { OE42 }
2273 { E3 } { OE43 }
2274 { E4 } { OE44 }
2275 { E5 } { OE45 }
2276 { E6 } { OE46 }
2277 { E7 } { OE47 }
2278 { E8 } { OE48 }
2279 { E9 } { OE49 }
2280 { EA } { OE4A }
2281 { EB } { OE4B }
2282 { EC } { OE4C }
2283 { ED } { OE4D }
2284 { EE } { OE4E }
2285 { EF } { OE4F }
2286 { FO } { OE50 }
2287 { F1 } { OE51 }
2288 { F2 } { OE52 }

```

2289 { F3 } { 0E53 }
2290 { F4 } { 0E54 }
2291 { F5 } { 0E55 }
2292 { F6 } { 0E56 }
2293 { F7 } { 0E57 }
2294 { F8 } { 0E58 }
2295 { F9 } { 0E59 }
2296 { FA } { 0E5A }
2297 { FB } { 0E5B }
2298 }
2299 {
2300 { DB }
2301 { DC }
2302 { DD }
2303 { DE }
2304 }
2305 </iso885911>
2306 {*iso885913}
2307 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso885913 }
2308 {
2309 { A1 } { 201D }
2310 { A5 } { 201E }
2311 { A8 } { 00D8 }
2312 { AA } { 0156 }
2313 { AF } { 00C6 }
2314 { B4 } { 201C }
2315 { B8 } { 00F8 }
2316 { BA } { 0157 }
2317 { BF } { 00E6 }
2318 { CO } { 0104 }
2319 { C1 } { 012E }
2320 { C2 } { 0100 }
2321 { C3 } { 0106 }
2322 { C6 } { 0118 }
2323 { C7 } { 0112 }
2324 { C8 } { 010C }
2325 { CA } { 0179 }
2326 { CB } { 0116 }
2327 { CC } { 0122 }
2328 { CD } { 0136 }
2329 { CE } { 012A }
2330 { CF } { 013B }
2331 { DO } { 0160 }
2332 { D1 } { 0143 }
2333 { D2 } { 0145 }
2334 { D4 } { 014C }
2335 { D8 } { 0172 }
2336 { D9 } { 0141 }
2337 { DA } { 015A }

```

```

2338 { DB } { 016A }
2339 { DD } { 017B }
2340 { DE } { 017D }
2341 { EO } { 0105 }
2342 { E1 } { 012F }
2343 { E2 } { 0101 }
2344 { E3 } { 0107 }
2345 { E6 } { 0119 }
2346 { E7 } { 0113 }
2347 { E8 } { 010D }
2348 { EA } { 017A }
2349 { EB } { 0117 }
2350 { EC } { 0123 }
2351 { ED } { 0137 }
2352 { EE } { 012B }
2353 { EF } { 013C }
2354 { FO } { 0161 }
2355 { F1 } { 0144 }
2356 { F2 } { 0146 }
2357 { F4 } { 014D }
2358 { F8 } { 0173 }
2359 { F9 } { 0142 }
2360 { FA } { 015B }
2361 { FB } { 016B }
2362 { FD } { 017C }
2363 { FE } { 017E }
2364 { FF } { 2019 }
2365 }
2366 {
2367 }
2368 </iso885913>
2369 <*iso885914>
2370 \__str_declare_eight_bit_encoding:nnn { iso885914 }
2371 {
2372 { A1 } { 1E02 }
2373 { A2 } { 1E03 }
2374 { A4 } { 010A }
2375 { A5 } { 010B }
2376 { A6 } { 1EOA }
2377 { A8 } { 1E80 }
2378 { AA } { 1E82 }
2379 { AB } { 1EOB }
2380 { AC } { 1EF2 }
2381 { AF } { 0178 }
2382 { BO } { 1E1E }
2383 { B1 } { 1E1F }
2384 { B2 } { 0120 }
2385 { B3 } { 0121 }
2386 { B4 } { 1E40 }

```

```

2387 { B5 } { 1E41 }
2388 { B7 } { 1E56 }
2389 { B8 } { 1E81 }
2390 { B9 } { 1E57 }
2391 { BA } { 1E83 }
2392 { BB } { 1E60 }
2393 { BC } { 1EF3 }
2394 { BD } { 1E84 }
2395 { BE } { 1E85 }
2396 { BF } { 1E61 }
2397 { D0 } { 0174 }
2398 { D7 } { 1E6A }
2399 { DE } { 0176 }
2400 { F0 } { 0175 }
2401 { F7 } { 1E6B }
2402 { FE } { 0177 }
2403 }
2404 {
2405 }
2406 </iso885914>
2407 (*iso885915)
2408 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso885915 }
2409 {
2410 { A4 } { 20AC }
2411 { A6 } { 0160 }
2412 { A8 } { 0161 }
2413 { B4 } { 017D }
2414 { B8 } { 017E }
2415 { BC } { 0152 }
2416 { BD } { 0153 }
2417 { BE } { 0178 }
2418 }
2419 {
2420 }
2421 </iso885915>
2422 (*iso885916)
2423 \_\_str\_declare\_eight\_bit\_encoding:nnn { iso885916 }
2424 {
2425 { A1 } { 0104 }
2426 { A2 } { 0105 }
2427 { A3 } { 0141 }
2428 { A4 } { 20AC }
2429 { A5 } { 201E }
2430 { A6 } { 0160 }
2431 { A8 } { 0161 }
2432 { AA } { 0218 }
2433 { AC } { 0179 }
2434 { AE } { 017A }
2435 { AF } { 017B }

```

```

2436 { B2 } { 010C }
2437 { B3 } { 0142 }
2438 { B4 } { 017D }
2439 { B5 } { 201D }
2440 { B8 } { 017E }
2441 { B9 } { 010D }
2442 { BA } { 0219 }
2443 { BC } { 0152 }
2444 { BD } { 0153 }
2445 { BE } { 0178 }
2446 { BF } { 017C }
2447 { C3 } { 0102 }
2448 { C5 } { 0106 }
2449 { D0 } { 0110 }
2450 { D1 } { 0143 }
2451 { D5 } { 0150 }
2452 { D7 } { 015A }
2453 { D8 } { 0170 }
2454 { DD } { 0118 }
2455 { DE } { 021A }
2456 { E3 } { 0103 }
2457 { E5 } { 0107 }
2458 { F0 } { 0111 }
2459 { F1 } { 0144 }
2460 { F5 } { 0151 }
2461 { F7 } { 015B }
2462 { F8 } { 0171 }
2463 { FD } { 0119 }
2464 { FE } { 021B }
2465 }
2466 {
2467 }
2468 ⟨/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	
\#	42, 622
\\$	620
\%	44, 630
\	621
*	65, 88, 971
\\"	39, 617, 839, 840, 843, 971, 1184, 1187, 1188, 1189, 1190, 1195, 1201,
\{	1206, 1213, 1371, 1374, 1375, 1376, 1378, 1384, 1389, 1394, 1543, 1550
\}	40, 618, 844
\^	41, 619, 844
__int_eval:w	199, 203, 208, 219, 234, 259, 299, 300, 316, 378, 423, 431, 712, 1246, 1261, 1291, 1439

```

\__int_eval_end: ..... 380, 434, 712, 1246, 1263, 1291, 1439
\__int_value:w ..... 260, 298, 670, 750, 772, 1165, 1238
\__msg_kernel_error:nmx ..... 486, 545
\__msg_kernel_error:nmx ..... 577
\__msg_kernel_new:nnn ..... 819, 821, 830
\__msg_kernel_new:nnnn ..... 691, 823, 834, 848, 854, 906, 953, 1048, 1173, 1354, 1361, 1531
\__prg_break: ..... 457, 473, 591, 639, 705, 734, 738, 780, 876, 925, 985, 991, 1227, 1310, 1480, 1620
\__prg_break:n ..... 346, 352, 364
\__prg_break_point: ... 347, 353, 458, 474, 592, 640, 706, 735, 739, 781, 877, 926, 986, 992, 1228, 1429, 1584
\__str_collect_delimit_by_q_stop:w . .... 216, 216, 315
\__str_collect_end:nnnnnnnnw ..... 216, 239, 243
\__str_collect_end:wn ..... 216, 227, 237
\__str_collect_loop:wn ..... 216, 218, 222, 233
\__str_collect_loop:wnNNNNNNN ..... 216, 225, 231
\__str_convert>NNnN ..... 533, 538, 541
\__str_convert:nNNnnn ..... 493, 494, 496, 501, 510, 515
\__str_convert:nnm ..... 536, 537, 551, 551
\__str_convert:nnnn ..... 551, 555, 560
\__str_convert:wwnnn .. 520, 525, 533, 533
\__str_convert_decode:_ .. 524, 667, 667
\__str_convert_decode_eight_bit:n .. 722, 728, 728
\__str_convert_decode_utf16: .. 1399
\__str_convert_decode_utf16be: .. 1399
\__str_convert_decode_utf16le: .. 1399
\__str_convert_decode_utf32: .. 1555
\__str_convert_decode_utf32be: .. 1555
\__str_convert_decode_utf32le: .. 1555
\__str_convert_decode_utf8: .. 1217
\__str_convert_encode:_ 529, 671, 675, 700
\__str_convert_encode_eight_bit:n .. 724, 774, 774
\__str_convert_encode_utf16: .. 1317
\__str_convert_encode_utf16be: .. 1317
\__str_convert_encode_utf16le: .. 1317
\__str_convert_encode_utf32: .. 1498
\__str_convert_encode_utf32be: .. 1498
\__str_convert_encode_utf32le: .. 1498
\__str_convert_encode_utf8: ..... 1147
\__str_convert_escape: .. 665, 665, 666
\__str_convert_escape_bytes: .. 665, 666
\__str_convert_escape_hex: .. 1058, 1059
\__str_convert_escape_name: .. 1064, 1067
\__str_convert_escape_string: 1089, 1092
\__str_convert_escape_url: .. 1123, 1124
\__str_convert_gmap:N ..... 451, 668, 741, 1060, 1068, 1093, 1125
\__str_convert_gmap_internal:N ..... 467, 467, 678, 783, 1148, 1330, 1500, 1504, 1506
\__str_convert_gmap_internal_loop:Nw ..... 467
\__str_convert_gmap_internal_loop:Nww ..... 471, 477, 481
\__str_convert_gmap_loop:NN ..... 451, 455, 461, 465
\__str_convert_lowercase_alphanum:n ..... 556, 588, 588
\__str_convert_lowercase_alphanum_loop:N ..... 588, 590, 594, 613
\__str_convert_unescape:_ ..... 653, 654, 656, 664
\__str_convert_unescape_bytes: 653, 664
\__str_convert_unescape_hex: .. 865, 866
\__str_convert_unescape_name: ..... 913
\__str_convert_unescape_string: 969, 976
\__str_convert_unescape_url: ..... 913
\__str_count:n ... 132, 135, 143, 148, 151
\__str_count_loop>NNNNNNNN ..... 132, 138, 144, 149, 162, 167
\__str_count_spaces_loop:wwwwwww .. 113, 119, 125, 130
\__str_count_unsafe:n .. 132, 141, 260, 298
\__str_declare_eight_bit_encoding:nnn ..... 718, 718, 1625, 1632, 1696, 1738, 1795, 1896, 1983, 2069, 2143, 2156, 2209, 2307, 2370, 2408, 2423
\__str_decode_eight_bit_char:N ..... 728, 741, 764
\__str_decode_eight_bit_load:nn .. 728, 732, 745, 752
\__str_decode_eight_bit_load_missing:n ..... 728, 736, 754, 762
\__str_decode_native_char:N 667, 668, 669
\__str_decode_utf_viii_aux:wNnnwN .. 1217, 1260, 1274
\__str_decode_utf_viii_continuation:wwN .. 1217, 1245, 1252, 1290

```

```

\__str_decode_utf_viii_end: .....
..... 1217, 1227, 1304
\__str_decode_utf_viii_overflow:w .....
..... 1217, 1288, 1297
\__str_decode_utf_viii_start:N .....
..... 1217, 1226,
1232, 1250, 1253, 1272, 1275, 1295
\__str_decode_utf_xvi:Nw .....
..... 1399,
1400, 1402, 1411, 1414, 1415, 1418
\__str_decode_utf_xvi_bom:NN .....
..... 1399, 1405, 1408
\__str_decode_utf_xvi_error:nNN .....
.. 1433, 1451, 1469, 1478, 1483, 1484
\__str_decode_utf_xvi_extra>NNw .....
..... 1433, 1441, 1482
\__str_decode_utf_xvi_pair:NN .....
.. 1427, 1433, 1433, 1445, 1448, 1471
\__str_decode_utf_xvi_pair_end:Nw .....
..... 1433, 1436, 1452, 1473
\__str_decode_utf_xvi_quad>NNNN ... .
..... 1433, 1440, 1447
\__str_decode_utf_xxxii:Nw ... .
1555,
1556, 1558, 1567, 1570, 1571, 1574
\__str_decode_utf_xxxii_bom>NNNN ... .
..... 1555, 1561, 1564
\__str_decode_utf_xxxii_end:w .....
..... 1555, 1591, 1611
\__str_decode_utf_xxxii_loop>NNNN ... .
..... 1555, 1582, 1588, 1609
\__str_encode_eight_bit_char:n .....
..... 774, 783, 797
\__str_encode_eight_bit_char_aux:n .. .
..... 774, 808, 811
\__str_encode_eight_bit_load:nn .....
..... 774, 778, 787, 795
\__str_encode_native_char:n 674, 678, 682
\__str_encode_native_filter:N .....
..... 699
\__str_encode_native_flush: .....
..... 699
\__str_encode_native_loop:w .....
..... 699, 704, 709, 714
\__str_encode_utf_viii_char:n .....
..... 1147, 1148, 1149
\__str_encode_utf_viii_loop:wwnnw ...
..... 1147, 1151, 1158, 1164
\__str_encode_utf_xvi_aux:N .....
..... 1317, 1319, 1323, 1325, 1326
\__str_encode_utf_xvi_char:n .....
..... 1317, 1330, 1333
\__str_encode_utf_xxxii_be:n .....
..... 1498, 1500, 1504, 1507
\__str_encode_utf_xxxii_be_aux:nn ...
..... 1498, 1509, 1512
\__str_encode_utf_xxxii_le:n .....
..... 1498, 1506, 1518
\__str_encode_utf_xxxii_le_aux:nn ...
..... 1498, 1520, 1523
\__str_escape_hex_char:N 1058, 1060, 1061
\__str_escape_name_char:N 1064, 1068, 1069
\__str_escape_string_char:N .....
..... 1089, 1093, 1094
\__str_escape_url_char:N 1123, 1125, 1126
\__str_filter_bytes:n .....
..... 633, 634, 636, 660, 934, 1001
\__str_filter_bytes_aux:N .....
..... 633, 638, 642, 650
\__str_gset_other:Nn .....
9, 87, 93, 519
\__str_gset_other_end:w ...
87, 105, 110
\__str_gset_other_loop:w ...
87, 97, 101, 108
\__str_head:w .....
169, 172, 176
\__str_hexadecimal_use:N .....
..... 373
\__str_hexadecimal_use:NTF .....
..... 9, 373, 886, 896, 937, 939
\__str_if_contains_char>NN .....
..... 343
\__str_if_contains_char>NNT ...
343, 1098
\__str_if_contains_char>NNTF .....
..... 343, 1077, 1083
\__str_if_contains_char:nN .....
..... 350
\__str_if_contains_char:nNTF .....
..... 343, 1134, 1140
\__str_if_contains_char_aux>NN .....
..... 343, 345, 352, 356, 361
\__str_if_contains_char_true: .....
..... 343, 359, 363
\__str_if_escape_name:N .....
..... 1074
\__str_if_escape_name:NTF ...
1064, 1071
\__str_if_escape_string:N .....
..... 1110
\__str_if_escape_string:NTF ...
1089, 1096
\__str_if_escape_url:N .....
..... 1131
\__str_if_escape_url:NTF ...
1123, 1128
\__str_if_flag_error:nnx 483, 483, 502,
511, 661, 679, 742, 784, 880, 928,
929, 994, 995, 1230, 1331, 1431, 1586
\__str_if_flag_no_error:nnx .....
..... 483, 489, 502, 511
\__str_if_flag_times:nT .....
..... 491, 491, 1176, 1177, 1178,
1179, 1364, 1365, 1366, 1534, 1535
\__str_item:ww .....
245, 258, 264
\__str_item_unsafe:nn ...
245, 250, 255, 256
\__str_load_catcodes: ...
571, 615, 615

```

__str_octal_use:N	365	__str_unescape_url_loop:wNN	913, 967
__str_octal_use:NTF	365, 1004, 1006, 1008	__tl_build_end:	707
__str_output_byte:n		__tl_build_one:n	713
.	9, 417, 417, 448, 449, 603,	__tl_gbuild_x:Nw	703
	685, 793, 814, 1161, 1167, 1516, 1525	\~	43, 625
__str_output_byte:w			
.	417, 418, 419, 873, 899, 936, 1003	\u	65, 88, 624, 844, 1187, 1188, 1189
__str_output_byte_pair:nnN			
.	435, 437, 442, 445	A	
__str_output_byte_pair_be:n		\A	66, 89
.	435, 435, 1319, 1323, 1515		
__str_output_byte_pair_le:n		B	
.	435, 440, 1325, 1526	\bool_gset_false:N	500, 509
__str_output_end:	417,	\bool_gset_true:N	490
	418, 426, 433, 878, 898, 950, 1037, 1041	\bool_if:NTF	504, 513
__str_output_hexadecimal:n		\bool_new:N	61
.	417, 425, 1062, 1072, 1129		
__str_output_hexadecimal:w	417, 426, 427	C	
__str_skip_c_zero:w		\c__str_byte_-1_tl	397
.	197, 197, 208, 272, 280, 317	\c__str_byte_0_tl	397
__str_skip_end:NNNNNNNN	197, 211, 215	\c__str_byte_1_tl	397
__str_skip_end:w	197, 202, 209	\c__str_byte_255_tl	397
__str_skip_loop:wNNNNNNNN	197, 200, 207	\c__str_escape_name_not_str	
__str_substr:nN	284, 301, 305	1064, 1065, 1077
__str_substr:nnw	284, 309, 313	\c__str_escape_name_str	1064, 1066, 1083
__str_substr:www	284, 297, 307	\c__str_escape_string_str	1089, 1090, 1099
__str_substr_normalize_range:nn		\c__str_replacement_char_int	38,
.	284, 310, 311, 319	38, 760, 1242, 1268, 1282, 1302,	
__str_substr_unsafe:nnn		1309, 1341, 1489, 1597, 1602, 1618	
.	284, 289, 294, 295	\c_backslash_str	2, 39,
__str_tail_auxi:w	183, 186, 190	39, 1022, 1051, 1053, 1091, 1100, 1104	
__str_tail_auxii:w	183, 193, 196	\c_eight	156, 199, 208, 1106, 1107
__str_tmp:w	25, 25, 914, 962,	\c_fifty_eight	29, 30
	966, 1329, 1336, 1341, 1343, 1346,	\c_five	157
	1347, 1424, 1439, 1444, 1455, 1458,	\c四十_eight	29, 29
	1463, 1464, 1579, 1594, 1599, 1605	\c_four	157, 1439
__str_to_other:n	64, 70, 251, 290	\c_hash_str	2, 39, 42, 962, 1072
__str_to_other_end:w	64, 79, 84	\c_lbrace_str	2, 39, 40
__str_to_other_loop:w	64, 72, 75, 81	\c_max_char_int	2, 36, 36
__str_unescape_hex_auxi:N		\c_max_int	301
.	865, 874, 883, 890, 899	\c_max_register_int	799
__str_unescape_hex_auxii:N		\c_minus_one	1151
.	865, 887, 893, 903	\c_nine	130, 167
__str_unescape_name_loop:wNN	913, 963	\c_ninety_one	29, 32, 597
__str_unescape_string_loop:wNNN		\c_ninety_seven	29, 33, 607
.	969, 990, 998, 1038, 1041	\c_one	158, 272, 316, 327, 367, 599,
__str_unescape_string_newlines:wN			751, 761, 794, 876, 1030, 1034, 1456
.	969, 984, 1042, 1046	\c_one_hundred_twenty_seven	29, 35
__str_unescape_string_repeat:NNNNNN		\c_one_hundred_twenty_three	29, 34, 606
.	969, 1013, 1015, 1017, 1040		

```

\c_percent_str ..... 2, 39, 44, 966, 1129
\c_rbrace_str ..... 2, 39, 41
\c_seven ..... 156, 224, 234
\c_six ..... 156
\c_sixteen ..... 1599
\c_sixty_five ..... 29, 31, 598
\c_thirty_two ..... 603
\c_three ..... 157, 1010
\c_tilde_str ..... 2, 39, 43
\c_two ..... 158, 375, 1167
\c_two_hundred_fifty_six ..... 684, 813
\c_zero ..... 158, 212, 215,
             240, 266, 299, 323, 325, 341, 398,
             401, 406, 943, 1151, 1152, 1456, 1594
\char_set_catcode_alignment:N ..... 621
\char_set_catcode_comment:N ..... 630
\char_set_catcode_escape:N ..... 617
\char_set_catcode_group_begin:N ... 618
\char_set_catcode_group_end:N ... 619
\char_set_catcode_ignore:N ..... 624
\char_set_catcode_letter:N ..... 627
\char_set_catcode_math_superscript:N
             ..... 623
\char_set_catcode_math_toggle:N ... 620
\char_set_catcode_other:N ..... 629,
             972, 973, 1315, 1316, 1495, 1496, 1497
\char_set_catcode_other:n ..... 398, 672
\char_set_catcode_parameter:N ..... 622
\char_set_catcode_space:N ..... 625
\char_set_lccode:nn ..... .
             65, 66, 88, 89, 401, 406, 971
\cs:w ..... 423, 431
\cs_end: ..... 434
\cs_generate_variant:Nn ..... 9, 23, 396
\cs_gset_eq:cc ..... 582, 586
\cs_if_exist:cF ..... 553, 562, 566, 580
\cs_if_exist:NF ..... 6
\cs_new:Npn ..... .
             7, 70, 75, 84, 101, 110, 113, 115,
             125, 133, 141, 146, 151, 178, 191,
             196, 197, 207, 209, 215, 216, 222,
             231, 237, 243, 246, 254, 256, 264,
             285, 293, 295, 305, 307, 313, 319,
             356, 417, 425, 435, 440, 445, 461,
             477, 588, 594, 636, 642, 669, 682,
             764, 797, 811, 883, 893, 932, 998,
             1040, 1042, 1061, 1069, 1094, 1126,
             1149, 1158, 1232, 1252, 1274, 1297,
             1333, 1433, 1447, 1473, 1482, 1484,
             1507, 1512, 1518, 1523, 1588, 1611
\cs_new_eq:NN ..... 8,
             394, 395, 634, 664, 666, 861, 862, 863
\cs_new_nopar:Npn ..... 132, 169,
             183, 245, 284, 363, 419, 427, 433, 1304
\cs_new_protected:cpn ..... 916
\cs_new_protected:cpx ..... 21
\cs_new_protected:Npn .. 93, 451, 467,
             483, 489, 491, 515, 533, 541, 551,
             560, 615, 709, 718, 728, 745, 754,
             774, 787, 1326, 1408, 1418, 1564, 1574
\cs_new_protected_nopar:cpn ..... .
             721, 723, 1147, 1217,
             1317, 1322, 1324, 1399, 1401, 1403,
             1498, 1503, 1505, 1555, 1557, 1559
\cs_new_protected_nopar:Npn .... 25,
             493, 495, 654, 656, 665, 667, 675,
             700, 866, 976, 1059, 1067, 1092, 1124
\cs_set:Npn ..... .
             .. 162, 170, 176, 184, 190, 1424, 1579
\cs_set_eq:NN ..... 502, 511, 1329
\cs_set_protected:Npn ..... 914
\cs_to_str:N ..... 39, 40, 41, 42, 43, 44

```

E

```

\else: 201, 212, 226, 240, 326, 329, 332,
             341, 369, 377, 387, 547, 602, 605,
             608, 647, 686, 801, 815, 1045, 1079,
             1082, 1114, 1117, 1136, 1139, 1239,
             1244, 1264, 1283, 1286, 1337, 1342,
             1345, 1457, 1468, 1477, 1598, 1603
\exp_after:wN ..... 72, 97, 119,
             128, 138, 149, 165, 172, 180, 186,
             193, 200, 202, 208, 211, 218, 225,
             227, 233, 239, 258, 259, 271, 279,
             297, 298, 299, 315, 316, 345, 359,
             379, 389, 421, 422, 429, 430, 455,
             456, 471, 472, 520, 525, 590, 704,
             750, 759, 805, 924, 962, 966, 984,
             990, 1162, 1164, 1226, 1245, 1246,
             1260, 1263, 1290, 1291, 1405, 1427,
             1440, 1441, 1467, 1561, 1582, 1591
\exp_args:Nc ..... 23
\exp_args:Ncc ..... 538
\exp_args:Nf .. 248, 250, 287, 289, 301,
             309, 437, 442, 1346, 1347, 1509, 1520
\exp_args>NNf ..... 792
\exp_args:No ..... 114, 132,
             169, 183, 245, 255, 284, 294, 400, 660
\exp_args:Nx ..... 555
\exp_last_unbraced:Nf 874, 1175, 1363, 1533

```

\exp_last_unbraced:Nx	732, 736, 778	\g__str_error_bool	
\exp_not:c	22	61, 61, 490, 500, 504, 509, 513
\exp_not:N	22	\g__str_result_tl	
\exp_stop_f:	212, 240, 324, 330,	28, 28, 453, 457, 469, 473,
	367, 375, 599, 645, 768, 803, 804,	519, 531, 659, 660, 703, 705, 871,	
	1076, 1080, 1112, 1115, 1133, 1137,	875, 922, 924, 982, 985, 988, 991,	
	1160, 1235, 1237, 1257, 1258, 1277,	1224, 1226, 1320, 1400, 1402, 1406,	
	1279, 1335, 1338, 1339, 1455, 1458	1425, 1501, 1556, 1558, 1561, 1580	
\ExplFileVersion	4	\g_tmpa_str	6
\ExplFileDescription	4	\g_tmpb_str	6
\ExplFileName	4	\group_begin:	64, 87, 397, 517, 570, 671,
\ExplFileVersion	4	730, 776, 868, 918, 970, 978, 1314, 1494	
		\group_end:	
		69, 92, 415, 530, 573, 717, 743,
		785, 881, 930, 996, 1056, 1491, 1622	
			I
\fi:	80, 84, 85, 106, 110, 111,	\if_case:w	212, 240, 378, 1438
	129, 166, 190, 204, 212, 215, 228,	\if_charcode:w	187, 358, 1045
	240, 241, 328, 334, 335, 341, 346,	\if_int_compare:w	
	352, 360, 371, 390, 392, 549, 601,	199, 212, 224, 240, 323, 324,
	604, 610, 611, 612, 649, 689, 770,	330, 340, 367, 375, 597, 598, 599,	
	771, 806, 807, 809, 817, 1012, 1045,	606, 607, 645, 684, 767, 768, 799,	
	1085, 1086, 1119, 1120, 1142, 1143,	802, 803, 813, 1010, 1076, 1080,	
	1163, 1243, 1247, 1257, 1269, 1285,	1112, 1115, 1133, 1137, 1160, 1235,	
	1289, 1292, 1297, 1299, 1344, 1348,	1237, 1256, 1257, 1277, 1279, 1335,	
	1349, 1437, 1442, 1453, 1459, 1470,	1338, 1339, 1454, 1455, 1594, 1599	
	1473, 1475, 1479, 1592, 1606, 1607		
\file_if_exist:nTF	568	\if_meaning:w	78, 104, 127,
\file_input:n	572		164, 543, 1287, 1435, 1450, 1476, 1590
\flag_clear:n	658, 677, 740,	\int_compare:nNnTF	266, 268, 276
	782, 869, 919, 920, 979, 980, 1219,	\int_const:Nn	
	1220, 1221, 1222, 1223, 1328, 1420,	29, 30, 31, 32, 33, 34, 35, 36, 38
	1421, 1422, 1423, 1576, 1577, 1578	\int_div_truncate:nn	438,
\flag_clear_new:n	1169, 1170, 1171,		443, 1105, 1106, 1165, 1346, 1510, 1521
	1172, 1351, 1352, 1353, 1529, 1530	\int_eval:n	
\flag_height:n	492	117, 153, 321, 410, 1444, 1461, 1604
\flag_if_raised:nT		\int_mod:n	1106, 1107, 1347
	\int_new:N	
	490, 492, 1193, 1199, 1204,	\int_set:Nn	27
	1211, 1382, 1387, 1392, 1541, 1548	\int_use:N	631, 870, 921, 981
\flag_if_raised:nTF	485		
\flag_new:n	62, 63		203, 208, 219, 234,
\flag_raise:n	648, 687, 800,		259, 299, 300, 316, 423, 431, 760,
	816, 889, 902, 942, 947, 1033, 1240,		1242, 1246, 1261, 1268, 1282, 1291,
	1241, 1266, 1267, 1280, 1281, 1300,		1302, 1309, 1489, 1597, 1602, 1618
	1301, 1307, 1308, 1340, 1486, 1487,		
	1595, 1596, 1600, 1601, 1615, 1616		
			L
		\l__str_end_flag	1351
		\l__str_extra_flag	1169, 1351

\l__str_internal_int	25, 27, 712, 731, 748, 749, 750, 751, 757, 758, 759, 761, 767, 777, 790, 791, 792, 794, 802	S
\l__str_internal_tl 25, 26, 399, 400, 402, 404, 564, 565, 566, 568, 572, 576, 583, 720	\s__stop . 1400, 1402, 1406, 1418, 1556, 1558, 1562, 1574, 1583, 1590, 1611
\l__str_missing_flag 1169, 1351	\s__tl 473, 477, 670, 705, 709, 766, 772, 1236, 1248, 1253, 1265, 1270, 1275, 1278, 1293, 1306, 1309, 1443, 1444, 1460, 1466, 1482, 1488, 1489, 1593, 1608, 1617, 1618
\l__str_overflow_flag 1169	\scan_stop: 188, 748, 749, 757, 758, 790, 791
\l__str_overlong_flag 1169	\str_byte 62
\l__tl_build_offset_int 702	\str_case:nnn 6, 343
\l_tmpa_str 6	\str_case:onn 343
\l_tmpb_str 6	\str_case_x:nnn 6, 343, 1020
O		\str_const:cn 10
\or: 213, 241, 382, 383, 384, 385, 386, 1440, 1441	\str_const:cx 10
P		\str_const:Nn 3, 10, 1065, 1066
\pdftex_if_engine:F 832	\str_const:Nx 10, 1090
\pdftex_if_engine:TF 37, 633, 653, 673	\str_count:N 3, 132, 132, 861
\pdftex_strcmp:D 340	\str_count:n 132, 132, 133, 862
\prg_do_nothing: 523, 1478	\str_count_ignore_spaces:n 3, 132, 146, 863
\prg_new_conditional:Npnn 338, 343, 350, 365, 373, 1074, 1110, 1131	\str_count_spaces:N 4, 113, 113
\prg_new_protected_conditional:Npnn 497, 506	\str_count_spaces:n 113, 114, 115, 137
\prg_return_false: 341, 348, 354, 370, 388, 504, 513, 1078, 1081, 1084, 1113, 1116, 1135, 1138, 1141	\str_end 1529
\prg_return_true: 341, 364, 368, 376, 391, 504, 513, 1078, 1084, 1118, 1135, 1141	\str_error 62
\prop_get:NnNF 564	\str_gput_left:cn 10
\prop_gput:Nnn 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60	\str_gput_left:cx 10
\prop_new:N 45	\str_gput_left:Nn 3, 10
\ProvidesExplPackage 3	\str_gput_left:Nx 10
Q		\str_gput_right:cn 10
\q_mark 73, 84	\str_gput_right:cx 10
\q_nil 1428, 1435, 1450, 1476	\str_gput_right:Nn 3, 10
\q_stop 73, 76, 82, 84, 98, 110, 122, 159, 174, 181, 188, 190, 194, 196, 243, 262, 303, 473, 479, 521, 526, 534, 705, 711, 734, 738, 747, 756, 780, 789, 1156, 1158, 1166, 1249, 1287	\str_gput_right:Nx 10
R		\str_gset:cn 10
\RequirePackage 5	\str_gset:cx 10
\reverse_if:N 187	\str_gset:Nn 3, 10
		\str_gset:Nx 10
		\str_gset_convert:Nnnn 8, 493, 495, 507
		\str_gset_convert:NnnnTF 8, 493
		\str_head:N 4, 169, 169
		\str_head:n 169, 169, 170
		\str_head_ignore_spaces:n 4, 169, 178
		\str_if_eq>NN 338
		\str_if_eq:NNTF 338
		\str_if_eq:nnTF 5, 338
		\str_if_eq_p:NN 338
		\str_if_eq_p:nn 5, 338
		\str_if_eq_x:nnTF 5, 338, 1410, 1413, 1566, 1569
		\str_if_eq_x_p:nn 5, 338

\str_item:Nn	4, 245, 245	\tex_toks:D	750, 759, 769, 792, 804
\str_item:nn	245, 245, 246	\tl_clear:N	576
\str_item_ignore_spaces:nn ..	4, 245, 254	\tl_const:cn	416, 725, 726
\str_length:N	861, 861	\tl_const:cx	409
\str_length:n	861, 862	\tl_const:Nx	39, 40, 41, 42, 43, 44
\str_length_ignore_spaces:n ..	861, 863	\tl_gput_left:Nx	1320, 1501
\str_new:c	8	\tl_gset:Nx	95, 453, 469, 659,
\str_new:N	3, 8, 8, 9	871, 922, 982, 988, 1224, 1425, 1580	
\str_overflow	1529	\tl_gset_eq:NN	496, 512
\str_put_left:cn	10	\tl_if_empty:nF	544, 1613
\str_put_left:cx	10	\tl_if_empty:nTF	306
\str_put_left:Nn	3, 10	\tl_map_function:nN	626, 628
\str_put_left:Nx	10	\tl_map_inline:Nn	402, 404
\str_put_right:cn	10	\tl_map_inline:nn	10, 400
\str_put_right:cx	10	\tl_new:N	8, 26, 28
\str_put_right:Nn	3, 10	\tl_set:Nn	565, 720
\str_put_right:Nx	10	\tl_set:Nx	399
\str_set:cn	10	\tl_set_eq:NN	494, 503
\str_set:cx	10	\tl_show:N	395
\str_set:Nn	3, 10	\tl_show:n	394
\str_set:Nx	10	\tl_to_lowercase:n ..	67, 90, 407, 713, 974
\str_set_convert:Nnnn ..	8, 493, 493, 498	\tl_to_str:N	3, 340, 875
\str_set_convert:NnnnTF	8, 493	\tl_to_str:n	22, 72, 97,
\str_show:c	394	120, 138, 149, 173, 181, 188, 194,	
\str_show:N	6, 394, 395, 396	248, 255, 287, 294, 399, 521, 526, 591	
\str_show:n	394, 394	\tl_use:c	733, 737, 779
\str_substr:Nnm	5, 284, 284	\token_to_str:N	367, 375, 379
\str_substr:nnn	284, 284, 285		
\str_substr_ignore_spaces:nnn	5, 284, 293		
\str_tail:N	4, 183, 183		
\str_tail:n	183, 183, 184		
\str_tail_ignore_spaces:n ..	4, 183, 191		

T

\tex_advance:D	751, 761, 794
\tex_dimen:D	748,
	757, 767, 768, 769, 790, 802, 803, 804
\tex_endlinechar:D	631
\tex_escapechar:D	870, 921, 981
\tex_lccode:D	712
\tex_roman numeral:D	272, 280, 317
\tex_skip:D	749, 758, 768, 791, 803
\tex_the:D	769, 804

U

\use:n	634
\use:nn	438
\use_i:nn	943, 948, 1034, 1038
\use_i:nnn	422, 951, 1467
\use_i_delimit_by_q_stop:nw	177, 180, 271, 279
\use_ii_i:nn	6, 6, 7, 443, 528
\use_none:n	364,
	389, 416, 430, 463, 596, 644, 885,
	895, 935, 1002, 1175, 1255, 1363, 1533
\use_none:nn	487, 805
\use_none_delimit_by_q_stop:w	128, 165, 269,
	277, 479, 711, 747, 756, 789, 1162, 1249