

The **I3regex** package: regular expressions in **T_EX**^{*}

The L^AT_EX3 Project[†]

Released 2012/01/19

1 I3regex documentation

The **I3regex** 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 we act on lists of 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 token list variable using `\regex_const:Nn`. For example,

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

stores in `\c_foo_regex_tl` 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 v3209, last revised 2012/01/19.

[†]E-mail: 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 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{\langle regex\rangle}` 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).

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 character in the subject string.

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

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

`[x-y]` Range (can be used with escaped characters).

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 } \l_tmpa_int yields 2, but replacing \G by ^ would result in \l_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;
- O 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<character or class> Matches a token with category code X (any of CBEMTPUDSLOA) if it is also matched by the *<character or class>*. For instance, \cL[A-Z] matches uppercase letters of category code letter, \cC. matches any control sequence, and \c0\d matches digits of category other.
- \c[XYZ]<character or class> Matches a token with category X, Y, or Z (each being any of CBEMTPUDSLOA), if it is also matched by the *<character or class>*. For instance, \c[LSO]. matches tokens of category letter, space, or other.
- \c[^XYZ]<character or class> Matches a token with category different from X, Y, or Z (each being any of CBEMTPUDSLOA), if it is also matched by the *<character or class>*. For instance, \c[LSO]. matches tokens of category letter, space, or other.

The category code tests can be used inside classes; for instance, [\c0\d \c[L0][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.

Options can be set with (?<option>) and unset with (?-<option>). Options are local to the group in which they are set, and revert to their previous setting upon reaching the closing parenthesis. For instance, in (?i)a(b(?-i)c|d)e, the i option applies to the letters a, b and e.

- (?i) and (?-i) Toggle to a case insensitive/sensitive mode. This only applies to ascii letters (mapping A-Z to a-z). For instance, (?i)[Y-\\" matches the characters Y, Z, [, \, and the lower case letters y and z, while (?i)[^aeiou] matches any character which is not a vowel.

In character classes, only ^, -,], \ and spaces are special, and should be escaped. Other non-alphanumeric characters can still be escaped without harm. The escape sequences \d, \D, etc. are also 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] is equivalent to [^6-9].

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 strings using for instance \regex_extract_once:nnNTF.

1.2 Syntax of in 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]?l|o) . } { \(\0\-\-\1\) } \l_my_tl
```

results in `\l_my_tl` holding `H(e11--e1)(o,--o) w(or--o)(1d--1)`!

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.

`\c{XY}` 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{\langle text\rangle}` Produces the control sequence with csname `\langle text\rangle`. The `\langle text\rangle` may contain references to the submatches `\0`, `\1` etc. As an experimental feature, `\c{\dots}` can be nested, but the behaviour is not yet fully specified. Use at own risks.

1.3 Precompiling 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 precompiled regular expression is stored as a token list variable. All of the `\l3regex` module's functions can be given their regular expression argument either as an explicit string or as a precompiled regular expression.

`\regex_set:Nn` `\regex_set:Nn <tl var> {\langle regex\rangle}`

Stores a precompiled version of the `\langle regular expression\rangle` in the `\langle tl var\rangle`. For instance, this function can be used as

```
\tl_new:N \l_my_regex_tl
\regex_set:Nn \l_my_regex_tl { 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 precompiled expressions which will never change.

`\regex_to_str:N` `\regex_to_str:N \langle regex var\rangle`

Converts the `\langle regex var\rangle`, previously defined using `\regex_(g)set:Nn`, to a string, which can be safely read back when the L^AT_EX3 syntax is active (as triggered by `\ExplSyntaxOn`). Line breaks are inserted at various places, to try and keep the line length short. For instance, if `\l_my_stream` is an open stream for writing, and `\l_my_regex_tl` is a regex variable, you can save its definition using

```
\iow_now:Nx \l_my_stream
{
  \tl_to_str:n { \tl_set:Nn \l_my_regex_tl } { \iow_newline:
    \regex_to_str:N \l_my_regex_tl
  }
}
```

1.4 Matching

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

`\regex_match:nnTF`
`\regex_match:NnTF`

`\regex_match:nnTF {<regex>} {<token list>} {<true code>} {<false code>}`

Tests whether the *<regular expression>* matches any part of the *<token list>*. For instance,

```
\regex_match:nntF { b [cde]* } { abecdcx } { 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 *<int var>* within the current TeX group level equal to the number of times *<regular expression>* appears in *<token list>*. 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 string: then we count one match between each pair of characters. For instance,

```
\int_new:N \l_foo_int
\regex_count:nnN { (b+|c) } { abbababcbb } \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 *<regular expression>* in the *<token list>*. If it exists, the match is stored as the zeroeth item of the *<seq var>*, and further items are the contents of capturing groups, in the order of their opening parenthesis. 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_once:nnNTF { \A(La)?TeX(!*)\Z } { LaTeX!!! } \l_foo_seq
{ true } { false }
```

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

```
\regex_extract_all:nnNTF
\regex_extract_all:NnNTF
```

```
\regex_extract_all:nnN {<regular expression>} {<token list>} <seq var>
\regex_extract_all:nnNTF {<regular expression>} {<token list>} <seq var> {<true code>} {<false 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.

```
\regex_split:nnNTF
\regex_split:NnNTF
```

```
\regex_split:nnN {<regular expression>} {<token list>} <seq var>
\regex_split:nnNTF {<regular expression>} {<token list>} <seq var> {<true code>}
{<false 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

```
\regex_replace_once:nnNTF
\regex_replace_once:NnNTF
```

```
\regex_replace_once:nnN {<regular expression>} {<replacement>} <tl var>
\regex_replace_once:nnNTF {<regular expression>} {<replacement>} <tl var> {<true code>} {<false 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.

```
\regex_replace_all:nnN {<regular expression>} {<replacement>} <tl var>
\regex_replace_all:nnNTF {<regular expression>} {<replacement>} <tl var> {<true
code>} {<false 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.

- The `{}` quantifiers are only partially implemented.
- Clean up the use of messages.
- Triple check magic numbers.
- Test `\par`.
- Add tests for every regex and replacement feature.
- Make sure the documentation is correct, and the code comments are as well.

Code improvements to come.

- `\regex_show:N` to show how a given regular expression is interpreted.
- Test for the maximum register `\c_max_register_int`.
- Use `\dimen` registers rather than `\l_regex_nesting_tl` to build `\regex_nesting:n`.
- Reduce the number of epsilon-transitions in alternatives.
- Improve nesting of `\c` in the replacement text.
- Move the “reconstruction” part of `I3regex` to `I3tl-analysis`.
- Optimize regexes for csnames when the regex is a simple string.
- Optimize simple strings: use less states (`abcade` should give two states, for `abc` and `ade`).
- Optimize groups with no alternative.
- Optimize the use of `\prg_stepwise_...` functions.

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

- 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 (partially implemented).
- `\K` for resetting the beginning of the match.
- 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, hence a lot of optimization ahead.

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

- `\ddd`, matching the character with code `ddd` in octal;
- POSIX character classes `[:alpha:]` etc., this is redundant;
- 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: Xe_TE_X and Lu_AT_E_X serve us characters directly, and splitting those into bytes is tricky, encoding dependent, and most likely not useful anyways.

2 l3regex implementation

```
<*package>
 1 \ProvidesExplPackage
 2   {\ExplFileName}{\ExplFileVersion}{\ExplFileDescription}
 3 \RequirePackage{l3str, l3tl-analysis, l3flag}
```

Most regex engines use backtracking. This allows to provide very powerful features (back-references come to mind first), but it is costly. Since _TE_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 build a non-deterministic finite automaton (NFA) with roughly n states, which accepts precisely those token lists matching that regular expression. We then run the token list through the NFA, and check the return value.

The code is structured as follows. Various helper functions are introduced in the next subsection, to limit the clutter in later parts. Then functions pertaining to parsing the regular expression are introduced: that part is rather long because of the many bells and whistles that we need to cater for. The next subsection takes care of running the NFA, and describes how the various _TE_X registers are (ab)used in this module. Finally, user functions.

2.1 Constants and variables

```
\regex_tmp:w
\l_regex_tma_tl
\l_regex_tmb_tl
\l_regex_tmc_tl
\l_regex_tma_int
\l_regex_tmb_int
 4 \cs_new:Npn \regex_tmp:w { }
 5 \tl_new:N \l_regex_tma_tl
 6 \tl_new:N \l_regex_tmb_tl
 7 \tl_new:N \l_regex_tmc_tl
 8 \int_new:N \l_regex_tma_int
 9 \int_new:N \l_regex_tmb_int
(End definition for \regex_tmp:w. This function is documented on page ??.)
```

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

```
10 \flag_new:N \l_regex_group_begin_flag
11 \flag_new:N \l_regex_group_end_flag
(End definition for \l_regex_group_begin_flag. This function is documented on page ??.)
```

2.1.1 Variables used while building

`\l_regex_build_mode_int` While building, ten modes are recognized, labelled $-63, -23, -6, -2, 0, 2, 3, 6, 23, 63$. See section 2.3.1.

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

`\l_regex_max_state_int` The end-point states of the last group (which any quantifier would repeat) are stored as `\l_regex_left/right_state_int`. For simple strings of characters, the left and right pointers only differ by one. 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.

```
13 \int_new:N \l_regex_max_state_int
14 \int_new:N \l_regex_left_state_int
15 \int_new:N \l_regex_right_state_int
(End definition for \l_regex_max_state_int. This function is documented on page ??.)
```

`\l_regex_left_state_seq` Alternatives are implemented by branching from a state into the various choices, then merging those into another state. We store information about those states in two sequences.

```
16 \seq_new:N \l_regex_left_state_seq
17 \seq_new:N \l_regex_right_state_seq
(End definition for \l_regex_left_state_seq and \l_regex_right_state_seq. These functions are documented on page ??.)
```

`\l_regex_end_group_seq` These sequences hold actions to be performed at the end of a group, and at the end of each branch of the alternation, respectively. Currently, `\l_regex_end_group_seq` is used to keep track of letter case, and `\l_regex_end_alternation_seq` is used for `(?|...)` groups.

```
18 \seq_new:N \l_regex_end_group_seq
19 \seq_new:N \l_regex_end_alternation_seq
(End definition for \l_regex_end_group_seq. This function is documented on page ??.)
```

`\l_regex_capturing_group_int` `\l_regex_capturing_group_int` is the ID number of the current capturing group, starting at 0 for a group enclosing the full regular expression, and counting in the order of their left parenthesis. This number is used when a branch of the alternation ends. Capturing groups can be arbitrarily nested, and we keep track of the stack of ID numbers in `\l_regex_capturing_group_seq`.

```
20 \int_new:N \l_regex_capturing_group_int
21 \seq_new:N \l_regex_capturing_group_seq
22 \int_new:N \l_regex_capturing_group_max_int
(End definition for \l_regex_capturing_group_int. This function is documented on page ??.)
```

`\l_regex_one_or_group_tl` When looking for quantifiers, this variable holds either “one” or “group” depending on whether the object to which the quantifier applies matches one character (*i.e.*, is a character or character class), or is a group.

```
23 \tl_new:N \l_regex_one_or_group_tl
```

(End definition for `\l_regex_one_or_group_t1`. This function is documented on page ??.)

```
\l_regex_catcodes_int
  \l_regex_catcodes_default_int
\c_regex_catcode_C_int
\c_regex_catcode_B_int
\c_regex_catcode_E_int
\c_regex_catcode_M_int
\c_regex_catcode_T_int
\c_regex_catcode_P_int
\c_regex_catcode_U_int
\c_regex_catcode_D_int
\c_regex_catcode_S_int
\c_regex_catcode_L_int
\c_regex_catcode_O_int
\c_regex_catcode_A_int
\c_regex_catcodes_all_int

We wish to allow constructions such as \c[^BE](..\cL[a-z]..), matching two tokens which are neither a begin-group nor an end-group token, followed by a token of category letter and character code in [a-z], followed by two more tokens which are neither begin-group nor end-group tokens. For this to work, we need to keep track of lists of allowed category codes: \l_regex_catcodes_int and \l_regex_catcodes_default_int are bitmaps, sums of  $4^c$ , for all allowed catcodes  $c$ . The latter is local to each capturing group, and we reset \l_regex_catcodes_int to that value after each character or class, changing it only when encountering a \c escape.

24 \int_new:N \l_regex_catcodes_int
25 \int_new:N \l_regex_catcodes_default_int
26 \int_const:Nn \c_regex_catcode_C_int { "1" }
27 \int_const:Nn \c_regex_catcode_B_int { "4" }
28 \int_const:Nn \c_regex_catcode_E_int { "10" }
29 \int_const:Nn \c_regex_catcode_M_int { "40" }
30 \int_const:Nn \c_regex_catcode_T_int { "100" }
31 \int_const:Nn \c_regex_catcode_P_int { "1000" }
32 \int_const:Nn \c_regex_catcode_U_int { "4000" }
33 \int_const:Nn \c_regex_catcode_D_int { "10000" }
34 \int_const:Nn \c_regex_catcode_S_int { "100000" }
35 \int_const:Nn \c_regex_catcode_L_int { "400000" }
36 \int_const:Nn \c_regex_catcode_O_int { "1000000" }
37 \int_const:Nn \c_regex_catcode_A_int { "4000000" }
38 \int_const:Nn \c_regex_catcodes_all_int { "5515155" }

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

`\l_regex_catcodes_bool` Controls whether the bitmap of category codes built should be inverted or not.

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

`\l_regex_tmpa_regex_t1` This holds a temporary pre-compiled regular expression when matching a control sequence name.

```
40 \tl_new:N \l_regex_tmpa_regex_t1
(End definition for \l_regex_tmpa_regex_t1. This function is documented on page ??.)
```

2.1.2 Character classes

`\l_regex_class_bool` is false for negative character classes. `\l_regex_class_t1` holds the tests which should be performed to decide whether the `\l_regex_current_char_int` matches that character class. In nested class, which can only occur as `[\c[...]]`, the two variables from the outer class must be saved while processing the inner class.

```
41 \bool_new:N \l_regex_class_bool
42 \tl_new:N \l_regex_class_t1
43 \bool_new:N \l_regex_class_saved_bool
44 \tl_new:N \l_regex_class_saved_t1
```

(End definition for `\l_regex_class_bool` and `\l_regex_class_t1`. These functions are documented on page ??.)

`\c_regex_d_t1` These constant token lists encode which characters are recognized by `\d`, `\D`, `\w`, etc. in regular expressions. Namely, `\d=[0-9]`, `\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 ranges appear is optimized for usual mostly lower case letter text.

```

45 \tl_const:Nn \c_regex_d_t1
46 {
47   \regex_item_range:nn { \c_forty_eight } { 57 } % 0--9
48 }
49 \tl_const:Nn \c_regex_D_t1
50 {
51   \regex_item_geq:n { \c_fifty_eight } % > '9
52   \regex_item_range:nn { \c_zero } { 47 } % 0
53 }
54 \tl_const:Nn \c_regex_h_t1
55 {
56   \regex_item_equal:n { \c_thirty_two } % space
57   \regex_item_equal:n { \c_nine } % tab
58 }
59 \tl_const:Nn \c_regex_H_t1
60 {
61   \regex_item_geq:n { 33 } % > space
62   \regex_item_range:nn { \c_ten } { 31 } % tab < ... < space
63   \regex_item_range:nn { \c_zero } { \c_eight } % < tab
64 }
65 \tl_const:Nn \c_regex_s_t1
66 {
67   \regex_item_equal:n { \c_thirty_two } % space
68   \regex_item_range:nn { \c_nine } { \c_ten } % tab, lf
69   \regex_item_range:nn { \c_twelve } { \c_thirteen } % ff, cr
70 }
71 \tl_const:Nn \c_regex_S_t1
72 {
73   \regex_item_geq:n { 33 } % > space
74   \regex_item_range:nn { \c_fourteen } { 31 } % tab < ... < space
75   \regex_item_range:nn { \c_zero } { \c_eight } % < tab
76   \regex_item_equal:n { \c_eleven } % vtab
77 }
78 \tl_const:Nn \c_regex_v_t1
79 {
80   \regex_item_range:nn { \c_ten } { \c_thirteen } % lf, vtab, ff, cr
81 }
82 \tl_const:Nn \c_regex_V_t1
83 {
84   \regex_item_geq:n { \c_fourteen } % >cr
85   \regex_item_range:nn { \c_zero } { \c_nine } % < lf

```

```

86    }
87 \tl_const:Nn \c_regex_w_tl
88 {
89     \regex_item_range:nn { \c_ninety_seven } { 122 } % a--z
90     \regex_item_range:nn { \c_sixty_five } { 90 } % A--Z
91     \regex_item_range:nn { \c_forty_eight } { 57 } % 0--9
92     \regex_item_equal:n { 95 } % _
93 }
94 \tl_const:Nn \c_regex_W_tl
95 {
96     \regex_item_range:nn { \c_zero } { 47 } % <'0
97     \regex_item_range:nn { \c_fifty_eight } { 64 } % ('9+1)--('A-1)
98     \regex_item_range:nn { \c_ninety_one } { 94 } % ('Z+1)--('_-1)
99     \regex_item_equal:n { 96 } % '
100    \regex_item_geq:n { \c_one_hundred_twenty_three } % z
101 }
102 \tl_const:Nn \c_regex_N_tl
103 {
104     \regex_item_geq:n { \c_eleven } % > lf
105     \regex_item_range:nn { \c_zero } { \c_nine } % < lf
106 }
(End definition for \c_regex_d_tl and \c_regex_D_tl. These functions are documented on page ??.)
```

2.1.3 Variables used when matching

\l_regex_nesting_int This integer is used to keep track of begin and end-group tokens.

```

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

\l_regex_min_step_int \l_regex_max_step_int This integer holds the value of the step corresponding to the left end of the token list, *i.e.*, when the token list is stored in \toks registers, the \toks register \l_regex_min_step_int holds the first token in the token list. In fact, this number is always one more than \l_regex_max_state_int, but it is more practical to give that a name.

```

108 \int_new:N \l_regex_min_step_int
109 \int_new:N \l_regex_max_step_int
(End definition for \l_regex_min_step_int. This function is documented on page ??.)
```

\l_regex_current_step_int \l_regex_start_step_int \l_regex_success_step_int While reading through the query token list, \l_regex_current_step_int is the position in the token list, starting at \l_regex_min_step_int for the left-most token. Each match begins at the position given by \l_regex_start_step_int. Whenever an execution thread succeeds, the corresponding step is stored into \l_regex_success_step_int, which will be the next starting step (except in the case of empty matches).

```

110 \int_new:N \l_regex_current_step_int
111 \int_new:N \l_regex_start_step_int
112 \int_new:N \l_regex_success_step_int
(End definition for \l_regex_current_step_int. This function is documented on page ??.)
```

`\l_regex_unique_id_int` In the case of repeated matches, `\l_regex_current_step_int` is reset to the end-position of the previous match. In contrast, `\l_regex_unique_id_int` is simply incremented to provide a unique number for each iteration of the matching loop. This is handy to attach each set of submatch information to a given iteration (and automatically discard it when it corresponds to a past iteration).

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

`\l_regex_current_char_int` The character codes of the character at the current position in the token list, and at the previous position, and the current character with its case changed (A-Z↔a-z). The `\l_regex_last_char_int` is used to test for word boundaries (\b and \B). The `\l_regex_case_changed_char_int` is only computed if the “case insensitive” option (?i) is used in the regex.

```
114 \int_new:N \l_regex_current_char_int
115 \int_new:N \l_regex_current_catcode_int
116 \tl_new:N \l_regex_current_token_tl
117 \int_new:N \l_regex_last_char_int
118 \int_new:N \l_regex_case_changed_char_int
(End definition for \l_regex_current_char_int. This function is documented on page ??.)
```

`\l_regex_caseless_bool` True if caseless matching is used within the regular expression. This controls whether `\l_regex_case_changed_char_int` is computed.

```
119 \bool_new:N \l_regex_caseless_bool
(End definition for \l_regex_caseless_bool. This function is 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. In the case of groups with quantifiers, `\l_regex_current_state_int` is shifted to a fake value for transitions to point to the correct states.

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

`\l_regex_current_submatches_prop` The submatches for the thread which lies at the `\l_regex_current_state_int` are stored in a property list variable. This property list is stored by `\regex_action_cost:n` into the `\toks` register for the target state of the transition. When a thread succeeds, this property list is copied to `\l_regex_success_submatches_prop` and only the last sucessful thread will remain there.

```
121 \prop_new:N \l_regex_current_submatches_prop
122 \prop_new:N \l_regex_success_submatches_prop
(End definition for \l_regex_current_submatches_prop. This function is documented on page ??.)
```

`\l_regex_max_index_int` All the currently active states are kept in order of precedence in the `\skip` registers, which for our purpose serve as an array: the *i*th item of the array is `\skipi`. The largest index used after treating the previous character is `\l_regex_max_index_int`. At the start of every step, the whole array is unpacked, so that the space can immediately be reused, and `\l_regex_max_index_int` reset to zero, effectively clearing the array.

```
123 \int_new:N \l_regex_max_index_int
```

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

\l_regex_every_match_t1 Every time a match is found, this token list is used. For single matching, the token list is set to removing the remainder of the query token list. For multiple matching, the token list is set to repeat the matching.

124 \tl_new:N \l_regex_every_match_t1

(End definition for \l_regex_every_match_t1. This function is documented on page ??.)

\g_regex_success_bool
\l_regex_saved_success_bool
\l_regex_success_match_bool The boolean \g_regex_success_bool is true if there was at least one match, and \l_regex_success_match_bool is true if the current match attempt was successful. The variable \g_regex_success_bool is the only global variable in this whole module. When nesting \regex functions internally, the value of \g_regex_success_bool is saved into \l_regex_saved_success_bool, which is local, hence not affected by the changes due to inner regex functions.

125 \bool_new:N \g_regex_success_bool

126 \bool_new:N \l_regex_saved_success_bool

127 \bool_new:N \l_regex_success_match_bool

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

\regex_last_match_empty:F
\regex_last_match_empty_no:F
\regex_last_match_empty_yes:F When doing multiple matches, we need to avoid infinite loops where each iteration matches the same empty token list. When we detect such a situation, the next match attempt is shifted by one character. Namely, an empty match is discarded if it follows an empty match at the same position. If the previous match was non-empty, \regex_last_match_empty:F is simply \use:n, and keeps the match. If it was empty, then we test whether the new match has moved or not: if it has not, then the success is discarded.

128 \cs_new_protected:Npn \regex_last_match_empty_no:F #1 {#1}

129 \cs_new_protected:Npn \regex_last_match_empty_yes:F

130 { \int_compare:nNnF \l_regex_start_step_int = \l_regex_current_step_int }

131 \cs_new_eq:NN \regex_last_match_empty:F \regex_last_match_empty_no:F

(End definition for \regex_last_match_empty:F. This function is documented on page ??.)

\l_regex_success_empty_bool
\l_regex_fresh_thread_bool When a match succeeds, \l_regex_success_empty_bool records whether it is empty. This information is used to initialize \regex_last_match_empty:F before starting the next match attempt. The boolean \l_regex_fresh_thread_bool is true when the current thread has started from the beginning of the regular expression at this character. This is probably suboptimal. Improvements welcome.

132 \bool_new:N \l_regex_success_empty_bool

133 \bool_new:N \l_regex_fresh_thread_bool

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

2.1.4 Variables used when building the replacement

\l_regex_replacement_int

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

\l_regex_replacement_csnames_int

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

\l_regex_submatch_int

\l_regex_submatch_start_int

```
136 \int_new:N \l_regex_submatch_int
137 \int_new:N \l_regex_submatch_start_int
(End definition for \l_regex_submatch_int. This function is documented on page ??.)
```

2.1.5 Variables used 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.

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

2.2 Helpers

2.2.1 Toks

When performing the matching, the \toks registers hold submatch information, followed by the instruction for a given state of the NFA. The two parts are separated by \s_stop.

\regex_toks_put_left:Nx

During the building phase, every \toks register starts with \s_stop, and we wish to add x-expanded material to those registers. The expansion is done “by hand” for optimization (these operations are used quite a lot). When adding material to the left, we define \regex_tmp:w to remove the \s_stop marker and put it back to the left of the new material.

```
139 \cs_new_protected:Npn \regex_toks_put_left:Nx #1#2
140  {
141      \cs_set_nopar:Npx \regex_tmp:w \s_stop { \s_stop #2 }
142      \tex_toks:D #1 \exp_after:wN \exp_after:wN \exp_after:wN
143          { \exp_after:wN \regex_tmp:w \tex_the:D \tex_toks:D #1 }
144  }
145 \cs_new_protected:Npn \regex_toks_put_right:Nx #1#2
146  {
147      \cs_set_nopar:Npx \regex_tmp:w {#2}
148      \tex_toks:D #1 \exp_after:wN
149          { \tex_the:D \tex_toks:D \exp_after:wN #1 \regex_tmp:w }
150  }
(End definition for \regex_toks_put_left:Nx. This function is documented on page ??.)
```

2.2.2 Extracting parts of the query token list

\regex_query_substr:nn
\regex_query_submatch:nn
\regex_query_submatch:w

The token list is stored in \toks registers: the first is \l_regex_min_step_int, the last is (we don't care yet).

```

151 \cs_new:Npn \regex_query_substr:nn #1#2
152 {
153     \if_num:w \int_eval:w #1 < \int_eval:w #2 \int_eval_end:
154         \str_aux_toks_range:nn {#1} {#2}
155     \fi:
156 }
157 \cs_new:Npn \regex_query_submatch:nn #1#2
158 {
159     \if_num:w #1 < \l_regex_capturing_group_int
160         \exp_after:wN \regex_query_submatch:w
161         \int_use:N \int_eval:w #1 + #2 ;
162     \fi:
163 }
164 \cs_new:Npn \regex_query_submatch:w #1 ;
165 {
166     \regex_query_substr:nn
167     { \tex_skip:D #1 }
168     { \etex_gluestretch:D \tex_skip:D #1 }
169 }
```

(End definition for \regex_query_substr:nn. This function is documented on page ??.)

2.2.3 Sequences

\regex_seq_pop_int:NN
\regex_seq_get_int:NN
\regex_seq_push_int:NN

When building the regular expression, we keep track of some integers (pointers to various states) without help from TeX's grouping. Here are variants of \seq_pop:NN and \seq_get:NN which assign using \int_set:Nn rather than \tl_set:Nn.

```

170 \cs_new_protected:Npn \regex_seq_pop_int:NN #1#2
171 {
172     \seq_pop:NN #1 \l_regex_tmpa_tl
173     \int_set:Nn #2 \l_regex_tmpa_tl
174 }
175 \cs_new_protected:Npn \regex_seq_get_int:NN #1#2
176 {
177     \seq_get:NN #1 \l_regex_tmpa_tl
178     \int_set:Nn #2 \l_regex_tmpa_tl
179 }
180 \cs_new_protected:Npn \regex_seq_push_int:NN #1#2
181 { \seq_push:No #1 { \int_use:N #2 } }
```

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

\regex_seq_pop_use:N
\regex_seq_get_use:N

When building the regular expression, some settings are kept local to capturing groups without any help from TeX's grouping. This is done "by hand", in sequences whose items should be run immediately.

```

182 \cs_new_protected:Npn \regex_seq_pop_use:N #1
183 {
```

```

184     \seq_pop>NN #1 \l_regex_tmpa_t1
185     \l_regex_tmpa_t1
186   }
187 \cs_new_protected:Npn \regex_seq_get_use:N #1
188   {
189     \seq_get>NN #1 \l_regex_tmpa_t1
190     \l_regex_tmpa_t1
191   }
(End definition for \regex_seq_pop_use:N. This function is documented on page ??.)
```

2.2.4 Testing characters

`\regex_break_point:TF` When testing whether a character of the query token list matches a given character class 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⟩ ... ⟨testn⟩
\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.

```

192 \cs_new_protected:Npn \regex_break_true:w
193   #1 \regex_break_point:TF #2 #3 {#2}
194 \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_dot:` The dot meta-character matches any character, except the end marker.

```

195 \cs_new_protected_nopar:Npn \regex_item_dot:
196   {
197     \if_num:w \l_regex_current_char_int > - \c_two
198       \exp_after:wN \regex_break_true:w
199     \fi:
200   }
(End definition for \regex_item_dot:. This function is documented on page ??.)
```

`\regex_item_caseful_equal:n` Simple comparisons triggering `\regex_break_true:w` when true.

```

201 \cs_new_protected:Npn \regex_item_caseful_equal:n #1
202   {
203     \if_num:w #1 = \l_regex_current_char_int
204       \exp_after:wN \regex_break_true:w
205     \fi:
206   }
207 \cs_new_protected:Npn \regex_item_caseful_range:nn #1 #2
208   {
209     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
210       \reverse_if:N \if_num:w #2 < \l_regex_current_char_int
211         \exp_after:wN \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
212     \fi:
```

```

213     \fi:
214   }
215 \cs_new_protected:Npn \regex_item_caseful_geq:n #1
216   {
217     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
218       \exp_after:wN \regex_break_true:w
219     \fi:
220   }
(End definition for \regex_item_caseful_equal:n. This function is documented on page ??.)
```

\regex_item_caseless_equal:n For caseless matching, we perform the test both on \l_regex_current_char_int and on \l_regex_case_changed_char_int.

```

221 \cs_new_protected:Npn \regex_item_caseless_equal:n #1
222   {
223     \if_num:w #1 = \l_regex_current_char_int
224       \exp_after:wN \regex_break_true:w
225     \fi:
226     \if_num:w #1 = \l_regex_case_changed_char_int
227       \exp_after:wN \regex_break_true:w
228     \fi:
229   }
230 \cs_new_protected:Npn \regex_item_caseless_range:nn #1 #2
231   {
232     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
233       \reverse_if:N \if_num:w #2 < \l_regex_current_char_int
234         \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
235       \fi:
236     \fi:
237     \reverse_if:N \if_num:w #1 > \l_regex_case_changed_char_int
238       \reverse_if:N \if_num:w #2 < \l_regex_case_changed_char_int
239         \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
240       \fi:
241     \fi:
242   }
243 \cs_new_protected:Npn \regex_item_caseless_geq:n #1
244   {
245     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
246       \exp_after:wN \regex_break_true:w
247     \fi:
248     \reverse_if:N \if_num:w #1 > \l_regex_case_changed_char_int
249       \exp_after:wN \regex_break_true:w
250     \fi:
251   }
(End definition for \regex_item_caseless_equal:n. This function is documented on page ??.)
```

\regex_item_equal:n By default, matching takes the letter case into account. Note that those functions are not protected: they will expand at the building step, hard-coding which states take care of caseless versus caseful matching.

```
252 \cs_new:Npn \regex_item_equal:n { \regex_item_caseful_equal:n }
```

```

253 \cs_new:Npn \regex_item_range:nn { \regex_item_caseful_range:nn }
254 \cs_new:Npn \regex_item_geq:n    { \regex_item_caseful_geq:n }
(End definition for \regex_item_equal:n. This function is documented on page ??.)
```

\regex_build_caseless: Switch between caseful and caseless matching. This is only done during the building step.

```

255 \cs_new_protected_nopar:Npn \regex_build_caseless:
256 {
257     \bool_set_true:N \l_regex_caseless_bool
258     \cs_set:Npn \regex_item_equal:n { \regex_item_caseless_equal:n }
259     \cs_set:Npn \regex_item_range:nn { \regex_item_caseless_range:nn }
260     \cs_set:Npn \regex_item_geq:n { \regex_item_caseless_geq:n }
261 }
262 \cs_new_protected_nopar:Npn \regex_build_caseful:
263 {
264     \bool_set_false:N \l_regex_caseless_bool
265     \cs_set:Npn \regex_item_equal:n { \regex_item_caseful_equal:n }
266     \cs_set:Npn \regex_item_range:nn { \regex_item_caseful_range:nn }
267     \cs_set:Npn \regex_item_geq:n { \regex_item_caseful_geq:n }
268 }
```

(End definition for \regex_build_caseless: and \regex_build_caseful:. These functions are documented on page ??.)

\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.

```

269 \cs_new_protected:Npn \regex_item_catcode:nT #1
270 {
271     \if_int_odd:w \int_eval:w #1 /
272         \if_case:w \l_regex_current_catcode_int
273             1          \or: 4          \or: 10          \or: 40
274             \or: 100        \or:           \or: 1000        \or: 4000
275             \or: 10000      \or:           \or: 100000      \or: 400000
276             \or: 1000000    \or: 4000000 \else: 1*\c_zero
277             \fi:
278         \int_eval_end:
279         \exp_after:wN \use:n
280     \else:
281         \exp_after:wN \use_none:n
282     \fi:
283 }
```

(End definition for \regex_item_catcode:nT. This function is documented on page ??.)

\regex_item_cs:n Match a control sequence (the argument is a pre-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 \l_regex_current_token_tl (of the form \exp_not:n {\<control sequence>}) to <control sequence>.

```

284 \cs_new_protected:Npn \regex_item_cs:n #1
285   {
286     \int_compare:nNnT \l_regex_current_catcode_int = \c_zero
287     {
288       \tl_set:Nn \l_regex_tmpa_regex_tl {#1}
289       \bool_set_eq:NN \l_regex_saved_success_bool \g_regex_success_bool
290       \exp_args:NNx \regex_match:NnTF \l_regex_tmpa_regex_tl
291       {
292         \exp_after:wN \exp_after:wN
293         \exp_after:wN \cs_to_str:N \l_regex_current_token_tl
294       }
295       {
296         \bool_gset_eq:NN \g_regex_success_bool \l_regex_saved_success_bool
297         \regex_break_true:w
298       }
299       { \bool_gset_eq:NN \g_regex_success_bool \l_regex_saved_success_bool }
300     }
301   }
302
(End definition for \regex_item_cs:n. This function is documented on page ??.)
```

2.2.5 Grabbing digits

\regex_get_digits:nw Grabs digits (of category code other), skipping any intervening space, until encountering a non-digit, and places the result in a brace group after #1. This is used when parsing the { quantifier.

```

302 \cs_new_protected:Npn \regex_get_digits:nw #1
303   {
304     \tex_afterassignment:D \regex_tmp:w
305     \cs_set_nopar:Npx \regex_tmp:w
306     {
307       \exp_not:n {#1}
308       { \if_false: } } \fi:
309       \regex_get_digits_aux:NN
310     }
311 \cs_new:Npn \regex_get_digits_aux:NN #1#2
312   {
313     \if_meaning:w \regex_build_raw:N #1
314     \if_charcode:w \c_space_token \exp_not:N #2
315     \else:
316       \if_num:w 9 < 1 \exp_not:N #2 \exp_stop_f:
317         #2
318       \else:
319         \regex_get_digits_end:w #1 #2
320       \fi:
321     \fi:
322   \else:
323     \regex_get_digits_end:w #1 #2
324   \fi:
325   \regex_get_digits_aux:NN
```

```

326    }
327 \cs_new:Npn \regex_get_digits_end:w #1 \fi: #2 \regex_get_digits_aux:NN
328  {
329    \fi: #2
330    \if_false: { { \fi: } }
331    #1
332  }
(End definition for \regex_get_digits:nw. This function is documented on page ??.)
```

2.2.6 More character testing

\regex_token_if_other_digit:N In the replacement text, `\g{i}{int}` denotes the `{int}`-th submatch. Parsing this construction robustly requires a test of whether a token is a digit or not.

```

333 \prg_new_conditional:Npnn \regex_token_if_other_digit:N #1 { TF }
334  {
335    \if_num:w \c_nine < 1 \exp_not:N #1 \exp_stop_f:
336      \prg_return_true: \else: \prg_return_false: \fi:
337  }
(End definition for \regex_token_if_other_digit:N. This function is documented on page ??.)
```

\regex_aux_char_if_alphanumeric: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,

- alphanumerics 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).

```

338 \prg_new_conditional:Npnn \regex_aux_char_if_special:N #1 { TF }
339  {
340    \if_num:w '#1 < \c_ninety_one
341      \if_num:w '#1 < \c_fifty_eight
342        \if_num:w '#1 < \c_forty_eight
343          \if_num:w '#1 < \c_thirty_two
344            \prg_return_false: \else: \prg_return_true: \fi:
345          \else: \prg_return_false: \fi:
346        \else:
347          \if_num:w '#1 < \c_sixty_five
348            \prg_return_true: \else: \prg_return_false: \fi:
349          \fi:
```

```

350     \else:
351         \if_num:w '#1 < \c_one_hundred_twenty_three
352             \if_num:w '#1 < \c_ninety_seven
353                 \prg_return_true: \else: \prg_return_false: \fi:
354             \else:
355                 \if_num:w '#1 < \c_one_hundred_twenty_seven
356                     \prg_return_true: \else: \prg_return_false: \fi:
357                 \fi:
358             \fi:
359     }
360 \prg_new_conditional:Npnn \regex_aux_char_if_alphanumeric:N #1 { TF }
361 {
362     \if_num:w '#1 < \c_ninety_one
363         \if_num:w '#1 < \c_fifty_eight
364             \if_num:w '#1 < \c_forty_eight
365                 \prg_return_false: \else: \prg_return_true: \fi:
366             \else:
367                 \if_num:w '#1 < \c_sixty_five
368                     \prg_return_false: \else: \prg_return_true: \fi:
369                 \fi:
370             \else:
371                 \if_num:w '#1 < \c_one_hundred_twenty_three
372                     \if_num:w '#1 < \c_ninety_seven
373                         \prg_return_false: \else: \prg_return_true: \fi:
374                     \else:
375                         \prg_return_false:
376                     \fi:
377                 \fi:
378 }
(End definition for \regex_aux_char_if_alphanumeric:NTF. This function is documented on page ??.)
```

2.3 Building

2.3.1 Build mode

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

- 6 $[\backslash c\{ \dots \}]$ control sequence in a class,
- 2 $\backslash c\{ \dots \}$ control sequence,
- 0 ... outer,
- 2 $\backslash c\dots$ catcode test,
- 6 $[\backslash c\dots]$ catcode test in a class,
- 63 $[\backslash c\{[\dots]\}]$ 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_build_if_in_class:TF Test whether we are currently in a character class (at the inner-most 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.

```

379 \cs_new_nopar:Npn \regex_build_if_in_class:TF
380 {
381   \if_int_odd:w \l_regex_build_mode_int
382     \exp_after:wN \use_i:nn
383   \else:
384     \exp_after:wN \use_ii:nn
385   \fi:
386 }

```

(End definition for \regex_build_if_in_class:TF. This function is documented on page ??.)

2.3.2 Helpers for building an NFA

\regex_build_new_state: Here, we add a new state to the NFA. At the end of the building phase, we want every \toks register to start with \s_stop, hence initialize the new register appropriately. Then set \l_regex_left/right_state_int to their new values.

```

387 \cs_new_protected:Npn \regex_build_new_state:
388 {
389     \tex_toks:D \l_regex_max_state_int { \s_stop }
390     \int_set_eq:NN \l_regex_left_state_int \l_regex_right_state_int
391     \int_set_eq:NN \l_regex_right_state_int \l_regex_max_state_int
392     \int_incr:N \l_regex_max_state_int
393 }
(End definition for \regex_build_new_state:. This function is documented on page ??.)
```

\regex_build_transition>NN These functions create a new state, and put one or two transitions starting from the old current state.

```

394 \cs_new_protected:Npn \regex_build_transition:NN #1#
395 {
396     \regex_build_new_state:
397     \regex_toks_put_right:Nx \l_regex_left_state_int
398     { #1 { \int_eval:n { #2 - \l_regex_left_state_int } } }
399 }
400 \cs_new_protected:Npn \regex_build_transitions:NNNN #1#2#3#4
401 {
402     \regex_build_new_state:
403     \regex_toks_put_right:Nx \l_regex_left_state_int
404     {
405         #1 { \int_eval:n { #2 - \l_regex_left_state_int } }
406         #3 { \int_eval:n { #4 - \l_regex_left_state_int } }
407     }
408 }
(End definition for \regex_build_transition:NN. This function is documented on page ??.)
```

2.3.3 From regex to NFA: framework

\regex_build:n First, reset a few variables. Then use the generic framework defined in \l3str to parse the regular expression once, recognizing which characters are raw characters, and which have special meanings. The search is not anchored: to achieve that, we insert state(s) responsible for repeating the match attempt on every token in the token list. The trailing \prg_do_nothing: ensure that the look-ahead done by some of the operations is harmless. Finally, \regex_build_end: adds the finishing code (checking that parentheses are properly nested, for instance).

```

409 \cs_new_protected:Npn \regex_build:n #1
410 {
411     \int_set_eq:NN \l_regex_build_mode_int \c_zero
412     \regex_build:w
413     \regex_build_new_state:
414     \regex_toks_put_right:Nx \l_regex_left_state_int
```

```

415     { \regex_action_start_wildcard: }
416   \regex_build_open_aux:
417     \str_escape_use:NNNn
418       \regex_build_i_unescaped:N
419       \regex_build_i_escaped:N
420       \regex_build_i_raw:N
421       { #1 }
422     \prg_do_nothing: \prg_do_nothing:
423     \prg_do_nothing: \prg_do_nothing:
424   \regex_build_end:
425 }

```

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

`\regex_build_i_unescaped:N` The `\I3str` function `\str_escape_use:NNNn` goes through the regular expression and finds the `\a`, `\e`, `\f`, `\n`, `\r`, `\t`, and `\x` escape sequences, then distinguishes three cases: non-escaped characters, escaped characters, and “raw” characters coming from one of the escape sequences. In the particular case of regular expressions, escaped alphanumerics and non-escaped non-alphanumeric printable ascii characters may have special meanings, while everything else should be treated as a raw character.

```

426 \cs_new:Npn \regex_build_i_unescaped:N #1
427 {
428   \regex_aux_char_if_special:NTF #1
429   { \exp_not:N \regex_build_special:N #1 }
430   { \exp_not:N \regex_build_raw:N #1 }
431 }
432 \cs_new:Npn \regex_build_i_escaped:N #1
433 {
434   \regex_aux_char_if_alphanumeric:NTF #1
435   { \exp_not:N \regex_build_escaped:N #1 }
436   { \exp_not:N \regex_build_raw:N #1 }
437 }
438 \cs_new:Npn \regex_build_i_raw:N #1
439 { \exp_not:N \regex_build_raw:N #1 }

```

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

`\regex_build_special:N` If the control character 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 in character classes, unknown escaped alphanumeric characters raise an error, while special characters are silently converted to raw characters.

```

440 \cs_new_protected:Npn \regex_build_special:N #1
441 {
442   \cs_if_exist_use:cF { \regex_build_#1: }
443   { \regex_build_raw:N #1 }
444 }
445 \cs_new_protected:Npn \regex_build_escaped:N #1
446 {
447   \cs_if_exist_use:cF { \regex_build_/#1: }

```

```

448     { \regex_build_raw:N #1 }
449   }
(End definition for \regex_build_special:N. This function is documented on page ??.)
```

\regex_build:w Hopefully, we didn't forget to initialize anything here.

```

450 \cs_new_protected_nopar:Npn \regex_build:w
451   {
452     \int_set_eq:NN \l_regex_catcodes_default_int \c_regex_catcodes_all_int
453     \int_set_eq:NN \l_regex_capturing_group_int \c_zero
454     \int_zero:N \l_regex_max_state_int
455     \regex_build_new_state:
456     \tl_clear:N \l_regex_class_tl
457     \bool_set_true:N \l_regex_class_bool
458     \int_set_eq:NN \l_regex_catcodes_int \l_regex_catcodes_default_int
459   }
(End definition for \regex_build:w. This function is documented on page ??.)
```

\regex_build_end: If parentheses are not nested properly, an error is raised, and the correct number of parentheses is closed. After that, we insert an instruction for the match to succeed.

```

460 \cs_new_protected_nopar:Npn \regex_build_end:
461   {
462     \regex_seq_push_int:NN \l_regex_capturing_group_seq \c_zero
463     \regex_build_close_aux: \regex_build_group_:
464     \seq_if_empty:NF \l_regex_capturing_group_seq
465     {
466       \msg_kernel_error:nnx { regex } { missing-rparen }
467       { \seq_length:N \l_regex_capturing_group_seq }
468       \prg_replicate:nn
469       { \seq_length:N \l_regex_capturing_group_seq }
470       { \regex_build_close_aux: \regex_build_group_: }
471     }
472     \regex_toks_put_right:Nx \l_regex_right_state_int
473     { \regex_action_success: }
474     \int_incr:N \l_regex_capturing_group_int
475   }
(End definition for \regex_build_end:. This function is documented on page ??.)
```

\regex_build_one:n In a class, add the argument to the current class. Outside a class, this argument is the whole “class”, and we look for quantifiers.

```

476 \cs_new_protected:Npn \regex_build_one:n #1
477   { \regex_build_one:x { \exp_not:n {#1} } }
478 \cs_new_protected:Npn \regex_build_one:x #1
479   {
480     \tl_put_right:Nx \l_regex_class_tl
481     {
482       \if_num:w \l_regex_catcodes_int = \c_regex_catcodes_all_int
483         \exp_after:wN \use:n
484       \else:
485         \regex_item_catcode:nT { \int_use:N \l_regex_catcodes_int }
```

```

486           \fi:
487           {#1}
488       }
489   \if_num:w \l_regex_build_mode_int = \c_two
490     \l_regex_build_mode_int = \c_zero
491   \else:
492     \if_num:w \l_regex_build_mode_int = \c_six
493       \l_regex_build_mode_int = \c_three
494     \fi:
495   \fi:
496   \if_int_odd:w \l_regex_build_mode_int \else:
497     \exp_after:wN \regex_build_one_quantifier:
498   \fi:
499 }
(End definition for \regex_build_one:n and \regex_build_one:x. These functions are documented on page ??.)
```

\regex_build_tmp_class:n The argument is the target state if the test succeeds.

```

500 \cs_new:Npn \regex_build_tmp_class:n #1
501 {
502   \exp_not:o \l_regex_class_tl
503   \bool_if:NTF \l_regex_class_bool
504     { \regex_break_point:TF { \regex_action_cost:n {#1} } { } }
505     { \regex_break_point:TF { } { \regex_action_cost:n {#1} } }
506 }
(End definition for \regex_build_tmp_class:n. This function is documented on page ??.)
```

\regex_class_put:N This function is used to rapidly put a character in a character class being built, avoiding some tests implied by \regex_build_one:x.

```

507 \cs_new_protected:Npn \regex_class_put:N #1
508 {
509   \tl_put_right:Nx \l_regex_class_tl
510   { \regex_item_equal:n { \int_value:w '#1 ~ } }
511 }
512 \cs_new_protected:Npn \regex_class_put:NN #1#2
513 {
514   \tl_put_right:Nx \l_regex_class_tl
515   { \regex_item_range:nn { \int_value:w '#1 ~ } { \int_value:w '#2 ~ } }
516 }
(End definition for \regex_class_put:N. This function is documented on page ??.)
```

\regex_build_raw_alphanum_is_error:N Within character classes, some special character and escape sequences do not have any meaning. Special characters are interpreted as their “raw” counterpart, but escape sequences without a meaning, e.g., \A, are most likely an error: this should be A.

```

517 \cs_new_protected:Npn \regex_build_raw_alphanum_is_error:N #1
518 {
519   \regex_aux_char_if_alphanumeric:NTF #1
520   { \msg_kernel_error:nnx { regex } { class-bad-escape } { #1 } }
521   { \regex_build_raw:N #1 }
522 }
```

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

\regex_build_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.

```
523 \cs_new_protected:Npn \regex_build_raw:N #1#2#3
  {
  525   \regex_build_if_in_class:TF
    { \str_if_eq:nnTF {#2#3} { \regex_build_special:N - } }
    { \use_i:nn }
    { \regex_class_range:Nw #1 }
    {
      \regex_build_one:x { \regex_item_equal:n { \int_value:w '#1 ~ } }
      #2 #3
    }
  }
  533 }
```

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

\regex_class_range:Nw We have just read a raw character followed by a dash; this should be followed by an end-point for the range. If the following character is an escaped alphanumeric, or if it is an unescaped right bracket, then we have an error, so we put the initial character and the dash back, as raw characters. Otherwise, build the range, checking that it is in the right order, and optimizing for equal end-points.

```
534 \cs_new_protected:Npn \regex_class_range:Nw #1#2#3
  {
  535   \token_if_eq_meaning:NNTF #2 \regex_build_escaped:N % [
    { \use_i:nn } { \str_if_eq:nnTF {#2#3} { \regex_build_special:N } }
    {
      \msg_kernel_warning:nnx { regex } { end-range-missing } % [
        { #1 - \if_meaning:w #3 } \else: \c_backslash_str \fi: #3
      \regex_class_put:N #1
      \regex_class_put:N -
      #2#3
    }
  545   {
    \if_num:w '#1 > '#3 \exp_stop_f:
      \msg_kernel_error:nnxx { regex } { backwards-range } {#1} {#3}
    \else:
      \if_num:w '#1 = '#3 \exp_stop_f:
        \regex_class_put:N #3
      \else:
        \regex_class_put>NN #1#3
      \fi:
    \fi:
  555   }
  556 }
```

(End definition for \regex_class_range:Nw. This function is documented on page ??.)

2.3.4 Character properties

\regex_build_:: In a class, the dot has no special meaning. Outside, insert \regex_item_dot:, which matches any character or control sequence, and refuses -2, which marks the end of the token list.

```

557 \cs_new_protected_nopar:cpn { regex_build_:: }
558 {
559     \regex_build_if_in_class:TF
560     { \regex_build_raw:N . }
561     { \regex_build_one:x \regex_item_dot: }
562 }
```

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

\regex_build_/d: The constants \c_regex_d_tl, 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.

```

563 \tl_map_inline:nn { dDhHsSvVwWN }
564 {
565     \cs_new_protected_nopar:cpn { regex_build_/#1: }
566     { \regex_build_one:n \exp_not:c { c_regex_#1_tl } }
567 }
```

(End definition for \regex_build_/d: and \regex_build_/D:. These functions are documented on page ??.)

\regex_build_/W:

\regex_build_/N:

2.3.5 Anchoring and simple assertions

\regex_build_simple_assertion:Nn Assertions are not allowed within character classes: the raw character #1 is inserted instead. Otherwise, use the test #2; if the assertion is successful, move from the left state to the right state.

```

568 \cs_new_protected:Npn \regex_build_simple_assertion:Nn #1#2
569 {
570     \regex_build_if_in_class:TF
571     { \regex_build_raw_alphanum_is_error:N #1 }
572     {
573         \regex_build_new_state:
574         \regex_toks_put_right:Nx \l_regex_left_state_int
575         {
576             \exp_not:n {#2}
577             {
578                 \regex_action_free:n
579                 {
580                     \int_eval:n
581                     { \l_regex_right_state_int - \l_regex_left_state_int }
582                 }
583             }
584         }
585         %^^A \regex_assertion_quantifier:
586     }
587 }
```

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

\regex_build_^\!: Anchoring at the start corresponds to checking that the current character is the first in
\regex_build_/A\: Anchoring to the beginning of the match attempt uses \l_regex_start_step_int instead of \c_zero. End anchors match the end of the token list, marked by
\regex_build_\$\: a character code of -2.
\regex_build_/Z:
588 \cs_new_protected_nopar:cpn { regex_build_^\!: }
\regex_build/_z:
589 {
590 \regex_build_simple_assertion:Nn ^
591 { \int_compare:nNnT \l_regex_min_step_int = \l_regex_current_step_int }
592 }
593 \cs_new_protected_nopar:cpn { regex_build_/A\: }
594 {
595 \regex_build_simple_assertion:Nn A
596 { \int_compare:nNnT \l_regex_min_step_int = \l_regex_current_step_int }
597 }
598 \cs_new_protected_nopar:cpn { regex_build_/G\: }
599 {
600 \regex_build_simple_assertion:Nn G
601 { \int_compare:nNnT \l_regex_start_step_int = \l_regex_current_step_int }
602 }
603 \cs_new_protected_nopar:cpn { regex_build_\$\: } % \$
604 {
605 \regex_build_simple_assertion:Nn \$ % \$
606 { \int_compare:nNnT \l_regex_current_char_int < \c_minus_one }
607 }
608 \cs_new_protected_nopar:cpn { regex_build_/Z\: }
609 {
610 \regex_build_simple_assertion:Nn Z
611 { \int_compare:nNnT \l_regex_current_char_int < \c_minus_one }
612 }
613 \cs_new_protected_nopar:cpn { regex_build/_z\: }
614 {
615 \regex_build_simple_assertion:Nn z
616 { \int_compare:nNnT \l_regex_current_char_int < \c_minus_one }
617 }
(End definition for \regex_build_^\!: This function is documented on page ??.)

\regex_build_/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. Case sensitivity does not change word boundaries.
\regex_build_/B:
\regex_if_word_boundary:TF

618 \cs_new_protected_nopar:cpn { regex_build_/b\: }
619 {
620 \regex_build_simple_assertion:Nn b
621 { \regex_if_word_boundary:TF \use:n \use_none:n }
622 }
623 \cs_new_protected_nopar:cpn { regex_build_/B\: }
624 {
625 \regex_build_simple_assertion:Nn B

```

626          { \regex_if_word_boundary:TF { } }
627      }
628 \cs_new_protected_nopar:Npn \regex_if_word_boundary:TF
629  {
630     \group_begin:
631         \int_set_eq:NN \l_regex_current_char_int \l_regex_last_char_int
632         \c_regex_w_tl
633         \regex_break_point:TF
634         { \group_end: \c_regex_W_tl \regex_item_equal:n { -2 } }
635         { \group_end: \c_regex_w_tl }
636         \regex_break_point:TF
637     }
638
(End definition for \regex_build_/:. This function is documented on page ??.)
```

2.3.6 Entering and exiting character classes

\regex_build_[:]: In a class, left brackets mean nothing. Outside a class, this starts a class, whose first characters may have a special meaning.

```

638 \cs_new_protected_nopar:cpn { regex_build_[:] }
639  {
640      \regex_build_if_in_class:TF
641      { \regex_build_raw:N [ }
642      { \regex_class_first:NNNN }
643  }
644
(End definition for \regex_build_[:]. This function is documented on page ??.)
```

\regex_build_[]: Outside a class, right brackets have no meaning. In a class, change the mode ($m \rightarrow (m - 15)/13$, truncated) to reflect the fact that we are leaving the class. If we are still in a class after leaving one, then this is the case [...\cL[...]]..., and we insert the relevant closing material in \l_regex_class_tl. Otherwise look for quantifiers.

```

644 \cs_new_protected:cpn { regex_build_[]: }
645  {
646      \regex_build_if_in_class:TF
647      {
648          \if_num:w \l_regex_build_mode_int > \c_sixteen
649              \tl_set:Nx \l_regex_class_tl
650              {
651                  \exp_not:o \l_regex_class_saved_tl
652                  \if_num:w \l_regex_catcodes_int < \c_regex_catcodes_all_int
653                      \regex_item_catcode:nT { \int_use:N \l_regex_catcodes_int }
654                  \else:
655                      \exp_after:wN \use:n
656                  \fi:
657                  {
658                      \exp_not:o \l_regex_class_tl
659                      \bool_if:NF \l_regex_class_bool
660                          { \regex_break_point:TF { } { \regex_break_true:w } }
661                  }
662  }
```

```

663           \bool_set_eq:NN \l_regex_class_bool \l_regex_class_saved_bool
664   \fi:
665   \tex_advance:D \l_regex_build_mode_int - \c_fifteen
666   \tex_divide:D \l_regex_build_mode_int \c_thirteen
667   \if_int_odd:w \l_regex_build_mode_int \else:
668     \exp_after:wN \regex_build_one_quantifier:
669   \fi:
670 }
671 { \regex_build_raw:N ]
672 }

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

```

\regex_class_first>NNNN This starts a class. Change the mode by appending 3 to it, and reset the variables \l_regex_class_tl and \l_regex_bool_tl. In the special case of mode 63 ([\c[...]]), we open a group, to avoid overriding the setting of \l_regex_class_bool and \l_regex_class_tl; the group ends at the matching right bracket. If the first character is ^, then the class is inverted. We keep track of this in \l_regex_class_bool. If the next character is a right bracket, then it should be changed to a raw one (dirty hack here; the F argument of \str_if_eq:nnTF is the trailing #3).

```

673 \cs_new_protected:Npn \regex_class_first>NNNN #1#2#3#4
674 {
675   \l_regex_build_mode_int = \int_value:w \l_regex_build_mode_int 3 ~ %
676   \if_num:w \l_regex_build_mode_int > \c_sixteen
677     \tl_set_eq:NN \l_regex_class_saved_tl \l_regex_class_tl
678     \bool_set_eq:NN \l_regex_class_saved_bool \l_regex_class_bool
679   \fi:
680   \token_if_eq_meaning:NNTF #1 \regex_build_special:N
681   {
682     \token_if_eq_charcode:NNTF #2 ^
683     {
684       \bool_set_false:N \l_regex_class_bool % [
685       \str_if_eq:nnTF {#3#4} { \regex_build_special:N ] }
686       { \regex_build_raw:N }
687     }
688     { % [
689       \token_if_eq_charcode:NNTF #2 ]
690       { \regex_build_raw:N #2 }
691       { #1 #2 }
692     }
693   }
694   { #1 #2 }
695   #3 #4
696 }

(End definition for \regex_class_first>NNNN. This function is documented on page ??.)

```

2.3.7 Catcodes and csnames

\regex_build_/_c: The \c escape sequence is only allowed in modes 0 and 3, *i.e.*, not within any other \c escape sequence.

```

697 \cs_new_protected:cpn { regex_build_/c: }
698   {
699     \if_case:w \l_regex_build_mode_int
700       \exp_after:wN \regex_build_c_aux:wNN
701     \or: \or: \or: \exp_after:wN \regex_build_c_aux:wNN
702     \fi:
703     \msg_kernel_error:nn { regex } { c-bad-mode }
704     \s_stop
705   }

```

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

`\regex_build_c_aux:wNN` 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.

```

706 \cs_new_protected:Npn \regex_build_c_aux:wNN #1 \s_stop #2#3
707   {
708     \token_if_eq_meaning:NNTF #2 \regex_build_raw:N
709     {
710       \cs_if_exist:cTF { c_regex_catcode_#3_int }
711       {
712         \int_set_eq:Nc \l_regex_catcodes_int { c_regex_catcode_#3_int }
713         \l_regex_build_mode_int
714         = \if_case:w \l_regex_build_mode_int \c_two \else: \c_six \fi:
715       }
716     }
717     { \cs_if_exist_use:cF { regex_build_c_#3: } }
718     {
719       \msg_kernel_error:nnxx { regex } { c-command }
720       { regular-expression } { #3 }
721       #2 #3
722     }
723   }

```

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

`\regex_build_c_{}`: The case of a left brace is easy, based on what we have done so far: in a group, build the regular expression, after changing the mode to forbid nesting `\c`.

```

724 \cs_new_protected:cpn { regex_build_c_ \c_lbrace_str : }
725   {
726     \group_begin:
727     \l_regex_build_mode_int
728     = - \if_case:w \l_regex_build_mode_int \c_two \else: \c_six \fi:
729     \regex_build:w
730   }

```

(End definition for `\regex_build_c_{}`. This function is documented on page ??.)

```

\regex_build_c_[:]
\regex_build_c_lbrack_loop:NN
\regex_build_c_lbrack_end:
\regex_build_add:N
731 \cs_new_protected:cpn { regex_build_c_[:] } #1#2
732   {
733     \int_zero:N \l_regex_catcodes_int

```

```

734 \str_if_eq:nnTF { #1 #2 } { \regex_build_special:N ^ }
735 {
736     \bool_set_false:N \l_regex_catcodes_bool
737     \regex_build_c_lbrack_loop:NN
738 }
739 {
740     \bool_set_true:N \l_regex_catcodes_bool
741     \regex_build_c_lbrack_loop:NN
742     #1#2
743 }
744 }
745 \cs_new_protected:Npn \regex_build_c_lbrack_loop:NN #1#2
746 {
747     \token_if_eq_meaning:NNTF #1 \regex_build_raw:N
748     {
749         \cs_if_exist:cTF { c_regex_catcode_#2_int }
750         {
751             \exp_args:Nc \regex_build_c_lbrack_add:N
752             { c_regex_catcode_#2_int }
753             \regex_build_c_lbrack_loop:NN
754         }
755     }
756     { % [
757         \token_if_eq_charcode:NNTF #2 ]
758         { \regex_build_c_lbrack_end: }
759     }
760     {
761         \msg_kernel_error:nnxx { regex } { c-command }
762         { regular-expression } { #2 }
763         \regex_build_c_lbrack_end:
764     }
765 }
766 \cs_new_protected_nopar:Npn \regex_build_c_lbrack_end:
767 {
768     \l_regex_build_mode_int
769     = \if_case:w \l_regex_build_mode_int \c_two \else: \c_six \fi:
770     \if_meaning:w \c_false_bool \l_regex_catcodes_bool
771     \int_set:Nn \l_regex_catcodes_int
772     { \c_regex_catcodes_all_int - \l_regex_catcodes_int }
773     \fi:
774 }
775 \cs_new_protected:Npn \regex_build_c_lbrack_add:N #1
776 {
777     \if_int_odd:w \int_eval:w \l_regex_catcodes_int / #1 \int_eval_end:
778     \else:
779         \tex_advance:D \l_regex_catcodes_int #1
780         \fi:
781 }
(End definition for \regex_build_c_[:. This function is documented on page ??.)
```

\regex_build}: Non-escaped right braces are only special if they appear when building the regular expression for a csname. Otherwise, replace the brace with an escaped brace.

```

782 \cs_new_protected:cpn { regex_build_ \c_rbrace_str : }
783   { %^A incorrect recovery for \c{...(...[...] with missing ) or ]}.
784     \int_compare:nNnTF \l_regex_build_mode_int < \c_zero
785     {
786       \use:c { regex_build_/Z: }
787       \regex_toks_put_right:Nx \l_regex_right_state_int
788         { \regex_action_success: }
789       \int_incr:N \l_regex_capturing_group_int
790       \use:x
791       {
792         \group_end:
793         \regex_build_one:n
794           { \regex_item_cs:n { \regex_set_aux:N ? } }
795       }
796     }
797   { \exp_after:wN \regex_build_raw:N \c_rbrace_str }
798 }
```

(End definition for \regex_build}.. This function is documented on page ??.)

2.3.8 Quantifiers

\regex_build_quantifier:w This looks ahead and finds any quantifier (control character equal to either of ?+*{}). When all characters for the quantifier are found, the corresponding function is called.

```

799 \cs_new_protected:Npn \regex_build_quantifier:w #1#2
800   {
801     \token_if_eq_meaning:NNTF #1 \regex_build_special:N
802     {
803       \cs_if_exist_use:cF { regex_build_quantifier_#2:w }
804       {
805         \regex_build_quantifier_end:nn { } { }
806         #1 #2
807       }
808     }
809   {
810     \regex_build_quantifier_end:nn { } { }
811     #1 #2
812   }
813 }
```

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

\regex_build_quantifier_?:w For each “basic” quantifier, ?, *, +, feed the correct arguments to \regex_build_quantifier_aux:nnNN.

\regex_build_quantifier_*:w

```

814 \cs_new_protected_nopar:cpn { regex_build_quantifier_?:w }
815   { \regex_build_quantifier_aux:nnNN { } { ? } }
816 \cs_new_protected_nopar:cpn { regex_build_quantifier_*:w }
817   { \regex_build_quantifier_aux:nnNN { } { * } }
```

```

818 \cs_new_protected_nopar:cpn { regex_build_quantifier_+:+w }
819   { \regex_build_quantifier_aux:nnNN { } { + } }
(End definition for \regex_build_quantifier_?:w. This function is documented on page ??.)
```

\regex_build_quantifier_aux:nnNN Once the “main” quantifier (?*, + or a braced construction) is found, we check whether it is lazy (followed by a question mark), and calls the appropriate function. Here #1 holds some extra arguments that the final function needs in the case of braced constructions, and is empty otherwise.

```

820 \cs_new_protected:Npn \regex_build_quantifier_aux:nnNN #1#2#3#4
821   {
822     \str_if_eq:nnTF { #3 #4 } { \regex_build_special:N ? }
823     { \regex_build_quantifier_end:nn { #2 #4 } {#1} }
824     {
825       \regex_build_quantifier_end:nn { #2 } {#1}
826       #3 #4
827     }
828   }
(End definition for \regex_build_quantifier_aux:nnNN. This function is documented on page ??.)
```

\regex_build_quantifier_{:w} Three possible syntaxes: {<int>}, {<int>,}, or {<int>,<int>}.

```

829 \cs_new_protected_nopar:cpn { regex_build_quantifier_ \c_lbrace_str :w }
830   { \regex_get_digits:nw { \regex_build_quantifier_lbrace:n } }
831 \cs_new_protected:Npn \regex_build_quantifier_lbrace:n #
832   {
833     \tl_if_empty:nTF {#1}
834     {
835       \regex_build_quantifier_end:nn { } { }
836       \exp_after:wN \regex_build_raw:N \c_lbrace_str
837     }
838     { \regex_build_quantifier_lbrace:nw {#1} }
839   }
840 \cs_new_protected:Npx \regex_build_quantifier_lbrace:nw #1#2#3
841   {
842     \exp_not:N \prg_case_str:nnn { #2 #3 }
843     {
844       { \exp_not:N \regex_build_special:N , }
845       {
846         \exp_not:N \regex_get_digits:nw
847         { \exp_not:N \regex_build_quantifier_lbrace:nnw {#1} }
848       }
849       { \exp_not:N \regex_build_special:N \c_rbrace_str }
850       { \exp_not:N \regex_build_quantifier_end:nn {n} { {#1} } }
851     }
852     {
853       \exp_not:N \regex_build_quantifier_end:nn { } { }
854       \exp_not:N \regex_build_raw:N \c_lbrace_str #1#2
855     }
856   }
857 \cs_new_protected:Npn \regex_build_quantifier_lbrace:nnw #1#2#3#4
```

```

858   {
859     \str_if_eq:xxTF
860       { \exp_not:n {#3#4} }
861       { \exp_not:N \regex_build_special:N \c_rbrace_str }
862     {
863       \tl_if_empty:nTF {#2}
864         { \regex_build_quantifier_aux:nnNN { {#1} } { n* } }
865         { \regex_build_quantifier_aux:nnNN { {#1} {#2} } { nn } }
866     } %^A todo: catch the case m>n
867   {
868     \regex_build_quantifier_end:nn { } { }
869     \use:x
870     {
871       \exp_args:No \tl_map_function:nN
872         { \c_lbrace_str #1 , #2 }
873         \regex_build_raw:N
874     }
875     #3 #4
876   }
877 }

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

\regex_build_quantifier_end:nn When all quantifiers are found, we will call the relevant \regex_build_one/group-⟨quantifiers⟩: function.

```

878 \cs_new_protected:Npn \regex_build_quantifier_end:nn #1#2
879   {
880     \use:c { regex_build_ \l_regex_one_or_group_tl _ #1 : } #2
881     \tl_clear:N \l_regex_class_tl
882     \bool_set_true:N \l_regex_class_bool
883     \int_set_eq:NN \l_regex_catcodes_int \l_regex_catcodes_default_int
884   }

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

2.3.9 Quantifiers for one character or character class

\regex_build_one_quantifier: Used for one single character, or a character class. Contrarily to \regex_build_group-⟨quantifier⟩:, we don’t need to keep track of submatches, and no thread can be created within one repetition, so things are relatively easy.

```

885 \cs_new_protected_nopar:Npn \regex_build_one_quantifier:
886   {
887     \tl_set:Nx \l_regex_one_or_group_tl { one }
888     \regex_build_quantifier:w
889   }

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

\regex_build_one_: If no quantifier is found, then the character or character class should just be built into a transition from the current “right” state to a new state.

```
890 \cs_new_protected_nopar:Npn \regex_build_one_:
```

```

891   {
892     \regex_build_transition:NN
893       \regex_build_tmp_class:n \l_regex_right_state_int
894   }
(End definition for \regex_build_one_:. This function is documented on page ??.)
```

\regex_build_one_?: The two transitions are a costly transition controlled by the character class, and a free transition, both going to a common new state. The only difference between the greedy and lazy operators is the order of transitions.

```

895 \cs_new_protected_nopar:cpn { regex_build_one_?: }
896   {
897     \regex_build_transitions:NNNN
898       \regex_build_tmp_class:n \l_regex_right_state_int
899       \regex_action_free:n      \l_regex_right_state_int
900   }
901 \cs_new_protected_nopar:cpn { regex_build_one_???: }
902   {
903     \regex_build_transitions:NNNN
904       \regex_action_free:n      \l_regex_right_state_int
905       \regex_build_tmp_class:n \l_regex_right_state_int
906   }
(End definition for \regex_build_one_???: This function is documented on page ??.)
```

\regex_build_one_*: Build a costly transition going from the current state to itself, and a free transition moving to a new state.

```

907 \cs_new_protected_nopar:cpn { regex_build_one_*: }
908   {
909     \regex_build_transitions:NNNN
910       \regex_build_tmp_class:n \l_regex_left_state_int
911       \regex_action_free:n      \l_regex_right_state_int
912   }
913 \cs_new_protected_nopar:cpn { regex_build_one_*?: }
914   {
915     \regex_build_transitions:NNNN
916       \regex_action_free:n      \l_regex_right_state_int
917       \regex_build_tmp_class:n \l_regex_left_state_int
918   }
(End definition for \regex_build_one_*?: This function is documented on page ??.)
```

\regex_build_one_+: Build a transition from the current state to a new state, controlled by the character class, then build two transitions from this new state to the original state (for repetition) and to another new state (to move on to the rest of the pattern).

```

919 \cs_new_protected_nopar:cpn { regex_build_one_+: }
920   {
921     \regex_build_one_:
922       \int_set_eq:NN \l_regex_tmpr_int \l_regex_left_state_int
923       \regex_build_transitions:NNNN
924         \regex_action_free:n \l_regex_tmpr_int
```

```

925     \regex_action_free:n \l_regex_right_state_int
926   }
927 \cs_new_protected_nopar:cpn { regex_build_one_+?: }
928 {
929   \regex_build_one_:
930   \int_set_eq:NN \l_regex_tmpr_int \l_regex_left_state_int
931   \regex_build_transitions:NNNN
932   \regex_action_free:n \l_regex_right_state_int
933   \regex_action_free:n \l_regex_tmpr_int
934 }
(End definition for \regex_build_one_+?: This function is documented on page ??.)
```

\regex_build_one_n: This function is called in case the syntax is `{<int>}`. Greedy and lazy operators are identical, since the number of repetitions is fixed. Simply repeat #1 times the effect of \regex_build_one_+.

```

935 \cs_new_protected:Npn \regex_build_one_n: #1
936   { \prg_replicate:nn {#1} { \regex_build_one_+?: } }
937 \cs_new_eq:cN { regex_build_one_n?: } \regex_build_one_n:
(End definition for \regex_build_one_n:. This function is documented on page ??.)
```

\regex_build_one_n*: This function is called in case the syntax is `{<int>,}`.

```

938 \cs_new_protected:cpx { regex_build_one_n*: } #1
939 {
940   \exp_not:N \prg_replicate:nn {#1} { \exp_not:N \regex_build_one_+?: }
941   \exp_not:c { regex_build_one_*?: }
942 }
943 \cs_new_protected:cpx { regex_build_one_n*?: } #1
944 {
945   \exp_not:N \prg_replicate:nn {#1} { \exp_not:N \regex_build_one_+?: }
946   \exp_not:c { regex_build_one_*?: }
947 }
(End definition for \regex_build_one_n*:. This function is documented on page ??.)
```

\regex_build_one_nn: This function is called when the syntax is `{<int>,<int>}`.

```

948 \cs_new_protected:Npn \regex_build_one_nn_aux:Nnn #1#2#3
949 {
950   \prg_replicate:nn {#2} { \regex_build_one_+?: }
951   \prg_replicate:nn {#3-#2} {#1}
952 }
953 \cs_new_protected_nopar:Npx \regex_build_one_nn:
954   { \regex_build_one_nn_aux:Nnn \exp_not:c { regex_build_one_??: } }
955 \cs_new_protected_nopar:cpx { regex_build_one_nn??: }
956   { \regex_build_one_nn_aux:Nnn \exp_not:c { regex_build_one_???: } }
(End definition for \regex_build_one_nn:. This function is documented on page ??.)
```

2.3.10 Groups and alternation

We support the syntax $\langle \text{expr}_1 \rangle \dots \langle \text{expr}_n \rangle$ for alternations.

```

\regex_build_(: Grouping and alternation go together.
\regex_build_:
\regex_build_open_aux:
\regex_build_|:
\regex_build_begin_alternation:
\regex_build_end_alternation:

• Allocate the next available number for the end vertex of the alternation/group and
  store it on a stack (so that nested alternations work).

• Put free transitions to separate all cases of the alternation.

• Build each branch separately, and merge them to the common end-node.

• Test for a quantifier, and if needed, transfer the initial vertex to a new vertex.

957 \cs_new_protected:cpn { regex_build_(:) #1#2
958   {
959     \regex_build_if_in_class:TF
960     {
961       \regex_build_raw:N (
962         #1 #2
963     }
964     {
965       \str_if_eq:nnTF { #1 #2 } { \regex_build_special:N ? }
966       { \regex_build_special_group:NN }
967       {
968         \int_incr:N \l_regex_capturing_group_int
969         \regex_seq_push_int:NN
970           \l_regex_capturing_group_seq \l_regex_capturing_group_int
971         \regex_build_open_aux:
972           #1 #2
973       }
974     }
975   }
976 \cs_new_protected_nopar:Npn \regex_build_open_aux:
977   {
978     \regex_build_new_state:
979     \regex_seq_push_int:NN \l_regex_left_state_seq \l_regex_left_state_int
980     \regex_seq_push_int:NN \l_regex_right_state_seq \l_regex_right_state_int
981     \bool_if:NTF \l_regex_caseless_bool
982       { \seq_push:Nn \l_regex_end_group_seq \regex_build_caseless: }
983       { \seq_push:Nn \l_regex_end_group_seq \regex_build_caseful: }
984     \seq_push:Nn \l_regex_end_alternation_seq { }
985     \regex_build_begin_alternation:
986   }
987 \cs_new_protected_nopar:cpn { regex_build_|: }
988   {
989     \regex_build_if_in_class:TF { \regex_build_raw:N | }
990     {
991       \regex_build_end_alternation:
992       \regex_build_begin_alternation:

```

```

993     }
994   }
995 \cs_new_protected_nopar:cpn { regex_build_}: }
996   {
997     \regex_build_if_in_class:TF { \regex_build_raw:N } }
998   {
999     \seq_if_empty:NTF \l_regex_capturing_group_seq
1000       { \msg_kernel_error:n { regex } { extra-rparen } }
1001     {
1002       \regex_build_close_aux:
1003       \regex_build_group_quantifier:
1004     }
1005   }
1006 }
1007 \cs_new_protected_nopar:Npn \regex_build_close_aux:
1008   {
1009     \regex_build_end_alternation:
1010     \regex_seq_pop_int:NN \l_regex_left_state_seq \l_regex_left_state_int
1011     \regex_seq_pop_int:NN \l_regex_right_state_seq \l_regex_right_state_int
1012     \regex_seq_pop_use:N \l_regex_end_group_seq
1013     \seq_pop:NN \l_regex_end_alternation_seq \l_regex_tmpa_tl
1014   }

```

Building each branch.

```

1015 \cs_new_protected_nopar:Npn \regex_build_begin_alternation:
1016   {
1017     \regex_build_new_state:
1018     \regex_seq_get_int:NN \l_regex_left_state_seq \l_regex_left_state_int
1019     \regex_toks_put_right:Nx \l_regex_left_state_int
1020     {
1021       \regex_action_free:n
1022       {
1023         \int_eval:n
1024           { \l_regex_right_state_int - \l_regex_left_state_int }
1025       }
1026     }
1027   }
1028 \cs_new_protected_nopar:Npn \regex_build_end_alternation:
1029   {
1030     \int_set_eq:NN \l_regex_left_state_int \l_regex_right_state_int
1031     \regex_seq_get_int:NN \l_regex_right_state_seq \l_regex_right_state_int
1032     \regex_toks_put_right:Nx \l_regex_left_state_int
1033     {
1034       \regex_action_free:n
1035       {
1036         \int_eval:n
1037           { \l_regex_right_state_int - \l_regex_left_state_int }
1038       }
1039     }
1040     \regex_seq_get_use:N \l_regex_end_alternation_seq

```

```

1041    }
(End definition for \regex_build_(: and \regex_build_). These functions are documented on
page ??.)
```

\regex_build_special_group:NN Same method as elsewhere: if the combination (?#1 is known, then use that. Otherwise, treat the question mark as if it had been escaped.

```

1042 \cs_new_protected:Npn \regex_build_special_group:NN #1#2
1043 {
1044     \cs_if_exist_use:cF { regex_build_special_group_\token_to_str:N #2 : }
1045     {
1046         \msg_kernel_error:nnx { regex } { unsupported }
1047         { ( ? \token_to_str:N #2 ) % )
1048         \regex_build_special:N ( % )
1049         \regex_build_raw:N ?
1050         #1 #2
1051     }
1052 }
```

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

\regex_build_special_group:: Non-capturing groups are like capturing groups, except that we set the group id to -1, which will then inhibit submatching in \regex_build_group_submatches:NN. The group number is not increased.

```

1053 \cs_new_protected_nopar:cpn { regex_build_special_group:: } :
1054 {
1055     \regex_seq_push_int:NN \l_regex_capturing_group_seq \c_minus_one
1056     \regex_build_open_aux:
1057 }
```

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

\regex_build_special_group_|: The special group (?|...|...) is non-capturing (hence we set the `capturing_group` to -1), and resets the group number in each branch of the alternation. We use a variant of \regex_build_open_aux:, adding some code to be performed at every alternation, and at the end of the group. Namely, we keep track of the maximal value that \l_regex_capturing_group_int takes, and restore that value when the group end, and in every branch, we reset the capturing group number.

```

1058 \cs_new_protected_nopar:cpn { regex_build_special_group_|: } :
1059 {
1060     \regex_seq_push_int:NN \l_regex_capturing_group_seq \c_minus_one
1061     \regex_build_new_state:
1062     \regex_seq_push_int:NN \l_regex_left_state_seq \l_regex_left_state_int
1063     \regex_seq_push_int:NN \l_regex_right_state_seq \l_regex_right_state_int
1064     \seq_push:Nx \l_regex_end_alternation_seq
1065     {
1066         \exp_not:N \int_compare:nNnT
1067         \l_regex_capturing_group_int
1068         > \l_regex_capturing_group_max_int
1069         {
1070             \int_set_eq:NN
```

```

1071           \l_regex_capturing_group_max_int
1072           \l_regex_capturing_group_int
1073       }
1074   \int_set:Nn \l_regex_capturing_group_int
1075       { \int_use:N \l_regex_capturing_group_int }
1076   }
1077 \seq_push:Nx \l_regex_end_group_seq
1078   {
1079     \bool_if:NTF \l_regex_caseless_bool
1080         \regex_build_caseless:
1081         \regex_build_caseful:
1082     \int_set_eq:NN
1083         \l_regex_capturing_group_int
1084         \l_regex_capturing_group_max_int
1085     }
1086     \regex_build_begin_alternation:
1087   }
(End definition for \regex_build_special_group_1/. This function is documented on page ??.)
```

\regex_build_special_group_i: The match can be made case-insensitive by setting the option with (?i).

```

1088 \cs_new_protected_nopar:Npn \regex_build_special_group_i:
1089   {
1090     \regex_build_options:NNN +
1091     \regex_build_raw:N i
1092   }
1093 \cs_new_protected_nopar:cpn { regex_build_special_group_-: }
1094   {
1095     \regex_build_options:NNN -
1096   }
1097 \cs_new_protected:Npn \regex_build_options:NNN #1#2#3
1098   {
1099     \token_if_eq_meaning:NNTF \regex_build_raw:N #2
1100     {
1101       \cs_if_exist_use:cF { regex_build_option_#1#3: }
1102       { \msg_error:nnx { regex } { unknown-option } { #3 } }
1103       \regex_build_options:NNN #1
1104     }
1105     {
1106       \prg_case_str:nnn { #3 }
1107       { % (
1108         { ) } { }
1109         { - } { \regex_build_options:NNN - }
1110       }
1111       { \msg_error:nnx { regex } { invalid-in-option } { #3 } }
1112     }
1113   }
1114 \cs_new_protected_nopar:cpn { regex_build_option_+i: }
1115   {
1116     \regex_build_caseless:
1117     \cs_set_eq:NN \regex_match_loop_case_hook:
```

```

1118     \regex_match_loop_caseless_hook:
1119   }
1120 \cs_new_protected_nopar:cpn { regex_build_option_-i: }
1121 { \regex_build_caseful: }
(End definition for \regex_build_special_group_i:. This function is documented on page ??.)
```

2.3.11 Quantifiers for groups

\regex_build_group_quantifier: Used for one group. We need to keep track of submatches, threads can be created within one repetition, so things are hard. The code for the group that was just built starts at \l_regex_left_state_int and ends at \l_regex_right_state_int.

```

1122 \cs_new_protected_nopar:Npn \regex_build_group_quantifier:
1123 {
1124   \tl_set:Nn \l_regex_one_or_group_tl { group }
1125   \regex_build_quantifier:w
1126 }
(End definition for \regex_build_group_quantifier:. This function is documented on page ??.)
```

\regex_build_group_submatches:NN Once the quantifier is found by \regex_build_quantifier:w, we insert the code for tracking submatches.

```

1127 \cs_new_protected:Npn \regex_build_group_submatches:NN #1#2
1128 {
1129   \seq_pop:NN \l_regex_capturing_group_seq \l_regex_tmpa_tl
1130   \int_compare:nNnF { \l_regex_tmpa_tl } < \c_zero
1131   {
1132     \regex_toks_put_left:Nx #1
1133     { \regex_action_submatch:n { \l_regex_tmpa_tl < } }
1134     \regex_toks_put_left:Nx #2
1135     { \regex_action_submatch:n { \l_regex_tmpa_tl > } }
1136   }
1137 }
(End definition for \regex_build_group_submatches:NN. This function is documented on page ??.)
```

\regex_build_group_: When there is no quantifier, the group is simply inserted as is, and we only need to track submatches, and move to a new state.

```

1138 \cs_new_protected_nopar:Npn \regex_build_group_:
1139 {
1140   \regex_build_group_submatches:NN
1141   \l_regex_left_state_int \l_regex_right_state_int
1142   \regex_build_transition:NN
1143   \regex_action_free:n \l_regex_right_state_int
1144 }
```

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

\regex_build_group_shift:N Most quantifiers require to add an extra state before the group. This is done by shifting the current contents of the \tex_toks:D \l_regex_tmpa_int to a new state.

```

1145 \cs_new_protected:Npn \regex_build_group_shift:N #1
1146 {
```

```

1147   \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
1148   \regex_build_new_state:
1149   \tex_toks:D \l_regex_right_state_int = \tex_toks:D \l_regex_tmpa_int
1150   \regex_toks_put_left:Nx \l_regex_right_state_int
1151   {
1152     \int_set:Nn \l_regex_current_state_int
1153     { \int_use:N \l_regex_tmpa_int } % ^~A here we lie!
1154   }
1155   \cs_set:Npx \regex_tmp:w
1156   {
1157     \tex_toks:D \l_regex_tmpa_int
1158     {
1159       \s_stop
1160       #1 { \int_eval:n { \l_regex_right_state_int - \l_regex_tmpa_int } }
1161     }
1162   }
1163   \regex_tmp:w
1164   \regex_build_group_submatches>NN
1165   \l_regex_right_state_int \l_regex_left_state_int
1166 }
(End definition for \regex_build_group_shift:N. This function is documented on page ??.)
```

\regex_build_group_qs_aux:NN Shift the state at which the group begins using \regex_build_group_shift:N, then add two transitions. The first transition is taken once the group has been traversed: in the case of ? and ??, we should exit by going to \l_regex_right_state_int, while for * and *? we loop by going to \l_regex_tmpa_int. The second transition corresponds to skipping the group; it has lower priority (put_right) for greedy operators, and higher priority (put_left) for lazy operators.

```

1167 \cs_new_protected:Npn \regex_build_group_qs_aux>NN #1#2
1168 {
1169   \regex_build_group_shift:N \regex_action_free:n
1170   \int_set_eq:NN \l_regex_right_state_int \l_regex_left_state_int
1171   \regex_build_transition>NN \regex_action_free:n #1
1172   #2 \l_regex_tmpa_int
1173   {
1174     \regex_action_free:n
1175     { \int_eval:n { \l_regex_right_state_int - \l_regex_tmpa_int } }
1176   }
1177 }
1178 \cs_new_protected_nopar:cpn { regex_build_group_?: }
1179 {
1180   \regex_build_group_qs_aux>NN
1181   \l_regex_right_state_int \regex_toks_put_right:Nx
1182 }
1183 \cs_new_protected_nopar:cpn { regex_build_group_??: }
1184 {
1185   \regex_build_group_qs_aux>NN
1186   \l_regex_right_state_int \regex_toks_put_left:Nx
1187 }
```

```

1188 \cs_new_protected_nopar:cpn { regex_build_group_*: }
1189   {
1190     \regex_build_group_qs_aux:NN
1191     \l_regex_tmpa_int \regex_toks_put_right:Nx
1192   }
1193 \cs_new_protected_nopar:cpn { regex_build_group_*?: }
1194   {
1195     \regex_build_group_qs_aux:NN
1196     \l_regex_tmpa_int \regex_toks_put_left:Nx
1197   }
(End definition for \regex_build_group_qs_aux:NN. This function is documented on page ??.)
```

- \regex_build_group_+: Insert the submatch tracking code, then add two transitions from the current state to the left end of the group (repeating the group), and to a new state (to carry on with the rest of the regular expression).

```

1198 \cs_new_protected_nopar:cpn { regex_build_group_+: }
1199   {
1200     \regex_build_group_submatches:NN
1201     \l_regex_left_state_int \l_regex_right_state_int
1202     \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
1203     \regex_build_transitions:NNNN
1204     \regex_action_free:n \l_regex_tmpa_int
1205     \regex_action_free:n \l_regex_right_state_int
1206   }
1207 \cs_new_protected_nopar:cpn { regex_build_group_+?: }
1208   {
1209     \regex_build_group_submatches:NN
1210     \l_regex_left_state_int \l_regex_right_state_int
1211     \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
1212     \regex_build_transitions:NNNN
1213     \regex_action_free:n \l_regex_right_state_int
1214     \regex_action_free:n \l_regex_tmpa_int
1215   }
(End definition for \regex_build_group_+:. This function is documented on page ??.)
```

- \regex_build_group_n_aux:n The braced quantifiers rely on replicating the states corresponding to the group that has just been built, and joining the right state of each copy to the left state of the next copy. Once this function has been run, \l_regex_tmpa_int points to the last copy of the initial left-most state, \l_regex_left_state_int has its initial value, and \l_regex_right_state_int points to the last copy of the initial right-most state. Furthermore, \l_regex_max_state_int is set appropriately to the largest allocated \toks register.

```

1216 \cs_new_protected:Npn \regex_build_group_n_aux:n #1
1217   {
1218     \regex_toks_put_right:Nx \l_regex_right_state_int
1219     {
1220       \regex_action_free:n
1221       {
1222         \int_eval:n %^^A todo: document why that value.
1223         { \l_regex_max_state_int - \c_one - \l_regex_left_state_int }
```

```

1224     }
1225   }
1226   \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
1227   \int_set_eq:NN \l_regex_tmpb_int \l_regex_max_state_int
1228   \int_set:Nn \l_regex_max_state_int
1229   {
1230     \l_regex_left_state_int
1231     + #1 * ( \l_regex_max_state_int - \l_regex_left_state_int )
1232   }
1233   \int_while_do:nNnn \l_regex_tmpb_int < \l_regex_max_state_int
1234   {
1235     \tex_toks:D \l_regex_tmpb_int = \tex_toks:D \l_regex_tmpa_int
1236     \int_incr:N \l_regex_tmpa_int
1237     \int_incr:N \l_regex_tmpb_int
1238   }
1239 }

```

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

`\regex_build_group_n:` These functions are called in case the syntax is `{(int)}`. Greedy and lazy operators are identical, since the number of repetitions is fixed. We only record the submatch information at the last repetition.

```

1240 \cs_new_protected:Npn \regex_build_group_n: #1
1241   { % ^~A todo: catch case #1 <= 0.
1242     \regex_build_group_n_aux:n {#1}
1243     \regex_build_transition>NN
1244     \regex_action_free:n \l_regex_right_state_int
1245     \regex_build_group_submatches>NN
1246     \l_regex_tmpa_int \l_regex_left_state_int
1247   }
1248 \cs_new_eq:cN { regex_build_group_n?: } \regex_build_group_n:

```

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

`\regex_build_group_n*:` These functions are called in case the syntax is `{(int)}`. They are somewhat hybrid between the `{(int)}` and the `*` quantifiers. Contrarily to the `*` quantifier, for which we had to be careful not to overwrite the submatch information in case no iteration was made, here, we know that the submatch information is overwritten in any case.

```

1249 \cs_new_protected:cpn { regex_build_group_n*: } #1
1250   { % ^~A todo: catch case #1 <= 0.
1251     \regex_build_group_n_aux:n {#1}
1252     \regex_build_transitions>NNNN
1253     \regex_action_free:n \l_regex_tmpa_int
1254     \regex_action_free:n \l_regex_right_state_int
1255     \regex_build_group_submatches>NN
1256     \l_regex_tmpa_int \l_regex_left_state_int
1257   }
1258 \cs_new_protected:cpn { regex_build_group_n*?: } #1
1259   { % ^~A todo: catch case #1 <= 0.
1260     \regex_build_group_n_aux:n {#1}
1261     \regex_build_transitions>NNNN

```

```

1262     \regex_action_free:n \l_regex_right_state_int
1263     \regex_action_free:n \l_regex_tmpa_int
1264     \regex_build_group_submatches>NN
1265     \l_regex_tmpa_int \l_regex_left_state_int
1266 }

```

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

`\regex_build_group_nn`: These functions are called when the syntax is either `{<int>,}` or `{<int>,<int>}`.

```

\regex_build_group_nn?:
1267 \cs_new_protected:Npn \regex_build_group_nn: #1#2
1268   { % ^^A Not Implemented Yet!
1269     \msg_expandable_error:n { Quantifier~{m,n}~not~implemented~yet }
1270     \use:c { regex_build_group_n*: } {#1}
1271   }
1272 \cs_new_protected:cpn { regex_build_group_nn?: } #1#2
1273   { % ^^A Not Implemented Yet!
1274     \msg_expandable_error:n { Quantifier~{m,n}~not~implemented~yet }
1275     \use:c { regex_build_group_n*?: } {#1}
1276   }

```

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

2.4 Matching

2.4.1 Use of **T_EX** registers when matching

The first step in matching a regular expression is to build the corresponding NFA and store its states in the `\toks` registers. Then loop through the query token list one character (one “step”) at a time, exploring in parallel every possible path through the NFA. We keep track of an array of the states currently “active”. More precisely, `\skip` registers hold the state numbers to be considered when the next token is read.

At every step, we unpack that array of active states and empty it. Then loop over all active states, and perform the instruction at that state of the NFA. This can involve “free” transitions to other states, or transitions which “consume” the current character. For free transitions, the instruction at the new state of the NFA is performed. When a transition consumes a character, the new state is put in the array of `\skip` registers: it will be active again when the next character is read.

If two paths through the NFA “collide” in the sense that they reach the same state when reading a given character, then any future execution will be identical for both. Hence, it is indeed enough to keep track of which states are active. [In the presence of back-references, the future execution is affected by how the previous match took place; this is why we cannot support those non-regular features.]

Many of the functions require extracting the submatches for the “best” match. Execution paths through the NFA are ordered by precedence: for instance, the regular expression `a?` creates two paths, matching either an empty token list or a single `a`; the path matching an `a` has higher precedence. When two paths collide, the path with the highest precedence is kept, and the other one is discarded. The submatch information for a given path is stored at the start of the `\toks` register which holds the state at which that path currently is.

Deciding to store the submatch information in `\toks` registers alongside with states of the NFA unfortunately implies some shuffling around. The two other options are to store the submatch information in one control sequence per path, which wastes csnames, or to store all of the submatch information in one property list, which turns out to be too slow. A tricky aspect of submatch tracking is to know when to get rid of submatch information. This naturally happens when submatch information is stored in `\toks` registers: if the information is not moved, it will be overwritten later.

The presence of ϵ -transitions (transitions which consume no character) leads to potential infinite loops; for instance the regular expression $(a??)^*$ could lead to an infinite recursion, where $a??$ matches no character, $*$ loops back to the start of the group, and $a??$ matches no character again. Therefore, we need to keep track of the states of the NFA visited at the current step. More precisely, a state is marked as “visited” if the instructions for that state have been inserted in the input stream, by setting the corresponding `\dimen` register to a value which uniquely identifies at which step it was last inserted.

The current approach means that stretch and shrink components of `\skip` registers, as well as all `\muskip` registers are unused. It could seem that `\count` registers are also free for use, but we still want to be able to safely use integers, which are implemented as `\count` registers.

2.4.2 Helpers for running the NFA

`\regex_store_state:n` Put the given state in the array of `\skip` registers. This is done by increasing the pointer `\l_regex_max_index_int`, and converting the integer to a dimension (suitable for a `\skip` assignment) in scaled points.

```

1277 \cs_new_protected:Npn \regex_store_state:n #1
1278 {
1279   \int_incr:N \l_regex_max_index_int
1280   \tex_skip:D \l_regex_max_index_int #1 sp \scan_stop:
1281   \regex_store_submatches:n {#1}
1282 }
```

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

`\regex_state_use:` Use a given program instruction, unless it has already been executed at this step. The `\toks` registers begin with some submatch information, ignored by `\regex_state_use:`, but not by `\regex_state_use_with_submatches:`. A state is free if it is not marker as taken, namely if the corresponding `\dimen` register is not `\l_regex_unique_id_int` in `sp`. The primitive conditional is ended before unpacking the `\toks` register.

```

1283 \cs_new_protected_nopar:Npn \regex_state_use_with_submatches:
1284   { \regex_state_use_aux:n { } }
1285 \cs_new_protected_nopar:Npn \regex_state_use:
1286   { \regex_state_use_aux:n { \exp_after:wN \use_none_delimit_by_s_stop:w } }
1287 \cs_new_protected:Npn \regex_state_use_aux:n #1
1288 {
1289   \if_num:w \tex_dimen:D \l_regex_current_state_int
1290     < \l_regex_unique_id_int
1291   \tex_dimen:D \l_regex_current_state_int
```

```

1292     = \l_regex_unique_id_int sp \scan_stop:
1293     #1 \tex_the:D \tex_toks:D \exp_after:wN \l_regex_current_state_int
1294     \fi:
1295     \scan_stop:
1296   }

```

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

2.4.3 Submatch tracking when running the NFA

`\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.

```

1297 \cs_new_protected_nopar:Npn \regex_disable_submatches:
1298   {
1299     \cs_set_eq:NN \regex_state_use_with_submatches: \regex_state_use:
1300     \cs_set_eq:NN \regex_store_submatches:n
1301       \regex_protected_use_none:n
1302     \cs_set_eq:NN \regex_action_submatches:n
1303       \regex_protected_use_none:n
1304   }
1305 \cs_new_protected:Npn \regex_protected_use_none:n #1 { }

```

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

`\regex_store_submatches:n`
`\regex_store_submatches_aux:w`
`\regex_store_submatches_aux_ii:Nnnw` The submatch information pertaining to one given thread is moved from state to state as we execute the NFA. We make sure that most of the `\toks` register is not read before being assigned again to that same register.

```

1306 \cs_new_protected:Npn \regex_store_submatches:n #1
1307   {
1308     \tex_toks:D #1 \exp_after:wN
1309     {
1310       \tex_roman numeral:D
1311       \exp_after:wN \regex_store_submatches_aux:w
1312       \tex_the:D \tex_toks:D #1
1313     }
1314   }
1315 \cs_new_protected:Npn \regex_store_submatches_aux:w #1 \s_stop
1316   {
1317     \regex_store_submatches_aux_ii:Nnnw
1318     #1
1319     \regex_state_submatches:nn \c_minus_one \q_prop
1320     \s_stop
1321   }
1322 \cs_new_protected:Npn \regex_store_submatches_aux_ii:Nnnw
1323   \regex_state_submatches:nn #1 #2 #3 \s_stop
1324   {
1325     \exp_after:wN \c_zero
1326     \exp_after:wN \regex_state_submatches:nn \exp_after:wN
1327     {
1328       \int_value:w \int_eval:w

```

```

1329           \l_regex_unique_id_int + \c_one
1330           \exp_after:wN
1331       }
1332       \exp_after:wN { \l_regex_current_submatches_prop }
1333       \regex_state_submatches:nn {#1} {#2}
1334       \s_stop
1335   }

```

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

`\regex_state_submatches:nn` This function is inserted by `\regex_store_submatches:n` in the `\toks` register holding a given state, and it is performed when the state is used.

```

1336 \cs_new_protected:Npn \regex_state_submatches:nn #1#2
1337 {
1338     \if_num:w #1 = \l_regex_unique_id_int
1339         \tl_set:Nn \l_regex_current_submatches_prop { #2 }
1340     \fi:
1341 }

```

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

2.4.4 Matching: framework

`\regex_match:n` Then reset a few variables which should be set only once, before the first match, even in the case of multiple matches. Then run the NFA (`\regex_match_once:` matches multiple times when appropriate).

```

1342 \cs_new_protected:Npn \regex_match:n #1
1343 {
1344     \tl_set_analysis:Nn \l_regex_tma_tl {#1}
1345
1346     \int_zero:N \l_regex_nesting_int
1347     \int_set_eq:NN \l_regex_current_step_int \l_regex_max_state_int
1348     \regex_query_set:nnn { } { -1 } { -2 }
1349     \int_set_eq:NN \l_regex_min_step_int \l_regex_current_step_int
1350     \exp_after:wN \regex_query_set_loop:ww
1351         \l_regex_tma_tl
1352         \s_tl { ? \prg_map_break: } \s_tl
1353     \prg_break_point:n { }
1354     \int_set_eq:NN \l_regex_max_step_int \l_regex_current_step_int
1355     \regex_query_set:nnn { } { -1 } { -2 }
1356
1357     \regex_match_initial_setup:
1358     \regex_match_once:
1359 }
1360 \cs_new:Npn \regex_query_set_loop:ww #1 \s_tl #2#3 \s_tl
1361 {
1362     \use_none:n #2
1363     \regex_query_set:nnn {#1} {"#2} {#3}
1364     \if_case:w "#2 \exp_stop_f:
1365     \or: \int_incr:N \l_regex_nesting_int
1366     \or: \int_decr:N \l_regex_nesting_int

```

```

1367     \fi:
1368     \regex_query_set_loop:ww
1369   }
1370 \cs_new_protected:Npn \regex_query_set:nnn #1#2#3
1371   {
1372     \tex_muskip:D \l_regex_current_step_int
1373     = \etex_gluetomu:D
1374     #3 sp
1375     plus #2 sp
1376     minus \l_regex_nesting_int sp
1377     \scan_stop:
1378     \tex_toks:D \l_regex_current_step_int {#1}
1379     \int_incr:N \l_regex_current_step_int
1380   }
1381 \cs_new_protected_nopar:Npn \regex_query_get:
1382   {
1383     \tl_set:Nx \l_regex_current_token_tl
1384     { \tex_the:D \tex_toks:D \l_regex_current_step_int }
1385     \l_regex_current_char_int
1386     = \etex_mutoglu:D \tex_muskip:D \l_regex_current_step_int
1387     \l_regex_current_catcode_int = \etex_gluestretch:D
1388     \etex_mutoglu:D \tex_muskip:D \l_regex_current_step_int
1389   }
(End definition for \regex_match:n. This function is documented on page ??.)
```

\regex_match_once: Set up more variables in \regex_match_setup:. If there was a match, use the token list \l_regex_every_match_tl, which may call \regex_match_once: again to achieve multiple matches.

```

1390 \cs_new_protected_nopar:Npn \regex_match_once:
1391   {
1392     \regex_match_setup:
1393     \regex_query_get:
1394     \regex_match_loop:
1395     \prg_break_point:n { }
1396     \bool_if:NT \l_regex_success_match_bool
1397     {
1398       \bool_gset_true:N \g_regex_success_bool
1399       \l_regex_every_match_tl
1400     }
1401   }
(End definition for \regex_match_once:. This function is documented on page ??.)
```

\regex_match_initial_setup: This function holds the setup that should be done only once for one given pattern matching. It is called only once for the whole token list. On the other hand, \regex_match_setup: is called for every match in the token list in case of repeated matches.

```

1402 \cs_new_protected_nopar:Npn \regex_match_initial_setup:
1403   {
1404     \prg_stepwise_inline:nnnn {0} {1} { \l_regex_max_state_int - \c_one }
1405     { \tex_dimen:D ##1 \c_minus_one sp \scan_stop: }
```

```

1406   \int_set_eq:NN \l_regex_unique_id_int  \c_minus_one
1407   \int_set:Nn \l_regex_start_step_int
1408     { \l_regex_min_step_int - \c_one }
1409   \int_set_eq:NN \l_regex_current_step_int \l_regex_min_step_int
1410   \int_set_eq:NN \l_regex_success_step_int \l_regex_min_step_int
1411   \int_set_eq:NN \l_regex_submatch_int \l_regex_max_state_int
1412   \bool_set_false:N \l_regex_success_empty_bool
1413   \bool_gset_false:N \g_regex_success_bool
1414 }

```

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

`\regex_match_setup`: Every time a match starts, `\regex_match_setup`: resets a few variables.

```

1415 \cs_new_protected_nopar:Npn \regex_match_setup:
1416 {
1417   \prop_clear:N \l_regex_current_submatches_prop
1418   \bool_if:NTF \l_regex_success_empty_bool
1419     { \cs_set_eq:NN \regex_last_match_empty:F \regex_last_match_empty_yes:F }
1420     { \cs_set_eq:NN \regex_last_match_empty:F \regex_last_match_empty_no:F }
1421   \int_set_eq:NN \l_regex_start_step_int \l_regex_success_step_int
1422   \int_set:Nn \l_regex_current_step_int
1423     { \l_regex_start_step_int - \c_one }
1424   \bool_set_false:N \l_regex_success_match_bool
1425   \int_zero:N \l_regex_max_index_int
1426   \regex_store_state:n {0} %^A _state_int!
1427 }

```

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

`\regex_match_loop`: Setup what needs to be reset at every character, then set `\l_regex_current_char_int`

`\regex_match_one_index:n` to the character code of the token that is read (and -1 for the end of the token list), and loop over the elements of the `\skip` array. Then repeat. There are a couple of tests to stop reading the token list when no active state is left, or when the end is reached. At every step in reading the token list, we store the character code of the current character in `\l_regex_current_char_int`, unless the end was reached: then we store -1 .

```

1428 \cs_new_protected_nopar:Npn \regex_match_loop:
1429 {
1430   \int_incr:N \l_regex_current_step_int
1431   \int_incr:N \l_regex_unique_id_int
1432   \bool_set_false:N \l_regex_fresh_thread_bool
1433   \int_set_eq:NN \l_regex_last_char_int \l_regex_current_char_int
1434   \regex_query_get:
1435   \regex_match_loop_case_hook:
1436   \cs_set_nopar:Npx \regex_tmp:w
1437   {
1438     \int_zero:N \l_regex_max_index_int
1439     \regex_match_one_index:w 1 ; \prg_break_point:n { }
1440     \exp_not:N \prg_break_point:n { }
1441   }
1442   \regex_tmp:w
1443   \if_num:w \l_regex_current_char_int < \c_minus_one

```

```

1444     \exp_after:wN \prg_map_break:
1445     \fi:
1446     \if_num:w \l_regex_max_index_int = \c_zero
1447         \exp_after:wN \prg_map_break:
1448     \fi:
1449     \regex_match_loop:
1450 }
1451 \cs_new:Npn \regex_match_one_index:w #1;
1452 {
1453     \if_num:w #1 > \l_regex_max_index_int
1454         \exp_after:wN \prg_map_break:
1455     \fi:
1456     \regex_match_one_index_aux:n
1457     { \int_value:w \tex_skip:D #1 }
1458     \exp_after:wN \regex_match_one_index:w
1459     \int_use:N \int_eval:w #1 + \c_one ;
1460 }
1461 \cs_new_protected:Npn \regex_match_one_index_aux:n #1
1462 {
1463     \int_set:Nn \l_regex_current_state_int {#1}
1464     \prop_clear:N \l_regex_current_submatches_prop
1465     \regex_state_use_with_submatches:
1466 }
(End definition for \regex_match_loop:. This function is documented on page ??.)
```

\regex_match_loop_case_hook: In the case where the regular expression contains caseless matching, the \regex_match_loop_case_hook: (normally empty) is redefined to set \l_regex_case_changed_char_int properly.

```

1467 \cs_new_protected_nopar:Npn \regex_match_loop_case_hook: { }
1468 \cs_new_protected_nopar:Npn \regex_match_loop_caseless_hook:
1469 {
1470     \int_set_eq:NN \l_regex_case_changed_char_int \l_regex_current_char_int
1471     \if_num:w \l_regex_current_char_int < \c_ninety_one
1472         \if_num:w \l_regex_current_char_int < \c_sixty_five
1473         \else:
1474             \int_add:Nn \l_regex_case_changed_char_int { \c_thirty_two }
1475         \fi:
1476     \else:
1477         \if_num:w \l_regex_current_char_int < \c_one_hundred_twenty_three
1478             \if_num:w \l_regex_current_char_int < \c_ninety_seven
1479             \else:
1480                 \int_sub:Nn \l_regex_case_changed_char_int { \c_thirty_two }
1481             \fi:
1482         \fi:
1483     \fi:
1484 }
```

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

2.4.5 Actions when matching

\regex_action_startWildcard: : the first state has a free transition to the second state, where the regular expression really begins, and a costly transition to itself, to try again at the next character. The search is made unanchored at the start by putting a free transition to the real start of the NFA, and a costly transition to the same state, waiting for the next token in the query. This combination could be reused (with some changes). We sometimes need to know that the match for a given thread starts at this character. For that, we use the boolean \l_regex_fresh_thread_bool.

```

1485 \cs_new_protected_nopar:Npn \regex_action_startWildcard:
1486 {
1487     \bool_set_true:N \l_regex_fresh_thread_bool
1488     \regex_action_free:n {1}
1489     \bool_set_false:N \l_regex_fresh_thread_bool
1490     \regex_action_cost:n {0}
1491 }
```

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

\regex_action_cost:n A transition which consumes the current character and moves to state #1.

```

1492 \cs_new_protected:Npn \regex_action_cost:n #1
1493 {
1494     \exp_args:Nf \regex_store_state:n
1495     { \int_eval:n { \l_regex_current_state_int + #1 } }
1496 }
```

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

\regex_action_success: There is a successful match when an execution path reaches the end of the regular expression. Then store the current step and submatches. The current step is then interrupted with \prg_map_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.

```

1497 \cs_new_protected_nopar:Npn \regex_action_success:
1498 {
1499     \regex_last_match_empty:F
1500     {
1501         \bool_set_true:N \l_regex_success_match_bool
1502         \bool_set_eq:NN \l_regex_success_empty_bool
1503             \l_regex_fresh_thread_bool
1504         \int_set_eq:NN \l_regex_success_step_int \l_regex_current_step_int
1505         \prop_set_eq:NN \l_regex_success_submatches_prop
1506             \l_regex_current_submatches_prop
1507         \prg_map_break:
1508     }
1509 }
```

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

\regex_action_free:n To copy a thread, check whether the program state has already been used at this character. 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.

```

1510 \cs_new_protected:Npn \regex_action_free:n #1
1511   {
1512     \cs_set_nopar:Npx \regex_tmp:w
1513     {
1514       \int_add:Nn \l_regex_current_state_int {#1}
1515       \regex_state_use:
1516       \int_set:Nn \l_regex_current_state_int
1517         { \int_use:N \l_regex_current_state_int }
1518       \tl_set:Nn \exp_not:N \l_regex_current_submatches_prop
1519         { \exp_not:o \l_regex_current_submatches_prop }
1520     }
1521     \regex_tmp:w
1522   }

```

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

`\regex_action_submatch:n` Update the current submatches with the information from the current step.

```

1523 \cs_new_protected:Npn \regex_action_submatch:n #1
1524   {
1525     \prop_put:Nno \l_regex_current_submatches_prop {#1}
1526       { \int_use:N \l_regex_current_step_int }
1527   }

```

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

2.5 Replacement

`\regex_submatch_nesting_aux:n`

```

1528 \cs_new_protected:Npn \regex_submatch_nesting_aux:n #1
1529   {
1530     + \etex_glueshrink:D \etex_mutoglu:D \etex_muexpr:D
1531       \tex_muskip:D \etex_gluestretch:D \tex_skip:D #1
1532       - \tex_muskip:D \tex_skip:D #1
1533     \scan_stop:
1534   }

```

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

`\regex_replacement:n` Our goal here is to analyse the replacement text. First take care of detecting escaped and non-escaped characters using `\str_escape_use:NNNn` with three protected arguments. This inserts in the input stream a token list of the form $\langle fn\ 1 \rangle \langle char\ 1 \rangle \dots \langle fn\ N \rangle \langle char\ N \rangle$, where $\langle fn\ i \rangle$ is one of the three functions, and $\langle char\ i \rangle$ a character in the string #1.

```

1535 \cs_new:Npn \regex_nesting:n #1 { } %^A move. Rename?
1536 \tl_new:N \l_regex_nesting_tl
1537 \cs_new_protected:Npn \regex_replacement:n #1
1538   {
1539     \int_zero:N \l_regex_replacement_int
1540     \int_zero:N \l_regex_nesting_int
1541     \tl_clear:N \l_regex_nesting_tl

```

```

1542   \str_escape_use>NNNn
1543     \regex_replacement_unescaped:N
1544     \regex_replacement_escaped:N
1545     \regex_replacement_raw:N
1546     {#1}
1547   \prg_do_nothing: \prg_do_nothing:
1548   \cs_set:Npx \regex_nesting:n ##1
1549   {
1550     + \int_use:N \l_regex_nesting_int
1551     \l_regex_nesting_tl
1552     - \regex_submatch_nesting_aux:n {##1}
1553   }
1554   \use:x
1555   {
1556     \exp_not:n { \cs_set:Npn \regex_replacement_tl:n ##1 }
1557     { \str_aux_toks_range:mn \c_zero \l_regex_replacement_int }
1558   }
1559   % ^A rename! Careful with \cP\#.
1560 }

```

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

`\regex_replacement_raw:N`

```

1561 \cs_new_protected:Npn \regex_replacement_raw:N #1
1562   { \regex_replacement_put:n {#1} }

```

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

`\regex_replacement_put:n` Raw characters are stored in a toks register.

```

1563 \cs_new_protected:Npn \regex_replacement_put:n #1
1564   {
1565     \tex_toks:D \l_regex_replacement_int {#1}
1566     \int_incr:N \l_regex_replacement_int
1567   }

```

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

`\regex_replacement_escaped:N`

```

1568 \cs_new_protected:Npn \regex_replacement_unescaped:N #1
1569   {
1570     \if_charcode:w \c_rbrace_str #1
1571       \if_num:w \l_regex_replacement_csnames_int > \c_zero
1572         \regex_replacement_put:n \cs_end:
1573       \else:
1574         \regex_replacement_put:n #1
1575       \fi:
1576     \else:
1577       \regex_replacement_put:n #1
1578     \fi:
1579   }
1580 \cs_new_protected:Npn \regex_replacement_escaped:N #1
1581   {

```

```

1582 \if_charcode:w c #1
1583   \exp_after:wN \regex_replacement_c:w
1584 \else:
1585   \if_charcode:w g #1
1586     \exp_after:wN \exp_after:wN \exp_after:wN \regex_replacement_g:w
1587 \else:
1588   \if_num:w \c_one < 1#1 \exp_stop_f:
1589     \regex_replacement_put_submatch:n {#1}
1590   \else:
1591     \regex_replacement_put:n #1
1592   \fi:
1593 \fi:
1594 \fi:
1595 }
1596 \cs_new_protected:Npn \regex_replacement_put_submatch:n #1
1597 {
1598   \regex_replacement_put:n
1599   { \regex_query_submatch:nn {#1} {##1} }
1600   \if_num:w \l_regex_replacement_csnames_int = \c_zero
1601     \tl_put_right:Nn \l_regex_nesting_tl
1602     {
1603       \exp_not:N \if_num:w #1 < \l_regex_capturing_group_int
1604         \regex_submatch_nesting_aux:n { \int_eval:w #1+##1 \int_eval_end: }
1605       \exp_not:N \fi:
1606     }
1607   \fi:
1608 }
1609 \cs_new_protected:Npn \regex_replacement_error:NNN #1#2#3
1610 {
1611   \msg_kernel_error:nnnx { regex } { #1-command }
1612   { replacement-text } {#3}
1613   #2 #3
1614 }
1615 \cs_new_protected:Npn \regex_replacement_g:w #1#2
1616 {
1617   \str_if_eq:xxTF
1618   { \exp_not:n { #1#2 } }
1619   { \regex_replacement_unescaped:N \c_lbrace_str }
1620   {
1621     \int_zero:N \l_regex_tmpa_int
1622     \regex_replacement_g_digits:NN
1623   }
1624   { \regex_replacement_error:NNN g #1 #2 }
1625 }
1626 \cs_new_protected:Npn \regex_replacement_g_digits:NN #1#2
1627 {
1628   \token_if_eq_meaning:NNTF #1 \regex_replacement_unescaped:N
1629   {
1630     \if_num:w \c_one < 1#2 \exp_stop_f:
1631     \int_set:Nn \l_regex_tmpa_int { \c_ten * \l_regex_tmpa_int + #2 }

```

```

1632           \exp_after:wN \use_i:nnn
1633           \exp_after:wN \regex_replacement_g_digits:NN
1634     \else:
1635       \if_charcode:w \c_rbrace_str #2
1636         \exp_args:No \regex_replacement_put_submatch:n
1637           { \int_use:N \l_regex_tmpa_int }
1638         \exp_after:wN \exp_after:wN \exp_after:wN \use_none:nn
1639     \else:
1640       \exp_after:wN \exp_after:wN
1641       \exp_after:wN \regex_replacement_error:NNN
1642       \exp_after:wN \exp_after:wN \exp_after:wN g
1643     \fi:
1644   \fi:
1645 }
1646   { \regex_replacement_error:NNN g }
1647 #1 #2
1648 }
1649 \cs_new_protected:Npn \regex_replacement_c:w #1#2
1650 {
1651   \token_if_eq_meaning:NNTF #1 \regex_replacement_unescaped:N
1652   {
1653     \cs_if_exist_use:cF { \regex_replacement_c_#2:w }
1654     { \regex_replacement_error:NNN c #1#2 }
1655   }
1656   { \regex_replacement_error:NNN c #1#2 }
1657 }
1658 \cs_new_protected_nopar:cpx { \regex_replacement_c_ \c_lbrace_str :w }
1659 {
1660   \int_incr:N \l_regex_replacement_csnames_int
1661   \regex_replacement_put:n
1662   { \exp_not:n { \exp_after:wN \regex_replacement_exp_not:N \cs:w } }
1663 }
1664 \cs_new_eq:cc { \regex_replacement_c_C:w }
1665 { \regex_replacement_c_ \c_lbrace_str :w }
1666 \cs_new:Npn \regex_replacement_exp_not:N #1 { \exp_not:n {#1} }
1667 \group_begin:
1668
1669   \char_set_catcode_math_superscript:N \^{\circ}
1670   \cs_new_protected_nopar:Npn \regex_replacement_c_U:w
1671   { \regex_replacement_char:nNN { ^{\circ} } }
1672
1673   \char_set_catcode_alignment:N \^{\circ}
1674   \cs_new_protected_nopar:Npn \regex_replacement_c_T:w
1675   { \regex_replacement_char:nNN { ^{\circ} } }
1676
1677 \cs_new_protected:Npn \regex_replacement_c_S:w #1#2
1678 {
1679   \int_compare:nNnTF { '#2 } = \c_zero
1680   { \regex_replacement_error:NNN c #1#2 }
1681 }
```

```

1682           \char_set_lccode:nn {32} { '#2 }
1683           \tl_to_lowercase:n { \regex_replacement_put:n {~} }
1684       }
1685   }
1686
1687   \char_set_catcode_parameter:N \^^@
1688   \cs_new_protected_nopar:Npn \regex_replacement_c_P:w
1689     { \regex_replacement_char:nNN { ^^@^^@^^@^^@^^@^^@ } }
1690
1691   \char_set_catcode_other:N \^^@
1692   \cs_new_protected_nopar:Npn \regex_replacement_c_0:w
1693     { \regex_replacement_char:nNN { ^^@ } }
1694
1695   \char_set_catcode_math_toggle:N \^^@
1696   \cs_new_protected_nopar:Npn \regex_replacement_c_M:w
1697     { \regex_replacement_char:nNN { ^^@ } }
1698
1699   \char_set_catcode_letter:N \^^@
1700   \cs_new_protected_nopar:Npn \regex_replacement_c_L:w
1701     { \regex_replacement_char:nNN { ^^@ } }
1702
1703   \char_set_catcode_group_end:N \^^@
1704   \cs_new_protected_nopar:Npn \regex_replacement_c_E:w
1705     {
1706       \int_decr:N \l_regex_nesting_int
1707       \regex_replacement_char:nNN { \if_false: { \fi: ^^^@ } }
1708     }
1709
1710   \char_set_catcode_math_subscript:N \^^@
1711   \cs_new_protected_nopar:Npn \regex_replacement_c_D:w
1712     { \regex_replacement_char:nNN { ^^@ } }
1713
1714   \char_set_catcode_group_begin:N \^^@
1715   \cs_new_protected_nopar:Npn \regex_replacement_c_B:w
1716     {
1717       \int_incr:N \l_regex_nesting_int
1718       \regex_replacement_char:nNN { \exp_after:wN ^^^@ \if_false: } \fi: }
1719     }
1720
1721   \char_set_catcode_active:N \^^@
1722   \cs_new_protected_nopar:Npn \regex_replacement_c_A:w
1723     { \regex_replacement_char:nNN { \exp_not:N ^^^@ } }
1724
1725 \group_end:
1726
1727 \cs_new_protected:Npn \regex_replacement_char:nNN #1#2#3
1728   {
1729     \char_set_lccode:nn \c_zero { '#3 }
1730     \tl_to_lowercase:n
1731       { \regex_replacement_put:n { \exp_not:n {#1} } }

```

```

1732    }
(End definition for \regex_replacement_escaped:N. This function is documented on page ??.)
```

2.6 User commands

2.6.1 Precompiled pattern

A given pattern is often reused to match many different query token lists. We thus give a means of storing the NFA corresponding to a given pattern in a token list variable of the form

```

\regex_nfa:Nw <variable name>
⟨assignments⟩
\tex_toks:D 0 { ⟨instruction0⟩ }
...
\tex_toks:D n { ⟨instructionn⟩ }
\s_stop
```

where n is the number of states in the NFA, and the various $\langle instruction_i \rangle$ control how the NFA behaves in state i . The `\regex_nfa:Nw` function removes the whole NFA from the input stream and produces an error: the $\langle nfa var \rangle$ should only be accessed through dedicated functions. This rather drastic approach is taken because assignments triggered by the contents of $\langle nfa var \rangle$ may overwrite data which is used elsewhere, unless everything is done carefully in a group.

<code>\regex_gset:Nn</code> <code>\regex_const:Nn</code> <code>\regex_set:Nn</code> <code>\regex_set_aux>NNn</code>	<p>The three user functions only differ with which function is used to assign the pre-compiled regular expression to the user's variable. Internally, they all first build the NFA corresponding to the regex, then store the contents of all the necessary <code>\toks</code> registers in the user's variable.</p>
---	--

```

1733 \cs_new_protected_nopar:Npn \regex_set:Nn
1734   { \regex_set_aux>NNn \tl_set:Nn }
1735 \cs_new_protected_nopar:Npn \regex_gset:Nn
1736   { \regex_set_aux>NNn \tl_gset:Nn }
1737 \cs_new_protected_nopar:Npn \regex_const:Nn
1738   { \regex_set_aux>NNn \tl_const:Nn }
1739 \cs_new_protected:Npn \regex_set_aux>NNn #1#2#3
1740   {
1741     \group_begin:
1742       \regex_build:n {#3}
1743       \cs_set_nopar:Npx \regex_tmp:w
1744         { #1 \exp_not:N #2 { \regex_set_aux:N #2 } }
1745       \exp_after:wN
1746       \group_end:
1747       \regex_tmp:w
1748 }
```

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

\regex_set_aux:N Within a group, build the NFA corresponding to the given regular expression, with submatch tracking. Then save the contents of all relevant \toks registers into the variable outside the group. The auxiliary \regex_nfa:Nw is not protected: this ensures that the NFA will properly be replaced by an error message in expansion contexts.

```

1749 \cs_new:Npn \regex_set_aux:N #1
1750 {
1751     \exp_not:n { \regex_nfa:Nw #1 }
1752     \l_regex_max_state_int
1753         = \int_use:N \l_regex_max_state_int
1754     \l_regex_capturing_group_int
1755         = \int_use:N \l_regex_capturing_group_int
1756     \token_if_eq_meaning:NNT
1757         \regex_match_loop_case_hook:
1758         \regex_match_loop_caseless_hook:
1759     {
1760         \cs_set_eq:NN \regex_match_loop_case_hook:
1761             \regex_match_loop_caseless_hook:
1762     }
1763     \prg_stepwise_function:nnnN
1764         {0} {1} {\l_regex_max_state_int - \c_one }
1765     \regex_set_aux:n
1766     \s_stop
1767 }
1768 \cs_new:Npn \regex_set_aux:n #1
1769 { \tex_toks:D #1 { \tex_the:D \tex_toks:D #1 } }
1770 \cs_new:Npn \regex_nfa:Nw #1
1771 {
1772     \msg_expandable_kernel_error:nnn { regex } { nfa-misused } { #1 }
1773     \use_none_delimit_by_s_stop:w
1774 }
```

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

\regex_use:N No error-checking.

```

1775 \cs_new_protected_nopar:Npn \regex_use:N
1776 { \exp_last_unbraced:No \use_none:nn }
(End definition for \regex_use:N. This function is documented on page ??.)
```

\regex_to_str:N Of course, we could simply set \regex_to_str:N equal to \tl_to_str:N. After all, regex variables are in particular token list variables. We do some more processing to start each line with \tex_toks:D instead of a single very long line.

```

1777 \cs_new:Npn \regex_to_str:N #1
1778 {
1779     \exp_after:wN \regex_to_str_aux:Nw #1
1780     \tex_toks:D \q_stop \prg_map_break: \tex_toks:D
1781     \prg_break_point:n { }
1782 }
1783 \cs_new:Npn \regex_to_str_aux:Nw #1#2\tex_toks:D
1784 {
```

```

1785     \use_none_delimit_by_q_stop:w #2 \q_stop
1786     \tl_to_str:n {#1#2} \iow_newline:
1787     \regex_to_str_aux:Nw \tex_toks:D
1788 }
(End definition for \regex_to_str:N. This function is documented on page ??.)
```

2.6.2 Generic auxiliary functions

Most of l3regex's work is done within a group.

- \regex_aux_return: This function triggers either \prg_return_false: or \prg_return_true: as appropriate to whether a match was found or not.

```

1789 \cs_new_protected:Npn \regex_aux_return:
1790 {
1791     \if_meaning:w \c_true_bool \g_regex_success_bool
1792         \prg_return_true:
1793     \else:
1794         \prg_return_false:
1795     \fi:
1796 }
(End definition for \regex_aux_return:. This function is documented on page ??.)
```

- \regex_aux_build_match:nn This auxiliary is used by user functions whose *(regex)* argument is given as an explicit regular expression within braces. In that case, we need to build the automaton corresponding to that regular expression, then perform the matching on the given token list #2.

```

1797 \cs_new_protected:Npn \regex_aux_build_match:nn #1#2
1798 {
1799     \regex_build:n {#1}
1800     \regex_match:n {#2}
1801 }
(End definition for \regex_aux_build_match:nn. This function is documented on page ??.)
```

- \regex_aux_use_match:Nn This auxiliary is used by user functions whose *(regex)* argument is given as a pre-compiled regex variable. We make sure that the token list variable indeed is an automaton (by testing the first token). If not, the match is deemed unsuccessful, after raising an error. If we have an automaton, “use” it, then perform the matching on the given token list #2.

```

1802 \cs_new_protected:Npn \regex_aux_use_match:Nn #1#2
1803 {
1804     \exp_args:No \tl_if_head_eq_meaning:nNTF {#1} \regex_nfa:Nw
1805     {
1806         \regex_use:N #1
1807         \regex_match:n {#2}
1808     }
1809     {
1810         \msg_kernel_error:nnx { regex } { not-nfa } { \token_to_str:N #1 }
1811         \bool_gset_false:N \g_regex_success_bool
1812     }
1813 }
```

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

`\regex_extract_once:nnN` We define here 40 user functions, following a common pattern in terms of an auxiliary
`\regex_extract_once:NnN` such as `\regex_extract_once_aux:NnnN` (those auxiliaries are defined in the coming
`\regex_extract_all:nnN` sections). The arguments handed to the auxiliary are `\regex_aux_build_match:nn` or
`\regex_extract_all:NnN` `\regex_aux_use_match:Nn`, followed by the three arguments of the user function. The
`\regex_replace_once:nnN` conditionals call `\regex_aux_return:` to return either `true` or `false` once matching has
`\regex_replace_once:NnN` been performed.

```

1814 \cs_set_protected:Npn \regex_tmp:w #1#2#3
1815 {
  \cs_new_protected_nopar:Npn #1
    { #3 \regex_aux_build_match:nn }
  \cs_new_protected_nopar:Npn #2
    { #3 \regex_aux_use_match:Nn }
  \prg_new_protected_conditional:Npnn #1 ##1##2##3 { T , F , TF }
  {
    #3 \regex_aux_build_match:nn {##1} {##2} ##3
    \regex_aux_return:
  }
  \prg_new_protected_conditional:Npnn #2 ##1##2##3 { T , F , TF }
  {
    #3 \regex_aux_use_match:Nn {##1} {##2} ##3
    \regex_aux_return:
  }
1821 }
1822 \tl_map_inline:nn
1823 {
  { extract_once } { extract_all }
  { replace_once } { replace_all }
  { split }
1824 }
1825 \exp_args:Nccc \regex_tmp:w
1826   { regex_#1:nnN } { regex_#1:NnN } { regex_#1_aux:NnnN }
1827 }
```

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

2.6.3 Submatches, once the correct match is found

```

\regex_extract:
\regex_extract_aux_b:wn
\regex_extract_aux_e:wn
1841 \cs_new_protected_nopar:Npn \regex_extract:
1842 {
  \int_set_eq:NN \l_regex_submatch_start_int \l_regex_submatch_int
  \if_meaning:w \c_true_bool \g_regex_success_bool
  \prg_replicate:nn \l_regex_capturing_group_int
  {
    \tex_skip:D \l_regex_submatch_int \c_zero sp \scan_stop:
    \int_incr:N \l_regex_submatch_int
```

```

1849     }
1850   \prop_map_inline:Nn \l_regex_success_submatches_prop
1851   {
1852     \if_num:w ##1 \c_max_int
1853       \exp_after:wN \regex_extract_aux_b:wn \int_use:N
1854     \else:
1855       \exp_after:wN \regex_extract_aux_e:wn \int_use:N
1856     \fi:
1857     \int_eval:w \l_regex_submatch_start_int + ##1 {##2}
1858   }
1859   \fi:
1860 }
1861 \cs_new_protected:Npn \regex_extract_aux_b:wn #1 < #2
1862 {
1863   \tex_skip:D #1 = #2 sp
1864   plus \etex_gluestretch:D \tex_skip:D #1 \scan_stop:
1865 }
1866 \cs_new_protected:Npn \regex_extract_aux_e:wn #1 > #2
1867 {
1868   \tex_skip:D #1
1869   = 1 \tex_skip:D #1 plus #2 sp \scan_stop:
1870 }

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

```

\regex_group_end_extract_seq:N
\regex_extract_seq_aux:n
\regex_extract_seq_aux:ww
1871 \cs_new_protected:Npn \regex_group_end_extract_seq:N #1
1872 {
1873   \cs_set_eq:NN \seq_item:n \scan_stop:
1874   \flag_lower:N \l_regex_group_begin_flag
1875   \flag_lower:N \l_regex_group_end_flag
1876   \tl_set:Nx \l_regex_tmpa_tl
1877   {
1878     \prg_stepwise_function:nnnN
1879       \l_regex_max_state_int
1880       \c_one
1881       { \l_regex_submatch_int - \c_one }
1882       \regex_extract_seq_aux:n
1883   }
1884   \flag_test:NF \l_regex_group_begin_flag
1885   { \flag_test:NF \l_regex_group_end_flag { \use_none:nnn } }
1886   \msg_kernel_error:nn { regex } { sequence-unbalanced }
1887   \tl_set:Nx \l_regex_tmpa_tl { \l_regex_tmpa_tl }
1888   \exp_args:NNNo \group_end:
1889   \tl_set:Nn #1 \l_regex_tmpa_tl
1890 }
1891 \cs_new:Npn \regex_extract_seq_aux:n #1
1892 {
1893   \seq_item:n
1894   {
1895     \exp_after:wN \regex_extract_seq_aux:ww

```

```

1896           \int_value:w \regex_submatch_nesting_aux:n {#1} ; #1;
1897       }
1898   }
1899 \cs_new:Npn \regex_extract_seq_aux:ww #1; #2;
1900 { %^^A todo: use flags
1901   \if_num:w #1 < \c_zero
1902     \flag_raise:N \l_regex_group_end_flag
1903     \prg_replicate:nn {-#1} { \exp_not:n { { \if_false: } \fi: } }
1904   \fi:
1905   \regex_query_submatch:w #2;
1906   \if_num:w #1 > \c_zero
1907     \flag_raise:N \l_regex_group_begin_flag
1908     \prg_replicate:nn {#1} { \exp_not:n { \if_false: { \fi: } } }
1909   \fi:
1910 }
(End definition for \regex_group_end_extract_seq:N. This function is documented on page ??.)
```

2.6.4 Matching

\regex_match:nn We don't track submatches. Then either build the NFA corresponding to the regular expression, or use a precompiled pattern. Then match, using the internal \regex_match:n. Finally return the result after closing the group.

```

1911 \prg_new_protected_conditional:Npnn \regex_match:nn #1#2 { T , F , TF }
1912 {
1913   \regex_match_aux:n
1914   { \regex_aux_build_match:nn {#1} {#2} }
1915 }
1916 \prg_new_protected_conditional:Npnn \regex_match:Nn #1#2 { T , F , TF }
1917 {
1918   \regex_match_aux:n
1919   { \regex_aux_use_match:Nn #1 {#2} }
1920 }
1921 \cs_new_protected:Npn \regex_match_aux:n #1
1922 {
1923   \group_begin:
1924     \tl_clear:N \l_regex_every_match_tl
1925     \regex_disable_submatches:
1926     #1
1927   \group_end:
1928   \regex_aux_return:
1929 }
```

(End definition for \regex_match:nn. This function is documented on page ??.)

\regex_count:nnN Instead of aborting once the first “longest match” is found, we repeat the search. The code is such that the search will not start on the same character, hence avoiding infinite loops.

```

1930 \cs_new_protected_nopar:Npn \regex_count:nnN
1931   { \regex_count_aux:NnnN \regex_aux_build_match:nn }
1932 \cs_new_protected_nopar:Npn \regex_count:NnN
```

```

1933 { \regex_count_aux:NnnN \regex_aux_use_match:Nn }
1934 \cs_new_protected:Npn \regex_count_aux:NnnN #1#2#3#4
1935 {
1936   \group_begin:
1937     \regex_disable_submatches:
1938     \int_zero:N \l_regex_match_count_int
1939     \tl_set:Nn \l_regex_every_match_tl
1940     {
1941       \int_incr:N \l_regex_match_count_int
1942       \regex_match_once:
1943     }
1944   #1 {#2} {#3}
1945   \exp_args:NNNo
1946   \group_end:
1947   \int_set:Nn #4 { \int_use:N \l_regex_match_count_int }
1948 }
```

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

2.6.5 Submatch extraction

`\regex_extract_once_aux:NnnN` As announced, here comes the auxiliary for extracting one match. Since we only want one match, `\l_regex_every_match_tl` is empty, and does not trigger the matching code again. After matching, `\regex_extract:` extracts submatches into various `\skip` registers, and those are then concatenated into a sequence by `\regex_group_end_extract_seq:N`. That function is also responsible for closing the group.

```

1949 \cs_new_protected:Npn \regex_extract_once_aux:NnnN #1#2#3#4
1950 {
1951   \group_begin:
1952   \tl_set:Nn \l_regex_every_match_tl { \regex_extract: }
1953   #1 {#2} {#3}
1954   \regex_group_end_extract_seq:N #4
1955 }
```

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

`\regex_extract_all_aux:NnnN` The set of submatches will be built progressively in `\l_regex_result_seq`. For each match, extract the submatches, and concatenate that to the right of the result sequence, then start matching again. Finally, copy the result in the user's sequence variable.

```

1956 \cs_new_protected:Npn \regex_extract_all_aux:NnnN #1#2#3#4
1957 {
1958   \group_begin:
1959   \tl_set:Nn \l_regex_every_match_tl
1960   { \regex_extract: \regex_match_once: }
1961   #1 {#2} {#3}
1962   \regex_group_end_extract_seq:N #4
1963 }
```

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

2.6.6 Splitting a token list by matches of a regex

\regex_split_aux:NnnN Recurse through the matches, and for each, do the following. Extract the submatches into various \skip registers, then replace the match \0, which should not be kept in the final result, and replace it by the part of the token list before the match. This process must be inhibited to avoid creating empty items if the regex matched an empty token list at the place where the match attempt started. After the last successful match, we need to add to the result the part of the token list after the last match, unless the last match was empty and at the very end. Finally, \regex_group_end_extract_seq:N builds a sequence from all the \skip registers, and assigns it to #4 after closing the group.

```

1964 \cs_new_protected:Npn \regex_split_aux:NnnN #1#2#3#4
1965 {
1966   \group_begin:
1967     \tl_set:Nn \l_regex_every_match_tl
1968     {
1969       \if_num:w \l_regex_start_step_int < \l_regex_success_step_int
1970         \regex_extract:
1971         \tex_skip:D \l_regex_submatch_start_int
1972           = \l_regex_start_step_int sp
1973             plus \tex_skip:D \l_regex_submatch_start_int \scan_stop:
1974           \fi:
1975           \regex_match_once:
1976         }
1977       #1 {#2} {#3}
1978       \tex_skip:D \l_regex_submatch_int
1979         = \l_regex_start_step_int sp
1980           plus \l_regex_current_step_int sp \scan_stop:
1981       \int_incr:N \l_regex_submatch_int
1982       \if_num:w \l_regex_start_step_int = \l_regex_current_step_int
1983         \if_meaning:w \c_true_bool \l_regex_success_empty_bool
1984           \int_decr:N \l_regex_submatch_int
1985         \fi:
1986       \fi:
1987       \regex_group_end_extract_seq:N #4
1988     }

```

(End definition for \regex_split_aux:NnnN. This function is documented on page ??.)

2.6.7 Replacement

\regex_replace_once_aux:NnnN The replacement text is analysed by \regex_replacement:n, which defines \regex_replacement_tl:n to expand to the replaced token list, assuming that submatches are stored in various \skip registers, as done by \regex_extract:. If there is a match, we grab the parts before and after it, and get the result by x-expanding twice.

```

1989 \cs_new_protected:Npn \regex_replace_once_aux:NnnN #1#2#3#4
1990 {
1991   \group_begin:
1992     \tl_clear:N \l_regex_every_match_tl
1993     \regex_replacement:n {#3}

```

```

1994 \exp_args:Nno #1 {#2} #4
1995 \if_meaning:w \c_true_bool \g_regex_success_bool
1996   \regex_extract:
1997   \int_set:Nn \l_regex_tmpa_int
1998     { \regex_nesting:n { \l_regex_submatch_start_int } }
1999   \if_num:w \l_regex_tmpa_int = \c_zero
2000   \else:
2001     \msg_kernel_error:nnx { regex } { replace-unbalanced }
2002     { \l_regex_tmpa_int }
2003   \fi:
2004   \tl_set:Nx \l_regex_tmpa_tl
2005   {
2006     \if_num:w \l_regex_tmpa_int < \c_zero
2007       \prg_replicate:nn { - \l_regex_tmpa_int }
2008         { \exp_not:n { { \if_false: } \fi: } }
2009   \fi:
2010   \regex_query_substr:nn
2011     { \l_regex_min_step_int }
2012     { \tex_skip:D \l_regex_submatch_start_int }
2013   \regex_replacement_tl:n { \l_regex_submatch_start_int }
2014   \regex_query_substr:nn
2015     { \etex_gluestretch:D \tex_skip:D \l_regex_submatch_start_int }
2016     { \l_regex_max_step_int }
2017   \if_num:w \l_regex_tmpa_int > \c_zero
2018     \prg_replicate:nn { \l_regex_tmpa_int }
2019       { \exp_not:n { \if_false: { \fi: } } }
2020   \fi:
2021   }
2022   \tl_set:Nx \l_regex_tmpa_tl { \l_regex_tmpa_tl }
2023   \exp_args:NNNo \group_end:
2024   \tl_set:Nn #4 \l_regex_tmpa_tl
2025 \else:
2026   \group_end:
2027 \fi:
2028 }
2029 \cs_new_protected:Npn \regex_replace_once_aux:w #1;
2030 {
2031 }

```

(End definition for \regex_replace_once_aux:NnnN. This function is documented on page ??.)

\regex_replace_all_aux:NnnN For every match, extract submatches, and add the part before the beginning of the match, as well as the replacement, to the result. After the last match, extract the end of the token list, and add it to the replaced token list.

```

2032 \cs_new_protected:Npn \regex_replace_all_aux:NnnN #1#2#3#4
2033 {
2034   \group_begin:
2035   \tl_set:Nn \l_regex_every_match_tl
2036   {
2037     \regex_extract:

```

```

2038   \tex_skip:D \l_regex_submatch_start_int
2039   = \tex_the:D \tex_skip:D \l_regex_submatch_start_int
2040   minus \l_regex_start_step_int sp \scan_stop:
2041   \regex_match_once:
2042   }
2043   \regex_replacement:n {#3}
2044   \exp_args:Nno #1 {#2} #4
2045   \int_set:Nn \l_regex_tmpa_int
2046   {
2047   0
2048   \prg_stepwise_function:nnnN
2049   \l_regex_max_state_int
2050   \l_regex_capturing_group_int
2051   { \l_regex_submatch_int - \c_one }
2052   \regex_nesting:n
2053   }
2054   \if_num:w \l_regex_tmpa_int = \c_zero
2055   \else:
2056   \msg_kernel_error:nnx { regex } { replace-unbalanced }
2057   { \l_regex_tmpa_int }
2058   \fi:
2059   \tl_set:Nx \l_regex_tmpa_tl
2060   {
2061   \if_num:w \l_regex_tmpa_int < \c_zero
2062   \prg_replicate:nn { - \l_regex_tmpa_int }
2063   { \exp_not:n { { \if_false: } \fi: } }
2064   \fi:
2065   \prg_stepwise_function:nnnN
2066   \l_regex_max_state_int
2067   \l_regex_capturing_group_int
2068   { \l_regex_submatch_int - \c_one }
2069   \regex_replace_all_aux:n
2070   \regex_query_substr:nn
2071   \l_regex_start_step_int \l_regex_max_step_int
2072   \if_num:w \l_regex_tmpa_int > \c_zero
2073   \prg_replicate:nn { \l_regex_tmpa_int }
2074   { \exp_not:n { \if_false: { \fi: } } }
2075   \fi:
2076   }
2077   \tl_set:Nx \l_regex_tmpa_tl { \l_regex_tmpa_tl }
2078   \exp_args:NNNo \group_end:
2079   \tl_set:Nn #4 \l_regex_tmpa_tl
2080   }
2081 \cs_new:Npn \regex_replace_all_aux:n #1
2082   {
2083   \regex_query_substr:nn
2084   { \etex_glueshrink:D \tex_skip:D #1 } { \tex_skip:D #1 }
2085   \regex_replacement_tl:n {#1}
2086   }

```

(End definition for \regex_replace_all_aux:NnnN. This function is documented on page ??.)

2.7 Messages

```

2087 \msg_kernel_new:nnnn { regex } { sequence-unbalanced }
2088 {
2089   Missing-
2090   \flag_test:NTF \l_regex_group_end_flag
2091   {
2092     left~
2093     \flag_test:NTF \l_regex_group_begin_flag
2094     { and-right-braces } { brace }
2095   }
2096   { right-brace }
2097   \ inserted-in-extracted-match.
2098 }
2099 {
2100   LaTeX-was-asked-to-extract-submatches-or-split-a-token-list-
2101   according-to-a-given-regular-expression,-but-some-of-the-resulting-
2102   items-were-not-balanced.
2103 }
2104 \msg_kernel_new:nnnn { regex } { replace-unbalanced }
2105 {
2106   The-result-of-a-replacement-does-not-have-balanced-braces. }
2107 {
2108   LaTeX-was-asked-to-do-some-regular-expression-replacement,-
2109   and-the-resulting-token-list-would-not-have-the-same-number-
2110   of-begin-group-and-end-group-tokens. \\ \\
2111 \prg_case_int:nnn {#1}
2112 {
2113   { -1 } { A-left-brace-was }
2114   { 1 } { A-right-brace-was }
2115 }
2116 {
2117   \int_abs:n {#1} ~
2118   \int_compare:nNnTF {#1} < \c_zero { left } { right } ~
2119   braces ~ were
2120 }
2121 \ inserted.
2122 }

2123 \msg_kernel_new:nnnn { regex } { missing-rparen }
2124 {
2125   Missing-right-parenthesis-added-in-regular-expression. }
2126 {
2127   LaTeX-was-given-a-regular-expression-with-\int_eval:n{#1}-
2128   more-left-parenthes\int_compare:nTF{#1=1}{i}{e}s-than-right-
2129   parenthes\int_compare:nTF{#1=1}{i}{e}s.
2130 }
2131 \msg_kernel_new:nnnn { regex } { extra-rparen }
2132 {
2133   Extra-right-parenthesis-ignored-in-regular-expression. }
2134 {
2135   LaTeX-came-across-a-closing-parenthesis-when-no-submatch-group-
2136   was-open.-The-parenthesis-will-be-ignored.

```

```

2135    }
2136 \msg_kernel_new:nnnn { regex } { backwards-range }
2137   { Range-[#1-#2]-out-of-order-in-character-class. }
2138   {
2139     In-ranges-of-characters-[x-y]-appearing-in-character-classes,-
2140     the-first-character-code-must-not-be-larger-than-the-second.-.
2141     Here, ~#1-has-character-code-\int_eval:n {'#1},~while~#2-has-
2142     character-code-\int_eval:n {'#2}.
2143   }
2144 \msg_kernel_new:nnnn { regex } { unsupported }
2145   { Unsupported-construction-'#1'. }
2146   {
2147     The-construction-'#1'-is-not-supported-by-the-LaTeX-
2148     regular-expression-module.~Perhaps-some-character-should-
2149     have-been-escaped?
2150   }
2151 \msg_kernel_new:nnnn { regex } { not-nfa }
2152   { This-is-not-a-regular-expression-variable. }
2153   {
2154     LaTeX-was-expecting-'#1'-to-be-a-regular-expression-variable.\ \
2155     This-control-sequence-is-not-a-regex-variable.~It's-current-meaning-
2156     is-\ \\
2157     \ \ \ \ \ \token_to_str:N #1 = \token_to_meaning:N #1 .
2158   }
2159 % \msg_new:nnn { regex } { 1 } { \iow_char:N\at-end-of-pattern }
2160 % \msg_new:nnn { regex } { 2 } { \iow_char:N\c-at-end-of-pattern }
2161 % \msg_new:nnn { regex } { 4 }
2162 %   { Numbers-out-of-order-in-\iow_char:N\{\iow_char\}-quantifier. }
2163 % \msg_new:nnn { regex } { 6 }
2164 %   { Missing-terminating-\iow_char:N]-for-character-class }
2165 % \msg_new:nnn { regex } { 7 }
2166 %   { Invalid-escape-sequence-in-character-class }
2167 % \msg_new:nnn { regex } { 34 }
2168 %   { Character-value-in-\iow_char:N\x{...}-sequence-is-too-large }
2169 % \msg_new:nnn { regex } { 44 } { Invalid-UTF-8-string }
2170 % \msg_new:nnn { regex } { 46 }
2171 %   { Malformed-\iow_char:N\P-or\iow_char:N\p-sequence }
2172 % \msg_new:nnn { regex } { 47 }
2173 %   { Unknown-property-after-\iow_char:N\P-or\iow_char:N\p }
2174 % \msg_new:nnn { regex } { 68 }
2175 %   { \iow_char:N\c-must-be-followed-by-an-ASCII-character }

2176 \msg_kernel_new:nnnn { regex } { c-command }
2177   { Misused-\iow_char:N\c-or\iow_char:N\C-command-in-a-#1. }
2178   {
2179     In-a-#1,~the-\iow_char:N\C-escape-sequence-
2180     can-be-followed-by-one-of-the-letters-ABCDELMOPSTU-
2181     or-a-brace-group,~not-by-'#2'.
2182   }
2183 \msg_kernel_new:nnnn { regex } { unknown-option }

```

```

2184 { Unknown-option-'#1'~for~regular~expressions. }
2185 {
2186   LaTeX~came~across~something~like~'(?#1)'~in~a~regular~expression,~
2187   but~the~option~'#1'~is~not~known.~It~will~be~ignored.
2188 }
2189 \msg_kernel_new:nnnn { regex } { invalid-in-option }
2190   { Invalid-character~in~option~of~a~regular~expression. }
2191   {
2192     The~character~or~escape~sequence~'#1'~is~not~defined~
2193     as~an~option~within~regular~expressions.
2194   }
2195 \msg_kernel_new:nnnn { regex } { g-command }
2196   { Missing-brace~for~the~\iow_char:N\\g~construction~in~a~replacement~text. }
2197   {
2198     In~the~replacement~text~for~a~regular~expression~search,~
2199     submatches~are~represented~either~as~\iow_char:N \\g{dd..d},~
2200     or~\d,~where~'d'~are~single~digits.~Here,~a~brace~is~missing.
2201   }
2202 \msg_kernel_new:nnn { regex } { nfa-misused }
2203   { Automaton~#1~used~incorrectly. }
2204 
```