

# The `l3regex` package: regular expressions in $\text{\TeX}$ \*

The  $\text{\LaTeX}$ 3 Project<sup>†</sup>

Released 2012/07/09

## 1 `l3regex` documentation

The `l3regex` package provides regular expression testing, extraction of submatches, splitting, and replacement, all acting on token lists. The syntax of regular expressions is mostly a subset of the PCRE syntax (and very close to POSIX), with some additions due to the fact that  $\text{\TeX}$  manipulates tokens rather than characters. For performance reasons, only a limited set of features are implemented. Notably, back-references are not supported.

Let us give a few examples. After

```
\tl_set:Nn \l_my_tl { That~cat. }
\regex_replace_once:nnN { at } { is } \l_my_tl
```

the token list variable `\l_my_tl` holds the text “`This cat.`”, where the first occurrence of “`at`” was replaced by “`is`”. A more complicated example is a pattern to add a comma at the end of each word:

```
\regex_replace_all:nnN { \w+ } { \0 , } \l_my_tl
```

The `\w` sequence represents any “word” character, and `+` indicates that the `\w` sequence should be repeated as many times as possible (at least once), hence matching a word in the input token list. In the replacement text, `\0` denotes the full match (here, a word).

If a regular expression is to be used several times, it can be compiled once, and stored in a regex variable using `\regex_const:Nn`. For example,

```
\regex_const:Nn \c_foo_regex { \c{begin} \cB. (\c[~BE].*) \cE. }
```

stores in `\c_foo_regex` a regular expression which matches the starting marker for an environment: `\begin`, followed by a begin-group token (`\cB.`), then any number of tokens which are neither begin-group nor end-group character tokens (`\c[~BE].*`), ending with an end-group token (`\cE.`). As explained in the next section, the parentheses “capture” the result of `\c[~BE].*`, giving us access to the name of the environment when doing replacements.

---

\*This file describes v3940, last revised 2012/07/09.

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

## 1.1 Syntax of regular expressions

Most characters match exactly themselves, with an arbitrary category code. Some characters are special and must be escaped with a backslash (*e.g.*, `\*` matches a star character). Some escape sequences of the form backslash-letter also have a special meaning (for instance `\d` matches any digit). As a rule,

- every alphanumeric character (`A–Z`, `a–z`, `0–9`) matches exactly itself, and should not be escaped, because `\A`, `\B`, ... have special meanings;
- non-alphanumeric printable ascii characters can (and should) always be escaped: many of them have special meanings (*e.g.*, use `\(`, `\)`, `\?`, `\.`);
- spaces should always be escaped (even in character classes);
- any other character may be escaped or not, without any effect: both versions will match exactly that character.

Note that these rules play nicely with the fact that many non-alphanumeric characters are difficult to input into  $\text{\TeX}$  under normal category codes. For instance, `\\abc\%` matches the characters `\abc%` (with arbitrary category codes), but does not match the control sequence `\abc` followed by a percent character. Matching control sequences can be done using the `\c{regex}` syntax (see below).

Any special character which appears at a place where its special behaviour cannot apply matches itself instead (for instance, a quantifier appearing at the beginning of a string), after raising a warning.

Characters.

`\x{hh...}` Character with hex code `hh...`

`\xhh` Character with hex code `hh`.

`\a` Alarm (hex 07).

`\e` Escape (hex 1B).

`\f` Form-feed (hex 0C).

`\n` New line (hex 0A).

`\r` Carriage return (hex 0D).

`\t` Horizontal tab (hex 09).

Character types.

`.` A single period matches any token.

`\d` Any decimal digit.

`\h` Any horizontal space character, equivalent to `[\ \^I]`: space and tab.

`\s` Any space character, equivalent to `[\ \^I\^J\^L\^M]`.

`\v` Any vertical space character, equivalent to `[\^J\^K\^L\^M]`. Note that `\^K` is a vertical space, but not a space, for compatibility with Perl.

`\w` Any word character, *i.e.*, alpha-numerics and underscore, equivalent to `[A-Za-z0-9\_]`.

`\D` Any token not matched by `\d`.

`\H` Any token not matched by `\h`.

`\N` Any token other than the `\n` character (hex 0A).

`\S` Any token not matched by `\s`.

`\V` Any token not matched by `\v`.

`\W` Any token not matched by `\w`.

Of those, `.`, `\D`, `\H`, `\N`, `\S`, `\V`, and `\W` will match arbitrary control sequences.

Character classes match exactly one token in the subject.

`[...]` Positive character class. Matches any of the specified tokens.

`[^...]` Negative character class. Matches any token other than the specified characters.

`x-y` Within a character class, this denotes a range (can be used with escaped characters).

`[:<name>:]` Within a character class (one more set of brackets), this denotes the POSIX character class *<name>*, which can be `alnum`, `alpha`, `ascii`, `blank`, `cntrl`, `digit`, `graph`, `lower`, `print`, `punct`, `space`, `upper`, `word`, or `xdigit`.

`[:^<name>:]` Negative POSIX character class.

For instance, `[a-oq-z\cC.]` matches any lowercase latin letter except `p`, as well as control sequences (see below for a description of `\c`).

Quantifiers (repetition).

`?` 0 or 1, greedy.

`??` 0 or 1, lazy.

`*` 0 or more, greedy.

`*?` 0 or more, lazy.

`+` 1 or more, greedy.

`+`? 1 or more, lazy.

`{n}` Exactly *n*.

`{n,}` *n* or more, greedy.

`{n,}?` *n* or more, lazy.

$\{n, m\}$  At least  $n$ , no more than  $m$ , greedy.

$\{n, m\}?$  At least  $n$ , no more than  $m$ , lazy.

anchors and simple assertions.

`\b` Word boundary: either the previous token is matched by `\w` and the next by `\W`, or the opposite. For this purpose, the ends of the token list are considered as `\W`.

`\B` Not a word boundary: between two `\w` tokens or two `\W` tokens (including the boundary).

`^` or `\A` Start of the subject token list.

`$`, `\Z` or `\z` End of the subject token list.

`\G` Start of the current match. This is only different from `^` in the case of multiple matches: for instance `\regex_count:nnN { \G a } { aaba } \1_tmpa_int` yields 2, but replacing `\G` by `^` would result in `\1_tmpa_int` holding the value 1.

Alternation and capturing groups.

`A|B|C` Either one of A, B, or C.

`(...)` Capturing group.

`(?:...)` Non-capturing group.

`(?|...)` Non-capturing group which resets the group number for capturing groups in each alternative. The following group will be numbered with the first unused group number.

The `\c` escape sequence allows to test the category code of tokens, and match control sequences. Each character category is represented by a single uppercase letter:

- C for control sequences;
- B for begin-group tokens;
- E for end-group tokens;
- M for math shift;
- T for alignment tab tokens;
- P for macro parameter tokens;
- U for superscript tokens (up);
- D for subscript tokens (down);
- S for spaces;
- L for letters;

- 0 for others; and
- A for active characters.

The `\c` escape sequence is used as follows.

`\c{<regex>}` A control sequence whose *csname* matches the *<regex>*, anchored at the beginning and end, so that `\c{begin}` matches exactly `\begin`, and nothing else.

`\cX` Applies to the next object, which can be a character, character property, class, or group, and forces this object to only match tokens with category *X* (any of CBEMTPUDSLOA). For instance, `\cL[A-Z\d]` matches uppercase letters and digits of category code letter, `\cC.` matches any control sequence, and `\cO(abc)` matches `abc` where each character has category other.

`\c[XYZ]` Applies to the next object, and forces it to only match tokens with category *X*, *Y*, or *Z* (each being any of CBEMTPUDSLOA). For instance, `\c[LSO](..)` matches two tokens of category letter, space, or other.

`\c[^XYZ]` Applies to the next object and prevents it from matching any token with category *X*, *Y*, or *Z* (each being any of CBEMTPUDSLOA). For instance, `\c[~0]\d` matches digits which have any category different from other.

The category code tests can be used inside classes; for instance, `[\c0\d \c[LO][A-F]]` matches what `TEX` considers as hexadecimal digits, namely digits with category other, or uppercase letters from *A* to *F* with category either letter or other. Within a group affected by a category code test, the outer test can be overridden by a nested test: for instance, `\cL(ab\c0\*cd)` matches `ab*cd` where all characters are of category letter, except `*` which has category other.

The `\u` escape sequence allows to insert the contents of a token list directly into a regular expression or a replacement, avoiding the need to escape special characters. Namely, `\u{<tl var name>}` matches the exact contents of the token list *<tl var>*. Within a `\c{...}` control sequence matching, the `\u` escape sequence only expands its argument once, in effect performing `\tl_to_str:v`. Quantifiers are not supported directly: use a group.

The option `(?i)` makes the match case insensitive (identifying *A-Z* with *a-z*; no Unicode support yet). This applies until the end of the group in which it appears, and can be reverted using `(?-i)`. For instance, in `(?i)(a(?-i)b|c)d`, the letters *a* and *d* are affected by the *i* option. Characters within ranges and classes are affected individually: `(?i)[Y-\]` is equivalent to `[YZ\[\lyz]`, and `(?i)[^aeiou]` matches any character which is not a vowel. Neither character properties, nor `\c{...}` nor `\u{...}` are affected by the *i* option.

In character classes, only `[`, `^`, `-`, `]`, `\` and spaces are special, and should be escaped. Other non-alphanumeric characters can still be escaped without harm. Any escape sequence which matches a single character (`\d`, `\D`, *etc.*) is supported in character classes. If the first character is `^`, then the meaning of the character class is inverted. Ranges of characters can be expressed using `-`, for instance, `[\D 0-5]` and `[~6-9]` are equivalent.

Capturing groups are a means of extracting information about the match. Parenthesized groups are labelled in the order of their opening parenthesis, starting at 1. The

contents of those groups corresponding to the “best” match (leftmost longest) can be extracted and stored in a sequence of token lists using for instance `\regex_extract_once:nnNTF`.

The `\K` escape sequence resets the beginning of the match to the current position in the token list. This only affects what is reported as the full match. For instance,

```
\regex_extract_all:nnN { a \K . } { a123aaxyz } \l_foo_seq
```

results in `\l_foo_seq` containing the items `{1}` and `{a}`: the true matches are `{a1}` and `{aa}`, but they are trimmed by the use of `\K`. The `\K` command does not affect capturing groups: for instance,

```
\regex_extract_once:nnN { (. \K c)+ \d } { acbc3 } \l_foo_seq
```

results in `\l_foo_seq` containing the items `{c3}` and `{bc}`: the true match is `{acbc3}`, with first submatch `{bc}`, but `\K` resets the beginning of the match to the last position where it appears.

## 1.2 Syntax of the replacement text

Most of the features described in regular expressions do not make sense within the replacement text. Escaped characters are supported as inside regular expressions. The whole match is accessed as `\0`, and the first 9 submatches are accessed as `\1`, ..., `\9`. Submatches with numbers higher than 9 are accessed as `\g{<number>}` instead.

For instance,

```
\tl_set:Nn \l_my_tl { Hello,~world! }
\regex_replace_all:nnN { ([er]?1|o) . } { \(\0\-\-\1\ ) } \l_my_tl
```

results in `\l_my_tl` holding `H(ell--el)(o,--o) w(or--o)(ld--l)!`

The characters inserted by the replacement have category code 12 (other) by default. The escape sequence `\c` allows to insert characters with arbitrary category codes, as well as control sequences.

`\cXY` Produces the character `Y` (which can be given as an escape sequence such as `\t` for tab) with category code `X`, which must be one of `CBEMTPUDSLOA`.

`\c{<text>}` Produces the control sequence with csname `<text>`. The `<text>` may contain references to the submatches `\0`, `\1` etc.

## 1.3 Pre-compiling regular expressions

If a regular expression is to be used several times, it is better to compile it once rather than doing it each time the regular expression is used. The compiled regular expression is stored in a variable. All of the `l3regex` module’s functions can be given their regular expression argument either as an explicit string or as a compiled regular expression.

---

```
\regex_new:N \regex_new:N <regex var>
```

---

Creates a new `<regex var>` or raises an error if the name is already taken. The declaration is global. The `<regex var>` will initially be such that it never matches.

---

<code>\regex_set:Nn</code>	<code>\regex_set:Nn &lt;regex var&gt; {(regex)}</code>
<code>\regex_gset:Nn</code>	Stores a compiled version of the <i>&lt;regular expression&gt;</i> in the <i>&lt;regex var&gt;</i> . For instance, this function can be used as
<code>\regex_const:Nn</code>	

---

```

\regex_new:N \l_my_regex
\regex_set:Nn \l_my_regex { my\ (simple\ )? reg(ex|ular\ expression) }

```

The assignment is local for `\regex_set:Nn` and global for `\regex_gset:Nn`. Use `\regex_const:Nn` for compiled expressions which will never change.

---

<code>\regex_show:n</code>	<code>\regex_show:n {(regex)}</code>
<code>\regex_show:N</code>	Shows how l3regex interprets the <i>&lt;regex&gt;</i> . For instance, <code>\regex_show:n {\A X Y}</code> shows

---

```

+-branch
  anchor at start (\A)
  char code 88
+-branch
  char code 89

```

indicating that the anchor `\A` only applies to the first branch: the second branch is not anchored to the beginning of the match.

## 1.4 Matching

All regular expression functions are available in both `:n` and `:N` variants. The former require a “standard” regular expression, while the later require a compiled expression as generated by `\regex_(g)set:Nn`.

---

<code>\regex_match:nnTF</code>	<code>\regex_match:nnTF {(regex)} {(token list)} {(true code)} {(false code)}</code>
<code>\regex_match:NnTF</code>	Tests whether the <i>&lt;regular expression&gt;</i> matches any part of the <i>&lt;token list&gt;</i> . For instance,

---

```

\regex_match:nnTF { b [cde]* } { abecdex } { TRUE } { FALSE }
\regex_match:nnTF { [b-dq-w] } { example } { TRUE } { FALSE }

```

leaves TRUE then FALSE in the input stream.

---

```
\regex_count:nnN
\regex_count:NnN
```

---

```
\regex_count:nnN {<regex>} {<token list>} {<int var>}
```

Sets  $\langle int\ var \rangle$  within the current  $\text{\TeX}$  group level equal to the number of times  $\langle regular\ expression \rangle$  appears in  $\langle token\ list \rangle$ . The search starts by finding the left-most longest match, respecting greedy and ungreedy operators. Then the search starts again from the character following the last character of the previous match, until reaching the end of the token list. Infinite loops are prevented in the case where the regular expression can match an empty token list: then we count one match between each pair of characters. For instance,

```
\int_new:N \l_foo_int
\regex_count:nnN { (b+|c) } { abbababcb } \l_foo_int
```

results in  $\l_foo\_int$  taking the value 5.

## 1.5 Submatch extraction

---

```
\regex_extract_once:nnNTF
\regex_extract_once:NnNTF
```

---

```
\regex_extract_once:nnN {<regex>} {<token list>} {<seq var>}
\regex_extract_once:nnNTF {<regex>} {<token list>} {<seq var>} {<true code>} {<false code>}
```

Finds the first match of the  $\langle regular\ expression \rangle$  in the  $\langle token\ list \rangle$ . If it exists, the match is stored as the zeroth item of the  $\langle seq\ var \rangle$ , and further items are the contents of capturing groups, in the order of their opening parenthesis. The  $\langle seq\ var \rangle$  is assigned locally. If there is no match, the  $\langle seq\ var \rangle$  is cleared. The testing versions insert the  $\langle true\ code \rangle$  into the input stream if a match was found, and the  $\langle false\ code \rangle$  otherwise. For instance, assume that you type

```
\regex_extract_once:nnNTF { \A(La)?TeX(!*)\Z } { LaTeX!!! } \l_foo_seq
{ true } { false }
```

Then the regular expression (anchored at the start with  $\backslash\text{A}$  and at the end with  $\backslash\text{Z}$ ) will match the whole token list. The first capturing group,  $(\text{La})?$ , matches  $\text{La}$ , and the second capturing group,  $(!*)$ , matches  $!!!$ . Thus,  $\l_foo\_seq$  will contain the items  $\{\text{LaTeX}!!!\}$ ,  $\{\text{La}\}$ , and  $\{!!!\}$ , and the `true` branch is left in the input stream.



---

<code>\regex_extract_all:nnNTF</code> <code>\regex_extract_all:NnNTF</code>	<code>\regex_extract_all:nnN {&lt;regex&gt;} {&lt;token list&gt;} &lt;seq var&gt;</code> <code>\regex_extract_all:nnNTF {&lt;regex&gt;} {&lt;token list&gt;} &lt;seq var&gt; {&lt;true code&gt;} {&lt;false code&gt;}</code>
--	---

---

Finds all matches of the *<regular expression>* in the *<token list>*, and stores all the sub-match information in a single sequence (concatenating the results of multiple `\regex_extract_once:nnN` calls). The *<seq var>* is assigned locally. If there is no match, the *<seq var>* is cleared. The testing versions insert the *<true code>* into the input stream if a match was found, and the *<false code>* otherwise. For instance, assume that you type

```
\regex_extract_all:nnNTF { \w+ } { Hello,~world! } \l_foo_seq
{ true } { false }
```

Then the regular expression will match twice, and the resulting sequence contains the two items `{Hello}` and `{world}`, and the `true` branch is left in the input stream.

---

<code>\regex_split:nnNTF</code> <code>\regex_split:NnNTF</code>	<code>\regex_split:nnN {&lt;regular expression&gt;} {&lt;token list&gt;} &lt;seq var&gt;</code> <code>\regex_split:nnNTF {&lt;regular expression&gt;} {&lt;token list&gt;} &lt;seq var&gt; {&lt;true code&gt;} {&lt;false code&gt;}</code>
--	---

---

Splits the *<token list>* into a sequence of parts, delimited by matches of the *<regular expression>*. If the *<regular expression>* has capturing groups, then the token lists that they match are stored as items of the sequence as well. The assignment to *<seq var>* is local. If no match is found the resulting *<seq var>* has the *<token list>* as its sole item. If the *<regular expression>* matches the empty token list, then the *<token list>* is split into single tokens. The testing versions insert the *<true code>* into the input stream if a match was found, and the *<false code>* otherwise. For example, after

```
\seq_new:N \l_path_seq
\regex_split:nnNTF { / } { the/path/for/this/file.tex } \l_path_seq
{ true } { false }
```

the sequence `\l_path_seq` contains the items `{the}`, `{path}`, `{for}`, `{this}`, and `{file.tex}`, and the `true` branch is left in the input stream.

## 1.6 Replacement

---

<code>\regex_replace_once:nnNTF</code> <code>\regex_replace_once:NnNTF</code>	<code>\regex_replace_once:nnN {&lt;regular expression&gt;} {&lt;replacement&gt;} &lt;tl var&gt;</code> <code>\regex_replace_once:nnNTF {&lt;regular expression&gt;} {&lt;replacement&gt;} &lt;tl var&gt; {&lt;true code&gt;} {&lt;false code&gt;}</code>
--	---

---

Searches for the *<regular expression>* in the *<token list>* and replaces the first match with the *<replacement>*. The result is assigned locally to *<tl var>*. In the *<replacement>*, `\0` represents the full match, `\1` represent the contents of the first capturing group, `\2` of the second, *etc.*

---

<code>\regex_replace_all:nnNTF</code> <code>\regex_replace_all:NnNTF</code>	<code>\regex_replace_all:nnN {&lt;regular expression&gt;} {&lt;replacement&gt;} &lt;tl var&gt;</code> <code>\regex_replace_all:nnNTF {&lt;regular expression&gt;} {&lt;replacement&gt;} &lt;tl var&gt; {&lt;true code&gt;} {&lt;false code&gt;}</code>
--	---

---

Replaces all occurrences of the `\regular expression` in the `<token list>` by the `<replacement>`, where `\0` represents the full match, `\1` represent the contents of the first capturing group, `\2` of the second, *etc.* Every match is treated independently, and matches cannot overlap. The result is assigned locally to `<tl var>`.

## 1.7 Bugs, misfeatures, future work, and other possibilities

The following need to be done now.

- Change user function names!
- Clean up the use of messages.
- Rewrite the documentation in a more ordered way, perhaps add a BNF?

Additional error-checking to come.

- Detect that a trailing `\c<category>` is an invalid regex.
- Currently, `a{x34}` is recognized as `a{4}`.
- Cleaner error reporting in the replacement phase.
- Add tracing information.
- Detect attempts to use back-references.
- Test for the maximum register `\c_max_register_int`.
- Find out whether the fact that `\W` and friends match the end-marker leads to bugs. Possibly update `\__regex_item_reverse:n`.
- Enforce that `\cC` can only be followed by a match-all dot.

Code improvements to come.

- Change `\skip` to `\dimen` for the array of active threads, and shift the array of submatch informations so that it starts at `\skip0`.
- Optimize `\c{abc}` for matching a specific control sequence.
- Only build `.,` once.
- Use `\skip` for the left and right state stacks when compiling a regex.
- Should `\__regex_action_free_group:n` only be used for greedy `{n,}` quantifier? (I think not.)
- Quantifiers for `\u` and assertions.

- Improve digit grabbing for the `\g` escape in replacement. Allow arbitrary integer expressions for all those numbers?
- When matching, keep track of an explicit stack of `current_state` and `current_submatches`.
- If possible, when a state is reused by the same thread, kill other subthreads.
- Use `\dimen` registers rather than `\l__regex_balance_tl` to build `\__regex_replacement_balance_one_match:n`.
- Reduce the number of epsilon-transitions in alternatives.
- Optimize simple strings: use less states (`abcade` should give two states, for `abc` and `ade`). [Does that really make sense?]
- Optimize groups with no alternative.
- Optimize states with a single `\__regex_action_free:n`.
- Optimize the use of `\__regex_action_success:` by inserting it in state 2 directly instead of having an extra transition.
- Optimize the use of `\int_step...` functions.
- Groups don't capture within regexes for csnames; optimize and document.
- Decide and document what `\c{\c{...}}` should do in the replacement text, similar questions for `\u`.
- Better “show” for anchors, properties, and catcode tests.
- Does `\K` really need a new state for itself?
- When compiling, use a boolean `in_cs` and less magic numbers.
- Instead of checking whether the character is special or alphanumeric using its character code, check if it is special in regexes with `\cs_if_exist` tests.

The following features are likely to be implemented at some point in the future.

- Allow `\cL(abc)` in replacement text.
- General look-ahead/behind assertions.
- Regex matching on external files.
- Conditional subpatterns with look ahead/behind: “if what follows is [...], then [...]”.
- `(*...)` and `(?...)` sequences to set some options.
- UTF-8 mode for pdfTeX.

- Newline conventions are not done. In particular, we should have an option for `.` not to match newlines. Also, `\A` should differ from `^`, and `\Z`, `\z` and `$` should differ.
- Unicode properties: `\p{...}` and `\P{...}`; `\X` which should match any “extended” Unicode sequence. This requires to manipulate a lot of data, probably using tree-boxes.

The following features of PCRE or Perl will probably not be implemented.

- `\ddd`, matching the character with octal code `ddd`;
- Callout with `(?C...)`, we cannot run arbitrary user code during the matching, because the regex code uses registers in an unsafe way;
- Conditional subpatterns (other than with a look-ahead or look-behind condition): this is non-regular, isn’t it?
- Named subpatterns:  $\TeX$  programmers have lived so far without any need for named macro parameters.

The following features of PCRE or Perl will definitely not be implemented.

- `\cx`, similar to  $\TeX$ ’s own `\^^x`;
- Comments:  $\TeX$  already has its own system for comments.
- `\Q... \E` escaping: this would require to read the argument verbatim, which is not in the scope of this module.
- Atomic grouping, possessive quantifiers: those tools, mostly meant to fix catastrophic backtracking, are unnecessary in a non-backtracking algorithm, and difficult to implement.
- Subroutine calls: this syntactic sugar is difficult to include in a non-backtracking algorithm, in particular because the corresponding group should be treated as atomic. Also, we cannot afford to run user code within the regular expression matching, because of our “misuse” of registers.
- Recursion: this is a non-regular feature.
- Back-references: non-regular feature, this requires backtracking, which is prohibitively slow.
- Backtracking control verbs: intrinsically tied to backtracking.
- `\C` single byte in UTF-8 mode:  $\XeTeX$  and  $\LuaTeX$  serve us characters directly, and splitting those into bytes is tricky, encoding dependent, and most likely not useful anyways.

## 2 l3regex implementation

```

1 <*initex | package>
2 <@@=regex>
3 <*package>
4 \ProvidesExplPackage
5   {\ExplFileName}{\ExplFileDate}{\ExplFileVersion}{\ExplFileDescription}
6 \RequirePackage{l3tl-build, l3tl-analysis, l3flag, l3str}
7 </package>

```

### 2.1 Plan of attack

Most regex engines use backtracking. This allows to provide very powerful features (back-references come to mind first), but it is costly, and raises the problem of catastrophic backtracking. Since T<sub>E</sub>X is not first and foremost a programming language, complicated code tends to run slowly, and we must use faster, albeit slightly more restrictive, techniques, coming from automata theory.

Given a regular expression of  $n$  characters, we do the following:

- (Compiling.) Analyse the regex, finding invalid input, and convert it to an internal representation.
- (Building.) Convert the compiled regex to a non-deterministic finite automaton (NFA) with roughly  $n$  states which accepts precisely token lists matching that regex.
- (Matching.) Loop through the query token list one token (one “position”) at a time, exploring in parallel every possible path (“active thread”) through the NFA, considering active threads in an order determined by the quantifiers’ greediness.

We use the following vocabulary in the code comments (and in variable names).

- *Group*: index of the capturing group,  $-1$  for non-capturing groups.
- *Position*: each token in the query is labelled by an integer  $\langle position \rangle$ , with  $\text{min\_pos} - 1 \leq \langle position \rangle \leq \text{max\_pos}$ . The lowest and highest positions correspond to imaginary begin and end markers (with inaccessible category code and character code).
- *Query*: the token list to which we apply the regular expression.
- *State*: each state of the NFA is labelled by an integer  $\langle state \rangle$  with  $\text{min\_state} \leq \langle state \rangle < \text{max\_state}$ .
- *Active thread*: state of the NFA that is reached when reading the query token list for the matching. Those threads are ordered according to the greediness of quantifiers.
- *Step*: used when matching, starts at 0, incremented every time a character is read, and is not reset when searching for repeated matches. The integer `\l__regex_step_int` is a unique id for all the steps of the matching algorithm.

To achieve a good performance, we abuse TeX's registers in two ways. We access registers directly by number rather than tying them to control sequence using `\int_new:N` and other allocation functions. And we store integers in `\dimen` registers in scaled points (`sp`), using TeX's implicit conversion from dimensions to integers in some contexts. Specifically, the registers are used as follows. When compiling, `\toks` registers are used under the hood by functions from the `l3tl-build` module. When building,

- `\toks<state>` holds the tests and actions to perform in the `<state>` of the NFA.
- (Not implemented yet.) `\skipi` has the form `<group id> plus <left state> minus <right state>`.

When matching,

- `\dimen<state>` is equal to the last `<step>` in which the `<state>` was active.
- (Currently, we use `\skip` instead of `\dimen`.) `\dimen<thread>`, with `min_active ≤ <thread> <max_active>`, is equal to the `<state>` in which the `<thread>` currently is. The `<threads>` are ordered starting from the best to the least preferred.
- `\toks<thread>` holds the submatch information for the `<thread>`, as the contents of a property list.
- `\muskip<position>` holds as its main and stretch components the character and category code of the token at this `<position>` in the query.
- `\toks<position>` holds `<tokens>` which o- and x-expand to the `<position>`-th token in the query.
- `\skip` registers hold the value of end-points of all submatches as would be extracted by the `\regex_extract` functions. Since smaller `\skip` registers are used, the minimum index is twice `max_state`, and the used registers go up to `\l__regex_submatch_int`. They are organized in blocks of `capturing_group`, each block corresponding to one match with all its submatches stored in consecutive `\skips`.

`\count` registers are not abused, which means that we can safely use named integers in this module. Note that `\box` registers are not abused either; maybe we could leverage those for some purpose.

The code is structured as follows. Variables are introduced in the relevant section. First we present some generic helper functions. Then comes the code for compiling a regular expression, and for showing the result of the compilation. The building phase converts a compiled regex to NFA states, and the automaton is run by the code in the following section. The only remaining brick is parsing the replacement text and performing the replacement. We are then ready for all the user functions. Finally, messages, and a little bit of tracing code.

## 2.2 Helpers

`\tl_to_str:V` A variant we need for the `\u` escape in the replacement text.

```
⋮ \cs_generate_variant:Nn \tl_to_str:n { V }
(End definition for \tl_to_str:V)
```

### 2.2.1 Constants and variables

`\__regex_tmp:w` Temporary function used for various short-term purposes.

```
9 \cs_new:Npn \__regex_tmp:w { }
(End definition for \__regex_tmp:w)
```

`\l__regex_internal_a_tl` Temporary variables used for various purposes.

```
\l__regex_internal_b_tl 10 \tl_new:N \l__regex_internal_a_tl
\l__regex_internal_a_int 11 \tl_new:N \l__regex_internal_b_tl
\l__regex_internal_b_int 12 \int_new:N \l__regex_internal_a_int
\l__regex_internal_c_int 13 \int_new:N \l__regex_internal_b_int
\l__regex_internal_bool 14 \int_new:N \l__regex_internal_c_int
\l__regex_internal_seq 15 \bool_new:N \l__regex_internal_bool
\g__regex_internal_tl 16 \seq_new:N \l__regex_internal_seq
17 \tl_new:N \g__regex_internal_tl
```

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

`\c__regex_no_match_regex` This regular expression matches nothing, but is still a valid regular expression. We could use a failing assertion, but I went for an empty class. It is used as the initial value for regular expressions declared using `\regex_new:N`.

```
18 \tl_const:Nn \c__regex_no_match_regex
19 {
20   \__regex_branch:n
21   { \__regex_class:NnnnN \c_true_bool { } { 1 } { 0 } \c_true_bool }
22 }
```

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

`\l__regex_balance_int` The first thing we do when matching is to go once through the query token list and store the information for each token as `\muskip` and `\toks` registers. During this phase, `\l__regex_balance_int` counts the balance of begin-group and end-group character tokens which appear before a given point in the token list, and we store it as the shrink component of each `\muskip` register. This variable is also used to keep track of the balance in the replacement text.

```
23 \int_new:N \l__regex_balance_int
```

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

### 2.2.2 Testing characters

`\__regex_break_point:TF` When testing whether a character of the query token list matches a given character class  
`\__regex_break_true:w` in the regular expression, we often have to test it against several ranges of characters, checking if any one of those matches. This is done with a structure like

```
<test1> ... <test_n>
\__regex_break_point:TF {<true code>} {<false code>}
```

If any of the tests succeeds, it calls `\__regex_break_true:w`, which cleans up and leaves *<true code>* in the input stream. Otherwise, `\__regex_break_point:TF` leaves the *<false code>* in the input stream.

```

24 \cs_new_protected:Npn \__regex_break_true:w
25   #1 \__regex_break_point:TF #2 #3 {#2}
26 \cs_new_protected:Npn \__regex_break_point:TF #1 #2 { #2 }
(End definition for \__regex_break_point:TF This function is documented on page ??.)

```

`\__regex_item_reverse:n` This function makes showing regular expressions easier, and lets us define `\D` in terms of `\d` for instance. There is a subtlety: the end of the query is marked by `-2`, and will thus match `\D` and other negated properties; this case is caught by another part of the code.

```

27 \cs_new_protected:Npn \__regex_item_reverse:n #1
28   {
29     #1
30     \__regex_break_point:TF { } \__regex_break_true:w
31   }
(End definition for \__regex_item_reverse:n)

```

`\_regex_item_caseful_equal:n` Simple comparisons triggering `\__regex_break_true:w` when true.

```

\_regex_item_caseful_range:nn
32 \cs_new_protected:Npn \__regex_item_caseful_equal:n #1
33   {
34     \if_int_compare:w #1 = \l__regex_current_char_int
35     \exp_after:wN \__regex_break_true:w
36     \fi:
37   }
38 \cs_new_protected:Npn \__regex_item_caseful_range:nn #1 #2
39   {
40     \reverse_if:N \if_int_compare:w #1 > \l__regex_current_char_int
41     \reverse_if:N \if_int_compare:w #2 < \l__regex_current_char_int
42     \exp_after:wN \exp_after:wN \exp_after:wN \__regex_break_true:w
43     \fi:
44     \fi:
45   }
(End definition for \__regex_item_caseful_equal:n and \__regex_item_caseful_range:nn)

```

`\_regex_item_caseless_equal:n` For caseless matching, we perform the test both on the `current_char` and on the `case_`  
`\_regex_item_caseless_range:nn` `changed_char`. Before doing the second set of tests, we make sure that `case_changed_`  
`char` has been computed.

```

46 \cs_new_protected:Npn \__regex_item_caseless_equal:n #1
47   {
48     \if_int_compare:w #1 = \l__regex_current_char_int
49     \exp_after:wN \__regex_break_true:w
50     \fi:
51     \if_int_compare:w \l__regex_case_changed_char_int = \c_max_int
52     \__regex_compute_case_changed_char:
53     \fi:
54     \if_int_compare:w #1 = \l__regex_case_changed_char_int
55     \exp_after:wN \__regex_break_true:w

```



```

56   \fi:
57 }
58 \cs_new_protected:Npn \__regex_item_caseless_range:nn #1 #2
59 {
60   \reverse_if:N \if_int_compare:w #1 > \l__regex_current_char_int
61   \reverse_if:N \if_int_compare:w #2 < \l__regex_current_char_int
62   \exp_after:wN \exp_after:wN \exp_after:wN \__regex_break_true:w
63   \fi:
64   \fi:
65   \if_int_compare:w \l__regex_case_changed_char_int = \c_max_int
66   \__regex_compute_case_changed_char:
67   \fi:
68   \reverse_if:N \if_int_compare:w #1 > \l__regex_case_changed_char_int
69   \reverse_if:N \if_int_compare:w #2 < \l__regex_case_changed_char_int
70   \exp_after:wN \exp_after:wN \exp_after:wN \__regex_break_true:w
71   \fi:
72   \fi:
73 }

```

(End definition for `\__regex_item_caseless_equal:n` and `\__regex_item_caseless_range:nn`)

`\__regex_compute_case_changed_char:` This function is called when `\l__regex_case_changed_char_int` has not yet been computed (or rather, when it is set to the marker value `\c_max_int`). If the current character code is in the range [65,90] (upper-case), then add 32, making it lowercase. If it is in the lower-case letter range [97,122], subtract 32.

```

74 \cs_new_protected_nopar:Npn \__regex_compute_case_changed_char:
75 {
76   \int_set_eq:NN \l__regex_case_changed_char_int \l__regex_current_char_int
77   \if_int_compare:w \l__regex_current_char_int < \c_ninety_one
78   \if_int_compare:w \l__regex_current_char_int < \c_sixty_five
79   \else:
80     \int_add:Nn \l__regex_case_changed_char_int { \c_thirty_two }
81   \fi:
82   \else:
83     \if_int_compare:w \l__regex_current_char_int < \c_one_hundred_twenty_three
84     \if_int_compare:w \l__regex_current_char_int < \c_ninety_seven
85     \else:
86       \int_sub:Nn \l__regex_case_changed_char_int { \c_thirty_two }
87     \fi:
88   \fi:
89   \fi:
90 }

```

(End definition for `\__regex_compute_case_changed_char:`)

`\__regex_item_equal:n` Those must always be defined to expand to a `caseful` (default) or `caseless` version, and not be protected: they must expand when compiling, to hard-code which tests are caseless or caseful.

```

91 \cs_new_eq:NN \__regex_item_equal:n ?
92 \cs_new_eq:NN \__regex_item_range:nn ?

```

(End definition for `\__regex_item_equal:n` and `\__regex_item_range:nn`)

`\_regex\_item\_catcode:nT` The argument is a sum of powers of 4 with exponents given by the allowed category codes (between 0 and 13). Dividing by a given power of 4 gives an odd result if and only if that category code is allowed. If the catcode does not match, then skip the character code tests which follow.

```

93 \cs_new_protected:Npn \_regex\_item\_catcode:
94 {
95   "
96   \if_case:w \l\_regex\_current\_catcode\_int
97     1      \or: 4      \or: 10     \or: 40
98     \or: 100   \or:      \or: 1000  \or: 4000
99     \or: 10000 \or:      \or: 100000 \or: 400000
100    \or: 1000000 \or: 4000000 \else: 1*\c\_zero
101    \fi:
102  }
103 \cs_new_protected:Npn \_regex\_item\_catcode:nT #1
104 {
105   \if_int_odd:w \_int\_eval:w #1 / \_regex\_item\_catcode: \_int\_eval\_end:
106   \exp\_after:wN \use:n
107   \else:
108   \exp\_after:wN \use\_none:n
109   \fi:
110 }
111 \cs_new_protected:Npn \_regex\_item\_catcode\_reverse:nT #1#2
112 { \_regex\_item\_catcode:nT {#1} { \_regex\_item\_reverse:n {#2} } }

```

(End definition for `\_regex\_item\_catcode:nT` and `\_regex\_item\_catcode\_reverse:nT` These functions are documented on page ??.)

`\_regex\_item\_exact:nn` This matches an exact  $\langle category \rangle$ - $\langle character\ code \rangle$  pair, or an exact control sequence.  
`\_regex\_item\_exact\_cs:c`

```

113 \cs_new_protected:Npn \_regex\_item\_exact:nn #1#2
114 {
115   \if_int_compare:w #1 = \l\_regex\_current\_catcode\_int
116   \if_int_compare:w #2 = \l\_regex\_current\_char\_int
117   \exp\_after:wN \exp\_after:wN \exp\_after:wN \_regex\_break\_true:w
118   \fi:
119   \fi:
120 }
121 \cs_new_protected:Npn \_regex\_item\_exact\_cs:c #1
122 {
123   \int_compare:nNnTF \l\_regex\_current\_catcode\_int = \c\_zero
124   {
125     \str_if_eq_x:nnTF
126     {
127       \exp\_after:wN \exp\_after:wN \exp\_after:wN \cs\_to\_str:N
128       \tex\_the:D \tex\_toks:D \l\_regex\_current\_pos\_int
129     }
130     { #1 }
131     { \_regex\_break\_true:w } { }
132   }
133   { }

```

```

134 }
(End definition for \_regex_item_exact:nn and \_regex_item_exact_cs:c)

```

`\_regex_item_cs:n` Match a control sequence (the argument is a compiled regex). First test the catcode of the current token to be zero. Then perform the matching test, and break if the csname indeed matches. The three `\exp_after:wN` expand the contents of the `\toks{current position}` (of the form `\exp_not:n {<control sequence>}`) to `<control sequence>`.

```

135 \cs_new_protected:Npn \_regex_item_cs:n #1
136 {
137   \int_compare:nNnT \l__regex_current_catcode_int = \c_zero
138   {
139     \group_begin:
140       \_regex_single_match:
141       \_regex_disable_submatches:
142       \_regex_build_for_cs:n {#1}
143       \bool_set_eq:NN \l__regex_saved_success_bool \g__regex_success_bool
144       \exp_args:Nx \_regex_match:n
145       {
146         \exp_after:wN \exp_after:wN
147         \exp_after:wN \cs_to_str:N
148         \tex_the:D \tex_toks:D \l__regex_current_pos_int
149       }
150       \if_meaning:w \c_true_bool \g__regex_success_bool
151       \group_insert_after:N \_regex_break_true:w
152       \fi:
153       \bool_gset_eq:NN \g__regex_success_bool \l__regex_saved_success_bool
154     \group_end:
155   }
156 }
(End definition for \_regex_item_cs:n)

```

### 2.2.3 Character property tests

`\_regex_prop_d:` Character property tests for `\d`, `\W`, *etc.* These character properties are not affected by the `(?i)` option. The characters recognized by each one are as follows: `\d=[0-9]`, `\_regex_prop_h:` `\w=[0-9A-Z_a-z]`, `\s=[\_\^\^I\^\^J\^\^L\^\^M]`, `\h=[\_\^\^I]`, `\v=[\^\^J-\^\^M]`, and the upper case counterparts match anything that the lower case does not match. The order in which the various tests appear is optimized for usual mostly lower case letter text.

```

\_regex_prop_N:
157 \cs_new_protected_nopar:Npn \_regex_prop_d:
158 { \_regex_item_caseful_range:nn \c_forty_eight { 57 } } % 0--9
159 \cs_new_protected_nopar:Npn \_regex_prop_h:
160 {
161   \_regex_item_caseful_equal:n \c_thirty_two % space
162   \_regex_item_caseful_equal:n \c_nine % tab
163 }
164 \cs_new_protected_nopar:Npn \_regex_prop_s:
165 {
166   \_regex_item_caseful_equal:n \c_thirty_two % space

```

```

167     \_regex_item_caseful_equal:n \c_nine      % tab
168     \_regex_item_caseful_equal:n \c_ten      % lf
169     \_regex_item_caseful_equal:n \c_twelve   % ff
170     \_regex_item_caseful_equal:n \c_thirteen % cr
171 }
172 \cs_new_protected_nopar:Npn \_regex_prop_v:
173 { \_regex_item_caseful_range:nn \c_ten \c_thirteen } % lf, vtab, ff, cr
174 \cs_new_protected_nopar:Npn \_regex_prop_w:
175 {
176     \_regex_item_caseful_range:nn \c_ninety_seven { 122 } % a--z
177     \_regex_item_caseful_range:nn \c_sixty_five { 90 } % A--Z
178     \_regex_item_caseful_range:nn \c_forty_eight { 57 } % 0--9
179     \_regex_item_caseful_equal:n { 95 } % _
180 }
181 \cs_new_protected_nopar:Npn \_regex_prop_N:
182 { \_regex_item_reverse:n { \_regex_item_caseful_equal:n \c_ten } }

```

(End definition for \\_regex\_prop\_d: and others.)

```

\_regex_posix_alnum: POSIX properties. No surprise.
\_regex_posix_alpha: 183 \cs_new_protected_nopar:Npn \_regex_posix_alnum:
\_regex_posix_ascii: 184 { \_regex_posix_alpha: \_regex_posix_digit: }
\_regex_posix_blank: 185 \cs_new_protected_nopar:Npn \_regex_posix_alpha:
\_regex_posix_cntrl: 186 { \_regex_posix_lower: \_regex_posix_upper: }
\_regex_posix_digit: 187 \cs_new_protected_nopar:Npn \_regex_posix_ascii:
\_regex_posix_graph: 188 { \_regex_item_caseful_range:nn \c_zero \c_one_hundred_twenty_seven }
\_regex_posix_lower: 189 \cs_new_eq:NN \_regex_posix_blank: \_regex_prop_h:
\_regex_posix_print: 190 \cs_new_protected_nopar:Npn \_regex_posix_cntrl:
\_regex_posix_punct: 191 {
\_regex_posix_space: 192     \_regex_item_caseful_range:nn \c_zero { 31 }
\_regex_posix_upper: 193     \_regex_item_caseful_equal:n \c_one_hundred_twenty_seven
194 }
    \_regex_posix_word: 195 \cs_new_eq:NN \_regex_posix_digit: \_regex_prop_d:
\_regex_posix_xdigit: 196 \cs_new_protected_nopar:Npn \_regex_posix_graph:
197 { \_regex_item_caseful_range:nn { 33 } { 126 } }
198 \cs_new_protected_nopar:Npn \_regex_posix_lower:
199 { \_regex_item_caseful_range:nn \c_ninety_seven { 122 } }
200 \cs_new_protected_nopar:Npn \_regex_posix_print:
201 { \_regex_item_caseful_range:nn \c_thirty_two { 126 } }
202 \cs_new_protected_nopar:Npn \_regex_posix_punct:
203 {
204     \_regex_item_caseful_range:nn { 33 } { 47 }
205     \_regex_item_caseful_range:nn { 58 } { 64 }
206     \_regex_item_caseful_range:nn { 91 } { 96 }
207     \_regex_item_caseful_range:nn { 123 } { 126 }
208 }
209 \cs_new_protected_nopar:Npn \_regex_posix_space:
210 {
211     \_regex_item_caseful_equal:n \c_thirty_two
212     \_regex_item_caseful_range:nn \c_nine \c_thirteen
213 }

```

```

214 \cs_new_protected_nopar:Npn \__regex_posix_upper:
215   { \__regex_item_caseful_range:nn \c_sixty_five { 90 } }
216 \cs_new_eq:NN \__regex_posix_word: \__regex_prop_w:
217 \cs_new_protected_nopar:Npn \__regex_posix_xdigit:
218   {
219     \__regex_posix_digit:
220     \__regex_item_caseful_range:nn \c_sixty_five { 70 }
221     \__regex_item_caseful_range:nn \c_ninety_seven { 102 }
222   }

```

(End definition for `\__regex_posix_alnum:` and others.)

## 2.2.4 Simple character escape

Before actually parsing the regular expression or the replacement text, we go through them once, converting `\n` to the character 10, *etc.* In this pass, we also convert any special character (`*`, `?`, `{`, *etc.*) or escaped alphanumeric character into a marker indicating that this was a special sequence, and replace escaped special characters and non-escaped alphanumeric characters by markers indicating that those were “raw” characters. The rest of the code can then avoid caring about escaping issues (those can become quite complex to handle in combination with ranges in character classes).

Usage: `\__regex_escape_use:nnnn` *<inline 1>* *<inline 2>* *<inline 3>* *{<token list>}*  
The *<token list>* is converted to a string, then read from left to right, interpreting backslashes as escaping the next character. Unescaped characters are fed to the function *<inline 1>*, and escaped characters are fed to the function *<inline 2>* within an *x*-expansion context (typically those functions perform some tests on their argument to decide how to output them). The escape sequences `\a`, `\e`, `\f`, `\n`, `\r`, `\t` and `\x` are recognized, and those are replaced by the corresponding character, then fed to *<inline 3>*. The result is then left in the input stream. Spaces are ignored unless escaped.

The conversion is mostly done within an *x*-expanding assignment, except for the `\x` escape sequence, which is not amenable to that in general. For this, we use the general framework of `\__tl_build:Nw`.

`\__regex_escape_use:nnnn` The result is built in `\l__regex_internal_a_tl`, which is then left in the input stream. Go through #4 once, applying #1, #2, or #3 as relevant to each character (after de-escaping it). Note that we cannot replace `\tl_set:Nx` and `\__tl_build_one:o` by a single call to `\__tl_build_one:x`, because the *x*-expanding assignment may be interrupted by `\x`.

```

223 \cs_new_protected:Npn \__regex_escape_use:nnnn #1#2#3#4
224   {
225     <trace> \trace_push:nnn { regex } { 1 } { __regex_escape_use:nnnn }
226     \__tl_build:Nw \l__regex_internal_a_tl
227     \cs_set_nopar:Npn \__regex_escape_unescaped:N ##1 { #1 }
228     \cs_set_nopar:Npn \__regex_escape_escaped:N ##1 { #2 }
229     \cs_set_nopar:Npn \__regex_escape_raw:N ##1 { #3 }
230     \int_set:Nn \tex_escapechar:D { 92 }
231     \__str_gset_other:Nn \g__regex_internal_tl { #4 }
232     \tl_set:Nx \l__regex_internal_b_tl
233       {
234         \exp_after:wN \__regex_escape_loop:N \g__regex_internal_tl

```

```

235         { break } \_prg_break_point:
236     }
237     \_tl_build_one:o \l\_regex\_internal\_b\_tl
238     \_tl_build_end:
239 <trace> \trace\_pop:nnn { regex } { 1 } { \_regex\_escape\_use:nnnn }
240     \l\_regex\_internal\_a\_tl
241 }
(End definition for \_regex\_escape\_use:nnnn)

```

\\_regex\\_escape\\_loop:N \\_regex\\_escape\\_loop:N reads one character: if it is special (space, backslash, or end-marker), perform the associated action, otherwise it is simply an unescaped character. After a backslash, the same is done, but unknown characters are “escaped”.

```

242 \cs\_new:Npn \_regex\_escape\_loop:N #1
243 {
244     \cs\_if\_exist\_use:cF { \_regex\_escape\_token\_to\_str:N #1:w }
245     { \_regex\_escape\_unescaped:N #1 }
246     \_regex\_escape\_loop:N
247 }
248 \cs\_new\_nopar:cpn { \_regex\_escape\_c\_backslash\_str :w }
249     \_regex\_escape\_loop:N #1
250 {
251     \cs\_if\_exist\_use:cF { \_regex\_escape\_token\_to\_str:N #1:w }
252     { \_regex\_escape\_escaped:N #1 }
253     \_regex\_escape\_loop:N
254 }

```

(End definition for \\_regex\\_escape\\_loop:N This function is documented on page ??.)

\\_regex\\_escape\\_unescaped:N Those functions are never called before being given a new meaning, so their definitions here don’t matter.

```

255 \cs\_new\_eq:NN \_regex\_escape\_unescaped:N ?
256 \cs\_new\_eq:NN \_regex\_escape\_escaped:N ?
257 \cs\_new\_eq:NN \_regex\_escape\_raw:N ?

```

(End definition for \\_regex\\_escape\\_unescaped:N, \\_regex\\_escape\\_escaped:N, and \\_regex\\_escape\\_raw:N)

\\_regex\\_escape\\_break:w The loop is ended upon seeing the end-marker “break”, with an error if the string ended in a backslash. Spaces are ignored, and \a, \e, \f, \n, \r, \t take their meaning here.

```

\_regex\_escape\_break:w
\_regex\_escape\_break:w
\_regex\_escape\_a:w
\_regex\_escape\_e:w
\_regex\_escape\_f:w
\_regex\_escape\_n:w
\_regex\_escape\_r:w
\_regex\_escape\_t:w
\_regex\_escape\_ :w
258 \cs\_new\_eq:NN \_regex\_escape\_break:w \_prg\_break:
259 \cs\_new\_nopar:cpn { \_regex\_escape\_break:w }
260 {
261     \if\_false: { \fi: }
262     \_msg\_kernel\_error:nn { regex } { trailing-backslash }
263     \exp\_after:wN \use\_none:n \exp\_after:wN { \if\_false: } \fi:
264 }
265 \cs\_new\_nopar:cpn { \_regex\_escape\_~:w } { }
266 \cs\_new\_nopar:cpx { \_regex\_escape\_a:w }
267     { \exp\_not:N \_regex\_escape\_raw:N \iow\_char:N \^^G }
268 \cs\_new\_nopar:cpx { \_regex\_escape\_t:w }
269     { \exp\_not:N \_regex\_escape\_raw:N \iow\_char:N \^^I }
270 \cs\_new\_nopar:cpx { \_regex\_escape\_n:w }

```

```

271 { \exp_not:N \__regex_escape_raw:N \iow_char:N \^^J }
272 \cs_new_nopar:cpx { __regex_escape_/f:w }
273 { \exp_not:N \__regex_escape_raw:N \iow_char:N \^^L }
274 \cs_new_nopar:cpx { __regex_escape_/r:w }
275 { \exp_not:N \__regex_escape_raw:N \iow_char:N \^^M }
276 \cs_new_nopar:cpx { __regex_escape_/e:w }
277 { \exp_not:N \__regex_escape_raw:N \iow_char:N \^^[ }

```

(End definition for \\_\_regex\_escape\_break:w and others. These functions are documented on page ??.)

\\_\_regex\_escape\_/x:w  
\\_\_regex\_escape\_x\_end:w  
\\_\_regex\_escape\_x\_large:n

When \x is encountered, \\_\_regex\_escape\_x\_test:N is responsible for grabbing some hexadecimal digits, and feeding the result to \\_\_regex\_escape\_x\_end:w. If the number is < 256, then it is turned into a byte and fed to \\_\_regex\_escape\_raw:N. Otherwise, interrupt the assignment, and either produce an error, or use a standard \lowercase trick depending on the precise value.

```

278 \cs_new:cpn { __regex_escape_/x:w } \__regex_escape_loop:N
279 {
280   \exp_after:wN \__regex_escape_x_end:w
281   \__int_value:w "0 \__regex_escape_x_test:N
282 }
283 \cs_new:Npn \__regex_escape_x_end:w #1 ;
284 {
285   \int_compare:nNnTF {#1} < \c_two_hundred_fifty_six
286   {
287     \exp_last_unbraced:Nf \__regex_escape_raw:N
288     { \__str_output_byte:n {#1} }
289   }
290   { \__regex_escape_x_large:n {#1} }
291 }
292 \group_begin:
293 \char_set_catcode_other:n { 0 }
294 \cs_new:Npn \__regex_escape_x_large:n #1
295 {
296   \if_false: { \fi: }
297   \__tl_build_one:o \l__regex_internal_b_tl
298   \int_compare:nNnTF {#1} > \c_max_char_int
299   {
300     \__msg_kernel_error:nxx { regex } { x-overflow } {#1}
301     \tl_set:Nx \l__regex_internal_b_tl
302     { \if_false: } \fi: \__regex_escape_loop:N
303   }
304   {
305     \char_set_lccode:nn { \c_zero } {#1}
306     \tl_to_lowercase:n
307     {
308       \tl_set:Nx \l__regex_internal_b_tl
309       { \if_false: } \fi:
310       \__regex_escape_raw:N \^^@
311       \__regex_escape_loop:N
312     }
313   }

```

```

313     }
314 }
315 \group_end:

```

(End definition for `\_regex_escape_/x:w` This function is documented on page ??.)

`\_regex_escape_x_test:N` Find out whether the first character is a left brace (allowing any number of hexadecimal digits), or not (allowing up to two hexadecimal digits). We need to check for the end-of-string marker. Eventually, call either `\_regex_escape_x_loop:N` or `\_regex_escape_x_ii:N`.

```

316 \cs_new:Npn \_regex_escape_x_test:N #1
317 {
318   \str_if_eq_x:nnTF {#1} { break } { ; }
319   {
320     \if_charcode:w \c_space_token #1
321     \exp_after:wN \_regex_escape_x_test:N
322     \else:
323     \exp_after:wN \_regex_escape_x_test_ii:N
324     \exp_after:wN #1
325     \fi:
326   }
327 }
328 \cs_new:Npn \_regex_escape_x_test_ii:N #1
329 {
330   \if_charcode:w \c_lbrace_str #1
331   \exp_after:wN \_regex_escape_x_loop:N
332   \else:
333   \_str_hexadecimal_use:NTF #1
334   { \exp_after:wN \_regex_escape_x_ii:N }
335   { ; \exp_after:wN \_regex_escape_loop:N \exp_after:wN #1 }
336   \fi:
337 }

```

(End definition for `\_regex_escape_x_test:N`)

`\_regex_escape_x_ii:N` This looks for the second digit in the unbraced case.

```

338 \cs_new:Npn \_regex_escape_x_ii:N #1
339 {
340   \str_if_eq_x:nnTF {#1} { break } { ; }
341   {
342     \_str_hexadecimal_use:NTF #1
343     { ; \_regex_escape_loop:N }
344     { ; \_regex_escape_loop:N #1 }
345   }
346 }

```

(End definition for `\_regex_escape_x_ii:N`)

`\_regex_escape_x_loop:N` Grab hexadecimal digits, skip spaces, and at the end, check that there is a right brace, otherwise raise an error outside the assignment.

```

347 \cs_new:Npn \_regex_escape_x_loop:N #1
348 {

```



```

349  \__str_hexadecimal_use:NTF #1
350  { \__regex_escape_x_loop:N }
351  {
352    \token_if_eq_charcode:NNTF \c_space_token #1
353    { \__regex_escape_x_loop:N }
354    {
355      ;
356      \exp_after:wN \token_if_eq_charcode:NNTF \c_rbrace_str #1
357      { \__regex_escape_loop:N }
358      {
359        \if_false: { \fi: }
360        \__tl_build_one:o \l__regex_internal_b_tl
361        \__msg_kernel_error:nm { regex } { x-missing-rbrace } {#1}
362        \tl_set:Nx \l__regex_internal_b_tl
363          { \if_false: } \fi: \__regex_escape_loop:N #1
364      }
365    }
366  }
367 }

```

(End definition for \\_\_regex\_escape\_x\_loop:N)

```

\__regex_char_if_alphanumeric:NTF
\__regex_char_if_special:NTF

```

These two tests are used in the first pass when parsing a regular expression. That pass is responsible for finding escaped and non-escaped characters, and recognizing which ones have special meanings and which should be interpreted as “raw” characters. Namely,

- alphanumeric characters are “raw” if they are not escaped, and may have a special meaning when escaped;
- non-alphanumeric printable ascii characters are “raw” if they are escaped, and may have a special meaning when not escaped;
- characters other than printable ascii are always “raw”.

The code is ugly, and highly based on magic numbers and the ascii codes of characters. This is mostly unavoidable for performance reasons: testing for instance with `\__str_if_contains_char:nN` would be much slower. Maybe the tests can be optimized a little bit more. Here, “alphanumeric” means 0–9, A–Z, a–z; “special” character means non-alphanumeric but printable ascii, from space (hex 20) to del (hex 7E).

```

368 \prg_new_conditional:Npnn \__regex_char_if_special:N #1 { TF }
369 {
370   \if_int_compare:w '#1 < \c_ninety_one
371   \if_int_compare:w '#1 < \c_fifty_eight
372   \if_int_compare:w '#1 < \c_forty_eight
373   \if_int_compare:w '#1 < \c_thirty_two
374   \prg_return_false: \else: \prg_return_true: \fi:
375   \else: \prg_return_false: \fi:
376   \else:
377     \if_int_compare:w '#1 < \c_sixty_five
378     \prg_return_true: \else: \prg_return_false: \fi:
379   \fi:

```

```

380 \else:
381   \if_int_compare:w '#1 < \c_one_hundred_twenty_three
382   \if_int_compare:w '#1 < \c_ninety_seven
383   \prg_return_true: \else: \prg_return_false: \fi:
384 \else:
385   \if_int_compare:w '#1 < \c_one_hundred_twenty_seven
386   \prg_return_true: \else: \prg_return_false: \fi:
387 \fi:
388 \fi:
389 }
390 \prg_new_conditional:Npnn \__regex_char_if_alphanumeric:N #1 { TF }
391 {
392   \if_int_compare:w '#1 < \c_ninety_one
393   \if_int_compare:w '#1 < \c_fifty_eight
394   \if_int_compare:w '#1 < \c_forty_eight
395   \prg_return_false: \else: \prg_return_true: \fi:
396 \else:
397   \if_int_compare:w '#1 < \c_sixty_five
398   \prg_return_false: \else: \prg_return_true: \fi:
399 \fi:
400 \else:
401   \if_int_compare:w '#1 < \c_one_hundred_twenty_three
402   \if_int_compare:w '#1 < \c_ninety_seven
403   \prg_return_false: \else: \prg_return_true: \fi:
404 \else:
405   \prg_return_false:
406 \fi:
407 \fi:
408 }

```

(End definition for \\_\_regex\_char\_if\_alphanumeric:NTF and \\_\_regex\_char\_if\_special:NTF)

## 2.3 Compiling

A regular expression starts its life as a string of characters. In this section, we convert it to internal instructions, resulting in a “compiled” regular expression. This compiled expression is then turned into states of an automaton in the building phase. Compiled regular expressions consist of the following:

- \\_\_regex\_class:NnnnN <boolean> {<tests>} {<min>} {<more>} <lazyness>
- \\_\_regex\_group:nnnN {<branches>} {<min>} {<more>} <lazyness>, also \\_\_regex\_group\_no\_capture:nnnN and \\_\_regex\_group\_resetting:nnnN with the same syntax.
- \\_\_regex\_branch:n {<contents>}
- \\_\_regex\_command\_K:
- \\_\_regex\_assertion:Nn <boolean> {<assertion test>}, where the <assertion test> is \\_\_regex\_b\_test: or { \\_\_regex\_anchor:N <integer> }

Tests can be the following:

- `\__regex_item_caseful_equal:n`  $\{\langle char\ code\rangle\}$
- `\__regex_item_caseless_equal:n`  $\{\langle char\ code\rangle\}$
- `\__regex_item_caseful_range:nn`  $\{\langle min\rangle\}\{\langle max\rangle\}$
- `\__regex_item_caseless_range:nn`  $\{\langle min\rangle\}\{\langle max\rangle\}$
- `\__regex_item_catcode:nT`  $\{\langle catcode\ bitmap\rangle\}\{\langle tests\rangle\}$
- `\__regex_item_catcode_reverse:nT`  $\{\langle catcode\ bitmap\rangle\}\{\langle tests\rangle\}$
- `\__regex_item_reverse:n`  $\{\langle tests\rangle\}$
- `\__regex_item_exact:nn`  $\{\langle catcode\rangle\}\{\langle char\ code\rangle\}$
- `\__regex_item_exact_cs:c`  $\{\langle csname\rangle\}$
- `\__regex_item_cs:n`  $\{\langle compiled\ regex\rangle\}$

### 2.3.1 Variables used when compiling

<code>\l__regex_group_level_int</code>	<p>We make sure to open the same number of groups as we close.</p> <pre> 409 \int_new:N \l__regex_group_level_int (End definition for \l__regex_group_level_int This variable is documented on page ??.) </pre>
<code>\l__regex_mode_int</code>	<p>While compiling, ten modes are recognized, labelled <math>-63, -23, -6, -2, 0, 2, 3, 6, 23, 63</math>. See section <a href="#">2.3.3</a>.</p> <pre> 410 \int_new:N \l__regex_mode_int (End definition for \l__regex_mode_int This variable is documented on page ??.) </pre>
<code>\l__regex_catcodes_int</code> <code>\l__regex_default_catcodes_int</code> <code>\l__regex_catcodes_bool</code>	<p>We wish to allow constructions such as <code>\c[<sup>^</sup>BE](. . \cL[a-z] . .)</code>, where the outer catcode test applies to the whole group, but is superseded by the inner catcode test. For this to work, we need to keep track of lists of allowed category codes: <code>\l__regex_catcodes_int</code> and <code>\l__regex_default_catcodes_int</code> are bitmaps, sums of <math>4^c</math>, for all allowed catcodes <math>c</math>. The latter is local to each capturing group, and we reset <code>\l__regex_catcodes_int</code> to that value after each character or class, changing it only when encountering a <code>\c</code> escape. The boolean records whether the list of categories of a catcode test has to be inverted: compare <code>\c[<sup>^</sup>BE]</code> and <code>\c[BE]</code>.</p> <pre> 411 \int_new:N \l__regex_catcodes_int 412 \int_new:N \l__regex_default_catcodes_int 413 \bool_new:N \l__regex_catcodes_bool (End definition for \l__regex_catcodes_int and \l__regex_default_catcodes_int These functions are documented on page ??.) </pre>

`\c__regex_catcode_C_int` Constants:  $4^c$  for each category, and the sum of all powers of 4.

```

414 \int_const:Nn \c__regex_catcode_C_int { "1 }
415 \int_const:Nn \c__regex_catcode_B_int { "4 }
416 \int_const:Nn \c__regex_catcode_E_int { "10 }
417 \int_const:Nn \c__regex_catcode_M_int { "40 }
418 \int_const:Nn \c__regex_catcode_T_int { "100 }
419 \int_const:Nn \c__regex_catcode_P_int { "1000 }
420 \int_const:Nn \c__regex_catcode_U_int { "4000 }
421 \int_const:Nn \c__regex_catcode_D_int { "10000 }
422 \int_const:Nn \c__regex_catcode_S_int { "100000 }
423 \int_const:Nn \c__regex_catcode_L_int { "400000 }
424 \int_const:Nn \c__regex_catcode_O_int { "1000000 }
425 \int_const:Nn \c__regex_catcode_A_int { "4000000 }
426 \int_const:Nn \c__regex_all_catcodes_int { "5515155 }

```

*(End definition for \c\_\_regex\_catcode\_C\_int and others. These functions are documented on page ??.)*

`\l__regex_internal_regex` The compilation step stores its result in this variable.

```

427 \cs_new_eq:NN \l__regex_internal_regex \c__regex_no_match_regex

```

*(End definition for \l\_\_regex\_internal\_regex This variable is documented on page ??.)*

`\l__regex_show_prefix_seq` This sequence holds the prefix that makes up the line displayed to the user. The various items must be removed from the right, which is tricky with a token list, hence we use a sequence.

```

428 \seq_new:N \l__regex_show_prefix_seq

```

*(End definition for \l\_\_regex\_show\_prefix\_seq This variable is documented on page ??.)*

`\l__regex_show_lines_int` A hack. To know whether a given class has a single item in it or not, we count the number of lines when showing the class.

```

429 \int_new:N \l__regex_show_lines_int

```

*(End definition for \l\_\_regex\_show\_lines\_int This variable is documented on page ??.)*

### 2.3.2 Generic helpers used when compiling

`\__regex_get_digits:NTFw` If followed by some raw digits, collect them one by one in the integer variable #1, and take the true branch. Otherwise, take the false branch.

```

430 \cs_new_protected:Npn \__regex_get_digits:NTFw #1#2#3#4#5
431 {
432   \__regex_if_raw_digit:NNTF #4 #5
433   { #1 = #5 \__regex_get_digits_loop:nw {#2} }
434   { #3 #4 #5 }
435 }
436 \cs_new:Npn \__regex_get_digits_loop:nw #1#2#3
437 {
438   \__regex_if_raw_digit:NNTF #2 #3
439   { #3 \__regex_get_digits_loop:nw {#1} }
440   { \scan_stop: #1 #2 #3 }
441 }

```

*(End definition for \\_\_regex\_get\_digits:NTFw This function is documented on page ??.)*

`\_regex_if_raw_digit:NNTF` Test used when grabbing digits for the `{m,n}` quantifier. It only accepts non-escaped digits.

```
442 \prg_new_conditional:Npnn \_regex_if_raw_digit:NN #1#2 { TF }
443 {
444   \if_meaning:w \_regex_compile_raw:N #1
445     \if_int_compare:w \c_one < 1 #2 \exp_stop_f:
446       \prg_return_true:
447     \else:
448       \prg_return_false:
449     \fi:
450   \else:
451     \prg_return_false:
452   \fi:
453 }
(End definition for \_regex_if_raw_digit:NNTF)
```

### 2.3.3 Mode

When compiling the NFA corresponding to a given regex string, we can be in ten distinct modes, which we label by some magic numbers:

- 6 `[\c{...}]` control sequence in a class,
- 2 `\c{...}` control sequence,
- 0 ... outer,
- 2 `\c...` catcode test,
- 6 `[\c...] catcode test in a class,`
- 63 `[\c{[...]}]` class inside mode -6,
- 23 `\c{[...]}` class inside mode -2,
- 3 `[...] class inside mode -3,`
- 23 `\c[...] class inside mode 2,`
- 63 `[\c[...] class inside mode 6.`

This list is exhaustive, because `\c` escape sequences cannot be nested, and character classes cannot be nested directly. The choice of numbers is such as to optimize the most useful tests, and make transitions from one mode to another as simple as possible.

- Even modes mean that we are not directly in a character class. In this case, a left bracket appends 3 to the mode. In a character class, a right bracket changes the mode as  $m \rightarrow (m - 15)/13$ , truncated.
- Grouping, assertion, and anchors are allowed in non-positive even modes (0, -2, -6), and do not change the mode. Otherwise, they trigger an error.

- A left bracket is special in even modes, appending 3 to the mode; in those modes, quantifiers and the dot are recognized, and the right bracket is normal. In odd modes (within classes), the left bracket is normal, but the right bracket ends the class, changing the mode from  $m$  to  $(m - 15)/13$ , truncated; also, ranges are recognized.
- In non-negative modes, left and right braces are normal. In negative modes, however, left braces trigger a warning; right braces end the control sequence, going from  $-2$  to  $0$  or  $-6$  to  $3$ , with error recovery for odd modes.
- Properties (such as the `\d` character class) can appear in any mode.

`\__regex_if_in_class:TF` Test whether we are directly in a character class (at the innermost level of nesting). There, many escape sequences are not recognized, and special characters are normal. Also, for every raw character, we must look ahead for a possible raw dash.

```

454 \cs_new_nopar:Npn \__regex_if_in_class:TF
455 {
456   \if_int_odd:w \l__regex_mode_int
457     \exp_after:wN \use_i:nn
458   \else:
459     \exp_after:wN \use_ii:nn
460   \fi:
461 }
```

(End definition for `\__regex_if_in_class:TF`)

`\__regex_if_in_cs:TF` Right braces are special only directly inside control sequences (at the inner-most level of nesting, not counting groups).

```

462 \cs_new_nopar:Npn \__regex_if_in_cs:TF
463 {
464   \if_int_odd:w \l__regex_mode_int
465     \exp_after:wN \use_ii:nn
466   \else:
467     \if_int_compare:w \l__regex_mode_int < \c_zero
468       \exp_after:wN \exp_after:wN \exp_after:wN \use_i:nn
469     \else:
470       \exp_after:wN \exp_after:wN \exp_after:wN \use_ii:nn
471     \fi:
472   \fi:
473 }
```

(End definition for `\__regex_if_in_cs:TF`)

`\__regex_if_in_class_or_catcode:TF` Assertions are only allowed in modes  $0$ ,  $-2$ , and  $-6$ , *i.e.*, even, non-positive modes.

```

474 \cs_new_nopar:Npn \__regex_if_in_class_or_catcode:TF
475 {
476   \if_int_odd:w \l__regex_mode_int
477     \exp_after:wN \use_i:nn
478   \else:
479     \if_int_compare:w \l__regex_mode_int > \c_zero
480       \exp_after:wN \exp_after:wN \exp_after:wN \use_i:nn
```

```

481     \else:
482         \exp_after:wN \exp_after:wN \exp_after:wN \use_ii:nn
483     \fi:
484 \fi:
485 }

```

(End definition for `\__regex_if_in_class_or_catcode:TF`)

`\__regex_if_within_catcode:TF` This test takes the true branch if we are in a catcode test, either immediately following it (modes 2 and 6) or in a class on which it applies (modes 23 and 63). This is used to tweak how left brackets behave in modes 2 and 6.

```

486 \cs_new_nopar:Npn \__regex_if_within_catcode:TF
487 {
488     \if_int_compare:w \l__regex_mode_int > \c_zero
489         \exp_after:wN \use_i:nn
490     \else:
491         \exp_after:wN \use_ii:nn
492     \fi:
493 }

```

(End definition for `\__regex_if_within_catcode:TF`)

`\__regex_chk_c_allowed:T` The `\c` escape sequence is only allowed in modes 0 and 3, *i.e.*, not within any other `\c` escape sequence.

```

494 \cs_new_protected:Npn \__regex_chk_c_allowed:T
495 {
496     \if_int_compare:w \l__regex_mode_int = \c_zero
497         \exp_after:wN \use:n
498     \else:
499         \if_int_compare:w \l__regex_mode_int = \c_three
500             \exp_after:wN \exp_after:wN \exp_after:wN \use:n
501         \else:
502             \__msg_kernel_error:nn { regex } { c-bad-mode }
503             \exp_after:wN \exp_after:wN \exp_after:wN \use_none:n
504         \fi:
505     \fi:
506 }

```

(End definition for `\__regex_chk_c_allowed:T`)

`\__regex_mode_quit_c:` This function changes the mode as it is needed just after a catcode test.

```

507 \cs_new_protected:Npn \__regex_mode_quit_c:
508 {
509     \if_int_compare:w \l__regex_mode_int = \c_two
510         \l__regex_mode_int = \c_zero
511     \else:
512         \if_int_compare:w \l__regex_mode_int = \c_six
513             \l__regex_mode_int = \c_three
514         \fi:
515     \fi:
516 }

```

(End definition for `\__regex_mode_quit_c:`)

### 2.3.4 Framework

`\__regex_compile:w` Used when compiling a user regex or a regex for the `\c{...}` escape sequence within another regex. Start building a token list within a group (with x-expansion at the outset), and set a few variables (group level, catcodes), then start the first branch. At the end, make sure there are no dangling classes nor groups, close the last branch: we are done building `\l__regex_internal_regex`.

```

517 \cs_new_protected_nopar:Npn \__regex_compile:w
518 {
519   \__tl_build_x:Nw \l__regex_internal_regex
520   \int_zero:N \l__regex_group_level_int
521   \int_set_eq:NN \l__regex_default_catcodes_int \c__regex_all_catcodes_int
522   \int_set_eq:NN \l__regex_catcodes_int \l__regex_default_catcodes_int
523   \cs_set_nopar:Npn \__regex_item_equal:n { \__regex_item_caseful_equal:n }
524   \cs_set_nopar:Npn \__regex_item_range:nn { \__regex_item_caseful_range:nn }
525   \__tl_build_one:n { \__regex_branch:n { \if_false: } \fi: }
526 }
527 \cs_new_protected_nopar:Npn \__regex_compile_end:
528 {
529   \__regex_if_in_class:TF
530   {
531     \__msg_kernel_error:nn { regex } { missing-rbrack }
532     \use:c { __regex_compile_]: }
533     \prg_do_nothing: \prg_do_nothing:
534   }
535   { }
536   \if_int_compare:w \l__regex_group_level_int > \c_zero
537   \__msg_kernel_error:nnx { regex } { missing-rparen }
538   { \int_use:N \l__regex_group_level_int }
539   \prg_replicate:nn
540   { \l__regex_group_level_int }
541   {
542     \__tl_build_one:n
543     {
544       \if_false: { \fi: }
545       \if_false: { \fi: } { 1 } { 0 } \c_true_bool
546     }
547     \__tl_build_end:
548     \__tl_build_one:o \l__regex_internal_regex
549   }
550   \fi:
551   \__tl_build_one:n { \if_false: { \fi: } }
552   \__tl_build_end:
553 }

```

(End definition for `\__regex_compile:w` and `\__regex_compile_end:`)

`\__regex_compile:n` The compilation is done between `\__regex_compile:w` and `\__regex_compile_end:`, starting in mode 0. Then `\__regex_escape_use:nnnn` distinguishes special characters, escaped alphanumerics, and raw characters, interpreting `\a`, `\x` and other sequences. The



4 trailing `\prg_do_nothing:` are needed because some functions defined later look up to 4 tokens ahead. Before ending, make sure that any `\c{...}` is properly closed.

```

554 \cs_new_protected:Npn \__regex_compile:n #1
555 {
556   \__regex_compile:w
557   \int_set:Nn \tex_escapechar:D { 92 }
558   \int_set_eq:NN \l__regex_mode_int \c_zero
559   \__regex_escape_use:nnnn
560   {
561     \__regex_char_if_special:NTF ##1
562     \__regex_compile_special:N \__regex_compile_raw:N ##1
563   }
564   {
565     \__regex_char_if_alphanumeric:NTF ##1
566     \__regex_compile_escaped:N \__regex_compile_raw:N ##1
567   }
568   { \__regex_compile_raw:N ##1 }
569   { #1 }
570   \prg_do_nothing: \prg_do_nothing:
571   \prg_do_nothing: \prg_do_nothing:
572   \int_compare:nNnT \l__regex_mode_int < \c_zero
573   {
574     \__msg_kernel_error:nn { regex } { c-missing-rbrace }
575     \__regex_compile_end:
576     \__regex_compile_one:x
577     { \__regex_item_cs:n { \exp_not:o \l__regex_internal_regex } }
578     \prg_do_nothing: \prg_do_nothing:
579     \prg_do_nothing: \prg_do_nothing:
580   }
581   \__regex_compile_end:
582 }

```

*(End definition for \\_\_regex\_compile:n)*

`\__regex_compile_escaped:N` If the special character or escaped alphanumeric has a particular meaning in regexes, the corresponding function is used. Otherwise, it is interpreted as a raw character. We distinguish special characters from escaped alphanumeric characters because they behave differently when appearing as an end-point of a range.

```

583 \cs_new_protected:Npn \__regex_compile_special:N #1
584 {
585   \cs_if_exist_use:cF { __regex_compile_#1: }
586   { \__regex_compile_raw:N #1 }
587 }
588 \cs_new_protected:Npn \__regex_compile_escaped:N #1
589 {
590   \cs_if_exist_use:cF { __regex_compile_/#1: }
591   { \__regex_compile_raw:N #1 }
592 }

```

*(End definition for \\_\_regex\_compile\_escaped:N and \\_\_regex\_compile\_special:N)*

`\__regex_compile_one:x` This is used after finding one “test”, such as `\d`, or a raw character. If that followed a catcode test (e.g., `\cL`), then restore the mode. If we are not in a class, then the test is “standalone”, and we need to add `\__regex_class:NnnnN` and search for quantifiers. In any case, insert the test, possibly together with a catcode test if appropriate.

```

593 \cs_new_protected:Npn \__regex_compile_one:x #1
594 {
595   \__regex_mode_quit_c:
596   \__regex_if_in_class:TF { }
597   {
598     \__tl_build_one:n
599     { \__regex_class:NnnnN \c_true_bool { \if_false: } \fi: }
600   }
601   \__tl_build_one:x
602   {
603     \if_int_compare:w \l__regex_catcodes_int < \c_regex_all_catcodes_int
604     \__regex_item_catcode:nT { \int_use:N \l__regex_catcodes_int }
605     { \exp_not:N \exp_not:n {#1} }
606     \else:
607     \exp_not:N \exp_not:n {#1}
608     \fi:
609   }
610   \int_set_eq:NN \l__regex_catcodes_int \l__regex_default_catcodes_int
611   \__regex_if_in_class:TF { } { \__regex_compile_quantifier:w }
612 }

```

(End definition for `\__regex_compile_one:x`)

`\__regex_compile_abort_tokens:n` This function places the collected tokens back in the input stream, each as a raw character.  
`\__regex_compile_abort_tokens:x` Spaces are not preserved.

```

613 \cs_new_protected:Npn \__regex_compile_abort_tokens:n #1
614 {
615   \use:x
616   {
617     \exp_args:No \tl_map_function:nN { \tl_to_str:n {#1} }
618     \__regex_compile_raw:N
619   }
620 }
621 \cs_generate_variant:Nn \__regex_compile_abort_tokens:n { x }

```

(End definition for `\__regex_compile_abort_tokens:n` and `\__regex_compile_abort_tokens:x`)

### 2.3.5 Quantifiers

`\__regex_compile_quantifier:w` This looks ahead and finds any quantifier (special character equal to either of `?+*{}`).

```

622 \cs_new_protected:Npn \__regex_compile_quantifier:w #1#2
623 {
624   \token_if_eq_meaning:NNTF #1 \__regex_compile_special:N
625   {
626     \cs_if_exist_use:cF { __regex_compile_quantifier_#2:w }
627     { \__regex_compile_quantifier_none: #1 #2 }

```

```

628     }
629     { \_regex_compile_quantifier_none: #1 #2 }
630 }
(End definition for \_regex_compile_quantifier:w)

```

\\_regex\_compile\_quantifier\_none: Those functions are called whenever there is no quantifier, or a braced construction is invalid (equivalent to no quantifier, and whatever characters were grabbed are left raw).  
\\_regex\_compile\_quantifier\_abort:xNN

```

631 \cs_new_protected:Npn \_regex_compile_quantifier_none:
632 { \_tl_build_one:n { \if_false: { \fi: } { 1 } { 0 } \c_false_bool } }
633 \cs_new_protected:Npn \_regex_compile_quantifier_abort:xNN #1#2#3
634 {
635   \_regex_compile_quantifier_none:
636   \_msg_kernel_warning:nxxx { regex } { invalid-quantifier } {#1} {#3}
637   \_regex_compile_abort_tokens:x {#1}
638   #2 #3
639 }
(End definition for \_regex_compile_quantifier_none: This function is documented on page ??.)

```

\\_regex\_compile\_quantifier\_lazy:nnNN Once the “main” quantifier (?, \*, + or a braced construction) is found, we check whether it is lazy (followed by a question mark). We then add to the compiled regex a closing brace (ending \\_regex\_class:NnnnN and friends), the start-point of the range, its end-point, and a boolean, true for lazy and false for greedy operators.

```

640 \cs_new_protected:Npn \_regex_compile_quantifier_lazy:nnNN #1#2#3#4
641 {
642   \str_if_eq:nnTF { #3 #4 } { \_regex_compile_special:N ? }
643   { \_tl_build_one:n { \if_false: { \fi: } { #1 } { #2 } \c_true_bool } }
644   {
645     \_tl_build_one:n { \if_false: { \fi: } { #1 } { #2 } \c_false_bool }
646     #3 #4
647   }
648 }
(End definition for \_regex_compile_quantifier_lazy:nnNN)

```

\\_regex\_compile\_quantifier?:w For each “basic” quantifier, ?, \*, +, feed the correct arguments to \\_regex\_compile\_quantifier\_lazy:nnNN, -1 means that there is no upper bound on the number of repetitions.  
\\_regex\_compile\_quantifier\*:w  
\\_regex\_compile\_quantifier+:w

```

649 \cs_new_protected_nopar:cpn { \_regex_compile_quantifier?:w }
650 { \_regex_compile_quantifier_lazy:nnNN { 0 } { 1 } }
651 \cs_new_protected_nopar:cpn { \_regex_compile_quantifier*:w }
652 { \_regex_compile_quantifier_lazy:nnNN { 0 } { -1 } }
653 \cs_new_protected_nopar:cpn { \_regex_compile_quantifier+:w }
654 { \_regex_compile_quantifier_lazy:nnNN { 1 } { -1 } }
(End definition for \_regex_compile_quantifier?:w, \_regex_compile_quantifier*:w, and \_regex_compile_quantifier+:w)

```

\\_regex\_compile\_quantifier{w Three possible syntaxes: {<int>}, {<int>}, or {<int>,<int>}. Any other syntax causes us to abort and put whatever we collected back in the input stream, as raw characters, including the opening brace. Grab a number into \l\_\_regex\_internal\_a\_int. If the number is followed by a right brace, the range is [a, a]. If followed by a comma, grab one

more number, and call the `_ii` or `_iii` auxiliary. Those auxiliaries check for a closing brace, leading to the range  $[a, \infty]$  or  $[a, b]$ , encoded as  $\{a\}\{-1\}$  and  $\{a\}\{b-a\}$ .

```

655 \cs_new_protected:cpn { __regex_compile_quantifier_ \c_lbrace_str :w }
656 {
657   \__regex_get_digits:NTFw \l__regex_internal_a_int
658   { \__regex_compile_quantifier_braced_i:w }
659   { \__regex_compile_quantifier_abort:xNN { \c_lbrace_str } }
660 }
661 \cs_new_protected:Npn \__regex_compile_quantifier_braced_i:w #1#2
662 {
663   \str_case_x:nnn { #1 #2 }
664   {
665     { \__regex_compile_special:N \c_rbrace_str }
666     {
667       \exp_args:No \__regex_compile_quantifier_lazyness:nnNN
668       { \int_use:N \l__regex_internal_a_int } { 0 }
669     }
670     { \__regex_compile_special:N , }
671     {
672       \__regex_get_digits:NTFw \l__regex_internal_b_int
673       { \__regex_compile_quantifier_braced_iii:w }
674       { \__regex_compile_quantifier_braced_ii:w }
675     }
676   }
677   {
678     \__regex_compile_quantifier_abort:xNN
679     { \c_lbrace_str \int_use:N \l__regex_internal_a_int }
680     #1 #2
681   }
682 }
683 \cs_new_protected:Npn \__regex_compile_quantifier_braced_ii:w #1#2
684 {
685   \str_if_eq_x:nnTF
686   { #1 #2 } { \__regex_compile_special:N \c_rbrace_str }
687   {
688     \exp_args:No \__regex_compile_quantifier_lazyness:nnNN
689     { \int_use:N \l__regex_internal_a_int } { -1 }
690   }
691   {
692     \__regex_compile_quantifier_abort:xNN
693     { \c_lbrace_str \int_use:N \l__regex_internal_a_int , }
694     #1 #2
695   }
696 }
697 \cs_new_protected:Npn \__regex_compile_quantifier_braced_iii:w #1#2
698 {
699   \str_if_eq_x:nnTF
700   { #1 #2 } { \__regex_compile_special:N \c_rbrace_str }
701   {

```

```

702     \if_int_compare:w \l__regex_internal_a_int > \l__regex_internal_b_int
703     \__msg_kernel_error:nxx { regex } { backwards-quantifier }
704     { \int_use:N \l__regex_internal_a_int }
705     { \int_use:N \l__regex_internal_b_int }
706     \int_zero:N \l__regex_internal_b_int
707 \else:
708     \int_sub:Nn \l__regex_internal_b_int \l__regex_internal_a_int
709 \fi:
710 \exp_args:Noo \__regex_compile_quantifier_lazyness:nnNN
711 { \int_use:N \l__regex_internal_a_int }
712 { \int_use:N \l__regex_internal_b_int }
713 }
714 {
715     \__regex_compile_quantifier_abort:xNN
716     {
717         \c_lbrace_str
718         \int_use:N \l__regex_internal_a_int ,
719         \int_use:N \l__regex_internal_b_int
720     }
721     #1 #2
722 }
723 }

```

(End definition for \\_\_regex\_compile\_quantifier\_{:w This function is documented on page ??.)

### 2.3.6 Raw characters

**\\_\_regex\_compile\_raw\_error:N** Within character classes, and following catcode tests, some escaped alphanumeric sequences such as \b do not have any meaning. They are replaced by a raw character, after spitting out an error.

```

724 \cs_new_protected:Npn \__regex_compile_raw_error:N #1
725 {
726     \__msg_kernel_error:nxx { regex } { bad-escape } {#1}
727     \__regex_compile_raw:N #1
728 }

```

(End definition for \\_\_regex\_compile\_raw\_error:N)

**\\_\_regex\_compile\_raw:N** If we are in a character class and the next character is an unescaped dash, this denotes a range. Otherwise, the current character #1 matches itself.

```

729 \cs_new_protected:Npn \__regex_compile_raw:N #1#2#3
730 {
731     \__regex_if_in_class:TF
732     {
733         \str_if_eq:nnTF {#2#3} { \__regex_compile_special:N - }
734         { \__regex_compile_range:Nw #1 }
735         {
736             \__regex_compile_one:x
737             { \__regex_item_equal:n { \__int_value:w '#1 ~ } }
738             #2 #3
739         }
740     }

```

```

740     }
741     {
742         \_regex_compile_one:x
743         { \_regex_item_equal:n { \_int_value:w '#1 ~ } }
744         #2 #3
745     }
746 }

```

(End definition for \\_regex\_compile\_raw:N)

\\_regex\_compile\_range:Nw  
\\_regex\_if\_end\_range:NNTF

We have just read a raw character followed by a dash; this should be followed by an end-point for the range. Valid end-points are: any raw character; any special character, except a right bracket. In particular, escaped characters are forbidden.

```

747 \prg_new_protected_conditional:Npnn \_regex_if_end_range:NN #1#2 { TF }
748 {
749     \if_meaning:w \_regex_compile_raw:N #1
750     \prg_return_true:
751     \else:
752         \if_meaning:w \_regex_compile_special:N #1
753         \if_charcode:w ] #2
754         \prg_return_false:
755         \else:
756             \prg_return_true:
757         \fi:
758     \else:
759         \prg_return_false:
760     \fi:
761 \fi:
762 }
763 \cs_new_protected:Npn \_regex_compile_range:Nw #1#2#3
764 {
765     \_regex_if_end_range:NNTF #2 #3
766     {
767         \if_int_compare:w '#1 > '#3 \exp_stop_f:
768         \_msg_kernel_error:nxxx { regex } { range-backwards } {#1} {#3}
769     \else:
770         \_tl_build_one:x
771         {
772             \if_int_compare:w '#1 = '#3 \exp_stop_f:
773             \_regex_item_equal:n
774         \else:
775             \_regex_item_range:nn { \_int_value:w '#1 ~ }
776         \fi:
777         { \_int_value:w '#3 ~ }
778     }
779     \fi:
780 }
781 {
782     \_msg_kernel_warning:nxxx { regex } { range-missing-end }
783     {#1} { \c_backslash_str #3 }

```

```

784     \__tl_build_one:x
785     {
786         \__regex_item_equal:n { \__int_value:w '#1 ~ }
787         \__regex_item_equal:n { \__int_value:w '- ~ }
788     }
789     #2#3
790 }
791 }

```

(End definition for \\_\_regex\_compile\_range:Nw and \\_\_regex\_if\_end\_range:NNTF)

### 2.3.7 Character properties

\\_\_regex\_compile\_.: In a class, the dot has no special meaning. Outside, insert \\_\_regex\_prop\_., which matches any character or control sequence, and refuses -2 (end-marker).

```

792 \cs_new_protected_nopar:cpx { \__regex_compile_.: }
793 {
794     \exp_not:N \__regex_if_in_class:TF
795     { \__regex_compile_raw:N . }
796     { \__regex_compile_one:x \exp_not:c { \__regex_prop_.: } }
797 }
798 \cs_new_protected_nopar:cpn { \__regex_prop_.: }
799 {
800     \if_int_compare:w \l__regex_current_char_int > - \c_two
801     \exp_after:wN \__regex_break_true:w
802     \fi:
803 }

```

(End definition for \\_\_regex\_compile\_.: and \\_\_regex\_prop\_.:)

\\_\_regex\_compile\_/d: The constants \\_\_regex\_prop\_d:, etc. hold a list of tests which match the corresponding character class, and jump to the \\_\_regex\_break\_point:TF marker. As for a normal character, we check for quantifiers.

```

804 \cs_set_protected:Npn \__regex_tmp:w #1#2
805 {
806     \cs_new_protected_nopar:cpx { \__regex_compile_/#1: }
807     { \__regex_compile_one:x \exp_not:c { \__regex_prop_#1: } }
808     \cs_new_protected_nopar:cpx { \__regex_compile_/#2: }
809     {
810         \__regex_compile_one:x
811         { \__regex_item_reverse:n \exp_not:c { \__regex_prop_#1: } }
812     }
813 }
814 \__regex_tmp:w d D
815 \__regex_tmp:w h H
816 \__regex_tmp:w s S
817 \__regex_tmp:w v V
818 \__regex_tmp:w w W
819 \cs_new_protected_nopar:cpn { \__regex_compile_/N: }
820 { \__regex_compile_one:x \__regex_prop_N: }

```

(End definition for \\_\_regex\_compile\_/d: and others.)

### 2.3.8 Anchoring and simple assertions

`__regex_compile_anchor:Nf` In modes where assertions are allowed, anchor to the start of the query, the start of the match, or the end of the query, depending on the integer #1. In other modes, #2 treats the character as raw, with an error for escaped letters (\$ is valid in a class, but \A is definitely a mistake on the user's part).

```

__regex_compile_~:
__regex_compile_/A:
__regex_compile_/G:
__regex_compile_/$:
__regex_compile_/Z:
__regex_compile_/z:
821 \cs_new_protected:Npn __regex_compile_anchor:Nf #1#2
822 {
823   __regex_if_in_class_or_catcode:TF {#2}
824   {
825     __tl_build_one:n
826     { __regex_assertion:Nn \c_true_bool { __regex_anchor:N #1 } }
827   }
828 }
829 \cs_set_protected:Npn __regex_tmp:w #1#2
830 {
831   \cs_new_protected_nopar:cpn { __regex_compile_/#1: }
832   { __regex_compile_anchor:Nf #2 { __regex_compile_raw_error:N #1 } }
833 }
834 __regex_tmp:w A \l__regex_min_pos_int
835 __regex_tmp:w G \l__regex_start_pos_int
836 __regex_tmp:w Z \l__regex_max_pos_int
837 __regex_tmp:w z \l__regex_max_pos_int
838 \cs_set_protected:Npn __regex_tmp:w #1#2
839 {
840   \cs_new_protected_nopar:cpn { __regex_compile_#1: }
841   { __regex_compile_anchor:Nf #2 { __regex_compile_raw:N #1 } }
842 }
843 \exp_args:Nx __regex_tmp:w { \iow_char:N ^ } \l__regex_min_pos_int
844 \exp_args:Nx __regex_tmp:w { \iow_char:N $ } \l__regex_max_pos_int

```

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

`__regex_compile_/b:` Contrarily to `~` and `$`, which could be implemented without really knowing what precedes in the token list, this requires more information, namely, the knowledge of the last character code.

```

845 \cs_new_protected_nopar:cpn { __regex_compile_/b: }
846 {
847   __regex_if_in_class_or_catcode:TF
848   { __regex_compile_raw_error:N b }
849   {
850     __tl_build_one:n
851     { __regex_assertion:Nn \c_true_bool { __regex_b_test: } }
852   }
853 }
854 \cs_new_protected_nopar:cpn { __regex_compile_/B: }
855 {
856   __regex_if_in_class_or_catcode:TF
857   { __regex_compile_raw_error:N B }
858   {

```





`\_regex_compile_class_normal:w` In the “normal” case, we will insert `\_regex_class:NnnnN` (*boolean*) in the compiled code. The *boolean* is true for positive classes, and false for negative classes, characterized by a leading `^`. The auxiliary `\_regex_compile_class:TFNN` also checks for a leading `]` which has a special meaning.

```

891 \cs_new_protected_nopar:Npn \_regex_compile_class_normal:w
892 {
893   \_regex_compile_class:TFNN
894   { \_regex_class:NnnnN \c_true_bool }
895   { \_regex_class:NnnnN \c_false_bool }
896 }

```

(End definition for `\_regex_compile_class_normal:w`)

`\_regex_compile_class_catcode:w` This function is called for a left bracket in modes 2 or 6 (catcode test, and catcode test within a class). In mode 2 the whole construction needs to be put in a class (like single character). Then determine if the class is positive or negative, inserting `\_regex_item_catcode:nT` or the `reverse` variant as appropriate, each with the current catcodes bitmap `#1` as an argument, and reset the catcodes.

```

897 \cs_new_protected:Npn \_regex_compile_class_catcode:w #1;
898 {
899   \if_int_compare:w \l__regex_mode_int = \c_two
900     \__tl_build_one:n
901     { \_regex_class:NnnnN \c_true_bool { \if_false: } \fi: }
902   \fi:
903   \int_set_eq:NN \l__regex_catcodes_int \l__regex_default_catcodes_int
904   \_regex_compile_class:TFNN
905   { \_regex_item_catcode:nT {#1} }
906   { \_regex_item_catcode_reverse:nT {#1} }
907 }

```

(End definition for `\_regex_compile_class_catcode:w`)

`\_regex_compile_class:TFNN` If the first character is `^`, then the class is negative (use `#2`), otherwise it is positive (use `#1`). If the next character is a right bracket, then it should be changed to a raw one.

`\_regex_compile_class_ii:NN`

```

908 \cs_new_protected:Npn \_regex_compile_class:TFNN #1#2#3#4
909 {
910   \l__regex_mode_int = \__int_value:w \l__regex_mode_int 3 \exp_stop_f:
911   \str_if_eq:nnTF { #3 #4 } { \_regex_compile_special:N ^ }
912   {
913     \__tl_build_one:n { #2 { \if_false: } \fi: }
914     \_regex_compile_class_ii:NN
915   }
916   {
917     \__tl_build_one:n { #1 { \if_false: } \fi: }
918     \_regex_compile_class_ii:NN #3 #4
919   }
920 }
921 \cs_new_protected:Npn \_regex_compile_class_ii:NN #1#2
922 {
923   \token_if_eq_charcode:NNTF #2 ]

```

```

924     { \_regex_compile_raw:N #2 }
925     { #1 #2 }
926   }
(End definition for \_regex_compile_class:TFNN and \_regex_compile_class:ii:NN)

```

Here we check for a syntax such `[:alpha:]`. We also detect `[=` and `[.` which have a meaning in POSIX regular expressions, but are not implemented in `l3regex`. In case we see `[:`, grab raw characters until hopefully reaching `:`. If that's missing, or the POSIX class is unknown, abort. If all is right, add the test to the current class, with an extra `\_regex_item_reverse:n` for negative classes.

```

927 \cs_new_protected:Npn \_regex_compile_class_posix_test:w #1#2
928 {
929   \token_if_eq_meaning:NNT \_regex_compile_special:N #1
930   {
931     \str_case:nnn { #2 }
932     {
933       : { \_regex_compile_class_posix:NNNNw }
934       = { \_msg_kernel_warning:nxx { regex } { posix-unsupported } { = } }
935       . { \_msg_kernel_warning:nxx { regex } { posix-unsupported } { . } }
936     }
937   { }
938 }
939 \_regex_compile_raw:N [ #1 #2
940 }
941 \cs_new_protected:Npn \_regex_compile_class_posix:NNNNw #1#2#3#4#5#6
942 {
943   \str_if_eq:nnTF { #5 #6 } { \_regex_compile_special:N ^ }
944   {
945     \bool_set_false:N \l_regex_internal_bool
946     \tl_set:Nx \l_regex_internal_a_tl { \if_false: } \fi:
947     \_regex_compile_class_posix_loop:w
948   }
949   {
950     \bool_set_true:N \l_regex_internal_bool
951     \tl_set:Nx \l_regex_internal_a_tl { \if_false: } \fi:
952     \_regex_compile_class_posix_loop:w #5 #6
953   }
954 }
955 \cs_new:Npn \_regex_compile_class_posix_loop:w #1#2
956 {
957   \token_if_eq_meaning:NNTF \_regex_compile_raw:N #1
958   { #2 \_regex_compile_class_posix_loop:w }
959   { \if_false: { \fi: } \_regex_compile_class_posix_end:w #1 #2 }
960 }
961 \cs_new_protected:Npn \_regex_compile_class_posix_end:w #1#2#3#4
962 {
963   \str_if_eq:nnTF { #1 #2 #3 #4 }
964   { \_regex_compile_special:N : \_regex_compile_special:N ] }
965   {

```

```

966 \cs_if_exist:cTF { __regex_posix_ \l__regex_internal_a_tl : }
967 {
968   \__regex_compile_one:x
969   {
970     \bool_if:NF \l__regex_internal_bool \__regex_item_reverse:n
971     \exp_not:c { __regex_posix_ \l__regex_internal_a_tl : }
972   }
973 }
974 {
975   \__msg_kernel_warning:nxx { regex } { posix-unknown }
976   { \l__regex_internal_a_tl }
977   \__regex_compile_abort_tokens:x
978   {
979     [: \bool_if:NF \l__regex_internal_bool { ^ }
980     \l__regex_internal_a_tl :]
981   }
982 }
983 }
984 {
985   \__msg_kernel_error:nxx { regex } { posix-missing-close }
986   { [: \l__regex_internal_a_tl } { #2 #4 }
987   \__regex_compile_abort_tokens:x { [: \l__regex_internal_a_tl }
988   #1 #2 #3 #4
989 }
990 }

```

(End definition for `\__regex_compile_class_posix_test:w` and others.)

### 2.3.10 Groups and alternations

`\__regex_compile_group_begin:N` The contents of a regex group are turned into compiled code in `\l__regex_internal_regex`, which ends up with items of the form `\__regex_branch:n {⟨concatenation⟩}`. `\__regex_compile_group_end:` This construction is done using `l3tl-build` within a `TeX` group, which automatically makes sure that options (case-sensitivity and default catcode) are reset at the end of the group. The argument `#1` is `\__regex_group:nnnN` or a variant thereof. A small subtlety to support `\cL(abc)` as a shorthand for `(\cLa\cLb\cLc)`: exit any pending catcode test, save the category code at the start of the group as the default catcode for that group, and make sure that the catcode is restored to the default outside the group.

```

991 \cs_new_protected:Npn \__regex_compile_group_begin:N #1
992 {
993   \__tl_build_one:n { #1 { \if_false: } \fi: }
994   \__regex_mode_quit_c:
995   \__tl_build:Nw \l__regex_internal_regex
996   \int_set_eq:NN \l__regex_default_catcodes_int \l__regex_catcodes_int
997   \int_incr:N \l__regex_group_level_int
998   \__tl_build_one:n { \__regex_branch:n { \if_false: } \fi: }
999 }
1000 \cs_new_protected:Npn \__regex_compile_group_end:
1001 {
1002   \if_int_compare:w \l__regex_group_level_int > \c_zero

```

```

1003     \_tl\_build\_one:n { \if\_false: { \fi: } }
1004     \_tl\_build\_end:
1005     \int\_set\_eq:NN \l\_regex\_catcodes\_int \l\_regex\_default\_catcodes\_int
1006     \_tl\_build\_one:o \l\_regex\_internal\_regex
1007     \exp\_after:wN \_regex\_compile\_quantifier:w
1008     \else:
1009     \_msg\_kernel\_warning:nn { regex } { extra-rparen }
1010     \exp\_after:wN \_regex\_compile\_raw:N \exp\_after:wN )
1011     \fi:
1012   }
(End definition for \_regex\_compile\_group\_begin:N and \_regex\_compile\_group\_end:)

```

`\_regex\_compile\_(:` In a class, parentheses are not special. Outside, check for a ?, denoting special groups, and run the code for the corresponding special group.

```

1013 \cs\_new\_protected\_nopar:cpn { \_regex\_compile\_(: }
1014 {
1015   \_regex\_if\_in\_class:TF { \_regex\_compile\_raw:N ( }
1016   { \_regex\_compile\_lparen:w }
1017 }
1018 \cs\_new\_protected:Npn \_regex\_compile\_lparen:w #1#2#3#4
1019 {
1020   \str\_if\_eq:nnTF { #1 #2 } { \_regex\_compile\_special:N ? }
1021   {
1022     \cs\_if\_exist\_use:cF
1023     { \_regex\_compile\_special\_group\_token\_to\_str:N #4 :w }
1024     {
1025       \_msg\_kernel\_warning:nnx { regex } { special-group-unknown }
1026       { (? \token\_to\_str:N #4 }
1027       \_regex\_compile\_group\_begin:N \_regex\_group:nnnN
1028       \_regex\_compile\_raw:N ? #3 #4
1029     }
1030   }
1031   {
1032     \_regex\_compile\_group\_begin:N \_regex\_group:nnnN
1033     #1 #2 #3 #4
1034   }
1035 }
(End definition for \_regex\_compile\_(:)

```

`\_regex\_compile\_|:` In a class, the pipe is not special. Otherwise, end the current branch and open another one.

```

1036 \cs\_new\_protected\_nopar:cpn { \_regex\_compile\_|: }
1037 {
1038   \_regex\_if\_in\_class:TF { \_regex\_compile\_raw:N | }
1039   {
1040     \_tl\_build\_one:n
1041     { \if\_false: { \fi: } \_regex\_branch:n { \if\_false: } \fi: }
1042   }
1043 }

```

(End definition for `\_regex_compile_!`.)

`\_regex_compile_!`: Within a class, parentheses are not special. Outside, close a group.

```
1044 \cs_new_protected_nopar:cpn { \_regex_compile_! }
1045 {
1046   \_regex_if_in_class:TF { \_regex_compile_raw:N }
1047   { \_regex_compile_group_end: }
1048 }
```

(End definition for `\_regex_compile_!`.)

`\_regex_compile_special_group::w` Non-capturing, and resetting groups are easy to take care of during compilation; for those  
`\_regex_compile_special_group|:w` groups, the harder parts will come when building.

```
1049 \cs_new_protected_nopar:cpn { \_regex_compile_special_group::w }
1050 { \_regex_compile_group_begin:N \_regex_group_no_capture:nnnN }
1051 \cs_new_protected_nopar:cpn { \_regex_compile_special_group|:w }
1052 { \_regex_compile_group_begin:N \_regex_group_resetting:nnnN }
```

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

`\_regex_compile_special_group_i:w` The match can be made case-insensitive by setting the option with `(?i)`; the original  
`\_regex_compile_special_group-:w` behaviour is restored by `(?-i)`. This is the only supported option.

```
1053 \cs_new_protected:Npn \_regex_compile_special_group_i:w #1#2
1054 {
1055   \str_if_eq:nnTF { #1 #2 } { \_regex_compile_special:N }
1056   {
1057     \cs_set_nopar:Npn \_regex_item_equal:n { \_regex_item_caseless_equal:n }
1058     \cs_set_nopar:Npn \_regex_item_range:nn { \_regex_item_caseless_range:nn }
1059   }
1060   {
1061     \__msg_kernel_warning:nnx { regex } { unknown-option } { (?i #2 }
1062     \_regex_compile_raw:N (
1063       \_regex_compile_raw:N ?
1064       \_regex_compile_raw:N i
1065       #1 #2
1066     )
1067   }
1068 \cs_new_protected_nopar:cpn { \_regex_compile_special_group-:w } #1#2#3#4
1069 {
1070   \str_if_eq:nnTF { #1 #2 #3 #4 }
1071   { \_regex_compile_raw:N i \_regex_compile_special:N }
1072   {
1073     \cs_set_nopar:Npn \_regex_item_equal:n { \_regex_item_caseful_equal:n }
1074     \cs_set_nopar:Npn \_regex_item_range:nn { \_regex_item_caseful_range:nn }
1075   }
1076   {
1077     \__msg_kernel_warning:nnx { regex } { unknown-option } { (?-#2#4 }
1078     \_regex_compile_raw:N (
1079       \_regex_compile_raw:N ?
1080       \_regex_compile_raw:N -
1081       #1 #2 #3 #4
```

```

1082     }
1083 }
(End definition for \_regex_compile_special_group_i:w and \_regex_compile_special_group -:w)

```

### 2.3.11 Catcodes and csnames

\\_regex\_compile\_/c: The \c escape sequence can be followed by a capital letter representing a character category, by a left bracket which starts a list of categories, or by a brace group holding a regular expression for a control sequence name. Otherwise, raise an error.

```

1084 \cs_new_protected:cpn { \_regex_compile_/c: }
1085 { \_regex_chk_c_allowed:T { \_regex_compile_c_test:NN } }
1086 \cs_new_protected:Npn \_regex_compile_c_test:NN #1#2
1087 {
1088   \token_if_eq_meaning:NNTF #1 \_regex_compile_raw:N
1089   {
1090     \int_if_exist:cTF { c\_regex_catcode_#2_int }
1091     {
1092       \int_set_eq:Nc \l\_regex_catcodes_int { c\_regex_catcode_#2_int }
1093       \l\_regex_mode_int
1094       = \if_case:w \l\_regex_mode_int \c_two \else: \c_six \fi:
1095     }
1096   }
1097   { \cs_if_exist_use:cF { \_regex_compile_c_#2:w } }
1098   {
1099     \_msg_kernel_error:nxx { regex } { c-missing-category } {#2}
1100     #1 #2
1101   }
1102 }
(End definition for \_regex_compile_/c: and \_regex_compile_c_test:NN)

```

\\_regex\_compile\_c[:w When encountering \c[, the task is to collect uppercase letters representing character categories. First check for ^ which negates the list of category codes.

```

\_regex_compile_c_lbrack_loop:NN
\_regex_compile_c_lbrack_add:N
\_regex_compile_c_lbrack_end:
1103 \cs_new_protected:cpn { \_regex_compile_c[:w } #1#2
1104 {
1105   \l\_regex_mode_int
1106   = \if_case:w \l\_regex_mode_int \c_two \else: \c_six \fi:
1107   \int_zero:N \l\_regex_catcodes_int
1108   \str_if_eq:nnTF { #1 #2 } { \_regex_compile_special:N ^ }
1109   {
1110     \bool_set_false:N \l\_regex_catcodes_bool
1111     \_regex_compile_c_lbrack_loop:NN
1112   }
1113   {
1114     \bool_set_true:N \l\_regex_catcodes_bool
1115     \_regex_compile_c_lbrack_loop:NN
1116     #1 #2
1117   }
1118 }
1119 \cs_new_protected:Npn \_regex_compile_c_lbrack_loop:NN #1#2

```

```

1120 {
1121   \token_if_eq_meaning:NNTF #1 \_regex_compile_raw:N
1122   {
1123     \int_if_exist:CTF { c\_regex_catcode\_#2\_int }
1124     {
1125       \exp_args:Nc \_regex_compile_c_lbrack_add:N
1126       { c\_regex_catcode\_#2\_int }
1127       \_regex_compile_c_lbrack_loop:NN
1128     }
1129   }
1130   {
1131     \token_if_eq_charcode:NNTF #2 ]
1132     { \_regex_compile_c_lbrack_end: }
1133   }
1134   {
1135     \_msg_kernel_error:nxx { regex } { c-missing-rbrack } {#2}
1136     \_regex_compile_c_lbrack_end:
1137     #1 #2
1138   }
1139 }
1140 \cs_new_protected:Npn \_regex_compile_c_lbrack_add:N #1
1141 {
1142   \if_int_odd:w \_int_eval:w \l\_regex_catcodes_int / #1 \_int_eval_end:
1143   \else:
1144     \tex_advance:D \l\_regex_catcodes_int #1
1145   \fi:
1146 }
1147 \cs_new_protected_nopar:Npn \_regex_compile_c_lbrack_end:
1148 {
1149   \if_meaning:w \c_false_bool \l\_regex_catcodes_bool
1150   \int_set:Nn \l\_regex_catcodes_int
1151   { \c\_regex_all_catcodes_int - \l\_regex_catcodes_int }
1152   \fi:
1153 }

```

*(End definition for \\_regex\_compile\_c\_[w and others.]*

`\_regex_compile_c_{`: The case of a left brace is easy, based on what we have done so far: in a group, compile the regular expression, after changing the mode to forbid nesting `\c`. Additionally, disable submatch tracking since groups don't escape the scope of `\c{...}`.

```

1154 \cs_new_protected_nopar:cpn { \_regex_compile_c \c_lbrace_str :w }
1155 {
1156   \_regex_compile:w
1157   \_regex_disable_submatches:
1158   \l\_regex_mode_int
1159   = - \if_case:w \l\_regex_mode_int \c_two \else: \c_six \fi:
1160 }

```

*(End definition for \\_regex\_compile\_c\_{: This function is documented on page ??.)*

`\_regex_compile_}`: Non-escaped right braces are only special if they appear when compiling the regular expression for a csname, but not within a class: `\c{[]{} }` matches the control sequences



\} and \{... Admittedly, that would be better done as \c{[{}]}]. So, end compiling the inner regex (this closes any dangling class or group). Then insert the corresponding test in the outer regex.

```

1161 \cs_new_protected:cpn { __regex_compile_ \c_rbrace_str : }
1162 {
1163   \__regex_if_in_cs:TF
1164   {
1165     \__regex_compile_end:
1166     \__regex_compile_one:x
1167     { \__regex_item_cs:n { \exp_not:o \l__regex_internal_regex } }
1168   }
1169   { \exp_after:wN \__regex_compile_raw:N \c_rbrace_str }
1170 }

```

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

### 2.3.12 Raw token lists with \u

```

__regex_compile_/u:
__regex_compile_u_loop:NN

```

The \u escape is invalid in classes and directly following a catcode test. Otherwise, it must be followed by a left brace. We then collect the characters for the argument of \u within an x-expanding assignment. In principle we could just wait to encounter a right brace, but this is unsafe: if the right brace is missing, then we will reach the end-markers of the regex, and continue, leading to obscure fatal errors. Instead, we only allow raw and special characters, and stop when encountering a special right brace, any escaped character, or the end-marker.

```

1171 \cs_new_protected:cpn { __regex_compile_/u: } #1#2
1172 {
1173   \__regex_if_in_class_or_catcode:TF
1174   { \__regex_compile_raw_error:N u #1 #2 }
1175   {
1176     \str_if_eq_x:nnTF {#1#2} { \__regex_compile_special:N \c_lbrace_str }
1177     {
1178       \tl_set:Nx \l__regex_internal_a_tl { \if_false: } \fi:
1179       \__regex_compile_u_loop:NN
1180     }
1181     {
1182       \__msg_kernel_error:nn { regex } { u-missing-lbrace }
1183       \__regex_compile_raw:N u #1 #2
1184     }
1185   }
1186 }
1187 \cs_new:Npn \__regex_compile_u_loop:NN #1#2
1188 {
1189   \token_if_eq_meaning:NNTF #1 \__regex_compile_raw:N
1190   { #2 \__regex_compile_u_loop:NN }
1191   {
1192     \token_if_eq_meaning:NNTF #1 \__regex_compile_special:N
1193     {
1194       \exp_after:wN \token_if_eq_charcode:NNTF \c_rbrace_str #2

```

```

1195         { \if_false: { \fi: } \__regex_compile_u_end: }
1196         { #2 \__regex_compile_u_loop:NN }
1197     }
1198     {
1199         \if_false: { \fi: }
1200         \__msg_kernel_error:nxx { regex } { u-missing-rbrace } {#2}
1201         \__regex_compile_u_end:
1202         #1 #2
1203     }
1204 }
1205 }

```

(End definition for \\_\_regex\_compile\_/u: This function is documented on page ??.)

`\__regex_compile_u_end:` Once we have extracted the variable's name, we store the contents of that variable in `\l__regex_internal_a_tl`. The behaviour of `\u` then depends on whether we are within a `\c{...}` escape (in this case, the variable is turned to a string), or not.

```

1206 \cs_new_protected:Npn \__regex_compile_u_end:
1207 {
1208     \tl_set:Nv \l__regex_internal_a_tl { \l__regex_internal_a_tl }
1209     \if_int_compare:w \l__regex_mode_int = \c_zero
1210         \__regex_compile_u_not_cs:
1211     \else:
1212         \__regex_compile_u_in_cs:
1213     \fi:
1214 }

```

(End definition for \\_\_regex\_compile\_u\_end:)

`\__regex_compile_u_in_cs:` When `\u` appears within a control sequence, we convert the variable to a string with escaped spaces. Then for each character insert a class matching exactly that character, once.

```

1215 \cs_new_protected:Npn \__regex_compile_u_in_cs:
1216 {
1217     \exp_args:NNo \__str_gset_other:Nn \g__regex_internal_tl
1218     { \l__regex_internal_a_tl }
1219     \__tl_build_one:x
1220     {
1221         \tl_map_function:NN \g__regex_internal_tl
1222         \__regex_compile_u_in_cs_aux:n
1223     }
1224 }
1225 \cs_new:Npn \__regex_compile_u_in_cs_aux:n #1
1226 {
1227     \__regex_class:NnnnN \c_true_bool
1228     { \__regex_item_caseful_equal:n { \__int_value:w '#1 } }
1229     { 1 } { 0 } \c_false_bool
1230 }

```

(End definition for \\_\_regex\_compile\_u\_in\_cs:)

`\_regex_compile_u_not_cs:` In mode 0, the `\u` escape adds one state to the NFA for each token in `\l\_regex\_internal\_a\_tl`. If a given *<token>* is a control sequence, then insert a string comparison test, otherwise, `\_regex_item_exact:nn` which compares catcode and character code.

```

1231 \cs_new_protected:Npn \_regex_compile_u_not_cs:
1232 {
1233   \exp_args:No \_tl_analysis_map_inline:nn { \l\_regex_internal\_a\_tl }
1234   {
1235     \_tl_build_one:n
1236     {
1237       \_regex_class:NnnnN \c_true_bool
1238       {
1239         \if_int_compare:w "##2 = \c_zero
1240         \_regex_item_exact_cs:c { \exp_after:wN \cs_to_str:N ##1 }
1241         \else:
1242         \_regex_item_exact:nn { \_int_value:w "##2 } { ##3 }
1243         \fi:
1244       }
1245       { 1 } { 0 } \c_false_bool
1246     }
1247   }
1248 }

```

(End definition for `\_regex_compile_u_not_cs:`)

### 2.3.13 Other

`\_regex_compile_/K:` The `\K` control sequence is currently the only “command”, which performs some action, rather than matching something. It is allowed in the same contexts as `\b`. At the compilation stage, we leave it as a single control sequence, defined later.

```

1249 \cs_new_protected_nopar:cpn { \_regex_compile_/K: }
1250 {
1251   \int_compare:nNnTF \l\_regex_mode_int = \c_zero
1252   { \_tl_build_one:n { \_regex_command_K: } }
1253   { \_regex_compile_raw_error:N K }
1254 }

```

(End definition for `\_regex_compile_/K:`)

### 2.3.14 Showing regexes

`\_regex_show:Nx` Within a `\_tl_build:Nw ... \_tl_build_end:` group, we redefine all the function that can appear in a compiled regex, then run the regex. The result is then shown.

```

1255 \cs_new_protected:Npn \_regex_show:Nx #1#2
1256 {
1257   \_tl_build:Nw \l\_regex_internal\_a\_tl
1258   \cs_set_protected_nopar:Npn \_regex_branch:n
1259   {
1260     \seq_pop_right:NN \l\_regex_show_prefix_seq \l\_regex_internal\_a\_tl
1261     \_regex_show_one:n { +-branch }
1262     \seq_put_right:No \l\_regex_show_prefix_seq \l\_regex_internal\_a\_tl

```

```

1263         \use:n
1264     }
1265     \cs_set_protected_nopar:Npn \__regex_group:nnnN
1266     { \__regex_show_group_aux:nnnnN { } }
1267     \cs_set_protected_nopar:Npn \__regex_group_no_capture:nnnN
1268     { \__regex_show_group_aux:nnnnN { ~(no~capture) } }
1269     \cs_set_protected_nopar:Npn \__regex_group_resetting:nnnN
1270     { \__regex_show_group_aux:nnnnN { ~(resetting) } }
1271     \cs_set_eq:NN \__regex_class:NnnnN \__regex_show_class:NnnnN
1272     \cs_set_protected_nopar:Npn \__regex_command_K:
1273     { \__regex_show_one:n { reset~match~start~(\iow_char:N\K) } }
1274     \cs_set_protected:Npn \__regex_assertion:Nn ##1##2
1275     { \__regex_show_one:n { \bool_if:NF ##1 { negative~ } assertion:~##2 } }
1276     \cs_set_nopar:Npn \__regex_b_test: { word~boundary }
1277     \cs_set_eq:NN \__regex_anchor:N \__regex_show_anchor_to_str:N
1278     \cs_set_protected:Npn \__regex_item_caseful_equal:n ##1
1279     { \__regex_show_one:n { char~code~\int_eval:n{##1} } }
1280     \cs_set_protected:Npn \__regex_item_caseful_range:nn ##1##2
1281     { \__regex_show_one:n { range~[\int_eval:n{##1}, \int_eval:n{##2}] } }
1282     \cs_set_protected:Npn \__regex_item_caseless_equal:n ##1
1283     { \__regex_show_one:n { char~code~\int_eval:n{##1}~(caseless) } }
1284     \cs_set_protected:Npn \__regex_item_caseless_range:nn ##1##2
1285     {
1286         \__regex_show_one:n
1287         { Range~[\int_eval:n{##1}, \int_eval:n{##2}]~(caseless) }
1288     }
1289     \cs_set_protected:Npn \__regex_item_catcode:nT
1290     { \__regex_show_item_catcode:NnT \c_true_bool }
1291     \cs_set_protected:Npn \__regex_item_catcode_reverse:nT
1292     { \__regex_show_item_catcode:NnT \c_false_bool }
1293     \cs_set_protected:Npn \__regex_item_reverse:n
1294     { \__regex_show_scope:nn { Reversed~match } }
1295     \cs_set_protected:Npn \__regex_item_exact:nn ##1##2
1296     { \__regex_show_one:n { char~##2,~catcode~##1 } }
1297     \cs_set_protected:Npn \__regex_item_exact_cs:c ##1
1298     { \__regex_show_one:n { control~sequence~\iow_char:N\##1 } }
1299     \cs_set_protected:Npn \__regex_item_cs:n
1300     { \__regex_show_scope:nn { control~sequence } }
1301     \cs_set:cpn { __regex_prop_.: } { \__regex_show_one:n { any~token } }
1302     \seq_clear:N \l__regex_show_prefix_seq
1303     \__regex_show_push:n { ~ }
1304     #1
1305     \__tl_build_end:
1306     \__msg_show_variable:x { > Compiled~regex~#2: \l__regex_internal_a_tl }
1307 }

```

(End definition for \\_\_regex\_show:Nx)

\\_\_regex\_show\_one:n Every part of the final message go through this function, which adds one line to the output, with the appropriate prefix.

```

1308 \cs_new_protected:Npn \__regex_show_one:n #1
1309 {
1310   \int_incr:N \l__regex_show_lines_int
1311   \__tl_build_one:x
1312   {
1313     \iow_newline:
1314     \seq_map_function:NN \l__regex_show_prefix_seq \use:n
1315     #1
1316   }
1317 }
(End definition for \__regex_show_one:n)

```

\\_\_regex\_show\_push:n Enter and exit levels of nesting. The `scope` function prints its first argument as an “introduction”, then performs its second argument in a deeper level of nesting.

```

\__regex_show_pop:
\__regex_show_scope:nn
1318 \cs_new_protected:Npn \__regex_show_push:n #1
1319 { \seq_put_right:Nx \l__regex_show_prefix_seq { #1 ~ } }
1320 \cs_new_protected:Npn \__regex_show_pop:
1321 { \seq_pop_right:NN \l__regex_show_prefix_seq \l__regex_internal_a_tl }
1322 \cs_new_protected:Npn \__regex_show_scope:nn #1#2
1323 {
1324   \__regex_show_one:n {#1}
1325   \__regex_show_push:n { ~ }
1326   #2
1327   \__regex_show_pop:
1328 }
(End definition for \__regex_show_push:n, \__regex_show_pop:, and \__regex_show_scope:nn)

```

\\_\_regex\_show\_group\_aux:nnnnN We display all groups in the same way, simply adding a message, (no capture) or (resetting), to special groups. The odd `\use_ii:nn` avoids printing a spurious `+-branch` for the first branch.

```

1329 \cs_new_protected:Npn \__regex_show_group_aux:nnnnN #1#2#3#4#5
1330 {
1331   \__regex_show_one:n { , -group-begin #1 }
1332   \__regex_show_push:n { | }
1333   \use_ii:nn #2
1334   \__regex_show_pop:
1335   \__regex_show_one:n
1336   { ‘-group-end \__regex_msg_repeated:nnN {#3} {#4} #5 }
1337 }
(End definition for \__regex_show_group_aux:nnnnN)

```

\\_\_regex\_show\_class:NnnnN I’m entirely unhappy about this function: I couldn’t find a way to test if a class is a single test. Instead, collect the representation of the tests in the class. If that had more than one line, write `Match` or `Don’t match` on its own line, with the repeating information if any. Then the various tests on lines of their own, and finally a line. Otherwise, we need to evaluate the representation of the tests again (since the prefix is incorrect). That’s clunky, but not too expensive, since it’s only one test.

```

1338 \cs_set:Npn \__regex_show_class:NnnnN #1#2#3#4#5

```

```

1339 {
1340   \__tl_build:Nw \l__regex_internal_a_tl
1341   \int_zero:N \l__regex_show_lines_int
1342   \__regex_show_push:n {~}
1343   #2
1344   \exp_last_unbraced:Nf
1345   \int_case:nnn { \l__regex_show_lines_int }
1346   {
1347     {0}
1348     {
1349       \__tl_build_end:
1350       \__regex_show_one:n { \bool_if:NTF #1 { Fail } { Pass } }
1351     }
1352     {1}
1353     {
1354       \__tl_build_end:
1355       \bool_if:NTF #1
1356       {
1357         #2
1358         \__tl_build_one:n { \__regex_msg_repeated:nnN {#3} {#4} #5 }
1359       }
1360       {
1361         \__regex_show_one:n
1362         { Don't~match~\__regex_msg_repeated:nnN {#3} {#4} #5 }
1363         \__tl_build_one:o \l__regex_internal_a_tl
1364       }
1365     }
1366   }
1367   {
1368     \__tl_build_end:
1369     \__regex_show_one:n
1370     {
1371       \bool_if:NTF #1 { M } { Don't~m } atch
1372       \__regex_msg_repeated:nnN {#3} {#4} #5
1373     }
1374     \__tl_build_one:o \l__regex_internal_a_tl
1375   }
1376 }

```

(End definition for \\_\_regex\_show\_class:NnnnN)

\\_\_regex\_show\_anchor\_to\_str:N The argument is an integer telling us where the anchor is. We convert that to the relevant info.

```

1377 \cs_new:Npn \__regex_show_anchor_to_str:N #1
1378 {
1379   anchor~at~
1380   \str_case:nnn { #1 }
1381   {
1382     { \l__regex_min_pos_int } { start~(\iow_char:N\\A) }
1383     { \l__regex_start_pos_int } { start~of~match~(\iow_char:N\\G) }

```

```

1384         { \l__regex_max_pos_int   } { end~(\iow_char:N\Z) }
1385     }
1386     { <error:~'~'#1'~not~recognized> }
1387 }

```

(End definition for \\_\_regex\_show\_anchor\_to\_str:N)

\\_\_regex\_show\_item\_catcode:NnT

Produce a sequence of categories which the catcode bitmap #2 contains, and show it, indenting the tests on which this catcode constraint applies.

```

1388 \cs_new_protected:Npn \__regex_show_item_catcode:NnT #1#2
1389 {
1390     \seq_set_split:Nnn \l__regex_internal_seq { } { CBEMTPUDSLOA }
1391     \seq_set_filter:NNn \l__regex_internal_seq \l__regex_internal_seq
1392     { \int_if_odd_p:n { #2 / \int_use:c { c__regex_catcode_##1_int } } }
1393     \__regex_show_scope:nn
1394     {
1395         categories~
1396         \seq_map_function:NN \l__regex_internal_seq \use:n
1397         , ~
1398         \bool_if:NF #1 { negative~ } class
1399     }
1400 }

```

(End definition for \\_\_regex\_show\_item\_catcode:NnT)

## 2.4 Building

### 2.4.1 Variables used while building

\l\_\_regex\_min\_state\_int  
\l\_\_regex\_max\_state\_int

The last state that was allocated is `\l__regex_max_state_int - 1`, so that `\l__regex_max_state_int` always points to a free state. The `min_state` variable is always 0, but is included to avoid hard-coding this value.

```

1401 \int_new:N \l__regex_min_state_int
1402 \int_new:N \l__regex_max_state_int

```

(End definition for \l\_\_regex\_min\_state\_int and \l\_\_regex\_max\_state\_int These variables are documented on page ??.)

\l\_\_regex\_left\_state\_int  
\l\_\_regex\_right\_state\_int  
\l\_\_regex\_left\_state\_seq  
\l\_\_regex\_right\_state\_seq

Alternatives are implemented by branching from a `left` state into the various choices, then merging those into a `right` state. We store information about those states in two sequences. Those states are also used to implement group quantifiers. Most often, the left and right pointers only differ by 1.

```

1403 \int_new:N \l__regex_left_state_int
1404 \int_new:N \l__regex_right_state_int
1405 \seq_new:N \l__regex_left_state_seq
1406 \seq_new:N \l__regex_right_state_seq

```

(End definition for \l\_\_regex\_left\_state\_int and \l\_\_regex\_right\_state\_int These functions are documented on page ??.)

`\1_regex_capturing_group_int` `\1_regex_capturing_group_int` is the ID number that will be assigned to a capturing group if one was opened now. This starts at 0 for the group enclosing the full regular expression, and groups are counted in the order of their left parenthesis, except when encountering **resetting** groups.

<sup>1407</sup> `\int_new:N \1_regex_capturing_group_int`  
*(End definition for \1\_regex\_capturing\_group\_int This variable is documented on page ??.)*

### 2.4.2 Framework

This phase is about going from a compiled regex to an NFA. Each state of the NFA is stored in a `\toks`. The operations which can appear in the `\toks` are

- `\_regex_action_start_wildcard`: inserted at the start of the regular expression to make it unanchored.
- `\_regex_action_success`: marks the exit state of the NFA.
- `\_regex_action_cost:n {\langle shift \rangle}` is a transition from the current  $\langle state \rangle$  to  $\langle state \rangle + \langle shift \rangle$ , which consumes the current character: the target state is saved and will be considered again when matching at the next position.
- `\_regex_action_free:n {\langle shift \rangle}`, and `\_regex_action_free_group:n {\langle shift \rangle}` are free transitions, which immediately perform the actions for the state  $\langle state \rangle + \langle shift \rangle$  of the NFA. They differ in how they detect and avoid infinite loops. For now, we just need to know that the **group** variant must be used for transitions back to the start of a group.
- `\_regex_action_submatch:n {\langle key \rangle}` where the  $\langle key \rangle$  is a group number followed by `<` or `>` for the beginning or end of group. This causes the current position in the query to be stored as the  $\langle key \rangle$  submatch boundary.

We strive to preserve the following properties while building.

- The current capturing group is `capturing_group - 1`, and if a group is opened now, it will be labelled `capturing_group`.
- The last allocated state is `max_state - 1`, so `max_state` is a free state.
- The `left_state` points to a state to the left of the current group or of the last class.
- The `right_state` points to a newly created, empty state, with some transitions leading to it.
- The `left/right` sequences hold a list of the corresponding end-points of nested groups.



`\_regex_build:n` The `n`-type function first compiles its argument. Reset some variables. Allocate two states, and put a wildcard in state 0 (transitions to state 1 and 0 state). Then build the regex within a (capturing) group, which will be numbered 0 (current value of `capturing_group`). Finally, if the match reaches the last state, it is successful.

```

1408 \cs_new_protected:Npn \_regex_build:n #1
1409 {
1410   \_regex_compile:n {#1}
1411   \_regex_build:N \l__regex_internal_regex
1412 }
1413 \cs_new_protected:Npn \_regex_build:N #1
1414 {
1415   <trace> \trace_push:nnn { regex } { 1 } { __regex_build }
1416   \int_set:Nn \tex_escapechar:D { 92 }
1417   \int_zero:N \l__regex_capturing_group_int
1418   \int_set_eq:NN \l__regex_max_state_int \l__regex_min_state_int
1419   \_regex_build_new_state:
1420   \_regex_build_new_state:
1421   \_regex_toks_put_right:Nn \l__regex_left_state_int
1422   { \_regex_action_start_wildcard: }
1423   \_regex_group:nnnN {#1} { 1 } { 0 } \c_false_bool
1424   \_regex_toks_put_right:Nn \l__regex_right_state_int
1425   { \_regex_action_success: }
1426   <trace> \_regex_trace_states:n { 2 }
1427   <trace> \trace_pop:nnn { regex } { 1 } { __regex_build }
1428 }

```

*(End definition for \\_regex\_build:n and \\_regex\_build:N)*

`\_regex_build_for_cs:n` When using a regex to match a cs, we don't insert a wildcard, we anchor at the end, and since we ignore submatches, there is no need to surround the expression with a group. However, for branches to work properly at the outer level, we need to put the appropriate `left` and `right` states in their sequence.

```

1429 \cs_new_protected:Npn \_regex_build_for_cs:n #1
1430 {
1431   <trace> \trace_push:nnn { regex } { 1 } { __regex_build_for_cs }
1432   \int_set_eq:NN \l__regex_max_state_int \l__regex_min_state_int
1433   \_regex_build_new_state:
1434   \_regex_build_new_state:
1435   \_regex_push_lr_states:
1436   #1
1437   \_regex_pop_lr_states:
1438   \_regex_toks_put_right:Nn \l__regex_right_state_int
1439   {
1440     \if_int_compare:w \l__regex_current_pos_int = \l__regex_max_pos_int
1441     \exp_after:wN \_regex_action_success:
1442     \fi:
1443   }
1444   <trace> \_regex_trace_states:n { 2 }
1445   <trace> \trace_pop:nnn { regex } { 1 } { __regex_build_for_cs }
1446 }

```

(End definition for `\_regex_build_for_cs:n`)

### 2.4.3 Helpers for building an nfa

`\_regex_push_lr_states:` When building the regular expression, we keep track of pointers to the left-end and  
`\_regex_pop_lr_states:` right-end of each group without help from T<sub>E</sub>X’s grouping.

```

1447 \cs_new_protected_nopar:Npn \_regex_push_lr_states:
1448 {
1449   \seq_push:No \l__regex_left_state_seq
1450   { \int_use:N \l__regex_left_state_int }
1451   \seq_push:No \l__regex_right_state_seq
1452   { \int_use:N \l__regex_right_state_int }
1453 }
1454 \cs_new_protected_nopar:Npn \_regex_pop_lr_states:
1455 {
1456   \seq_pop:NN \l__regex_left_state_seq \l__regex_internal_a_tl
1457   \int_set:Nn \l__regex_left_state_int \l__regex_internal_a_tl
1458   \seq_pop:NN \l__regex_right_state_seq \l__regex_internal_a_tl
1459   \int_set:Nn \l__regex_right_state_int \l__regex_internal_a_tl
1460 }

```

(End definition for `\_regex_push_lr_states:` and `\_regex_pop_lr_states:`)

`\_regex_toks_put_left:Nx` During the building phase we wish to add x-expanded material to `\toks`, either to the left  
`\_regex_toks_put_right:Nx` or to the right. The expansion is done “by hand” for optimization (these operations are  
`\_regex_toks_put_right:Nn` used quite a lot). The `Nn` version of `\_regex_toks_put_right:Nx` is provided because  
it is more efficient than x-expanding with `\exp_not:n`.

```

1461 \cs_new_protected:Npn \_regex_toks_put_left:Nx #1#2
1462 {
1463   \cs_set_nopar:Npx \_regex_tmp:w { #2 }
1464   \tex_toks:D #1 \exp_after:wN \exp_after:wN \exp_after:wN
1465   { \exp_after:wN \_regex_tmp:w \tex_the:D \tex_toks:D #1 }
1466 }
1467 \cs_new_protected:Npn \_regex_toks_put_right:Nx #1#2
1468 {
1469   \cs_set_nopar:Npx \_regex_tmp:w {#2}
1470   \tex_toks:D #1 \exp_after:wN
1471   { \tex_the:D \tex_toks:D \exp_after:wN #1 \_regex_tmp:w }
1472 }
1473 \cs_new_protected:Npn \_regex_toks_put_right:Nn #1#2
1474 { \tex_toks:D #1 \exp_after:wN { \tex_the:D \tex_toks:D #1 #2 } }

```

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

`\_regex_build_transition_left:NNN` Add a transition from #2 to #3 using the function #1. The `left` function is used for  
`\_regex_build_transition_right:nNn` higher priority transitions, and the `right` function for lower priority transitions (which  
should be performed later). The signatures differ to reflect the differing usage later on.  
Both functions could be optimized.

```

1475 \cs_new_protected:Npn \_regex_build_transition_left:NNN #1#2#3
1476 { \_regex_toks_put_left:Nx #2 { #1 { \int_eval:n { #3 - #2 } } } }

```

```

1477 \cs_new_protected:Npn \__regex_build_transition_right:nNn #1#2#3
1478   { \__regex_toks_put_right:Nx #2 { #1 { \int_eval:n { #3 - #2 } } } }
(End definition for \__regex_build_transition_left:NNN and \__regex_build_transition_right:nNn)

```

**\\_\_regex\_build\_new\_state:** Add a new empty state to the NFA. Then update the left, right, and max states, so that the right state is the new empty state, and the left state points to the previously “current” state.

```

1479 \cs_new_protected_nopar:Npn \__regex_build_new_state:
1480   {
1481   (*trace)
1482     \trace:nxx { regex } { 2 }
1483     {
1484       regex~new~state~
1485       L=\int_use:N \l__regex_left_state_int ~ -> ~
1486       R=\int_use:N \l__regex_right_state_int ~ -> ~
1487       M=\int_use:N \l__regex_max_state_int ~ -> ~
1488       \int_eval:n { \l__regex_max_state_int + \c_one }
1489     }
1490   (/trace)
1491     \tex_toks:D \l__regex_max_state_int { }
1492     \int_set_eq:NN \l__regex_left_state_int \l__regex_right_state_int
1493     \int_set_eq:NN \l__regex_right_state_int \l__regex_max_state_int
1494     \int_incr:N \l__regex_max_state_int
1495   }
(End definition for \__regex_build_new_state:)

```

**\\_\_regex\_build\_transitions\_lazyness:NNNN** This function creates a new state, and puts two transitions starting from the old current state. The order of the transitions is controlled by #1, true for lazy quantifiers, and false for greedy quantifiers.

```

1496 \cs_new_protected:Npn \__regex_build_transitions_lazyness:NNNN #1#2#3#4#5
1497   {
1498     \__regex_build_new_state:
1499     \__regex_toks_put_right:Nx \l__regex_left_state_int
1500     {
1501       \if_meaning:w \c_true_bool #1
1502         #2 { \int_eval:n { #3 - \l__regex_left_state_int } }
1503         #4 { \int_eval:n { #5 - \l__regex_left_state_int } }
1504       \else:
1505         #4 { \int_eval:n { #5 - \l__regex_left_state_int } }
1506         #2 { \int_eval:n { #3 - \l__regex_left_state_int } }
1507       \fi:
1508     }
1509   }
(End definition for \__regex_build_transitions_lazyness:NNNN)

```

## 2.4.4 Building classes

`__regex_class:NnnnN` The arguments are:  $\langle \text{boolean} \rangle$   $\{ \langle \text{tests} \rangle \}$   $\{ \langle \text{min} \rangle \}$   $\{ \langle \text{more} \rangle \}$   $\langle \text{lazyness} \rangle$ . First store the tests with a trailing `__regex_action_cost:n`, in the true branch of `__regex_break_point:TF` for positive classes, or the false branch for negative classes. The integer  $\langle \text{more} \rangle$  is 0 for fixed repetitions,  $-1$  for unbounded repetitions, and  $\langle \text{max} \rangle - \langle \text{min} \rangle$  for a range of repetitions.

```

1510 \cs_new_protected:Npn __regex_class:NnnnN #1#2#3#4#5
1511 {
1512   \cs_set_nopar:Npx __regex_tests_action_cost:n ##1
1513   {
1514     \exp_not:n { \exp_not:n {#2} }
1515     \bool_if:NTF #1
1516     { __regex_break_point:TF { __regex_action_cost:n {##1} } { } }
1517     { __regex_break_point:TF { } { __regex_action_cost:n {##1} } }
1518   }
1519   \if_case:w - #4 \exp_stop_f:
1520     __regex_class_repeat:n {#3}
1521   \or: __regex_class_repeat:nN {#3} #5
1522   \else: __regex_class_repeat:nnN {#3} {#4} #5
1523   \fi:
1524 }
1525 \cs_new:Npn __regex_tests_action_cost:n { __regex_action_cost:n }

```

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

`__regex_class_repeat:n` This is used for a fixed number of repetitions. Build one state for each repetition, with a transition controlled by the tests that we have collected. That works just fine for  $\#1 = 0$  repetitions: nothing is built.

```

1526 \cs_new_protected:Npn __regex_class_repeat:n #1
1527 {
1528   \prg_replicate:nn {#1}
1529   {
1530     __regex_build_new_state:
1531     __regex_build_transition_right:nNn __regex_tests_action_cost:n
1532     \l__regex_left_state_int \l__regex_right_state_int
1533   }
1534 }

```

(End definition for `__regex_class_repeat:n`)

`__regex_class_repeat:nN` This implements unbounded repetitions of a single class (e.g. the  $*$  and  $+$  quantifiers). If the minimum number  $\#1$  of repetitions is 0, then build a transition from the current state to itself governed by the tests, and a free transition to a new state (hence skipping the tests). Otherwise, call `__regex_class_repeat:n` for the code to match  $\#1$  repetitions, and add free transitions from the last state to the previous one, and to a new one. In both cases, the order of transitions is controlled by the lazyness boolean  $\#2$ .

```

1535 \cs_new_protected:Npn __regex_class_repeat:nN #1#2
1536 {
1537   \if_int_compare:w #1 = \c_zero

```

```

1538     \_regex_build_transitions_lazyness:NNNN #2
1539     \_regex_action_free:n      \l__regex_right_state_int
1540     \_regex_tests_action_cost:n \l__regex_left_state_int
1541   \else:
1542     \_regex_class_repeat:n {#1}
1543     \int_set_eq:NN \l__regex_internal_a_int \l__regex_left_state_int
1544     \_regex_build_transitions_lazyness:NNNN #2
1545     \_regex_action_free:n \l__regex_right_state_int
1546     \_regex_action_free:n \l__regex_internal_a_int
1547   \fi:
1548 }
(End definition for \_regex_class_repeat:nN)

```

\\_regex\_class\_repeat:nnN

We want to build the code to match from #1 to #1 + #2 repetitions. Match #1 repetitions (can be 0). Compute the final state of the next construction as `a`. Build #2 > 0 states, each with a transition to the next state governed by the tests, and a transition to the final state `a`. The computation of `a` is safe because states are allocated in order, starting from `max_state`.

```

1549 \cs_new_protected:Npn \_regex_class_repeat:nnN #1#2#3
1550 {
1551   \_regex_class_repeat:n {#1}
1552   \int_set:Nn \l__regex_internal_a_int
1553     { \l__regex_max_state_int + #2 - \c_one }
1554   \prg_replicate:nn { #2 }
1555   {
1556     \_regex_build_transitions_lazyness:NNNN #3
1557     \_regex_action_free:n      \l__regex_internal_a_int
1558     \_regex_tests_action_cost:n \l__regex_right_state_int
1559   }
1560 }
(End definition for \_regex_class_repeat:nnN)

```

## 2.4.5 Building groups

\\_regex\_group\_aux:nnnnN

Arguments:  $\{\langle label \rangle\} \{\langle contents \rangle\} \{\langle min \rangle\} \{\langle more \rangle\} \langle lazyness \rangle$ . If  $\langle min \rangle$  is 0, we need to add a state before building the group, so that the thread which skips the group does not also set the start-point of the submatch. After adding one more state, the `left_state` is the left end of the group, from which all branches will stem, and the `right_state` is the right end of the group, and all branches end their course in that state. We store those two integers to be queried for each branch, we build the NFA states for the contents #2 of the group, and we forget about the two integers. Once this is done, perform the repetition: either exactly #3 times, or #3 or more times, or between #3 and #3 + #4 times, with lazyness #5. The  $\langle label \rangle$  #1 is used for submatch tracking. Each of the three auxiliaries expects `left_state` and `right_state` to be set properly.

```

1561 \cs_new_protected:Npn \_regex_group_aux:nnnnN #1#2#3#4#5
1562 {
1563   \trace { \trace_push:nnn { regex } { 1 } { \_regex_group } }
1564   \if_int_compare:w #3 = \c_zero

```

```

1565     \_regex_build_new_state:
1566 <assert>\assert_int:n { \l__regex_max_state_int = \l__regex_right_state_int + 1 }
1567     \_regex_build_transition_right:nNn \_regex_action_free_group:n
1568     \l__regex_left_state_int \l__regex_right_state_int
1569     \fi:
1570     \_regex_build_new_state:
1571     \_regex_push_lr_states:
1572     #2
1573     \_regex_pop_lr_states:
1574     \if_case:w - #4 \exp_stop_f:
1575         \_regex_group_repeat:nn {#1} {#3}
1576     \or: \_regex_group_repeat:nnN {#1} {#3} #5
1577     \else: \_regex_group_repeat:nnnN {#1} {#3} {#4} #5
1578     \fi:
1579 <trace> \trace_pop:nnn { regex } { 1 } { __regex_group }
1580     }

```

(End definition for \\_regex\_group\_aux:nnnnN)

\\_regex\_group:nnnN Hand to \\_regex\_group\_aux:nnnnN the label of that group (expanded), and the group itself, with some extra commands to perform.

```

1581 \cs_new_protected:Npn \_regex_group:nnnN #1
1582 {
1583     \exp_args:No \_regex_group_aux:nnnnN
1584     { \int_use:N \l__regex_capturing_group_int }
1585     {
1586         \int_incr:N \l__regex_capturing_group_int
1587         #1
1588     }
1589 }
1590 \cs_new_protected_nopar:Npn \_regex_group_no_capture:nnnN
1591 { \_regex_group_aux:nnnnN { -1 } }

```

(End definition for \\_regex\_group:nnnN and \\_regex\_group\_no\_capture:nnnN)

\\_regex\_group\_resetting:nnnN Again, hand the label -1 to \\_regex\_group\_aux:nnnnN, but this time we work a little  
\\_regex\_group\_resetting\_loop:nnNn bit harder to keep track of the maximum group label at the end of any branch, and to reset the group number at each branch. This relies on the fact that a compiled regex always is a sequence of items of the form \\_regex\_branch:n {<branch>}.

```

1592 \cs_new_protected:Npn \_regex_group_resetting:nnnN #1
1593 {
1594     \_regex_group_aux:nnnnN { -1 }
1595     {
1596         \exp_args:Noo \_regex_group_resetting_loop:nnNn
1597         { \int_use:N \l__regex_capturing_group_int }
1598         { \int_use:N \l__regex_capturing_group_int }
1599         #1
1600         { ?? \_prg_break:n } { }
1601         \_prg_break_point:
1602     }
1603 }

```

```

1604 \cs_new_protected:Npn \__regex_group_resetting_loop:nnNn #1#2#3#4
1605 {
1606   \use_none:nn #3 { \int_set:Nn \l__regex_capturing_group_int {#1} }
1607   \int_set:Nn \l__regex_capturing_group_int {#2}
1608   #3 {#4}
1609   \exp_args:Nf \__regex_group_resetting_loop:nnNn
1610     { \int_max:nn {#1} { \l__regex_capturing_group_int } }
1611     {#2}
1612 }

```

(End definition for \\_\_regex\_group\_resetting:nnnN This function is documented on page ??.)

**\\_\_regex\_branch:n** Add a free transition from the left state of the current group to a brand new state, starting point of this branch. Once the branch is built, add a transition from its last state to the right state of the group. The left and right states of the group are extracted from the relevant sequences.

```

1613 \cs_new_protected:Npn \__regex_branch:n #1
1614 {
1615   <trace> \trace_push:nnn { regex } { 1 } { __regex_branch }
1616   \__regex_build_new_state:
1617   \seq_get:NN \l__regex_left_state_seq \l__regex_internal_a_tl
1618   \int_set:Nn \l__regex_left_state_int \l__regex_internal_a_tl
1619   \__regex_build_transition_right:nNn \__regex_action_free:n
1620     \l__regex_left_state_int \l__regex_right_state_int
1621   #1
1622   \seq_get:NN \l__regex_right_state_seq \l__regex_internal_a_tl
1623   \__regex_build_transition_right:nNn \__regex_action_free:n
1624     \l__regex_right_state_int \l__regex_internal_a_tl
1625   <trace> \trace_pop:nnn { regex } { 1 } { __regex_branch }
1626 }

```

(End definition for \\_\_regex\_branch:n)

**\\_\_regex\_group\_repeat:nn** This function is called to repeat a group a fixed number of times #2; if this is 0 we remove the group altogether (but don't reset the `capturing_group` label). Otherwise, the auxiliary `\__regex_group_repeat_aux:n` copies #2 times the `\toks` for the group, and leaves `internal_a` pointing to the left end of the last repetition. We only record the submatch information at the last repetition. Finally, add a state at the end (the transition to it has been taken care of by the replicating auxiliary).

```

1627 \cs_new_protected:Npn \__regex_group_repeat:nn #1#2
1628 {
1629   \if_int_compare:w #2 = \c_zero
1630     \int_set:Nn \l__regex_max_state_int
1631       { \l__regex_left_state_int - \c_one }
1632     \__regex_build_new_state:
1633   \else:
1634     \__regex_group_repeat_aux:n {#2}
1635     \__regex_group_submatches:nNn {#1}
1636     \l__regex_internal_a_int \l__regex_right_state_int
1637     \__regex_build_new_state:

```

```

1638     \fi:
1639   }
(End definition for \_regex_group_repeat:nn)

```

\\_regex\_group\_submatches:nnN This inserts in states #2 and #3 the code for tracking submatches of the group #1, unless inhibited by a label of -1.

```

1640 \cs_new_protected:Npn \_regex_group_submatches:nnN #1#2#3
1641 {
1642   \if_int_compare:w #1 > \c_minus_one
1643     \_regex_toks_put_left:Nx #2 { \_regex_action_submatch:n { #1 < } }
1644     \_regex_toks_put_left:Nx #3 { \_regex_action_submatch:n { #1 > } }
1645   \fi:
1646 }
(End definition for \_regex_group_submatches:nnN)

```

\\_regex\_group\_repeat\_aux:n Here we repeat \toks ranging from left\_state to max\_state, #1 > 0 times. First add a transition so that the copies will “chain” properly. Compute the shift c between the original copy and the last copy we want. Shift the right\_state and max\_state to their final values. We then want to perform c copy operations. At the end, b is equal to the max\_state, and a points to the left of the last copy of the group.

```

1647 \cs_new_protected:Npn \_regex_group_repeat_aux:n #1
1648 {
1649   \_regex_build_transition_right:nnN \_regex_action_free:n
1650   \l__regex_right_state_int \l__regex_max_state_int
1651   \int_set_eq:NN \l__regex_internal_a_int \l__regex_left_state_int
1652   \int_set_eq:NN \l__regex_internal_b_int \l__regex_max_state_int
1653   \if_int_compare:w \_int_eval:w #1 > \c_one
1654     \int_set:Nn \l__regex_internal_c_int
1655     {
1656       ( #1 - \c_one )
1657       * ( \l__regex_internal_b_int - \l__regex_internal_a_int )
1658     }
1659     \tex_advance:D \l__regex_right_state_int \l__regex_internal_c_int
1660     \tex_advance:D \l__regex_max_state_int \l__regex_internal_c_int
1661     \prg_replicate:nn \l__regex_internal_c_int
1662     {
1663       \tex_toks:D \l__regex_internal_b_int
1664       = \tex_toks:D \l__regex_internal_a_int
1665       \tex_advance:D \l__regex_internal_a_int \c_one
1666       \tex_advance:D \l__regex_internal_b_int \c_one
1667     }
1668   \fi:
1669 }
(End definition for \_regex_group_repeat_aux:n)

```

\\_regex\_group\_repeat:nnN This function is called to repeat a group at least n times; the case n = 0 is very different from n > 0. Assume first that n = 0. Insert submatch tracking information at the start and end of the group, add a free transition from the right end to the “true” left state a



(remember: in this case we had added an extra state before the left state). This forms the loop, which we break away from by adding a free transition from a to a new state.

Now consider the case  $n > 0$ . Repeat the group  $n$  times, chaining various copies with a free transition. Add submatch tracking only to the last copy, then add a free transition from the right end back to the left end of the last copy, either before or after the transition to move on towards the rest of the NFA. This transition can end up before submatch tracking, but that is irrelevant since it only does so when going again through the group, recording new matches. Finally, add a state; we already have a transition pointing to it from `__regex_group_repeat_aux:n`.

```

1670 \cs_new_protected:Npn __regex_group_repeat:nnN #1#2#3
1671 {
1672   \if_int_compare:w #2 = \c_zero
1673     \__regex_group_submatches:nnN {#1}
1674     \l__regex_left_state_int \l__regex_right_state_int
1675     \int_set:Nn \l__regex_internal_a_int
1676       { \l__regex_left_state_int - \c_one }
1677     \__regex_build_transition_right:nNn __regex_action_free:n
1678     \l__regex_right_state_int \l__regex_internal_a_int
1679     \__regex_build_new_state:
1680     \if_meaning:w \c_true_bool #3
1681       \__regex_build_transition_left:NNN __regex_action_free:n
1682       \l__regex_internal_a_int \l__regex_right_state_int
1683     \else:
1684       \__regex_build_transition_right:nNn __regex_action_free:n
1685       \l__regex_internal_a_int \l__regex_right_state_int
1686     \fi:
1687   \else:
1688     \__regex_group_repeat_aux:n {#2}
1689     \__regex_group_submatches:nnN {#1}
1690     \l__regex_internal_a_int \l__regex_right_state_int
1691     \if_meaning:w \c_true_bool #3
1692       \__regex_build_transition_right:nNn __regex_action_free_group:n
1693       \l__regex_right_state_int \l__regex_internal_a_int
1694     \else:
1695       \__regex_build_transition_left:NNN __regex_action_free_group:n
1696       \l__regex_right_state_int \l__regex_internal_a_int
1697     \fi:
1698     \__regex_build_new_state:
1699   \fi:
1700 }

```

(End definition for `__regex_group_repeat:nnN`)

`__regex_group_repeat:nnnN`

We wish to repeat the group between `#2` and `#2 + #3` times, with a laziness controlled by `#4`. We insert submatch tracking up front: in principle, we could avoid recording submatches for the first `#2` copies of the group, but that forces us to treat specially the case `#2 = 0`. Repeat that group with submatch tracking `#2 + #3` times (the maximum number of repetitions). Then our goal is to add `#3` transitions from the end of the `#2`-th group, and each subsequent groups, to the end. For a lazy quantifier, we add those

transitions to the left states, before submatch tracking. For the greedy case, we add the transitions to the right states, after submatch tracking and the transitions which go on with more repetitions. In the greedy case with #2 = 0, the transition which skips over all copies of the group must be added separately, because its starting state does not follow the normal pattern: we had to add it “by hand” earlier.

```

1701 \cs_new_protected:Npn \__regex_group_repeat:nnnN #1#2#3#4
1702 {
1703   \__regex_group_submatches:nNN {#1}
1704   \l__regex_left_state_int \l__regex_right_state_int
1705   \__regex_group_repeat_aux:n { #2 + #3 }
1706   \if_meaning:w \c_true_bool #4
1707     \int_set_eq:NN \l__regex_left_state_int \l__regex_max_state_int
1708     \prg_replicate:nn { #3 }
1709     {
1710       \int_sub:Nn \l__regex_left_state_int
1711       { \l__regex_internal_b_int - \l__regex_internal_a_int }
1712       \__regex_build_transition_left:NNN \__regex_action_free:n
1713       \l__regex_left_state_int \l__regex_max_state_int
1714     }
1715   \else:
1716     \prg_replicate:nn { #3 - \c_one }
1717     {
1718       \int_sub:Nn \l__regex_right_state_int
1719       { \l__regex_internal_b_int - \l__regex_internal_a_int }
1720       \__regex_build_transition_right:nNn \__regex_action_free:n
1721       \l__regex_right_state_int \l__regex_max_state_int
1722     }
1723     \if_int_compare:w #2 = \c_zero
1724       \int_set:Nn \l__regex_right_state_int
1725       { \l__regex_left_state_int - \c_one }
1726     \else:
1727       \int_sub:Nn \l__regex_right_state_int
1728       { \l__regex_internal_b_int - \l__regex_internal_a_int }
1729     \fi:
1730     \__regex_build_transition_right:nNn \__regex_action_free:n
1731     \l__regex_right_state_int \l__regex_max_state_int
1732   \fi:
1733   \__regex_build_new_state:
1734 }

```

(End definition for \\_\_regex\_group\_repeat:nnnN)

## 2.4.6 Others

<code>\__regex_assertion:Nn</code> <code>\__regex_b_test:</code> <code>\__regex_anchor:N</code>	<p>Usage: <code>\__regex_assertion:Nn &lt;boolean&gt; {&lt;test&gt;}</code>, where the <code>&lt;test&gt;</code> is either of the two other functions. Add a free transition to a new state, conditionally to the assertion test.</p> <p>The <code>\__regex_b_test:</code> test is used by the <code>\b</code> and <code>\B</code> escape: check if the last character was a word character or not, and do the same to the current character. The boundary-markers of the string are non-word characters for this purpose. Anchors at the start or</p>
---	--

end of match use `\__regex_anchor:N`, with a position controlled by the integer #1.

```

1735 \cs_new_protected:Npn \__regex_assertion:Nn #1#2
1736 {
1737   \__regex_build_new_state:
1738   \__regex_toks_put_right:Nx \l__regex_left_state_int
1739   {
1740     \exp_not:n {#2}
1741     \__regex_break_point:TF
1742     \bool_if:NF #1 { { } }
1743     {
1744       \__regex_action_free:n
1745       {
1746         \int_eval:n
1747         { \l__regex_right_state_int - \l__regex_left_state_int }
1748       }
1749     }
1750     \bool_if:NT #1 { { } }
1751   }
1752 }
1753 \cs_new_protected:Npn \__regex_anchor:N #1
1754 {
1755   \if_int_compare:w #1 = \l__regex_current_pos_int
1756   \exp_after:wN \__regex_break_true:w
1757   \fi:
1758 }
1759 \cs_new_protected_nopar:Npn \__regex_b_test:
1760 {
1761   \group_begin:
1762   \int_set_eq:NN \l__regex_current_char_int \l__regex_last_char_int
1763   \__regex_prop_w:
1764   \__regex_break_point:TF
1765   { \group_end: \__regex_item_reverse:n \__regex_prop_w: }
1766   { \group_end: \__regex_prop_w: }
1767 }

```

(End definition for `\__regex_assertion:Nn`, `\__regex_b_test:`, and `\__regex_anchor:N`)

`\__regex_command_K:` Change the starting point of the 0-th submatch (full match), and transition to a new state, pretending that this is a fresh thread.

```

1768 \cs_new_protected_nopar:Npn \__regex_command_K:
1769 {
1770   \__regex_build_new_state:
1771   \__regex_toks_put_right:Nx \l__regex_left_state_int
1772   {
1773     \__regex_action_submatch:n { 0< }
1774     \bool_set_true:N \l__regex_fresh_thread_bool
1775     \__regex_action_free:n
1776     { \int_eval:n { \l__regex_right_state_int - \l__regex_left_state_int } }
1777     \bool_set_false:N \l__regex_fresh_thread_bool
1778   }

```

```

1779     }
(End definition for \_regex_command_K:)

```

## 2.5 Matching

We search for matches by running all the execution threads through the NFA in parallel, reading one token of the query at each step. The NFA contains “free” transitions to other states, and transitions which “consume” the current token. For free transitions, the instruction at the new state of the NFA is performed immediately. When a transition consumes a character, the new state is appended to a list of “active states”, stored in `\skip` registers: this thread will be active again when the next token is read from the query. At every step (for each token in the query), we unpack that list of active states and the corresponding submatch props, and empty those.

If two paths through the NFA “collide” in the sense that they reach the same state after reading a given token, then they only differ in how they previously matched, and the future execution will be identical for both. (Note that this would be wrong in the presence of back-references.) Hence, we only need to keep one of the two threads: the thread with the highest priority. Our NFA is built in such a way that higher priority actions always come before lower priority actions, which makes things work.

The explanation in the previous paragraph may make us think that we simply need to keep track of which states were visited at a given step: after all, the loop generated when matching `(a?)*` against `a` is broken, isn’t it? No. The group first matches `a`, as it should, then repeats; it attempts to match `a` again but fails; it skips `a`, and finds out that this state has already been seen at this position in the query: the match stops. The capturing group is (wrongly) `a`. What went wrong is that a thread collided with itself, and the later version, which has gone through the group one more times with an empty match, should have a higher priority than not going through the group.

We solve this by distinguishing “normal” free transitions `\_regex_action_free:n` from transitions `\_regex_action_free_group:n` which go back to the start of the group. The former will keep threads unless they have been visited by a “completed” thread, while the latter kind of transition also prevents going back to a state visited by the current thread.

### 2.5.1 Variables used when matching

<pre> \l__regex_min_pos_int \l__regex_max_pos_int \l__regex_current_pos_int \l__regex_start_pos_int \l__regex_success_pos_int </pre>	<p>The tokens in the query are indexed from <code>min_pos</code> for the first to <code>max_pos - 1</code> for the last, and their information is stored in <code>\muskip</code> and <code>\toks</code> registers with those numbers. We don’t start from 0 because the <code>\toks</code> registers with low numbers are used to hold the states of the NFA. We match without backtracking, keeping all threads in lockstep at the <code>current_pos</code> in the query. The starting point of the current match attempt is <code>start_pos</code>, and <code>success_pos</code>, updated whenever a thread succeeds, is used as the next starting position.</p>
--	--

```

1780 \int_new:N \l__regex_min_pos_int
1781 \int_new:N \l__regex_max_pos_int
1782 \int_new:N \l__regex_current_pos_int
1783 \int_new:N \l__regex_start_pos_int

```

1784 `\int_new:N \l__regex_success_pos_int`  
*(End definition for \l\_\_regex\_min\_pos\_int and others. These variables are documented on page ??.)*

`\l__regex_current_char_int`    The character and category codes of the token at the current position; the character code  
`\l__regex_current_catcode_int`    of the token at the previous position; and the character code of the result of changing the  
`\l__regex_last_char_int`    case of the current token (A-Z↔a-z). This last integer is only computed when necessary,  
`\l__regex_case_changed_char_int`    and is otherwise `\c_max_int`. The `current_char` variable is also used in various other  
phases to hold a character code.

1785 `\int_new:N \l__regex_current_char_int`  
1786 `\int_new:N \l__regex_current_catcode_int`  
1787 `\int_new:N \l__regex_last_char_int`  
1788 `\int_new:N \l__regex_case_changed_char_int`  
*(End definition for \l\_\_regex\_current\_char\_int and others. These variables are documented on page ??.)*

`\l__regex_current_state_int`    For every character in the token list, each of the active states is considered in turn. The  
variable `\l__regex_current_state_int` holds the state of the NFA which is currently  
considered: transitions are then given as shifts relative to the current state.

1789 `\int_new:N \l__regex_current_state_int`  
*(End definition for \l\_\_regex\_current\_state\_int This variable is documented on page ??.)*

`\l__regex_current_submatches_prop`    The submatches for the thread which is currently active are stored in the `current_-`  
`\l__regex_success_submatches_prop`    `submatches` property list variable. This property list is stored by `\__regex_action-`  
`cost:n` into the `\toks` register for the target state of the transition, to be retrieved when  
matching at the next position. When a thread succeeds, this property list is copied to  
`\l__regex_success_submatches_prop`: only the last successful thread will remain there.

1790 `\prop_new:N \l__regex_current_submatches_prop`  
1791 `\prop_new:N \l__regex_success_submatches_prop`  
*(End definition for \l\_\_regex\_current\_submatches\_prop and \l\_\_regex\_success\_submatches\_prop  
These variables are documented on page ??.)*

`\l__regex_step_int`    This integer, always even, is increased every time a character in the query is read, and not  
reset when doing multiple matches. For each  $\langle state \rangle$  in the NFA we store in  $\dimen\langle state \rangle$   
the last step in which this state was encountered. This lets us break infinite loops by not  
visiting the same state twice in the same step. In fact,  $\dimen\langle state \rangle$  is equal `step` when  
we have started performing the operations of  $\toks\langle state \rangle$ , but not finished yet. However,  
once we finish, we set  $\dimen\langle state \rangle$  to `step + 1`. This is needed to track submatches  
properly (see building phase). The `step` is also used to attach each set of submatch  
information to a given iteration (and automatically discard it when it corresponds to a  
past step).

1792 `\int_new:N \l__regex_step_int`  
*(End definition for \l\_\_regex\_step\_int This variable is documented on page ??.)*

`\l__regex_min_active_int` All the currently active states are kept in order of precedence in the `\skip` registers, and  
`\l__regex_max_active_int` the corresponding submatches in the `\toks`. For our purposes, those serve as an array, indexed from `min_active` (inclusive) to `max_active` (excluded). At the start of every step, the whole array is unpacked, so that the space can immediately be reused, and `max_active` is reset to `min_active`, effectively clearing the array.

1793 `\int_new:N \l__regex_min_active_int`

1794 `\int_new:N \l__regex_max_active_int`

*(End definition for `\l__regex_min_active_int` and `\l__regex_max_active_int` These variables are documented on page ??.)*

`\l__regex_every_match_tl` Every time a match is found, this token list is used. For single matching, the token list is empty. For multiple matching, the token list is set to repeat the matching, after performing some operation which depends on the user function. See `\__regex_single-match:` and `\__regex_multi_match:n`.

1795 `\tl_new:N \l__regex_every_match_tl`

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

`\l__regex_fresh_thread_bool` When doing multiple matches, we need to avoid infinite loops where each iteration  
`\l__regex_empty_success_bool` matches the same empty token list. When an empty token list is matched, the next  
`\__regex_if_two_empty_matches:F` successful match of the same empty token list is suppressed. We detect empty matches by setting `\l__regex_fresh_thread_bool` to `true` for threads which directly come from the start of the regex or from the `\K` command, and testing that boolean whenever a thread succeeds. The function `\__regex_if_two_empty_matches:F` is redefined at every match attempt, depending on whether the previous match was empty or not: if it was, then the function must cancel a purported success if it is empty and at the same spot as the previous match; otherwise, we definitely don't have two identical empty matches, so the function is `\use:n`.

1796 `\bool_new:N \l__regex_fresh_thread_bool`

1797 `\bool_new:N \l__regex_empty_success_bool`

1798 `\cs_new_eq:NN \__regex_if_two_empty_matches:F \use:n`

*(End definition for `\l__regex_fresh_thread_bool` and `\l__regex_empty_success_bool` These functions are documented on page ??.)*

`\g__regex_success_bool` The boolean `\l__regex_match_success_bool` is true if the current match attempt was  
`\l__regex_saved_success_bool` successful, and `\g__regex_success_bool` is true if there was at least one successful  
`\l__regex_match_success_bool` match. This is the only global variable in this whole module, but we would need it to be local when matching a control sequence with `\c{...}`. This is done by saving the global variable into `\l__regex_saved_success_bool`, which is local, hence not affected by the changes due to inner regex functions.

1799 `\bool_new:N \g__regex_success_bool`

1800 `\bool_new:N \l__regex_saved_success_bool`

1801 `\bool_new:N \l__regex_match_success_bool`

*(End definition for `\g__regex_success_bool`, `\l__regex_saved_success_bool`, and `\l__regex_match_success_bool` These variables are documented on page ??.)*

## 2.5.2 Matching: framework

`\__regex_match:n` First store the query into `\toks` and `\muskip` registers (see `\__regex_query_set:nnn`). Then initialize the variables that should be set once for each user function (even for multiple matches). Namely, the overall matching is not yet successful; none of the states should be marked as visited (`\dimen` registers), and we start at step 0; we pretend that there was a previous match ending at the start of the query, which was not empty (to avoid smothering an empty match at the start). Once all this is set up, we are ready for the ride. Find the first match.

```

1802 \cs_new_protected:Npn \__regex_match:n #1
1803 {
1804   \trace_push:nnx { regex } { 1 } { __regex_match }
1805   \trace:nnx { regex } { 1 } { analyzing-query-token-list }
1806   \int_zero:N \l__regex_balance_int
1807   \int_set:Nn \l__regex_current_pos_int { \c_two * \l__regex_max_state_int }
1808   \__regex_query_set:nnn { } { -1 } { -2 }
1809   \int_set_eq:NN \l__regex_min_pos_int \l__regex_current_pos_int
1810   \__tl_analysis_map_inline:nn {#1}
1811     { \__regex_query_set:nnn {##1} {##2} {##3} }
1812   \int_set_eq:NN \l__regex_max_pos_int \l__regex_current_pos_int
1813   \__regex_query_set:nnn { } { -1 } { -2 }
1814   \trace:nnx { regex } { 1 } { initializing }
1815   \bool_gset_false:N \g__regex_success_bool
1816   \int_step_inline:nnnn
1817     \l__regex_min_state_int \c_one { \l__regex_max_state_int - \c_one }
1818     { \tex_dimen:D ##1 \c_one sp \scan_stop: }
1819   \int_set_eq:NN \l__regex_min_active_int \l__regex_max_state_int
1820   \int_set_eq:NN \l__regex_step_int \c_zero
1821   \int_set_eq:NN \l__regex_success_pos_int \l__regex_min_pos_int
1822   \int_set:Nn \l__regex_submatch_int
1823     { \c_two * \l__regex_max_state_int }
1824   \bool_set_false:N \l__regex_empty_success_bool
1825   \__regex_match_once:
1826   \trace_pop:nnx { regex } { 1 } { __regex_match }
1827 }

```

(End definition for `\__regex_match:n`)

`\__regex_match_once:` This function finds one match, then does some action defined by the `every_match` token list, which may recursively call `\__regex_match_once:`. First initialize some variables: set the conditional which detects identical empty matches; this match attempt starts at the previous `success_pos`, is not yet successful, and has no submatches yet; clear the array of active threads, and put the starting state 0 in it. We are then almost ready to read our first token in the query, but we actually start one position earlier than the start, and `get` that token, so that the `last_char` will be set properly for word boundaries. Then call `\__regex_match_loop:`, which runs through the query until the end or until a successful match breaks early.

```

1828 \cs_new_protected_nopar:Npn \__regex_match_once:
1829 {

```

```

1830 \if_meaning:w \c_true_bool \l__regex_empty_success_bool
1831 \cs_set_nopar:Npn \__regex_if_two_empty_matches:F
1832 { \int_compare:nNnF \l__regex_start_pos_int = \l__regex_current_pos_int }
1833 \else:
1834 \cs_set_eq:NN \__regex_if_two_empty_matches:F \use:n
1835 \fi:
1836 \int_set_eq:NN \l__regex_start_pos_int \l__regex_success_pos_int
1837 \bool_set_false:N \l__regex_match_success_bool
1838 \prop_clear:N \l__regex_current_submatches_prop
1839 \int_set_eq:NN \l__regex_max_active_int \l__regex_min_active_int
1840 \__regex_store_state:n { \l__regex_min_state_int }
1841 \int_set:Nn \l__regex_current_pos_int
1842 { \l__regex_start_pos_int - \c_one }
1843 \__regex_query_get:
1844 \__regex_match_loop:
1845 \l__regex_every_match_tl
1846 }

```

(End definition for \\_\_regex\_match\_once:)

**\\_\_regex\_single\_match:** For a single match, the overall success is determined by whether the only match attempt  
**\\_\_regex\_multi\_match:n** is a success. When doing multiple matches, the overall matching is successful as soon as  
any match succeeds. Perform the action #1, then find the next match.

```

1847 \cs_new_protected_nopar:Npn \__regex_single_match:
1848 {
1849 \tl_set:Nn \l__regex_every_match_tl
1850 { \bool_gset_eq:NN \g__regex_success_bool \l__regex_match_success_bool }
1851 }
1852 \cs_new_protected:Npn \__regex_multi_match:n #1
1853 {
1854 \tl_set:Nn \l__regex_every_match_tl
1855 {
1856 \if_meaning:w \c_true_bool \l__regex_match_success_bool
1857 \bool_gset_true:N \g__regex_success_bool
1858 #1
1859 \exp_after:wN \__regex_match_once:
1860 \fi:
1861 }
1862 }

```

(End definition for \\_\_regex\_single\_match: and \\_\_regex\_multi\_match:n)

**\\_\_regex\_match\_loop:** At each new position, set some variables and get the new character and category from  
**\\_\_regex\_match\_one\_active:w** the query. Then unpack the array of active threads, and clear it by resetting its length  
(max\_active). This results in a sequence of \\_\_regex\_use\_state\_and\_submatches:nn  
{<state>} {<prop>}, and we consider those states one by one in order. As soon as a thread  
succeeds, exit the step, and, if there are threads to consider at the next position, and  
we have not reached the end of the string, repeat the loop. Otherwise, the last thread  
that succeeded is what \\_\_regex\_match\_once: matches. We explain the **fresh\_thread**  
business when describing \\_\_regex\_action\_wildcard:.



```

1863 \cs_new_protected_nopar:Npn \__regex_match_loop:
1864 {
1865   \tex_advance:D \l__regex_step_int \c_two
1866   \int_incr:N \l__regex_current_pos_int
1867   \int_set_eq:NN \l__regex_last_char_int \l__regex_current_char_int
1868   \int_set_eq:NN \l__regex_case_changed_char_int \c_max_int
1869   \__regex_query_get:
1870   \use:x
1871   {
1872     \int_set_eq:NN \l__regex_max_active_int \l__regex_min_active_int
1873     \exp_after:wN \__regex_match_one_active:w
1874     \int_use:N \l__regex_min_active_int ;
1875   }
1876   \__prg_break_point:
1877   \bool_set_false:N \l__regex_fresh_thread_bool %^^A was arg of break_point:n
1878   \if_int_compare:w \l__regex_max_active_int > \l__regex_min_active_int
1879     \if_int_compare:w \l__regex_current_pos_int < \l__regex_max_pos_int
1880       \exp_after:wN \exp_after:wN \exp_after:wN \__regex_match_loop:
1881     \fi:
1882   \fi:
1883 }
1884 \cs_new:Npn \__regex_match_one_active:w #1;
1885 {
1886   \if_int_compare:w #1 < \l__regex_max_active_int
1887     \__regex_use_state_and_submatches:nn
1888     { \__int_value:w \tex_skip:D #1 }
1889     { \tex_the:D \tex_toks:D #1 }
1890     \exp_after:wN \__regex_match_one_active:w
1891     \int_use:N \__int_eval:w #1 + \c_one \exp_after:wN ;
1892   \fi:
1893 }

```

(End definition for \\_\_regex\_match\_loop: This function is documented on page ??.)

\\_\_regex\_query\_set:nnn The arguments are: tokens that o and x expand to one token of the query, the catcode, and the character code. Store those, and the current brace balance (used later to check for overall brace balance) in a \muskip register and a \toks, then update the balance.

```

1894 \cs_new_protected:Npn \__regex_query_set:nnn #1#2#3
1895 {
1896   \tex_muskip:D \l__regex_current_pos_int
1897   = \etex_glue_tomu:D
1898   #3 sp
1899   plus #2 sp
1900   minus \l__regex_balance_int sp
1901   \scan_stop:
1902   \tex_toks:D \l__regex_current_pos_int {#1}
1903   \int_incr:N \l__regex_current_pos_int
1904   \if_case:w #2 \exp_stop_f:
1905   \or: \int_incr:N \l__regex_balance_int
1906   \or: \int_decr:N \l__regex_balance_int

```

```

1907   \fi:
1908 }
(End definition for \_regex_query_set:nnn)

```

`\_regex_query_get:` Extract the current character and category codes from the `\muskip` register of the current position: those are the main and the stretch components, and we need a conversion to avoid TeX’s “incompatible glue units” error.

```

1909 \cs_new_protected_nopar:Npn \_regex_query_get:
1910 {
1911   \l__regex_current_char_int
1912   = \etex_mutoglua:D \tex_muskip:D \l__regex_current_pos_int
1913   \l__regex_current_catcode_int = \etex_gluestretch:D
1914   \etex_mutoglua:D \tex_muskip:D \l__regex_current_pos_int
1915 }
(End definition for \_regex_query_get:)

```

### 2.5.3 Using states of the nfa

`\_regex_use_state:` Use the current NFA instruction. The state is initially marked as belonging to the current **step**: this allows normal free transition to repeat, but group-repeating transitions won’t. Once we are done exploring all the branches it spawned, the state is marked as **step + 1**: any thread hitting it at that point will be terminated.

```

1916 \cs_new_protected_nopar:Npn \_regex_use_state:
1917 {
1918   <trace>
1919   \trace:nnx { regex } { 2 } { state-\int_use:N \l__regex_current_state_int }
1920   </trace>
1921   \tex_dimen:D \l__regex_current_state_int
1922   = \l__regex_step_int sp \scan_stop:
1923   \tex_the:D \tex_toks:D \l__regex_current_state_int
1924   \tex_dimen:D \l__regex_current_state_int
1925   = \__int_eval:w \l__regex_step_int + \c_one \__int_eval_end: sp \scan_stop:
1926 }
(End definition for \_regex_use_state:)

```

`\_regex_use_state_and_submatches:nn` This function is called as one item in the array of active threads after that array has been unpacked for a new step. Update the `current_state` and `current_submatches` and use the state if it has not yet been encountered at this step.

```

1927 \cs_new_protected:Npn \_regex_use_state_and_submatches:nn #1 #2
1928 {
1929   \int_set:Nn \l__regex_current_state_int {#1}
1930   \if_int_compare:w \tex_dimen:D \l__regex_current_state_int
1931     < \l__regex_step_int
1932     \tl_set:Nn \l__regex_current_submatches_prop {#2}
1933     \exp_after:wN \_regex_use_state:
1934   \fi:
1935   \scan_stop:
1936 }
(End definition for \_regex_use_state_and_submatches:nn)

```

## 2.5.4 Actions when matching

`\_regex_action_start_wildcard:` For an unanchored match, state 0 has a free transition to the next and a costly one to itself, to repeat at the next position. To catch repeated identical empty matches, we need to know if a successful thread corresponds to an empty match. The instruction resetting `\l\_regex_fresh_thread_bool` may be skipped by a successful thread, hence we had to add it to `\_regex_match_loop:` too.

```

1937 \cs_new_protected_nopar:Npn \_regex_action_start_wildcard:
1938 {
1939   \bool_set_true:N \l\_regex_fresh_thread_bool
1940   \_regex_action_free:n {1}
1941   \bool_set_false:N \l\_regex_fresh_thread_bool
1942   \_regex_action_cost:n {0}
1943 }

```

*(End definition for \\_regex\_action\_start\_wildcard:)*

`\_regex_action_free:n` These functions copy a thread after checking that the NFA state has not already been used at this position. If not, store submatches in the new state, and insert the instructions for that state in the input stream. Then restore the old value of `\l\_regex_current_state_int` and of the current submatches. The two types of free transitions differ by how they test that the state has not been encountered yet: the **group** version is stricter, and will not use a state if it was used earlier in the current thread, hence forcefully breaking the loop, while the “normal” version will revisit a state when within the thread itself.

```

1944 \cs_new_protected_nopar:Npn \_regex_action_free:n
1945 { \_regex_action_free_aux:nn { > \l\_regex_step_int \else: } }
1946 \cs_new_protected_nopar:Npn \_regex_action_free_group:n
1947 { \_regex_action_free_aux:nn { < \l\_regex_step_int } }
1948 \cs_new_protected:Npn \_regex_action_free_aux:nn #1#2
1949 {
1950   \use:x
1951   {
1952     \int_add:Nn \l\_regex_current_state_int {#2}
1953     \exp_not:n
1954     {
1955       \if_int_compare:w \tex_dimen:D \l\_regex_current_state_int #1
1956       \exp_after:wN \_regex_use_state:
1957       \fi:
1958     }
1959     \int_set:Nn \l\_regex_current_state_int
1960     { \int_use:N \l\_regex_current_state_int }
1961     \tl_set:Nn \exp_not:N \l\_regex_current_submatches_prop
1962     { \exp_not:o \l\_regex_current_submatches_prop }
1963   }
1964 }

```

*(End definition for \\_regex\_action\_free:n and \\_regex\_action\_free\_group:n These functions are documented on page ??.)*

`\__regex_action_cost:n` A transition which consumes the current character and shifts the state by #1. The resulting state is stored in the `\skip` array for use at the next position, and we also store the current submatches.

```

1965 \cs_new_protected:Npn \__regex_action_cost:n #1
1966 {
1967     \exp_args:No \__regex_store_state:n
1968     { \int_use:N \__int_eval:w \l__regex_current_state_int + #1 }
1969 }
(End definition for \__regex_action_cost:n)

```

`\__regex_store_state:n` Put the given state in the array of `\skip` registers (converted to a dimension in scaled points), and increment the length of the array. Then store the current submatch in the `\__regex_store_submatches:`. This is done by increasing the pointer `\l__regex_max_active_int`, and converting the integer to a dimension (suitable for a `\skip` assignment) in scaled points.

```

1970 \cs_new_protected:Npn \__regex_store_state:n #1
1971 {
1972     \__regex_store_submatches:
1973     \tex_skip:D \l__regex_max_active_int = #1 sp \scan_stop:
1974     \int_incr:N \l__regex_max_active_int
1975 }
1976 \cs_new_protected_nopar:Npn \__regex_store_submatches:
1977 {
1978     \tex_toks:D \l__regex_max_active_int \exp_after:wN
1979     { \l__regex_current_submatches_prop }
1980 }
(End definition for \__regex_store_state:n This function is documented on page ??.)

```

`\__regex_disable_submatches:` Some user functions don't require tracking submatches. We get a performance improvement by simply defining the relevant functions to remove their argument and do nothing with it.

```

1981 \cs_new_protected_nopar:Npn \__regex_disable_submatches:
1982 {
1983     \cs_set_protected_nopar:Npn \__regex_store_submatches: { }
1984     \cs_set_protected:Npn \__regex_action_submatch:n ##1 { }
1985 }
(End definition for \__regex_disable_submatches:)

```

`\__regex_action_submatch:n` Update the current submatches with the information from the current position. Maybe a bottleneck.

```

1986 \cs_new_protected:Npn \__regex_action_submatch:n #1
1987 {
1988     \prop_put:Nno \l__regex_current_submatches_prop {#1}
1989     { \int_use:N \l__regex_current_pos_int }
1990 }
(End definition for \__regex_action_submatch:n)

```

**\\_\_regex\_action\_success:** There is a successful match when an execution path reaches the last state in the NFA, unless this marks a second identical empty match. Then mark that there was a successful match; it is empty if it is “fresh”; and we store the current position and submatches. The current step is then interrupted with **\\_\_prg\_break:**, and only paths with higher precedence are pursued further. The values stored here may be overwritten by a later success of a path with higher precedence.

```

1991 \cs_new_protected_nopar:Npn \__regex_action_success:
1992 {
1993   \__regex_if_two_empty_matches:F
1994   {
1995     \bool_set_true:N \l__regex_match_success_bool
1996     \bool_set_eq:NN \l__regex_empty_success_bool
1997     \l__regex_fresh_thread_bool
1998     \int_set_eq:NN \l__regex_success_pos_int \l__regex_current_pos_int
1999     \prop_set_eq:NN \l__regex_success_submatches_prop
2000     \l__regex_current_submatches_prop
2001     \__prg_break:
2002   }
2003 }
(End definition for \__regex_action_success:)

```

## 2.6 Replacement

### 2.6.1 Variables and helpers used in replacement

**\l\_\_regex\_replacement\_csnames\_int** The behaviour of closing braces inside a replacement text depends on whether a sequences **\c{}** or **\u{}** has been encountered. The number of “open” such sequences that should be closed by **}** is stored in **\l\_\_regex\_replacement\_csnames\_int**, and decreased by 1 by each **}**.

```

2004 \int_new:N \l__regex_replacement_csnames_int
(End definition for \l__regex_replacement_csnames_int This variable is documented on page ??.)

```

**\l\_\_regex\_balance\_tl** This token list holds the replacement text for **\\_\_regex\_replacement\_balance\_one\_match:n** while it is being built incrementally.

```

2005 \tl_new:N \l__regex_balance_tl
(End definition for \l__regex_balance_tl This variable is documented on page ??.)

```

**\\_\_regex\_replacement\_balance\_one\_match:n** This expects as an argument the first index of a range of **\skip** registers which hold the submatch information for a given match. It can be used within an integer expression to obtain the brace balance incurred by performing the replacement on that match. This combines the braces lost by removing the match, braces added by all the submatches appearing in the replacement, and braces appearing explicitly in the replacement. Even though it is always redefined before use, we initialize it as for an empty replacement. An important property is that concatenating several calls to that function must result in a valid integer expression (hence a leading **+** in the actual definition).

```

2006 \cs_new:Npn \__regex_replacement_balance_one_match:n #1
2007 { - \__regex_submatch_balance:n {#1} }

```

(End definition for `\_regex\_replacement\_balance\_one\_match:n`)

`\_regex\_replacement\_do\_one\_match:n` The input is the same as `\_regex\_replacement\_balance\_one\_match:n`. This function is redefined to expand to the part of the token list from the end of the previous match to a given match, followed by the replacement text. Hence concatenating the result of this function with all possible arguments (one call for each match), as well as the range from the end of the last match to the end of the string, will produce the fully replaced token list. The initialization does not matter, but we set it as for an empty replacement.

```

2008 \cs_new:Npn \_regex\_replacement\_do\_one\_match:n #1
2009 {
2010   \_regex\_query\_range:nn
2011   { \etex\_glueshrink:D \tex\_skip:D #1 }
2012   { \tex\_skip:D #1 }
2013 }

```

(End definition for `\_regex\_replacement\_do\_one\_match:n`)

`\_regex\_replacement\_exp\_not:N` This function lets us navigate around the fact that the primitive `\exp\_not:n` requires a braced argument. As far as I can tell, it is only needed if the user tries to include in the replacement text a control sequence set equal to a macro parameter character, such as `\c\_parameter\_token`. Indeed, within an x-expanding assignment, `\exp\_not:N #` behaves as a single `#`, whereas `\exp\_not:n {#}` behaves as a doubled `##`.

```

2014 \cs_new:Npn \_regex\_replacement\_exp\_not:N #1 { \exp\_not:n {#1} }

```

(End definition for `\_regex\_replacement\_exp\_not:N`)

## 2.6.2 Query and brace balance

`\_regex\_query\_range:nn` When it is time to extract submatches from the token list, the various tokens are stored in `\toks` registers numbered from `\l\_regex\_min\_pos\_int` inclusive to `\l\_regex\_max\_pos\_int` exclusive. The function `\_regex\_query\_range:nn {<min>} {<max>}` unpacks registers from the position `<min>` to the position `<max> - 1` included. Once this is expanded, a second x-expansion will result in the actual tokens from the query. That second expansion is only done by user functions at the very end of their operation, after checking (and correcting) the brace balance first.

`\_regex\_query\_range\_loop:ww`

```

2015 \cs_new:Npn \_regex\_query\_range:nn #1#2
2016 {
2017   \exp\_after:wN \_regex\_query\_range\_loop:ww
2018   \int\_use:N \_int\_eval:w #1 \exp\_after:wN ;
2019   \int\_use:N \_int\_eval:w #2 ;
2020   \_prg\_break\_point:
2021 }
2022 \cs_new:Npn \_regex\_query\_range\_loop:ww #1 ; #2 ;
2023 {
2024   \if\_int\_compare:w #1 < #2 \exp\_stop\_f:
2025   \else:
2026     \exp\_after:wN \_prg\_break:
2027   \fi:
2028   \tex\_the:D \tex\_toks:D #1 \exp\_stop\_f:

```

```

2029 \exp_after:wN \__regex_query_range_loop:ww
2030 \int_use:N \__int_eval:w #1 + \c_one ; #2 ;
2031 }

```

(End definition for \\_\_regex\_query\_range:nn This function is documented on page ??.)

**\\_\_regex\_query\_submatch:n** When this function is called, \skipi holds the start and end positions for the *i*-th overall submatch as its main and stretch components. In the case of repeated matches, submatches from all the matches are put one after the other in blocks of \l\_\_regex\_capturing\_group\_int \skip registers.

```

2032 \cs_new:Npn \__regex_query_submatch:n #1
2033 {
2034   \__regex_query_range:nn
2035   { \tex_skip:D \__int_eval:w #1 }
2036   { \etex_gluestretch:D \tex_skip:D \__int_eval:w #1 }
2037 }

```

(End definition for \\_\_regex\_query\_submatch:n)

**\\_\_regex\_submatch\_balance:n** Every user function must result in a balanced token list (unbalanced token lists cannot be stored by TeX). When we unpacked the query, we kept track of the brace balance as the shrink component of \muskip registers, hence the contribution from a given range is the difference between the shrink components of \muskip⟨*max pos*⟩ and \muskip⟨*min pos*⟩. For the *i*-th submatch, the end-points of the range are the main and stretch components of \skipi. The trailing \scan\_stop: is gobbled by \etex\_muexpr:D, and the whole expression can be cast safely to an integer (no trailing expansion).

```

2038 \cs_new_protected:Npn \__regex_submatch_balance:n #1
2039 {
2040   \etex_glueshrink:D \etex_mutogluue:D \etex_muexpr:D
2041   \tex_muskip:D \etex_gluestretch:D \tex_skip:D #1
2042   - \tex_muskip:D \tex_skip:D #1
2043   \scan_stop:
2044 }

```

(End definition for \\_\_regex\_submatch\_balance:n This function is documented on page ??.)

### 2.6.3 Framework

**\\_\_regex\_replacement:n** The replacement text is built incrementally by abusing \toks within a group (see l3tl-build). We keep track in \l\_\_regex\_balance\_int of the balance of explicit begin- and end-group tokens and \l\_\_regex\_balance\_tl will consist of some code to compute the brace balance from submatches (see its description). Detect unescaped right braces, and escaped characters, with trailing \prg\_do\_nothing: because some of the later function look-ahead. Once the whole replacement text has been parsed, make sure that there is no open cname. Finally, define the balance\_one\_match and do\_one\_match functions.

```

2045 \cs_new_protected:Npn \__regex_replacement:n #1
2046 {
2047   <trace> \trace_push:nnn { regex } { 1 } { __regex_replacement:n }
2048   \__tl_build:Nw \l__regex_internal_a_tl
2049   \int_zero:N \l__regex_balance_int

```

```

2050 \tl_clear:N \l__regex_balance_tl
2051 \__regex_escape_use:nnnn
2052 {
2053   \if_charcode:w \c_rbrace_str ##1
2054     \__regex_replacement_rbrace:N \else: \__tl_build_one:n \fi: ##1
2055   }
2056   { \__regex_replacement_escaped:N ##1 }
2057   { \__tl_build_one:n ##1 }
2058   {#1}
2059 \prg_do_nothing: \prg_do_nothing:
2060 \if_int_compare:w \l__regex_replacement_csnames_int > \c_zero
2061   \__msg_kernel_error:nnx { regex } { replacement-missing-rbrace }
2062   { \int_use:N \l__regex_replacement_csnames_int }
2063   \__tl_build_one:x
2064   { \prg_replicate:nn \l__regex_replacement_csnames_int \cs_end: }
2065 \fi:
2066 \cs_gset:Npx \__regex_replacement_balance_one_match:n ##1
2067 {
2068   + \int_use:N \l__regex_balance_int
2069   \l__regex_balance_tl
2070   - \__regex_submatch_balance:n {##1}
2071 }
2072 \__tl_build_end:
2073 \exp_args:No \__regex_replacement_aux:n \l__regex_internal_a_tl
2074 <trace> \trace_pop:nnn { regex } { 1 } { __regex_replacement:n }
2075 }
2076 \cs_new_protected:Npn \__regex_replacement_aux:n #1
2077 {
2078   \cs_set:Npn \__regex_replacement_do_one_match:n ##1
2079   {
2080     \__regex_query_range:nn
2081     { \etex_glueshrink:D \tex_skip:D ##1 }
2082     { \tex_skip:D ##1 }
2083     #1
2084   }
2085 }

```

(End definition for \\_\_regex\_replacement:n This function is documented on page ??.)

\\_\_regex\_replacement\_escaped:N As in parsing a regular expression, we use an auxiliary built from #1 if defined. Otherwise, check for escaped digits (standing from submatches from 0 to 9): anything else is a raw character.

```

2086 \cs_new_protected:Npn \__regex_replacement_escaped:N #1
2087 {
2088   \cs_if_exist_use:cF { __regex_replacement_#1:w }
2089   {
2090     \if_int_compare:w \c_one < 1#1 \exp_stop_f:
2091       \__regex_replacement_put_submatch:n {#1}
2092     \else:
2093       \__tl_build_one:n #1

```



```

2094         \fi:
2095     }
2096 }
(End definition for \_regex_replacement_escaped:N)

```

## 2.6.4 Submatches

`\_regex_replacement_put_submatch:n` Insert a submatch in the replacement text. This is dropped if the submatch number is larger than the number of capturing groups. Unless the submatch appears inside a `\c{...}` or `\u{...}` construction, it must be taken into account in the brace balance. Here, `##1` will receive a pointer to the 0-th submatch for a given match. We cannot use `\int_eval:n` because it is expandable, and would be expanded too early (short of adding `\exp_not:N`, making the code messy again).

```

2097 \cs_new_protected:Npn \_regex_replacement_put_submatch:n #1
2098 {
2099     \if_int_compare:w #1 < \l__regex_capturing_group_int
2100         \__tl_build_one:n { \_regex_query_submatch:n { #1 + ##1 } }
2101         \if_int_compare:w \l__regex_replacement_csnames_int = \c_zero
2102             \tl_put_right:Nn \l__regex_balance_tl
2103                 { + \_regex_submatch_balance:n { \__int_eval:w #1+##1 \__int_eval_end: } }
2104         \fi:
2105     \fi:
2106 }
(End definition for \_regex_replacement_put_submatch:n)

```

`\_regex_replacement_g:w` An ugly method to grab digits for the `\g` escape sequence. At the end of the run of digits, `\_regex_replacement_g_digits:NN` check that it ends with a right brace.

```

2107 \cs_new_protected:Npn \_regex_replacement_g:w #1#2
2108 {
2109     \str_if_eq_x:nnTF { #1#2 } { \__tl_build_one:n \c_lbrace_str }
2110     {
2111         \int_zero:N \l__regex_internal_a_int
2112         \_regex_replacement_g_digits:NN
2113     }
2114     { \_regex_replacement_error:NNN g #1 #2 }
2115 }
2116 \cs_new_protected:Npn \_regex_replacement_g_digits:NN #1#2
2117 {
2118     \token_if_eq_meaning:NNTF #1 \__tl_build_one:n
2119     {
2120         \if_int_compare:w \c_one < 1#2 \exp_stop_f:
2121             \int_set:Nn \l__regex_internal_a_int
2122                 { \c_ten * \l__regex_internal_a_int + #2 }
2123             \exp_after:wN \use_i:nnn
2124             \exp_after:wN \_regex_replacement_g_digits:NN
2125         \else:
2126             \exp_after:wN \_regex_replacement_error:NNN
2127             \exp_after:wN g

```

```

2128     \fi:
2129   }
2130   {
2131     \if_meaning:w \__regex_replacement_rbrace:N #1
2132     \exp_args:No \__regex_replacement_put_submatch:n
2133       { \int_use:N \l__regex_internal_a_int }
2134     \exp_after:wN \use_none:nn
2135     \else:
2136       \exp_after:wN \__regex_replacement_error:NNN
2137       \exp_after:wN g
2138     \fi:
2139   }
2140   #1 #2
2141 }

```

(End definition for \\_\_regex\_replacement\_g:w and \\_\_regex\_replacement\_g\_digits:NN)

### 2.6.5 Csnames in replacement

\\_\_regex\_replacement\_c:w \\_\_regex\_replacement\_c\_{:w

\c can be followed by a left brace, or by a letter for which we have defined a way to produce that category of characters. The appropriate definitions for catcodes are introduced later. For control sequences, if we are within a control sequence, convert the token list to a string, otherwise simply prevent expansion, with a weird cross-over between \exp\_not:n and \exp\_not:N (see this helper's description for an explanation).

```

2142 \cs_new_protected:Npn \__regex_replacement_c:w #1#2
2143 {
2144   \token_if_eq_meaning:NNTF #1 \__tl_build_one:n
2145   {
2146     \cs_if_exist_use:cF { \__regex_replacement_c_#2:w }
2147     { \__regex_replacement_error:NNN c #1#2 }
2148   }
2149   { \__regex_replacement_error:NNN c #1#2 }
2150 }
2151 \cs_new_protected_nopar:cpn { \__regex_replacement_c_ \c_lbrace_str :w }
2152 {
2153   \if_case:w \l__regex_replacement_csnames_int
2154     \__tl_build_one:n
2155     { \exp_not:n { \exp_after:wN \__regex_replacement_exp_not:N \cs:w } }
2156   \else:
2157     \__tl_build_one:n { \exp_not:n { \exp_after:wN \tl_to_str:N \cs:w } }
2158   \fi:
2159   \int_incr:N \l__regex_replacement_csnames_int
2160 }

```

(End definition for \\_\_regex\_replacement\_c:w This function is documented on page ??.)

\\_\_regex\_replacement\_u:w

Check that \u is followed by a left brace. If so, start a control sequence with \cs:w, which is then unpacked either with \exp\_not:V or \tl\_to\_str:V depending on the current context.

```

2161 \cs_new_protected:Npn \__regex_replacement_u:w #1#2

```

```

2162 {
2163   \str_if_eq_x:nnTF { #1#2 } { \__tl_build_one:n \c_lbrace_str }
2164   {
2165     \if_case:w \l__regex_replacement_csnames_int
2166     \__tl_build_one:n { \exp_not:n { \exp_after:wN \exp_not:V \cs:w } }
2167     \else:
2168     \__tl_build_one:n { \exp_not:n { \exp_after:wN \tl_to_str:V \cs:w } }
2169     \fi:
2170     \int_incr:N \l__regex_replacement_csnames_int
2171   }
2172   { \__regex_replacement_error:NNN u #1#2 }
2173 }

```

(End definition for \\_\_regex\_replacement\_u:w)

\\_\_regex\_replacement\_rbrace:N Within a \c{...} or \u{...} construction, end the control sequence, and decrease the brace count. Otherwise, this is a raw right brace.

```

2174 \cs_new_protected:Npn \__regex_replacement_rbrace:N #1
2175 {
2176   \if_int_compare:w \l__regex_replacement_csnames_int > \c_zero
2177   \__tl_build_one:n \cs_end:
2178   \int_decr:N \l__regex_replacement_csnames_int
2179   \else:
2180   \__tl_build_one:n #1
2181   \fi:
2182 }

```

(End definition for \\_\_regex\_replacement\_rbrace:N)

## 2.6.6 Characters in replacement

We will need to change the category code of the null character many times, hence work in a group. The catcode-specific macros below are defined in alphabetical order; if you are trying to understand the code, start from the end of the alphabet as those categories are simpler than active or begin-group.

```

2183 \group_begin:

```

\\_\_regex\_replacement\_char:nnN The only way to produce an arbitrary character-catcode pair is to use the \lowercase or \uppercase primitives. This is a wrapper for our purposes. The first argument is the null character with various catcodes. The second and third arguments are grabbed from the input stream: #3 is the character whose character code to reproduce.

```

2184 \cs_new_protected:Npn \__regex_replacement_char:nnN #1#2#3
2185 {
2186   \if_meaning:w \prg_do_nothing: #3
2187   \__msg_kernel_error:nn { regex } { replacement-catcode-end }
2188   \else:
2189   \tex_lccode:D \c_zero = '#3 \scan_stop:
2190   \tl_to_lowercase:n { \__tl_build_one:n {#1} }
2191   \fi:
2192 }

```

(End definition for `\_regex_replacement_char:nNN`)

`\_regex_replacement_c_A:w` For an active character, expansion must be avoided, twice because we later do two x-expansions, to unpack `\toks` for the query, and to expand their contents to tokens of the query.

```
2193 \char_set_catcode_active:N \^^@
2194 \cs_new_protected_nopar:Npn \_regex_replacement_c_A:w
2195   { \_regex_replacement_char:nNN { \exp_not:n { \exp_not:N \^^@ } } }
```

(End definition for `\_regex_replacement_c_A:w`)

`\_regex_replacement_c_B:w` An explicit begin-group token increases the balance, unless within a `\c{...}` or `\u{...}` construction. Add the desired begin-group character, using the standard `\if_false:` trick. We eventually x-expand twice. The first time must yield a balanced token list, and the second one gives the bare begin-group token. The `\exp_after:wN` is not strictly needed, but is more consistent with l3tl-analysis.

```
2196 \char_set_catcode_group_begin:N \^^@
2197 \cs_new_protected_nopar:Npn \_regex_replacement_c_B:w
2198   {
2199     \if_int_compare:w \l__regex_replacement_csnames_int = \c_zero
2200       \int_incr:N \l__regex_balance_int
2201     \fi:
2202     \_regex_replacement_char:nNN
2203     { \exp_not:n { \exp_after:wN \^^@ \if_false: } \fi: } }
2204   }
```

(End definition for `\_regex_replacement_c_B:w`)

`\_regex_replacement_c_C:w` This is not quite catcode-related: when the user requests a character with category “control sequence”, the one-character control symbol is returned. As for the active character, we prepare for two x-expansions.

```
2205 \cs_new_protected:Npn \_regex_replacement_c_C:w #1#2
2206   { \_tl_build_one:n { \exp_not:N \exp_not:N \exp_not:c {#2} } }
```

(End definition for `\_regex_replacement_c_C:w`)

`\_regex_replacement_c_D:w` Subscripts fit the mould: `\lowercase` the null byte with the correct category.

```
2207 \char_set_catcode_math_subscript:N \^^@
2208 \cs_new_protected_nopar:Npn \_regex_replacement_c_D:w
2209   { \_regex_replacement_char:nNN { \^^@ } }
```

(End definition for `\_regex_replacement_c_D:w`)

`\_regex_replacement_c_E:w` Similar to the begin-group case, the second x-expansion produces the bare end-group token.

```
2210 \char_set_catcode_group_end:N \^^@
2211 \cs_new_protected_nopar:Npn \_regex_replacement_c_E:w
2212   {
2213     \if_int_compare:w \l__regex_replacement_csnames_int = \c_zero
2214       \int_decr:N \l__regex_balance_int
2215     \fi:
2216     \_regex_replacement_char:nNN
```

```

2217         { \exp_not:n { \if_false: { \fi: ^^@ } }
2218     }

```

(End definition for `\_regex_replacement_c_E:w`)

`\_regex_replacement_c_L:w` Simply `\lowercase` a letter null byte to produce an arbitrary letter.

```

2219     \char_set_catcode_letter:N ^^@
2220     \cs_new_protected_nopar:Npn \_regex_replacement_c_L:w
2221         { \_regex_replacement_char:nNN { ^^@ } }

```

(End definition for `\_regex_replacement_c_L:w`)

`\_regex_replacement_c_M:w` No surprise here, we lowercase the null math toggle.

```

2222     \char_set_catcode_math_toggle:N ^^@
2223     \cs_new_protected_nopar:Npn \_regex_replacement_c_M:w
2224         { \_regex_replacement_char:nNN { ^^@ } }

```

(End definition for `\_regex_replacement_c_M:w`)

`\_regex_replacement_c_O:w` Lowercase an other null byte.

```

2225     \char_set_catcode_other:N ^^@
2226     \cs_new_protected_nopar:Npn \_regex_replacement_c_O:w
2227         { \_regex_replacement_char:nNN { ^^@ } }

```

(End definition for `\_regex_replacement_c_O:w`)

`\_regex_replacement_c_P:w` For macro parameters, expansion is a tricky issue. We need to prepare for two x-expansions and passing through various macro definitions. Note that we cannot replace one `\exp_not:n` by doubling the macro parameter characters because this would misbehave if a mischievous user asks for `\c{\cP\#}`, since that macro parameter character would be doubled.

```

2228     \char_set_catcode_parameter:N ^^@
2229     \cs_new_protected_nopar:Npn \_regex_replacement_c_P:w
2230         {
2231             \_regex_replacement_char:nNN
2232             { \exp_not:n { \exp_not:n { ^^@^^@^^@^^@ } } }
2233         }

```

(End definition for `\_regex_replacement_c_P:w`)

`\_regex_replacement_c_S:w` Spaces are normalized on input by `TeX` to have character code 32. It is in fact impossible to get a token with character code 0 and category code 10. Hence we use 32 instead of 0 as our base character.

```

2234     \cs_new_protected:Npn \_regex_replacement_c_S:w #1#2
2235     {
2236         \if_meaning:w \prg_do_nothing: #2
2237         \_msg_kernel_error:nn { regex } { replacement-catcode-end }
2238     \else:
2239         \if_int_compare:w '#2 = \c_zero
2240         \_msg_kernel_error:nn { regex } { replacement-null-space }
2241     \fi:
2242     \tex_lccode:D 32 = '#2 \scan_stop:
2243     \tl_to_lowercase:n { \_tl_build_one:n {~} }

```

```

2244         \fi:
2245     }
(End definition for \_regex_replacement_c_S:w)

```

\\_regex\_replacement\_c\_T:w No surprise for alignment tabs here. Those are surrounded by the appropriate braces whenever necessary, hence they don't cause trouble in alignment settings.

```

2246     \char_set_catcode_alignment:N \^^@
2247     \cs_new_protected_nopar:Npn \_regex_replacement_c_T:w
2248         { \_regex_replacement_char:nnn { ^^@ } }
(End definition for \_regex_replacement_c_T:w)

```

\\_regex\_replacement\_c\_U:w Simple call to \\_regex\_replacement\_char:nnn which lowercases the math superscript ^^@.

```

2249     \char_set_catcode_math_superscript:N \^^@
2250     \cs_new_protected_nopar:Npn \_regex_replacement_c_U:w
2251         { \_regex_replacement_char:nnn { ^^@ } }
(End definition for \_regex_replacement_c_U:w)
    Restore the catcode of the null byte.
2252 \group_end:

```

### 2.6.7 An error

\\_regex\_replacement\_error:nnn Simple error reporting by calling one of the messages replacement-c, replacement-g, or replacement-u.

```

2253 \cs_new_protected:Npn \_regex_replacement_error:nnn #1#2#3
2254     {
2255         \_msg_kernel_error:nnx { regex } { replacement-#1 } {#3}
2256         #2 #3
2257     }
(End definition for \_regex_replacement_error:nnn)

```

## 2.7 User functions

**\regex\_new:N** Before being assigned a sensible value, a regex variable matches nothing.

```

2258 \cs_new_protected:Npn \regex_new:N #1
2259     { \cs_new_eq:NN #1 \c__regex_no_match_regex }
(End definition for \regex_new:N This function is documented on page 6.)

```

**\regex\_set:Nn** Compile, then store the result in the user variable with the appropriate assignment function.  
**\regex\_gset:Nn**  
**\regex\_const:Nn**

```

2260 \cs_new_protected_nopar:Npn \regex_set:Nn #1#2
2261     {
2262         \_regex_compile:n {#2}
2263         \tl_set_eq:NN #1 \l__regex_internal_regex
2264     }
2265 \cs_new_protected_nopar:Npn \regex_gset:Nn #1#2
2266     {

```

```

2267   \__regex_compile:n {#2}
2268   \tl_gset_eq:NN #1 \l__regex_internal_regex
2269   }
2270 \cs_new_protected_nopar:Npn \regex_const:Nn #1#2
2271 {
2272   \__regex_compile:n {#2}
2273   \tl_const:Nx #1 { \exp_not:o \l__regex_internal_regex }
2274 }

```

(End definition for `\regex_set:Nn`, `\regex_gset:Nn`, and `\regex_const:Nn` These functions are documented on page 7.)

`\regex_show:N` User functions: the `n` variant requires compilation first. Then show the variable with some appropriate text. The auxiliary `\__regex_show:Nx` is defined in a different section.

`\regex_show:n`

```

2275 \cs_new_protected:Npn \regex_show:n #1
2276 {
2277   \__regex_compile:n {#1}
2278   \__regex_show:Nx \l__regex_internal_regex
2279   { { \tl_to_str:n {#1} } }
2280 }
2281 \cs_new_protected:Npn \regex_show:N #1
2282 { \__regex_show:Nx #1 { variable~\token_to_str:N #1 } }

```

(End definition for `\regex_show:N` and `\regex_show:n` These functions are documented on page 7.)

`\regex_match:nnTF`

`\regex_match:NnTF`

Those conditionals are based on a common auxiliary defined later. Its first argument builds the NFA corresponding to the regex, and the second argument is the query token list. Once we have performed the match, convert the resulting boolean to `\prg_return_true:` or `false`.

```

2283 \prg_new_protected_conditional:Npnn \regex_match:nn #1#2 { T , F , TF }
2284 {
2285   \__regex_if_match:nn { \__regex_build:n {#1} } {#2}
2286   \__regex_return:
2287 }
2288 \prg_new_protected_conditional:Npnn \regex_match:Nn #1#2 { T , F , TF }
2289 {
2290   \__regex_if_match:nn { \__regex_build:N #1 } {#2}
2291   \__regex_return:
2292 }

```

(End definition for `\regex_match:nnTF` and `\regex_match:NnTF` These functions are documented on page ??.)

`\regex_count:nnN`

`\regex_count:NnN`

Again, use an auxiliary whose first argument builds the NFA.

```

2293 \cs_new_protected:Npn \regex_count:nnN #1
2294 { \__regex_count:nnN { \__regex_build:n {#1} } }
2295 \cs_new_protected:Npn \regex_count:NnN #1
2296 { \__regex_count:nnN { \__regex_build:N #1 } }

```

(End definition for `\regex_count:nnN` and `\regex_count:NnN` These functions are documented on page ??.)

`\regex_extract_once:nnN` We define here 40 user functions, following a common pattern in terms of `:nnN` auxiliaries, defined in the coming subsections. The auxiliary is handed `\__regex_build:n` or `\__-regex_build:N` with the appropriate regex argument, then all other necessary arguments (replacement text, token list, *etc.* The conditionals call `\__regex_return:` to return either true or false once matching has been performed.

```

2297 \cs_set_protected:Npn \__regex_tmp:w #1#2#3
2298 {
2299   \cs_new_protected:Npn #2 ##1 { #1 { \__regex_build:n {##1} } }
2300   \cs_new_protected:Npn #3 ##1 { #1 { \__regex_build:N ##1 } }
2301   \prg_new_protected_conditional:Npnn #2 ##1##2##3 { T , F , TF }
2302     { #1 { \__regex_build:n {##1} } {##2} ##3 \__regex_return: }
2303   \prg_new_protected_conditional:Npnn #3 ##1##2##3 { T , F , TF }
2304     { #1 { \__regex_build:N ##1 } {##2} ##3 \__regex_return: }
2305 }
2306 \__regex_tmp:w \__regex_extract_once:nnN
2307 \regex_extract_once:nnN \regex_extract_once:NnN
2308 \__regex_tmp:w \__regex_extract_all:nnN
2309 \regex_extract_all:nnN \regex_extract_all:NnN
2310 \__regex_tmp:w \__regex_replace_once:nnN
2311 \regex_replace_once:nnN \regex_replace_once:NnN
2312 \__regex_tmp:w \__regex_replace_all:nnN
2313 \regex_replace_all:nnN \regex_replace_all:NnN
2314 \__regex_tmp:w \__regex_split:nnN \regex_split:nnN \regex_split:NnN
  
```

*(End definition for \regex\_extract\_once:nnN and others. These functions are documented on page ??.)*

### 2.7.1 Variables and helpers for user functions

`\l__regex_match_count_int` The number of matches found so far is stored in `\l__regex_match_count_int`. This is only used in the `\regex_count:nnN` functions.

```

2315 \int_new:N \l__regex_match_count_int
  
```

*(End definition for \l\_\_regex\_match\_count\_int This variable is documented on page ??.)*

`__regex_begin` Those flags are raised to indicate extra begin-group or end-group tokens when extracting submatches.

```

2316 \flag_new:n { __regex_begin }
2317 \flag_new:n { __regex_end }
  
```

*(End definition for \_\_regex\_begin and \_\_regex\_end These variables are documented on page ??.)*

`\l__regex_submatch_int` The end-points of each submatch are stored as main and stretch components of `\skip⟨submatch⟩`, where `⟨submatch⟩` ranges from `\l__regex_max_state_int` (inclusive) to `\l__regex_submatch_int` (exclusive). Each successful match comes with a 0-th submatch (the full match), and one match for each capturing group: submatches corresponding to the last successful match are labelled starting at `zeroth_submatch`. Additionally, the shrink component of this 0-th submatch is the position at which that match attempt started: this is used for splitting and replacements.

```

2318 \int_new:N \l__regex_submatch_int
2319 \int_new:N \l__regex_zeroth_submatch_int
  
```



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

`\__regex_return:` This function triggers either `\prg_return_false:` or `\prg_return_true:` as appropriate to whether a match was found or not. It is used by all user conditionals.

```

2320 \cs_new_protected_nopar:Npn \__regex_return:
2321 {
2322   \if_meaning:w \c_true_bool \g__regex_success_bool
2323     \prg_return_true:
2324   \else:
2325     \prg_return_false:
2326   \fi:
2327 }
```

(End definition for `\__regex_return:`)

## 2.7.2 Matching

`\__regex_if_match:nn` We don't track submatches, and stop after a single match. Build the NFA with #1, and perform the match on the query #2.

```

2328 \cs_new_protected:Npn \__regex_if_match:nn #1#2
2329 {
2330   \group_begin:
2331     \__regex_disable_submatches:
2332     \__regex_single_match:
2333     #1
2334     \__regex_match:n {#2}
2335   \group_end:
2336 }
```

(End definition for `\__regex_if_match:nn`)

`\__regex_count:nnN` Again, we don't care about submatches. Instead of aborting after the first "longest match" is found, we search for multiple matches, incrementing `\l__regex_match_count_int` every time to record the number of matches. Build the NFA and match. At the end, store the result in the user's variable.

```

2337 \cs_new_protected:Npn \__regex_count:nnN #1#2#3
2338 {
2339   \group_begin:
2340     \__regex_disable_submatches:
2341     \int_zero:N \l__regex_match_count_int
2342     \__regex_multi_match:n { \int_incr:N \l__regex_match_count_int }
2343     #1
2344     \__regex_match:n {#2}
2345     \exp_args:NNNo
2346   \group_end:
2347   \int_set:Nn #3 { \int_use:N \l__regex_match_count_int }
2348 }
```

(End definition for `\__regex_count:nnN`)

### 2.7.3 Extracting submatches

`\__regex_extract_once:nnN` Match once or multiple times. After each match (or after the only match), extract the submatches using `\__regex_extract:.` At the end, store the sequence containing all the submatches into the user variable #3 after closing the group.

```

2349 \cs_new_protected:Npn \__regex_extract_once:nnN #1#2#3
2350 {
2351   \group_begin:
2352   \__regex_single_match:
2353   #1
2354   \__regex_match:n {#2}
2355   \__regex_extract:
2356   \__regex_group_end_extract_seq:N #3
2357 }
2358 \cs_new_protected:Npn \__regex_extract_all:nnN #1#2#3
2359 {
2360   \group_begin:
2361   \__regex_multi_match:n { \__regex_extract: }
2362   #1
2363   \__regex_match:n {#2}
2364   \__regex_group_end_extract_seq:N #3
2365 }

```

*(End definition for `\__regex_extract_once:nnN` and `\__regex_extract_all:nnN`)*

`\__regex_split:nnN` Splitting at submatches is a bit more tricky. For each match, extract all submatches, and replace the zeroth submatch by the part of the query between the start of the match attempt and the start of the zeroth submatch. This is inhibited if the delimiter matched an empty token list at the start of this match attempt. After the last match, store the last part of the token list, which ranges from the start of the match attempt to the end of the query. This step is inhibited if the last match was empty and at the very end: decrement `\l__regex_submatch_int`, which controls which `\skip` registers will be used.

```

2366 \cs_new_protected:Npn \__regex_split:nnN #1#2#3
2367 {
2368   \group_begin:
2369   \__regex_multi_match:n
2370   {
2371     \if_int_compare:w \l__regex_start_pos_int < \l__regex_success_pos_int
2372     \__regex_extract:
2373     \tex_skip:D \l__regex_zeroth_submatch_int
2374     = \l__regex_start_pos_int sp
2375     plus \tex_skip:D \l__regex_zeroth_submatch_int \scan_stop:
2376     \fi:
2377   }
2378   #1
2379   \__regex_match:n {#2}
2380   <assert>\assert_int:n { \l__regex_current_pos_int = \l__regex_max_pos_int }
2381   \tex_skip:D \l__regex_submatch_int
2382   = \l__regex_start_pos_int sp plus \l__regex_max_pos_int sp \scan_stop:
2383   \int_incr:N \l__regex_submatch_int

```

```

2384     \if_meaning:w \c_true_bool \l__regex_empty_success_bool
2385     \if_int_compare:w \l__regex_start_pos_int = \l__regex_max_pos_int
2386     \int_decr:N \l__regex_submatch_int
2387     \fi:
2388   \fi:
2389   \__regex_group_end_extract_seq:N #3
2390 }

```

(End definition for \\_\_regex\_split:nnN)

\\_\_regex\_group\_end\_extract\_seq:N

The end-points of submatches are stored as the main and stretch components of \skip registers from \l\_\_regex\_max\_state\_int to \l\_\_regex\_submatch\_int (exclusive). Extract the relevant ranges into \l\_\_regex\_internal\_a\_tl. We detect unbalanced results using the two flags @@\_begin and @@\_end, raised whenever we see too many begin-group or end-group tokens in a submatch. We disable \\_\_seq\_item:n to prevent two x-expansions.

```

2391 \cs_new_protected:Npn \__regex_group_end_extract_seq:N #1
2392 {
2393   \cs_set_eq:NN \__seq_item:n \scan_stop:
2394   \flag_clear:n { __regex_begin }
2395   \flag_clear:n { __regex_end }
2396   \tl_set:Nx \l__regex_internal_a_tl
2397   {
2398     \int_step_function:nnnN
2399     { \c_two * \l__regex_max_state_int }
2400     \c_one
2401     { \l__regex_submatch_int - \c_one }
2402     \__regex_extract_seq_aux:n
2403   }
2404   \int_compare:nNnF
2405   { \flag_height:n { __regex_begin } + \flag_height:n { __regex_end } }
2406   = \c_zero
2407   {
2408     \_msg_kernel_error:nnxxx { regex } { result-unbalanced }
2409     { splitting-or-extracting-submatches }
2410     { \flag_height:n { __regex_end } }
2411     { \flag_height:n { __regex_begin } }
2412   }
2413   \use:x
2414   {
2415     \group_end:
2416     \tl_set:Nn \exp_not:N #1 { \l__regex_internal_a_tl }
2417   }
2418 }

```

(End definition for \\_\_regex\_group\_end\_extract\_seq:N)

\\_\_regex\_extract\_seq\_aux:n  
 \\_\_regex\_extract\_seq\_aux:ww

The :n auxiliary builds one item of the sequence of submatches. First compute the brace balance of the submatch, then extract the submatch from the query, adding the appropriate braces and raising a flag if the submatch is not balanced.

```

2419 \cs_new:Npn \__regex_extract_seq_aux:n #1

```

```

2420 {
2421   \__seq_item:n
2422   {
2423     \exp_after:wN \__regex_extract_seq_aux:ww
2424     \__int_value:w \__regex_submatch_balance:n {#1} ; #1;
2425   }
2426 }
2427 \cs_new:Npn \__regex_extract_seq_aux:ww #1; #2;
2428 {
2429   \if_int_compare:w #1 < \c_zero
2430     \flag_raise:n { __regex_end }
2431     \prg_replicate:nn {-#1} { \exp_not:n { { \if_false: } \fi: } }
2432   \fi:
2433   \__regex_query_submatch:n {#2}
2434   \if_int_compare:w #1 > \c_zero
2435     \flag_raise:n { __regex_begin }
2436     \prg_replicate:nn {#1} { \exp_not:n { \if_false: { \fi: } } }
2437   \fi:
2438 }

```

(End definition for \\_\_regex\_extract\_seq\_aux:n and \\_\_regex\_extract\_seq\_aux:ww)

```

\__regex_extract:
\__regex_extract_b:wn
\__regex_extract_e:wn

```

Our task here is to extract from the property list \l\_\_regex\_success\_submatches\_prop the list of end-points of submatches, and store them in \skip registers, from \l\_\_regex\_zeroth\_submatch\_int upwards. We begin by emptying those \skip registers. Then for each  $\langle key \rangle$ – $\langle value \rangle$  pair in the property list update the appropriate \skip component. This is somewhat a hack: the  $\langle key \rangle$  is a non-negative integer followed by < or >, which we use in a comparison to  $-1$ . At the end, store the information about the position at which the match attempt started, as a shrink component.

```

2439 \cs_new_protected_nopar:Npn \__regex_extract:
2440 {
2441   \if_meaning:w \c_true_bool \g__regex_success_bool
2442     \int_set_eq:NN \l__regex_zeroth_submatch_int \l__regex_submatch_int
2443     \prg_replicate:nn \l__regex_capturing_group_int
2444     {
2445       \tex_skip:D \l__regex_submatch_int \c_zero sp \scan_stop:
2446       \int_incr:N \l__regex_submatch_int
2447     }
2448   \prop_map_inline:Nn \l__regex_success_submatches_prop
2449   {
2450     \if_int_compare:w ##1 \c_minus_one
2451       \exp_after:wN \__regex_extract_e:wn \__int_value:w
2452     \else:
2453       \exp_after:wN \__regex_extract_b:wn \__int_value:w
2454     \fi:
2455     \__int_eval:w \l__regex_zeroth_submatch_int + ##1 {##2}
2456   }
2457   \tex_skip:D \l__regex_zeroth_submatch_int
2458   = \tex_the:D \tex_skip:D \l__regex_zeroth_submatch_int
2459   minus \l__regex_start_pos_int sp \scan_stop:

```

```

2460   \fi:
2461 }
2462 \cs_new_protected:Npn \__regex_extract_b:wn #1 < #2
2463 {
2464   \tex_skip:D #1 = #2 sp
2465   plus \etex_gluestretch:D \tex_skip:D #1 \scan_stop:
2466 }
2467 \cs_new_protected:Npn \__regex_extract_e:wn #1 > #2
2468 {
2469   \tex_skip:D #1
2470   = 1 \tex_skip:D #1 plus #2 sp \scan_stop:
2471 }

```

(End definition for \\_\_regex\_extract:, \\_\_regex\_extract\_b:wn, and \\_\_regex\_extract\_e:wn)

## 2.7.4 Replacement

\\_\_regex\_replace\_once:nnN Build the NFA and the replacement functions, then find a single match. If the match failed, simply exit the group. Otherwise, we do the replacement. Extract submatches. Compute the brace balance corresponding to replacing this match by the replacement (this depends on submatches). Prepare the replaced token list: the replacement function produces the tokens from the start of the query to the start of the match and the replacement text for this match; we need to add the tokens from the end of the match to the end of the query. Finally, store the result in the user's variable after closing the group: this step involves an additional x-expansion, and checks that braces are balanced in the final result.

```

2472 \cs_new_protected:Npn \__regex_replace_once:nnN #1#2#3
2473 {
2474   \group_begin:
2475   \__regex_single_match:
2476   #1
2477   \__regex_replacement:n {#2}
2478   \exp_args:No \__regex_match:n { #3 }
2479   \if_meaning:w \c_false_bool \g__regex_success_bool
2480   \group_end:
2481   \else:
2482     \__regex_extract:
2483     \int_set:Nn \l__regex_balance_int
2484     {
2485       \__regex_replacement_balance_one_match:n
2486       { \l__regex_zeroth_submatch_int }
2487     }
2488     \tl_set:Nx \l__regex_internal_a_tl
2489     {
2490       \__regex_replacement_do_one_match:n { \l__regex_zeroth_submatch_int }
2491       \__regex_query_range:nn
2492       { \etex_gluestretch:D \tex_skip:D \l__regex_zeroth_submatch_int }
2493       { \l__regex_max_pos_int }
2494     }
2495     \__regex_group_end_replace:N #3
2496   \fi:

```

```

2497     }
(End definition for \_regex_replace_once:nnN)

```

\\_regex\_replace\_all:nnN Match multiple times, and for every match, extract submatches and additionally store the position at which the match attempt started (as the shrink component of a \skip register). The \skip registers from \l\\_regex\_max\_state\_int to \l\\_regex\_submatch\_int hold information about submatches of every match in order; each match corresponds to \l\\_regex\_capturing\_group\_int consecutive \skip registers. Compute the brace balance corresponding to doing all the replacements: this is the sum of brace balances for replacing each match. Join together the replacement texts for each match (including the part of the query before the match), and the end of the query.

```

2498 \cs_new_protected:Npn \_regex_replace_all:nnN #1#2#3
2499 {
2500   \group_begin:
2501     \_regex_multi_match:n { \_regex_extract: }
2502     #1
2503     \_regex_replacement:n {#2}
2504     \exp_args:No \_regex_match:n {#3}
2505     \int_set:Nn \l\_regex_balance_int
2506       {
2507         0
2508         \int_step_function:nnnN
2509           { \c_two * \l\_regex_max_state_int }
2510           \l\_regex_capturing_group_int
2511           { \l\_regex_submatch_int - \c_one }
2512           \_regex_replacement_balance_one_match:n
2513       }
2514     \tl_set:Nx \l\_regex_internal_a_tl
2515     {
2516       \int_step_function:nnnN
2517         { \c_two * \l\_regex_max_state_int }
2518         \l\_regex_capturing_group_int
2519         { \l\_regex_submatch_int - \c_one }
2520         \_regex_replacement_do_one_match:n
2521         \_regex_query_range:nn
2522         \l\_regex_start_pos_int \l\_regex_max_pos_int
2523     }
2524     \_regex_group_end_replace:N #3
2525   }

```

(End definition for \\_regex\_replace\_all:nnN)

\\_regex\_group\_end\_replace:N If the brace balance is not 0, raise an error. Then set the user's variable #1 to the x-expansion of \l\\_regex\_internal\_a\_tl, adding the appropriate braces to produce a balanced result. And end the group.

```

2526 \cs_new_protected_nopar:Npn \_regex_group_end_replace:N #1
2527 {
2528   \if_int_compare:w \l\_regex_balance_int = \c_zero
2529   \else:

```

```

2530     \_msg_kernel_error:nnxxx { regex } { result-unbalanced }
2531     { replacing }
2532     { \int_max:nn { - \l__regex_balance_int } { \c_zero } }
2533     { \int_max:nn { \l__regex_balance_int } { \c_zero } }
2534 \fi:
2535 \use:x
2536 {
2537     \group_end:
2538     \tl_set:Nn \exp_not:N #1
2539     {
2540         \if_int_compare:w \l__regex_balance_int < \c_zero
2541         \prg_replicate:nn { - \l__regex_balance_int }
2542         { { \if_false: } \fi: }
2543         \fi:
2544         \l__regex_internal_a_tl
2545         \if_int_compare:w \l__regex_balance_int > \c_zero
2546         \prg_replicate:nn { \l__regex_balance_int }
2547         { \if_false: { \fi: } }
2548         \fi:
2549     }
2550 }
2551 }

```

(End definition for \\_regex\_group\_end\_replace:N)

## 2.7.5 Storing and showing compiled patterns

## 2.8 Messages

Messages for the preparsing phase.

```

2552 \_msg_kernel_new:nnnn { regex } { trailing-backslash }
2553 { Trailing~escape~character~\iow_char:N\\. }
2554 {
2555     A~regular~expression~or~its~replacement~text~ends~with~
2556     the~escape~character~\iow_char:N\\.~It~will~be~ignored.
2557 }
2558 \_msg_kernel_new:nnnn { regex } { x-missing-rbrace }
2559 { Missing~closing~brace~in~\iow_char:N\\x~hexadecimal~sequence. }
2560 {
2561     You~wrote~something~like~
2562     '\iow_char:N\\x\{\int_to_hexadecimal:n{#1}\}'~
2563     The~closing~brace~is~missing.
2564 }
2565 \_msg_kernel_new:nnnn { regex } { x-overflow }
2566 { Character~code~'#1'~too~large~in~\iow_char:N\\x~hexadecimal~sequence. }
2567 {
2568     You~wrote~something~like~
2569     '\iow_char:N\\x\{\int_to_hexadecimal:n{#1}\}'~
2570     The~character~code~'#1'~is~larger~than~\int_use:N \c_max_char_int.
2571 }

```

Invalid quantifier.

```

2572 \__msg_kernel_new:nnnn { regex } { invalid-quantifier }
2573 { Braced~quantifier~'#1'~may~not~be~followed~by~'#2'. }
2574 {
2575   The~character~'#2'~is~invalid~in~the~braced~quantifier~'#1'.~
2576   The~only~valid~quantifiers~are~'*',~'?','+',~'<int>',~
2577   '{<min>},'~and~'<min>,<max>}',~followed~or~not~by~'?''.
2578 }

```

Messages for missing or extra closing brackets and parentheses, with some fancy singular/plural handling for the case of parentheses.

```

2579 \__msg_kernel_new:nnnn { regex } { missing-rbrack }
2580 { Missing~right~bracket~inserted~in~regular~expression. }
2581 {
2582   LaTeX~was~given~a~regular~expression~where~a~character~class~
2583   was~started~with~'['~but~the~matching~']'~is~missing.
2584 }
2585 \__msg_kernel_new:nnnn { regex } { missing-rparen }
2586 {
2587   Missing~right~parenthes\int_compare:nTF{#1=1}{i}{e}s~
2588   inserted~in~regular~expression.
2589 }
2590 {
2591   LaTeX~was~given~a~regular~expression~with~\int_eval:n{#1}~
2592   more~left~parenthes\int_compare:nTF{#1=1}{i}{e}s~than~right~
2593   parenthes\int_compare:nTF{#1=1}{i}{e}s.
2594 }
2595 \__msg_kernel_new:nnnn { regex } { extra-rparen }
2596 { Extra~right~parenthesis~ignored~in~regular~expression. }
2597 {
2598   LaTeX~came~across~a~closing~parenthesis~when~no~submatch~group~
2599   was~open.~The~parenthesis~will~be~ignored.
2600 }

```

Some escaped alphanumerics are not allowed everywhere.

```

2601 \__msg_kernel_new:nnnn { regex } { bad-escape }
2602 {
2603   Invalid~escape~\c_backslash_str #1~
2604   \__regex_if_in_cs:TF { within~a~control~sequence. }
2605   {
2606     \__regex_if_in_class:TF
2607     { in~a~character~class. }
2608     { following~a~category~test. }
2609   }
2610 }
2611 {
2612   The~escape~sequence~\iow_char:N\#1~may~not~appear~
2613   \__regex_if_in_cs:TF
2614   {
2615     within~a~control~sequence~test~introduced~by~

```



```

2616         \iow_char:N\c\iow_char:N\{.
2617     }
2618     {
2619         \__regex_if_in_class:TF
2620         { within-a-character-class~ }
2621         { following-a-category-test-such-as-\iow_char:N\cL~ }
2622         because-it-does-not-match-exactly-one-character.
2623     }
2624 }

Range errors.

2625 \__msg_kernel_new:nnnn { regex } { range-missing-end }
2626 { Invalid-end-point-for-range~'#1-#2'~in-character-class. }
2627 {
2628     The-end-point~'#2'~of-the-range~'#1-#2'~may-not-serve-as-an-
2629     end-point-for-a-range:~alphanumeric~characters~should-not-be-
2630     escaped,~and-non-alphanumeric~characters~should-be-escaped.
2631 }
2632 \__msg_kernel_new:nnnn { regex } { range-backwards }
2633 { Range~[#1-#2]~out-of-order-in-character-class. }
2634 {
2635     In-ranges-of-characters~[x-y]~appearing-in-character-classes,~
2636     the-first-character-code-must-not-be-larger-than-the-second.~
2637     Here,~#1~has-character-code~\int_eval:n {'#1},~while~#2~has-
2638     character-code~\int_eval:n {'#2}.
2639 }

Errors related to \c and \u.

2640 \__msg_kernel_new:nnnn { regex } { c-bad-mode }
2641 { Invalid-nested-\iow_char:N\c~escape-in-regular-expression. }
2642 {
2643     The-\iow_char:N\c~escape-cannot-be-used-within-
2644     a-control-sequence-test~'\iow_char:N\c{...}'~.~
2645     To-combine-several-category-tests,~use~'\iow_char:N\c[...]'~.
2646 }
2647 \__msg_kernel_new:nnnn { regex } { c-missing-rbrace }
2648 { Missing-right-brace-inserted-for-\iow_char:N\c~escape. }
2649 {
2650     LaTeX-was-given-a-regular-expression-where-a-
2651     '\iow_char:N\c\iow_char:N\{...}'~construction-was-not-ended-
2652     with-a-closing-brace~'\iow_char:N\}'~.
2653 }
2654 \__msg_kernel_new:nnnn { regex } { c-missing-rbrack }
2655 { Missing-right-bracket-inserted-for-\iow_char:N\c~escape. }
2656 {
2657     A-construction~'\iow_char:N\c[...]'~appears-in-a-
2658     regular-expression,~but-the-closing~']'~is-not-present.
2659 }
2660 \__msg_kernel_new:nnnn { regex } { c-missing-category }
2661 { Invalid-character~'#1'~following-\iow_char:N\c~escape. }
2662 {

```

```

2663 In-regular-expressions,~the~\iow_char:N\\c~escape-sequence~
2664 may~only~be~followed~by~a~left~brace,~a~left~bracket,~or~a~
2665 capital~letter~representing~a~character~category,~namely~
2666 one~of~ABCDELMOPSTU.
2667 }
2668 \__msg_kernel_new:nnnn { regex } { u-missing-lbrace }
2669 { Missing~left~brace~following~\iow_char:N\\u~escape. }
2670 {
2671 The~\iow_char:N\\u~escape-sequence~must~be~followed~by~
2672 a~brace~group~with~the~name~of~the~variable~to~use.
2673 }
2674 \__msg_kernel_new:nnnn { regex } { u-missing-rbrace }
2675 { Missing~right~brace~inserted~for~\iow_char:N\\u~escape. }
2676 {
2677 LaTeX~
2678 \str_if_eq_x:nnTF { } {#2}
2679 { reached~the~end~of~the~string~ }
2680 { encountered~an~escaped~alphanumeric~character '~\iow_char:N\\#2'~ }
2681 when~parsing~the~argument~of~an~'\iow_char:N\\u\iow_char:N\{...\}'~escape.
2682 }

```

Errors when encountering the POSIX syntax [:...:].

```

2683 \__msg_kernel_new:nnnn { regex } { posix-unsupported }
2684 { POSIX-collating-element~'[#1 ~ #1]~not~supported. }
2685 {
2686 The~[.foo.]~and~[=bar=]~syntaxes~have~a~special~meaning~in~POSIX~
2687 regular~expressions.~This~is~not~supported~by~LaTeX.~Maybe~you~
2688 forgot~to~escape~a~left~bracket~in~a~character~class?
2689 }
2690 \__msg_kernel_new:nnnn { regex } { posix-unknown }
2691 { POSIX~class~[:#1:]~unknown. }
2692 {
2693 [:#1:]~is~not~among~the~known~POSIX~classes~
2694 [:alnum:],~[:alpha:],~[:ascii:],~[:blank:],~
2695 [:cntrl:],~[:digit:],~[:graph:],~[:lower:],~
2696 [:print:],~[:punct:],~[:space:],~[:upper:],~
2697 [:word:],~and~[:xdigit:].
2698 }
2699 \__msg_kernel_new:nnnn { regex } { posix-missing-close }
2700 { Missing~closing~'~'~for~POSIX~class. }
2701 { The~POSIX~syntax~'#1'~must~be~followed~by~'~',~not~'#2'. }

```

In various cases, the result of a `l3regex` operation can leave us with an unbalanced token list, which we must re-balance by adding begin-group or end-group character tokens.

```

2702 \__msg_kernel_new:nnnn { regex } { result-unbalanced }
2703 { Missing~brace~inserted~when~#1. }
2704 {
2705 LaTeX~was~asked~to~do~some~regular~expression~operation,~
2706 and~the~resulting~token~list~would~not~have~the~same~number~

```

```

2707 of~begin~group~and~end~group~tokens.~Braces~were~inserted:~
2708 #2~left,~#3~right.
2709 }
Error message for unknown options.
2710 \_msg_kernel_new:nnnn { regex } { unknown-option }
2711 { Unknown~option~'#1'~for~regular~expressions. }
2712 {
2713 The~only~available~option~is~'case-insensitive',~toggled~by~
2714 '(?i)'~and~'(?-i)'.
2715 }
Errors in the replacement text.
2716 \_msg_kernel_new:nnnn { regex } { replacement-c }
2717 { Misused~\iow_char:N\\c~command~in~a~replacement~text. }
2718 {
2719 In~a~replacement~text,~the~\iow_char:N\\c~escape~sequence~
2720 can~be~followed~by~one~of~the~letters~ABCDELMOPTU~
2721 or~a~brace~group,~not~by~'#1'.
2722 }
2723 \_msg_kernel_new:nnnn { regex } { replacement-u }
2724 { Misused~\iow_char:N\\u~command~in~a~replacement~text. }
2725 {
2726 In~a~replacement~text,~the~\iow_char:N\\u~escape~sequence~
2727 must~be~followed~by~a~brace~group~holding~the~name~of~the~
2728 variable~to~use.
2729 }
2730 \_msg_kernel_new:nnnn { regex } { replacement-g }
2731 { Missing~brace~for~the~\iow_char:N\\g~construction~in~a~replacement~text. }
2732 {
2733 In~the~replacement~text~for~a~regular~expression~search,~
2734 submatches~are~represented~either~as~\iow_char:N \\g{dd..d},~
2735 or~\\d,~where~'d'~are~single~digits.~Here,~a~brace~is~missing.
2736 }
2737 \_msg_kernel_new:nnnn { regex } { replacement-catcode-end }
2738 {
2739 Missing~character~for~the~\iow_char:N\\c<category><character>~
2740 construction~in~a~replacement~text.
2741 }
2742 {
2743 In~a~replacement~text,~the~\iow_char:N\\c~escape~sequence~
2744 can~be~followed~by~one~of~the~letters~ABCDELMOPTU~representing~
2745 the~character~category.~Then,~a~character~must~follow.~LaTeX~
2746 reached~the~end~of~the~replacement~when~looking~for~that.
2747 }
2748 \_msg_kernel_new:nnnn { regex } { replacement-null-space }
2749 { TeX~cannot~build~a~space~token~with~character~code~0. }
2750 {
2751 You~asked~for~a~character~token~with~category~'space',~
2752 and~character~code~0,~for~instance~through~
2753 '\iow_char:N\\cS\iow_char:N\\x00'.~

```

```

2754 This~specific~case~is~impossible~and~will~be~replaced~
2755 by~a~normal~space.
2756 }
2757 \__msg_kernel_new:nnnn { regex } { replacement-missing-rbrace }
2758 { Missing~right~brace~inserted~in~replacement~text. }
2759 {
2760   There~were~\int_use:N \l__regex_replacement_csnames_int \
2761   missing~right~braces.
2762 }

```

`\__regex_msg_repeated:nnN` This is not technically a message, but seems related enough to go there. The arguments are: #1 is the minimum number of repetitions; #2 is the number of allowed extra repetitions (−1 for infinite number), and #3 tells us about lazyness.

```

2763 \cs_new:Npn \__regex_msg_repeated:nnN #1#2#3
2764 {
2765   \str_if_eq_x:nnF { #1 #2 } { 1 0 }
2766   {
2767     , ~ repeated ~
2768     \int_case:nnn {#2}
2769     {
2770       { -1 } { #1~or~more~times,~\bool_if:NTF #3 { lazy } { greedy } }
2771       { 0 } { #1~times }
2772     }
2773     {
2774       between~#1~and~\int_eval:n {#1+#2}~times,~
2775       \bool_if:NTF #3 { lazy } { greedy }
2776     }
2777   }
2778 }

```

(End definition for `\__regex_msg_repeated:nnN`)

## 2.9 Code for tracing

The tracing code is still very experimental, and is meant to be used with the `l3trace` package, currently in `l3trial`.

`\__regex_trace_states:n` This function lists the contents of all states of the NFA, stored in `\toks` from 0 to `\l__regex_max_state_int` (excluded).

```

2779 <*trace>
2780 \cs_new_protected:Npn \__regex_trace_states:n #1
2781 {
2782   \int_step_inline:nnnn
2783     \l__regex_min_state_int
2784     \c_one
2785     { \l__regex_max_state_int - 1 }
2786   {
2787     \trace:nnx { regex } { #1 }
2788     { \iow_char:N \toks ##1 = { \tex_the:D \tex_toks:D ##1 } }
2789   }

```



\\_regex_compile_):	.....	1	\\_regex_compile_group_begin:N	....	
\\_regex_compile_::	.....	1	.....	0, 0, 0, 0, 0, 1	
\\_regex_compile_/A:	.....	1	\\_regex_compile_group_end:	....	0, 0, 1
\\_regex_compile_/B:	.....	1	\\_regex_compile_lparen:w	.....	0, 0
\\_regex_compile_/D:	.....	1	\\_regex_compile_one:x	.....	
\\_regex_compile_/G:	.....	1	.....	0, 0, 0, 0, 0, 0, 0, 0, 0, 1	
\\_regex_compile_/H:	.....	1	\\_regex_compile_quantifier:w	....	
\\_regex_compile_/K:	.....	1	.....	0, 0, 0, 0, 1	
\\_regex_compile_/N:	.....	1	\\_regex_compile_quantifier_*:w	....	1
\\_regex_compile_/S:	.....	1	\\_regex_compile_quantifier_+:w	....	1
\\_regex_compile_/V:	.....	1	\\_regex_compile_quantifier_?:w	....	1
\\_regex_compile_/W:	.....	1	\\_regex_compile_quantifier_abort:xNN	.....	
\\_regex_compile_/Z:	.....	1	.....	0, 0, 0, 0, 0, 1	
\\_regex_compile_/b:	.....	1	\\_regex_compile_quantifier_braced_i:w	.....	
\\_regex_compile_/c:	.....	1	.....	0, 0, 1	
\\_regex_compile_/d:	.....	1	\\_regex_compile_quantifier_braced_ii:w	.....	
\\_regex_compile_/h:	.....	1	.....	0, 0, 1	
\\_regex_compile_/s:	.....	1	\\_regex_compile_quantifier_braced_iii:w	.....	
\\_regex_compile_/u:	.....	1	.....	0, 0, 1	
\\_regex_compile_/v:	.....	1	\\_regex_compile_quantifier_lazyness:nnNN	.....	
\\_regex_compile_/w:	.....	1	.....	0, 0, 0, 0, 0, 0, 0, 1	
\\_regex_compile_/z:	.....	1	\\_regex_compile_quantifier_none:	..	
\\_regex_compile_[:	.....	1	.....	0, 0, 0, 0, 1	
\\_regex_compile_]:	.....	1	\\_regex_compile_range:Nw	....	0, 0, 1
\\_regex_compile^:	.....	1	\\_regex_compile_raw:N	... 0, 0, 0, 0,	
\\_regex_compile_abort_tokens:n	0, 0, 1		0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,		
\\_regex_compile_abort_tokens:x	....		0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1		
.....	0, 0, 0, 1		\\_regex_compile_raw_error:N	.....	
\\_regex_compile_anchor:Nf	... 0, 0, 0, 1		.....	0, 0, 0, 0, 0, 0, 1	
\\_regex_compile_c[:w	.....	1	\\_regex_compile_special:N	0, 0, 0, 0, 0,	
\\_regex_compile_c_lbrack_add:N	0, 0, 1		0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1		
\\_regex_compile_c_lbrack_end:	0, 0, 0, 1		\\_regex_compile_special_group_~:w	..	1
\\_regex_compile_c_lbrack_loop:NN	..		\\_regex_compile_special_group_::w	..	1
.....	0, 0, 0, 0, 1		\\_regex_compile_special_group_i:w	0, 1	
\\_regex_compile_c_test:NN	....	0, 0, 1	\\_regex_compile_u_end:	....	0, 0, 0, 1
\\_regex_compile_class:TFNN	..	0, 0, 0, 1	\\_regex_compile_u_in_cs:	.....	0, 0, 1
\\_regex_compile_class_catcode:w	0, 0, 1		\\_regex_compile_u_in_cs_aux:n	... 0, 0	
\\_regex_compile_class_ii:NN	0, 0, 0, 1		\\_regex_compile_u_loop:NN	0, 0, 0, 0, 1	
\\_regex_compile_class_normal:w	0, 0, 1		\\_regex_compile_u_not_cs:	....	0, 0, 1
\\_regex_compile_class_posix:NNNNw	..		\\_regex_compute_case_changed_char:	.....	
.....	0, 0, 1		.....	0, 0, 0, 1	
\\_regex_compile_class_posix_end:w	..		\\_regex_count:nnN	.....	0, 0, 0, 1
.....	0, 0, 1		\\_regex_disable_submatches:	.....	
\\_regex_compile_class_posix_loop:w	.....	0, 0, 0, 0, 1	.....	0, 0, 0, 0, 0, 1	
\\_regex_compile_class_posix_test:w	.....	0, 0, 1	\\_regex_end	.....	1
.....	0, 0, 1		\\_regex_escape_\\:w	.....	1
\\_regex_compile_end:	....	0, 0, 0, 0, 1	\\_regex_escape_/a:w	.....	1
\\_regex_compile_escaped:N	....	0, 0, 1	\\_regex_escape_/break:w	.....	1
			\\_regex_escape_/e:w	.....	1
			\\_regex_escape_/f:w	.....	1

\\_regex_escape_/n:w	1	\\_regex_if_in_cs:TF	0, 0, 0, 0, 1
\\_regex_escape_/r:w	1	\\_regex_if_match:nn	0, 0, 0, 1
\\_regex_escape_/t:w	1	\\_regex_if_raw_digit:NN	0
\\_regex_escape_/x:w	1	\\_regex_if_raw_digit:NNTF	0, 0, 1
\\_regex_escape_\\:w	1	\\_regex_if_two_empty_matches:F	0, 0, 0, 0, 1
\\_regex_escape_break:w	0, 1	\\_regex_if_within_catcode:TF	0, 0, 1
\\_regex_escape_escaped:N	0, 0, 0, 1	\\_regex_item_caseful_equal:n	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\\_regex_escape_loop:N	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1	\\_regex_item_caseful_range:nn	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\\_regex_escape_raw:N	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1	\\_regex_item_caseless_equal:n	0, 0, 0, 1
\\_regex_escape_unescaped:N	0, 0, 0, 1	\\_regex_item_caseless_range:nn	0, 0, 0, 1
\\_regex_escape_use:nnnn	0, 0, 0, 1	\\_regex_item_catcode:	0, 0, 1
\\_regex_escape_x_end:w	0, 0, 1	\\_regex_item_catcode:nT	0, 0, 0, 0, 0, 1
\\_regex_escape_x_ii:N	0, 0, 1	\\_regex_item_catcode_reverse:nT	0, 0, 0, 1
\\_regex_escape_x_large:n	0, 0, 1	\\_regex_item_cs:n	0, 0, 0, 0, 1
\\_regex_escape_x_loop:N	0, 0, 0, 0, 1	\\_regex_item_equal:n	0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\\_regex_escape_x_test:N	0, 0, 0, 1	\\_regex_item_exact:nn	0, 0, 0, 1
\\_regex_escape_x_test_ii:N	0, 0	\\_regex_item_exact_cs:c	0, 0, 0, 1
\\_regex_extract:	0, 0, 0, 0, 0, 0, 1	\\_regex_item_range:nn	0, 0, 0, 0, 0, 1
\\_regex_extract_all:nnN	0, 0, 1	\\_regex_item_reverse:n	0, 0, 0, 0, 0, 1
\\_regex_extract_b:wn	0, 0, 1	\\_regex_match:n	0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\\_regex_extract_e:wn	0, 0, 1	\\_regex_match_loop:	0, 0, 0, 1
\\_regex_extract_once:nnN	0, 0, 1	\\_regex_match_once:	0, 0, 0, 1
\\_regex_extract_seq_aux:n	0, 0, 1	\\_regex_match_one_active:w	0, 0, 0, 1
\\_regex_extract_seq_aux:ww	0, 0, 1	\\_regex_mode_quit_c:	0, 0, 0, 1
\\_regex_get_digits:NNTFw	0, 0, 0, 1	\\_regex_msg_repeated:nnN	0, 0, 0, 0, 0, 1
\\_regex_get_digits_loop:nw	0, 0, 0	\\_regex_multi_match:n	0, 0, 0, 0, 0, 1
\\_regex_get_digits_loop:w	1	\\_regex_pop_lr_states:	0, 0, 0, 1
\\_regex_group:nnnN	0, 0, 0, 0, 0, 1	\\_regex_posix_alnum:	0, 1
\\_regex_group_aux:nnnnN	0, 0, 0, 0, 1	\\_regex_posix_alpha:	0, 0, 1
\\_regex_group_end_extract_seq:N	0, 0, 0, 0, 1	\\_regex_posix_ascii:	0, 1
\\_regex_group_end_replace:N	0, 0, 0, 1	\\_regex_posix_blank:	0, 1
\\_regex_group_no_capture:nnnN	0, 0, 0, 1	\\_regex_posix_cntrl:	0, 1
\\_regex_group_repeat:nn	0, 0, 1	\\_regex_posix_digit:	0, 0, 0, 1
\\_regex_group_repeat:nnN	0, 0, 1	\\_regex_posix_graph:	0, 1
\\_regex_group_repeat:nnnN	0, 0, 1	\\_regex_posix_lower:	0, 0, 1
\\_regex_group_repeat_aux:n	0, 0, 0, 0, 1	\\_regex_posix_print:	0, 1
\\_regex_group_resetting:nnnN	0, 0, 0, 1	\\_regex_posix_punct:	0, 1
\\_regex_group_resetting_loop:nnNn	0, 0, 0, 1	\\_regex_posix_space:	0, 1
\\_regex_group_submatches:nnN	0, 0, 0, 0, 0, 1	\\_regex_posix_upper:	0, 0, 1
\\_regex_if_end_range:NN	0	\\_regex_posix_word:	0, 1
\\_regex_if_end_range:NNTF	0, 1	\\_regex_posix_xdigit:	0, 1
\\_regex_if_in_class:TF	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1	\\_regex_prop_:	1
\\_regex_if_in_class_or_catcode:TF	0, 0, 0, 0, 0, 1		









```

\int_eval:n ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_if_exist:cTF ..... 0, 0
\int_if_odd_p:n ..... 0
\int_incr:N ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_max:nn ..... 0, 0, 0
\int_new:N ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_set:Nn ..... 0, 0, 0, 0, 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_set_eq:Nc ..... 0
\int_set_eq:NN ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_step_function:nnnN ..... 0, 0, 0
\int_step_inline:nnnn ..... 0, 0
\int_sub:Nn ..... 0, 0, 0, 0, 0
\int_to_hexadecimal:n ..... 0, 0
\int_use:c ..... 0
\int_use:N ..... 0, 0, 0, 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\int_zero:N ..... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\iow_char:N ..... 0, 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
\iow_newline: ..... 0

L
\l__regex_balance_int ..... 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_balance_tl ..... 0, 0, 0, 0, 1
\l__regex_capturing_group_int ... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_case_changed_char_int ....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_catcodes_bool ... 0, 0, 0, 0, 1
\l__regex_catcodes_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_current_catcode_int .....
    ..... 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_current_char_int .... 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_current_pos_int 0, 0, 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_current_state_int .....
    .... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_current_submatches_prop ..
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_default_catcodes_int ....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_empty_success_bool .....
    ..... 0, 0, 0, 0, 0, 0, 1
\l__regex_every_match_tl ... 0, 0, 0, 0, 1
\l__regex_fresh_thread_bool .....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_group_level_int .....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_a_int 0, 0, 0, 0, 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_a_tl ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_b_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_b_tl .....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_bool .. 0, 0, 0, 0, 0, 1
\l__regex_internal_c_int .. 0, 0, 0, 0, 0, 1
\l__regex_internal_regex .....
    .... 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_internal_seq .... 0, 0, 0, 0, 0, 1
\l__regex_last_char_int ..... 0, 0, 0, 1
\l__regex_left_state_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_left_state_seq ... 0, 0, 0, 0, 1
\l__regex_match_count_int .. 0, 0, 0, 0, 1
\l__regex_match_success_bool .....
    ..... 0, 0, 0, 0, 0, 0, 1
\l__regex_max_active_int .....
    ..... 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_max_pos_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_max_state_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_min_active_int 0, 0, 0, 0, 0, 0, 1
\l__regex_min_pos_int .. 0, 0, 0, 0, 0, 0, 1
\l__regex_min_state_int 0, 0, 0, 0, 0, 0, 1
\l__regex_mode_int ..... 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1
\l__regex_replacement_csnames_int 0,
    0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1

```



<code>\tex_escapechar:D</code> .....	0, 0, 0	<code>\tl_to_str:n</code> .....	0, 0, 0
<code>\tex_lccode:D</code> .....	0, 0	<code>\tl_to_str:V</code> .....	0, 1
<code>\tex_muskip:D</code> .....	0, 0, 0, 0, 0	<code>\token_if_eq_charcode:NNTF</code> ..	0, 0, 0, 0, 0
<code>\tex_skip:D</code> .....	0, 0, 0, 0, 0,	<code>\token_if_eq_meaning:NNT</code> .....	0
	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	<code>\token_if_eq_meaning:NNTF</code> .....	
<code>\tex_the:D</code> .....	0, 0, 0, 0, 0, 0, 0, 0, 0, 0		0, 0, 0, 0, 0, 0, 0, 0, 0
<code>\tex_toks:D</code> .....	0,	<code>\token_to_str:N</code> .....	0, 0, 0, 0, 0
	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	<code>\trace:nxx</code> .....	0, 0, 0, 0, 0
<code>\tl_clear:N</code> .....	0	<code>\trace_pop:nnn</code> .....	0, 0, 0, 0, 0, 0
<code>\tl_const:Nn</code> .....	0	<code>\trace_pop:nxx</code> .....	0
<code>\tl_const:Nx</code> .....	0	<code>\trace_push:nnn</code> .....	0, 0, 0, 0, 0, 0
<code>\tl_gset_eq:NN</code> .....	0	<code>\trace_push:nxx</code> .....	0
<code>\tl_map_function:NN</code> .....	0		
<code>\tl_map_function:nN</code> .....	0		
			U
<code>\tl_new:N</code> .....	0, 0, 0, 0, 0	<code>\use:c</code> .....	0
<code>\tl_put_right:Nn</code> .....	0	<code>\use:n</code> .....	0, 0, 0, 0, 0, 0, 0, 0
<code>\tl_set:Nn</code> .....	0, 0, 0, 0, 0, 0	<code>\use:x</code> .....	0, 0, 0, 0, 0
<code>\tl_set:Nv</code> .....	0	<code>\use_i:nn</code> .....	0, 0, 0, 0, 0
<code>\tl_set:Nx</code> .....	0, 0, 0, 0, 0, 0, 0, 0, 0, 0	<code>\use_i:nnn</code> .....	0
<code>\tl_set_eq:NN</code> .....	0	<code>\use_ii:nn</code> .....	0, 0, 0, 0, 0, 0
<code>\tl_to_lowercase:n</code> .....	0, 0, 0	<code>\use_none:n</code> .....	0, 0, 0
<code>\tl_to_str:N</code> .....	0	<code>\use_none:nn</code> .....	0, 0