

# The **I3regex** package: regular expressions in **T<sub>E</sub>X**<sup>\*</sup>

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

Released 2012/04/23

## 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 **T<sub>E</sub>X** manipulates tokens rather than characters. For performance reasons, only a limited set of features are implemented. Notably, back-references are not supported.

Let us give a few examples. After

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

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

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

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

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

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

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

---

<sup>\*</sup>This file describes v3570, last revised 2012/04/23.

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

## 1.1 Syntax of regular expressions

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

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

Note that these rules play nicely with the fact that many non-alphanumeric characters are difficult to input into 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), after raising a warning.

Characters.

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

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

`\a` Alarm (hex 07).

`\e` Escape (hex 1B).

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

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

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

`\t` Horizontal tab (hex 09).

Character types.

. A single period matches any token.

`\d` Any decimal digit.

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

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

- \v Any vertical space character, equivalent to `[\^J\^K\^L\^M]`. Note that `\^K` is a vertical space, but not a space, for compatibility with Perl.
- \w Any word character, *i.e.*, alpha-numerics and underscore, equivalent to `[A-Za-z0-9\_]`.
- \D Any token not matched by \d.
- \H Any token not matched by \h.
- \N Any token other than the \n character (hex 0A).
- \S Any token not matched by \s.
- \V Any token not matched by \v.
- \W Any token not matched by \w.

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

Character classes match exactly one token in the subject.

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

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

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

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

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

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

Quantifiers (repetition).

`?` 0 or 1, greedy.

`??` 0 or 1, lazy.

`*` 0 or more, greedy.

`*?` 0 or more, lazy.

`+` 1 or more, greedy.

`+?` 1 or more, lazy.

`{n}` Exactly *n*.

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

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

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

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

Anchors and simple assertions.

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

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

$\sim$  or  $\backslash A$  Start of the subject token list.

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

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

Alternation and capturing groups.

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

$(\dots)$  Capturing group.

$(?:\dots)$  Non-capturing group.

$(?|\dots)$  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  $\backslash c$  escape sequence allows to test the category code of tokens, and match control sequences. Each character category is represented by a single uppercase letter:

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

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

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

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

`\cX` Applies to the next object, which can be a character, character property, class, or group, and forces this object to only match tokens with category X (any of CBEMTPUDSLOA. For instance, `\cL[A-Z\d]` matches uppercase letters and digits of category code letter, `\cC.` matches any control sequence, and `\c0(abc)` matches `abc` where each character has category other.

`\c[XYZ]` Applies to the next object, and forces it to only match tokens with category X, Y, or Z (each being any of CBEMTPUDSLOA). For instance, `\c[LS0](..)` matches two tokens of category letter, space, or other.

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

The category code tests can be used inside classes; for instance, `[\c0\d \c[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. Within a group affected by a category code test, the outer test can be overridden by a nested test: for instance, `\cL(ab\c0\*cd)` matches `ab*cd` where all characters are of category letter, except \* which has category other.

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

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

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

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

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

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

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

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

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

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

## 1.2 Syntax of the replacement text

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

For instance,

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

results in `\l_my_tl` holding `H(e\~ll--e\~l)(o,\~--o) w(or--o)(ld--l)!`

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

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

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

## 1.3 Pre-compiling regular expressions

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

---

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

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

---

`\regex_set:Nn`  
`\regex_gset:Nn`  
`\regex_const:Nn`

---

`\regex_set:Nn <regex var> {<regex>}`

Stores a compiled version of the *<regular expression>* in the *<regex var>*. For instance, this function can be used as

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

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

---

`\regex_show:n`  
`\regex_show:N`

---

`\regex_show:n {<regex>}`

Shows how `\regex` interprets the *<regex>*. For instance, `\regex_show:n {\A X|Y}` shows

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

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

## 1.4 Matching

All regular expression functions are available in both :`n` and :`N` variants. The former require a “standard” regular expression, while the latter require a compiled 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]* } { abecdex } { TRUE } { FALSE }
\regex_match:nntf { [b-dq-w] } { example } { TRUE } { FALSE }
```

leaves `TRUE` then `FALSE` in the input stream.

---

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

---

`\regex_count:nnN {⟨regular expression⟩} {⟨token list⟩} <int var>`

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

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

results in `\l_foo_int` taking the value 5.

## 1.5 Submatch extraction

---

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

---

`\regex_extract_once:nnN {⟨regular expression⟩} {⟨token list⟩} <seq var>`

`\regex_extract_once:nnNTF {⟨regular expression⟩} {⟨token list⟩} <seq var> {⟨true code⟩} {⟨false code⟩}`

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

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

Then the regular expression (anchored at the start with `\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:nnNF
\regex_replace_all:NnNF
```

---

```
\regex_replace_all:nn {⟨regular expression⟩} {⟨replacement⟩} ⟨tl var⟩
\regex_replace_all:nnNF {⟨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.

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

Additional error-checking to come.

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

Code improvements to come.

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

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

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

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

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

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

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

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

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

## 2 l3regex implementation

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

### 2.1 Plan of attack

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

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

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

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

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

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

- $\text{\toks(state)}$  holds the tests and actions to perform in the  $\langle state \rangle$  of the NFA.
- (Not implemented yet.)  $\text{\skip}_i$  has the form  $\langle group\ id \rangle$  plus  $\langle left\ state \rangle$  minus  $\langle right\ state \rangle$ .

When matching,

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

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

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

## 2.2 Helpers

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

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

### 2.2.1 Constants and variables

\regex_tmp:w	Temporary function used for various short-term purposes. 5 \cs_new:Npn \regex_tmp:w { } (End definition for \regex_tmp:w.)
\l_regex_internal_a_tl \l_regex_internal_b_tl \l_regex_internal_a_int \l_regex_internal_b_int \l_regex_internal_c_int \l_regex_internal_bool \l_regex_internal_seq \g_regex_internal_tl	Temporary variables used for various purposes. 6 \tl_new:N \l_regex_internal_a_tl 7 \tl_new:N \l_regex_internal_b_tl 8 \int_new:N \l_regex_internal_a_int 9 \int_new:N \l_regex_internal_b_int 10 \int_new:N \l_regex_internal_c_int 11 \bool_new:N \l_regex_internal_bool 12 \seq_new:N \l_regex_internal_seq 13 \tl_new:N \g_regex_internal_tl (End definition for \l_regex_internal_a_tl and others. These variables are documented on page ??.)
\c_regex_no_match_regex	This regular expression matches nothing, but is still a valid regular expression. We could use a failing assertion, but I went for an empty class. It is used as the initial value for regular expressions declared using \regex_new:N. 14 \tl_const:Nn \c_regex_no_match_regex 15 { 16     \regex_branch:n 17     { \regex_class:NnnnN \c_true_bool { } { 1 } { 0 } \c_true_bool } 18 } (End definition for \c_regex_no_match_regex. This variable is documented on page ??.)
\l_regex_balance_int	The first thing we do when matching is to go once through the query token list and store the information for each token as \muskip and \toks registers. During this phase, \l_regex_balance_int counts the balance of begin-group and end-group character tokens which appear before a given point in the token list, and we store it as the shrink component of each \muskip register. This variable is also used to keep track of the balance in the replacement text. 19 \int_new:N \l_regex_balance_int (End definition for \l_regex_balance_int. This variable is documented on page ??.)

### 2.2.2 Testing characters

\regex_break_point:TF \regex_break_true:w	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
	$\langle test_1 \rangle \dots \langle test_n \rangle$ \regex_break_point:TF {\langle true code \rangle} {\langle false code \rangle}

If any of the tests succeeds, it calls `\regex_break_true:w`, which cleans up and leaves `\langle true code \rangle` in the input stream. Otherwise, `\regex_break_point:TF` leaves the `\langle false code \rangle` in the input stream.

```

20 \cs_new_protected:Npn \regex_break_true:w
21     #1 \regex_break_point:TF #2 #3 {#2}
22 \cs_new_protected:Npn \regex_break_point:TF #1 #2 { #2 }
(End definition for \regex_break_point:TF. This function is documented on page ??.)
```

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

```

23 \cs_new_protected:Npn \regex_item_reverse:n #1
24 {
25     #1
26     \regex_break_point:TF { } \regex_break_true:w
27 }
```

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

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

```

28 \cs_new_protected:Npn \regex_item_caseful_equal:n #1
29 {
30     \if_num:w #1 = \l_regex_current_char_int
31         \exp_after:wN \regex_break_true:w
32     \fi:
33 }
34 \cs_new_protected:Npn \regex_item_caseful_range:nn #1 #2
35 {
36     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
37         \reverse_if:N \if_num:w #2 < \l_regex_current_char_int
38             \exp_after:wN \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
39         \fi:
40     \fi:
41 }
```

(End definition for \regex\_item\_caseful\_equal:n and \regex\_item\_caseful\_range:nn. These functions are documented on page ??.)

`\regex_item_caseless_equal:n` For caseless matching, we perform the test both on the `current_char` and on the `case_changed_char`. Before doing the second set of tests, we make sure that `case_changed_char` has been computed.

```

42 \cs_new_protected:Npn \regex_item_caseless_equal:n #1
43 {
44     \if_num:w #1 = \l_regex_current_char_int
45         \exp_after:wN \regex_break_true:w
46     \fi:
47     \if_num:w \l_regex_case_changed_char_int = \c_max_int
48         \regex_compute_case_changed_char:
49     \fi:
50     \if_num:w #1 = \l_regex_case_changed_char_int
```

```

51      \exp_after:wN \regex_break_true:w
52      \fi:
53  }
54 \cs_new_protected:Npn \regex_item_caseless_range:nn #1 #2
55 {
56     \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
57     \reverse_if:N \if_num:w #2 < \l_regex_current_char_int
58     \exp_after:wN \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
59     \fi:
60 \fi:
61 \if_num:w \l_regex_case_changed_char_int = \c_max_int
62     \regex_compute_case_changed_char:
63 \fi:
64 \reverse_if:N \if_num:w #1 > \l_regex_case_changed_char_int
65     \reverse_if:N \if_num:w #2 < \l_regex_case_changed_char_int
66     \exp_after:wN \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
67     \fi:
68 \fi:
69 }

```

(End definition for `\regex_item_caseless_equal:n` and `\regex_item_caseless_range:nn`. These functions are documented on page ??.)

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

```

70 \cs_new_protected_nopar:Npn \regex_compute_case_changed_char:
71 {
72     \int_set_eq:NN \l_regex_case_changed_char_int \l_regex_current_char_int
73     \if_num:w \l_regex_current_char_int < \c_ninety_one
74     \if_num:w \l_regex_current_char_int < \c_sixty_five
75     \else:
76         \int_add:Nn \l_regex_case_changed_char_int { \c_thirty_two }
77     \fi:
78 \else:
79     \if_num:w \l_regex_current_char_int < \c_one_hundred_twenty_three
80     \if_num:w \l_regex_current_char_int < \c_ninety_seven
81     \else:
82         \int_sub:Nn \l_regex_case_changed_char_int { \c_thirty_two }
83     \fi:
84     \fi:
85 \fi:
86 }

```

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

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

```
87 \cs_new_eq:NN \regex_item_equal:n ?
```

```
88 \cs_new_eq:NN \regex_item_range:nn ?
```

(End definition for `\regex_item_equal:n` and `\regex_item_range:nn`. These functions are documented on page ??.)

`\regex_item_catcode:nT`

`\regex_item_catcode_reverse:nT`

`\regex_item_catcode_aux:`

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.

```
89 \cs_new_protected:Npn \regex_item_catcode_aux:
90 {
91   "
92   \if_case:w \l_regex_current_catcode_int
93     1          \or: 4          \or: 10         \or: 40
94     \or: 100      \or:          \or: 1000       \or: 4000
95     \or: 10000     \or:          \or: 100000     \or: 400000
96     \or: 1000000   \or: 4000000 \else: 1*\c_zero
97   \fi:
98 }
99 \cs_new_protected:Npn \regex_item_catcode:nT #1
100 {
101   \if_int_odd:w \int_eval:w #1 / \regex_item_catcode_aux: \int_eval_end:
102   \exp_after:wN \use:n
103   \else:
104   \exp_after:wN \use_none:n
105   \fi:
106 }
107 \cs_new_protected:Npn \regex_item_catcode_reverse:nT #1#2
108 { \regex_item_catcode:nT {#1} { \regex_item_reverse:n {#2} } }
```

(End definition for `\regex_item_catcode:nT` and `\regex_item_catcode_reverse:nT`. These functions are documented on page ??.)

`\regex_item_exact:nn`

This matches an exact `<category>-<character code>` pair, or an exact control sequence.

`\regex_item_exact_cs:c`

```
109 \cs_new_protected:Npn \regex_item_exact:nn #1#2
110 {
111   \if_num:w #1 = \l_regex_current_catcode_int
112   \if_num:w #2 = \l_regex_current_char_int
113   \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
114   \fi:
115   \fi:
116 }
117 \cs_new_protected:Npn \regex_item_exact_cs:c #1
118 {
119   \int_compare:nNnTF \l_regex_current_catcode_int = \c_zero
120   {
121     \str_if_eq:xxTF
122     {
123       \exp_after:wN \exp_after:wN \exp_after:wN \cs_to_str:N
124       \tex_the:D \tex_toks:D \l_regex_current_pos_int
125     }
126 }
```

```

126      { #1 }
127      { \regex_break_true:w } { }
128    }
129  { }
130 }

```

(End definition for `\regex_item_exact:nn` and `\regex_item_exact_cs:c`. These functions are documented on page ??.)

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

```

131 \cs_new_protected:Npn \regex_item_cs:n #1
132   {
133     \int_compare:nNnT \l_regex_current_catcode_int = \c_zero
134     {
135       \group_begin:
136         \regex_single_match:
137         \regex_disable_submatches:
138         \regex_build_for_cs:n {#1}
139         \bool_set_eq:NN \l_regex_saved_success_bool \g_regex_success_bool
140         \exp_args:Nx \regex_match:n
141         {
142           \exp_after:wN \exp_after:wN
143           \exp_after:wN \cs_to_str:N
144           \tex_the:D \tex_toks:D \l_regex_current_pos_int
145         }
146         \if_meaning:w \c_true_bool \g_regex_success_bool
147           \group_insert_after:N \regex_break_true:w
148         \fi:
149         \bool_gset_eq:NN \g_regex_success_bool \l_regex_saved_success_bool
150         \group_end:
151       }
152     }

```

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

### 2.2.3 Character property tests

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

```

153 \cs_new_protected_nopar:Npn \regex_prop_d:
154   { \regex_item_caseful_range:nn \c_forty_eight { 57 } } % 0--9
155 \cs_new_protected_nopar:Npn \regex_prop_h:
156   {
157     \regex_item_caseful_equal:n \c_thirty_two % space
158     \regex_item_caseful_equal:n \c_nine % tab

```

```

159    }
160 \cs_new_protected_nopar:Npn \regex_prop_s:
161 {
162     \regex_item_caseful_equal:n \c_thirty_two % space
163     \regex_item_caseful_equal:n \c_nine      % tab
164     \regex_item_caseful_equal:n \c_ten       % lf
165     \regex_item_caseful_equal:n \c_twelve     % ff
166     \regex_item_caseful_equal:n \c_thirteen   % cr
167 }
168 \cs_new_protected_nopar:Npn \regex_prop_v:
169 {
170     \regex_item_caseful_range:nn \c_ten \c_thirteen } % lf, vtab, ff, cr
171 \cs_new_protected_nopar:Npn \regex_prop_w:
172 {
173     \regex_item_caseful_range:nn \c_ninety_seven { 122 } % a--z
174     \regex_item_caseful_range:nn \c_sixty_five   { 90 } % A--Z
175     \regex_item_caseful_range:nn \c_forty_eight { 57 } % 0--9
176     \regex_item_caseful_equal:n { 95 } % _
177 \cs_new_protected_nopar:Npn \regex_prop_N:
178 {
179     \regex_item_reverse:n { \regex_item_caseful_equal:n \c_ten } }
(End definition for \regex_prop_d: and others.)

```

\regex\_posix\_alnum: POSIX properties. No surprise.

```

\regex_posix_alpha:
\regex_posix_ascii:
\regex_posix_blank:
\regex_posix_cntrl:
\regex_posix_digit:
\regex_posix_graph:
\regex_posix_lower:
\regex_posix_print:
\regex_posix_punct:
\regex_posix_space:
\regex_posix_upper:
\regex_posix_word:
\regex_posix_xdigit:

```

179 \cs\_new\_protected\_nopar:Npn \regex\_posix\_alnum:
180 {
181 \regex\_posix\_alpha: \regex\_posix\_digit: }
182 \cs\_new\_protected\_nopar:Npn \regex\_posix\_alpha:
183 {
184 \regex\_posix\_lower: \regex\_posix\_upper: }
185 \cs\_new\_protected\_nopar:Npn \regex\_posix\_ascii:
186 {
187 \regex\_item\_caseful\_range:nn \c\_zero \c\_one\_hundred\_twenty\_seven }
188 \cs\_new\_eq:NN \regex\_posix\_blank: \regex\_prop\_h:
189 \cs\_new\_protected\_nopar:Npn \regex\_posix\_cntrl:
190 {
191 \regex\_item\_caseful\_range:nn \c\_zero { 31 }
192 \regex\_item\_caseful\_equal:n \c\_one\_hundred\_twenty\_seven
193 }
194 \cs\_new\_eq:NN \regex\_posix\_digit: \regex\_prop\_d:
195 \cs\_new\_protected\_nopar:Npn \regex\_posix\_graph:
196 {
197 \regex\_item\_caseful\_range:nn { 33 } { 126 } }
198 \cs\_new\_protected\_nopar:Npn \regex\_posix\_lower:
199 {
200 \regex\_item\_caseful\_range:nn \c\_ninety\_seven { 122 } }
201 \cs\_new\_protected\_nopar:Npn \regex\_posix\_print:
202 {
203 \regex\_item\_caseful\_range:nn \c\_thirty\_two { 126 } }
204 \cs\_new\_protected\_nopar:Npn \regex\_posix\_punct:
205 {
206 \regex\_item\_caseful\_range:nn { 33 } { 47 }
207 \regex\_item\_caseful\_range:nn { 58 } { 64 }
208 \regex\_item\_caseful\_range:nn { 91 } { 96 }
209 \regex\_item\_caseful\_range:nn { 123 } { 126 }
210 }
211 \cs\_new\_protected\_nopar:Npn \regex\_posix\_space:

```

206  {
207    \regex_item_caseful_equal:n \c_thirty_two
208    \regex_item_caseful_range:nn \c_nine \c_thirteen
209  }
210 \cs_new_protected_nopar:Npn \regex_posix_upper:
211   { \regex_item_caseful_range:nn \c_sixty_five { 90 } }
212 \cs_new_eq:NN \regex_posix_word: \regex_prop_w:
213 \cs_new_protected_nopar:Npn \regex_posix_xdigit:
214   {
215     \regex_posix_digit:
216     \regex_item_caseful_range:nn \c_sixty_five { 70 }
217     \regex_item_caseful_range:nn \c_ninety_seven { 102 }
218   }

```

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

#### 2.2.4 Simple character escape

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

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

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

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

```

219 \cs_new_protected:Npn \regex_escape_use:nnnn #1#2#3#4
220   {
221   (trace)    \trace_push:nnn { regex } { 1 } { regex_escape_use:nnnn }
222     \tl_set_build:Nw \l_regex_internal_a_tl
223     \cs_set_nopar:Npn \regex_escape_unescaped:N ##1 { #1 }
224     \cs_set_nopar:Npn \regex_escape_escaped:N ##1 { #2 }
225     \cs_set_nopar:Npn \regex_escape_raw:N ##1 { #3 }
226     \int_set:Nn \tex_escapechar:D { 92 }

```

```

227   \str_gset_other:Nn \g_regex_internal_tl { #4 }
228   \tl_set:Nx \l_regex_internal_b_tl
229   {
230     \exp_after:wN \regex_escape_loop:N \g_regex_internal_tl
231     { break } \prg_break_point:n { }
232   }
233   \tl_build_one:o \l_regex_internal_b_tl
234   \tl_build_end:
235   <trace>   \trace_pop:nnn { regex } { 1 } { regex_escape_use:nnnn }
236   \l_regex_internal_a_tl
237 }

```

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

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

```

238 \cs_new:Npn \regex_escape_loop:N #1
239   {
240     \cs_if_exist_use:cF { regex_escape_`token_to_str:N #1:w }
241     { \regex_escape_unescaped:N #1 }
242     \regex_escape_loop:N
243   }
244 \cs_new_nopar:cpn { regex_escape_`c_backslash_str :w }
245   \regex_escape_loop:N #1
246   {
247     \cs_if_exist_use:cF { regex_escape_`token_to_str:N #1:w }
248     { \regex_escape_escaped:N #1 }
249     \regex_escape_loop:N
250   }

```

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

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

```

251 \cs_new_eq:NN \regex_escape_unescaped:N ?
252 \cs_new_eq:NN \regex_escape_escaped:N    ?
253 \cs_new_eq:NN \regex_escape_raw:N      ?

```

(End definition for `\regex_escape_unescaped:N`, `\regex_escape_escaped:N`, and `\regex_escape_raw:N`.)

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

```

254 \cs_new_eq:NN \regex_escape_break:w \prg_map_break:
255 \cs_new_nopar:cpn { regex_escape_`break:w }
256   {
257     \if_false: { \fi: }
258     \msg_kernel_error:nn { regex } { trailing-backslash }
259     \exp_after:wN \use_none:n \exp_after:wN { \if_false: } \fi:
260   }
261 \cs_new_nopar:cpn { regex_escape_`~:w } { }
262 \cs_new_nopar:cpx { regex_escape_`a:w }

```

```

263   { \exp_not:N \regex_escape_raw:N \iow_char:N \^G }
264 \cs_new_nopar:cpx { regex_escape/_t:w }
265   { \exp_not:N \regex_escape_raw:N \iow_char:N \^I }
266 \cs_new_nopar:cpx { regex_escape/_n:w }
267   { \exp_not:N \regex_escape_raw:N \iow_char:N \^J }
268 \cs_new_nopar:cpx { regex_escape/_f:w }
269   { \exp_not:N \regex_escape_raw:N \iow_char:N \^L }
270 \cs_new_nopar:cpx { regex_escape/_r:w }
271   { \exp_not:N \regex_escape_raw:N \iow_char:N \^M }
272 \cs_new_nopar:cpx { regex_escape/_e:w }
273   { \exp_not:N \regex_escape_raw:N \iow_char:N \^_[ }
(End definition for \regex_escape_break:w and others. These functions are documented on page ??.)
```

```
\regex_escape/_x:w
\regex_escape_x_end:w
\regex_escape_x_large:n
```

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

```

274 \cs_new:cpn { regex_escape/_x:w } \regex_escape_loop:N
275   {
276     \exp_after:wN \regex_escape_x_end:w
277     \int_value:w "0 \regex_escape_x_test:N
278   }
279 \cs_new:Npn \regex_escape_x_end:w #1 ;
280   {
281     \int_compare:nNnTF {#1} < \c_two_hundred_fifty_six
282     {
283       \exp_last_unbraced:Nf \regex_escape_raw:N
284         { \str_output_byte:n {#1} }
285     }
286     { \regex_escape_x_large:n {#1} }
287   }
288 \group_begin:
289   \char_set_catcode_other:n { 0 }
290 \cs_new:Npn \regex_escape_x_large:n #1
291   {
292     \if_false: { \fi: }
293     \tl_build_one:o \l_regex_internal_b_tl
294     \int_compare:nNnTF {#1} > \c_max_char_int
295     {
296       \msg_kernel_error:nnx { regex } { x-overflow } {#1}
297       \tl_set:Nx \l_regex_internal_b_tl
298         { \if_false: } \fi: \regex_escape_loop:N
299     }
300   {
301     \char_set_lccode:nn { \c_zero } {#1}
302     \tl_to_lowercase:n
303       {
304         \tl_set:Nx \l_regex_internal_b_tl
```

```

305           { \if_false: } \fi:
306           \regex_escape_raw:N ^@^
307           \regex_escape_loop:N
308       }
309   }
310 }
311 \group_end:
(End definition for \regex_escape_x:w. This function is documented on page ??.)
```

\regex\_escape\_x\_test:N Find out whether the first character is a left brace (allowing any number of hexadecimal digits), or not (allowing up to two hexadecimal digits). We need to check for the end-of-string marker. Eventually, call either \regex\_escape\_x\_loop:N or \regex\_escape\_x\_i:N.

```

312 \cs_new:Npn \regex_escape_x_test:N #1
313 {
314     \str_if_eq:xxTF {#1} { break } { ; }
315     {
316         \if_charcode:w \c_space_token #1
317         \exp_after:wN \regex_escape_x_test:N
318     \else:
319         \exp_after:wN \regex_escape_x_test_ii:N
320         \exp_after:wN #1
321     \fi:
322 }
323 }
324 \cs_new:Npn \regex_escape_x_test_ii:N #1
325 {
326     \if_charcode:w \c_lbrace_str #1
327     \exp_after:wN \regex_escape_x_loop:N
328 \else:
329     \str_aux_hexadecimal_use:NTF #1
330     { \exp_after:wN \regex_escape_x_ii:N }
331     { ; \exp_after:wN \regex_escape_loop:N \exp_after:wN #1 }
332 \fi:
333 }
```

(End definition for \regex\_escape\_x\_test:N.)

\regex\_escape\_x\_ii:N This looks for the second digit in the unbraced case.

```

334 \cs_new:Npn \regex_escape_x_ii:N #1
335 {
336     \str_if_eq:xxTF {#1} { break } { ; }
337     {
338         \str_aux_hexadecimal_use:NTF #1
339         { ; \regex_escape_loop:N }
340         { ; \regex_escape_loop:N #1 }
341     }
342 }
```

(End definition for \regex\_escape\_x\_ii:N.)

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

```

343 \cs_new:Npn \regex_escape_x_loop:N #1
344   {
345     \str_aux_hexadecimal_use:NTF #1
346     { \regex_escape_x_loop:N }
347     {
348       \token_if_eq_charcode:NNTF \c_space_token #1
349       { \regex_escape_x_loop:N }
350       {
351         ;
352         \exp_after:wN \token_if_eq_charcode:NNTF \c_rbrace_str #1
353         { \regex_escape_loop:N }
354         {
355           \if_false: { \fi: }
356           \tl_build_one:o \l_regex_internal_b_tl
357           \msg_kernel_error:nn { regex } { x-missing-rbrace } {#1}
358           \tl_set:Nx \l_regex_internal_b_tl
359           { \if_false: } \fi: \regex_escape_loop:N #1
360         }
361       }
362     }
363   }
(End definition for \regex_escape_x_loop:N.)
```

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

- 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).

```

364 \prg_new_conditional:Npnn \regex_char_if_special:N #1 { TF }
365   {
366     \if_num:w '#1 < \c_ninety_one
367       \if_num:w '#1 < \c_fifty_eight
368         \if_num:w '#1 < \c_forty_eight
369           \if_num:w '#1 < \c_thirty_two
370             \prg_return_false: \else: \prg_return_true: \fi:
```

```

371           \else: \prg_return_false: \fi:
372
373           \else:
374             \if_num:w '#1 < \c_sixty_five
375               \prg_return_true: \else: \prg_return_false: \fi:
376             \fi:
377           \else:
378             \if_num:w '#1 < \c_one_hundred_twenty_three
379               \if_num:w '#1 < \c_ninety_seven
380                 \prg_return_true: \else: \prg_return_false: \fi:
381               \else:
382                 \if_num:w '#1 < \c_one_hundred_twenty_seven
383                   \prg_return_true: \else: \prg_return_false: \fi:
384                   \fi:
385               \fi:
386           }
387 \prg_new_conditional:Npn \regex_char_if_alphanumeric:N #1 { TF }
388   {
389     \if_num:w '#1 < \c_ninety_one
390       \if_num:w '#1 < \c_fifty_eight
391         \if_num:w '#1 < \c_forty_eight
392           \prg_return_false: \else: \prg_return_true: \fi:
393         \else:
394           \if_num:w '#1 < \c_sixty_five
395             \prg_return_false: \else: \prg_return_true: \fi:
396             \fi:
397           \else:
398             \if_num:w '#1 < \c_one_hundred_twenty_three
399               \if_num:w '#1 < \c_ninety_seven
400                 \prg_return_false: \else: \prg_return_true: \fi:
401                 \else:
402                   \prg_return_false:
403                   \fi:
404               \fi:
405   }

```

(End definition for `\regex_char_if_alphanumeric:NTF` and `\regex_char_if_special:NTF`.)

## 2.3 Compiling

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

- `\regex_class:NnnnN`  $\langle \text{boolean} \rangle \{ \langle \text{tests} \rangle \} \{ \langle \text{min} \rangle \} \{ \langle \text{more} \rangle \} \langle \text{lazyness} \rangle$
- `\regex_group:nnnN`  $\{ \langle \text{branches} \rangle \} \{ \langle \text{min} \rangle \} \{ \langle \text{more} \rangle \} \langle \text{lazyness} \rangle$ , also `\regex_group_no_capture:nnnN` and `\regex_group_resetting:nnnN` with the same syntax.
- `\regex_branch:n`  $\{ \langle \text{contents} \rangle \}$

- `\regex_command_K`:
- `\regex_assertion:Nn <boolean> {<assertion test>}`, where the `<assertion test>` is `\regex_b_test:` or `\{ \regex_anchor:N <integer> }`

Tests can be the following:

- `\regex_item_caseful_equal:n {<char code>}`
- `\regex_item_caseless_equal:n {<char code>}`
- `\regex_item_caseful_range:nn {<min>} {<max>}`
- `\regex_item_caseless_range:nn {<min>} {<max>}`
- `\regex_item_catcode:nT {<catcode bitmap>} {<tests>}`
- `\regex_item_catcode_reverse:nT {<catcode bitmap>} {<tests>}`
- `\regex_item_reverse:n {<tests>}`
- `\regex_item_exact:nn {<catcode>} {<char code>}`
- `\regex_item_exact_cs:c {<csname>}`
- `\regex_item_cs:n {<compiled regex>}`

### 2.3.1 Variables used when compiling

`\l_regex_group_level_int` We make sure to open the same number of groups as we close.

405 `\int_new:N \l_regex_group_level_int`

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

`\l_regex_mode_int` While compiling, ten modes are recognized, labelled  $-63, -23, -6, -2, 0, 2, 3, 6, 23, 63$ . See section 2.3.3.

406 `\int_new:N \l_regex_mode_int`

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

`\l_regex_catcodes_int` We wish to allow constructions such as `\c[^BE] (.. \cL[a-z] ..)`, where the outer catcode test applies to the whole group, but is superseeded by the inner catcode test. For this to work, we need to keep track of lists of allowed category codes: `\l_regex_catcodes_int` and `\l_regex_default_catcodes_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. The boolean records whether the list of categories of a catcode test has to be inverted: compare `\c[^BE]` and `\c[BE]`.

407 `\int_new:N \l_regex_catcodes_int`

408 `\int_new:N \l_regex_default_catcodes_int`

409 `\bool_new:N \l_regex_catcodes_bool`

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

```

\c_regex_catcode_C_int      Constants: 4c for each category, and the sum of all powers of 4.
\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_all_catcodes_int

410 \int_const:Nn \c_regex_catcode_C_int { "1 }
411 \int_const:Nn \c_regex_catcode_B_int { "4 }
412 \int_const:Nn \c_regex_catcode_E_int { "10 }
413 \int_const:Nn \c_regex_catcode_M_int { "40 }
414 \int_const:Nn \c_regex_catcode_T_int { "100 }
415 \int_const:Nn \c_regex_catcode_P_int { "1000 }
416 \int_const:Nn \c_regex_catcode_U_int { "4000 }
417 \int_const:Nn \c_regex_catcode_D_int { "10000 }
418 \int_const:Nn \c_regex_catcode_S_int { "100000 }
419 \int_const:Nn \c_regex_catcode_L_int { "400000 }
420 \int_const:Nn \c_regex_catcode_O_int { "1000000 }
421 \int_const:Nn \c_regex_catcode_A_int { "4000000 }
422 \int_const:Nn \c_regex_all_catcodes_int { "5515155 }

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

\l\_regex\_internal\_regex The compilation step stores its result in this variable.

```
423 \cs_new_eq:NN \l_regex_internal_regex \c_regex_no_match_regex
```

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

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

```
424 \seq_new:N \l_regex_show_prefix_seq
```

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

\l\_regex\_show\_lines\_int A hack. To know whether a given class has a single item in it or not, we count the number of lines when showing the class.

```
425 \int_new:N \l_regex_show_lines_int
```

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

### 2.3.2 Generic helpers used when compiling

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

```

426 \cs_new_protected:Npn \regex_get_digits:NTFw #1#2#3#4#5
427   {
428     \regex_if_raw_digit:NNTF #4 #5
429     { #1 = #5 \regex_get_digits_loop:nw {#2} }
430     { #3 #4 #5 }
431   }
432 \cs_new:Npn \regex_get_digits_loop:nw #1#2#3
433   {
434     \regex_if_raw_digit:NNTF #2 #3
435     { #3 \regex_get_digits_loop:nw {#1} }
436     { \scan_stop: #1 #2 #3 }
437   }
```

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

\regex\_if\_raw\_digit:NNTF Test used when grabbing digits for the  $\{m,n\}$  quantifier. It only accepts non-escaped digits.

```

438 \prg_new_conditional:Npnn \regex_if_raw_digit:NN #1#2 { TF }
439   {
440     \if_meaning:w \regex_compile_raw:N #1
441       \if_int_compare:w \c_one < 1 #2 \exp_stop_f:
442         \prg_return_true:
443       \else:
444         \prg_return_false:
445       \fi:
446     \else:
447       \prg_return_false:
448     \fi:
449   }

```

(End definition for \regex\_if\_raw\_digit:NNTF.)

### 2.3.3 Mode

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

- 6  $[\backslash c\{\dots\}]$  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  $\backslash c\{\dots\}$  class inside mode -2,
- 3  $[\dots]$  class inside mode -3,
- 23  $\backslash c[\dots]$  class inside mode 2,
- 63  $[\backslash c[\dots]]$  class inside mode 6.

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

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

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

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

```

450 \cs_new_nopar:Npn \regex_if_in_class:TF
451 {
452   \if_int_odd:w \l_regex_mode_int
453     \exp_after:wN \use_i:nn
454   \else:
455     \exp_after:wN \use_ii:nn
456   \fi:
457 }
```

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

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

```

458 \cs_new_nopar:Npn \regex_if_in_cs:TF
459 {
460   \if_int_odd:w \l_regex_mode_int
461     \exp_after:wN \use_ii:nn
462   \else:
463     \if_int_compare:w \l_regex_mode_int < \c_zero
464       \exp_after:wN \exp_after:wN \exp_after:wN \use_i:nn
465     \else:
466       \exp_after:wN \exp_after:wN \exp_after:wN \use_ii:nn
467     \fi:
468   \fi:
469 }
```

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

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

```

470 \cs_new_nopar:Npn \regex_if_in_class_or_catcode:TF
471 {
472   \if_int_odd:w \l_regex_mode_int
473     \exp_after:wN \use_i:nn
474   \else:
475     \if_int_compare:w \l_regex_mode_int > \c_zero
476       \exp_after:wN \exp_after:wN \exp_after:wN \use_i:nn
477 
```

```

477     \else:
478         \exp_after:wN \exp_after:wN \exp_after:wN \use_i:nn
479     \fi:
480 \fi:
481 }

```

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

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

```

482 \cs_new_nopar:Npn \regex_if_within_catcode:TF
483 {
484     \if_int_compare:w \l_regex_mode_int > \c_zero
485         \exp_after:wN \use_i:nn
486     \else:
487         \exp_after:wN \use_i:nn
488     \fi:
489 }

```

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

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

```

490 \cs_new_protected:Npn \regex_chk_c_allowed:T
491 {
492     \if_num:w \l_regex_mode_int = \c_zero
493         \exp_after:wN \use:n
494     \else:
495         \if_num:w \l_regex_mode_int = \c_three
496             \exp_after:wN \exp_after:wN \exp_after:wN \use:n
497         \else:
498             \msg_kernel_error:nn { regex } { c-bad-mode }
499             \exp_after:wN \exp_after:wN \exp_after:wN \use:none:n
500         \fi:
501     \fi:
502 }

```

(End definition for `\regex_chk_c_allowed:T`.)

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

```

503 \cs_new_protected:Npn \regex_mode_quit_c:
504 {
505     \if_num:w \l_regex_mode_int = \c_two
506         \l_regex_mode_int = \c_zero
507     \else:
508         \if_num:w \l_regex_mode_int = \c_six
509             \l_regex_mode_int = \c_three
510         \fi:
511     \fi:
512 }

```

(End definition for `\regex_mode_quit_c`.)

### 2.3.4 Framework

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

```

513 \cs_new_protected_nopar:Npn \regex_compile:w
514 {
515     \tl_set_build_x:Nw \l_regex_internal_regex
516     \int_zero:N \l_regex_group_level_int
517     \int_set_eq:NN \l_regex_default_catcodes_int \c_regex_all_catcodes_int
518     \int_set_eq:NN \l_regex_catcodes_int \l_regex_default_catcodes_int
519     \cs_set_nopar:Npn \regex_item_equal:n { \regex_item_caseful_equal:n }
520     \cs_set_nopar:Npn \regex_item_range:nn { \regex_item_caseful_range:nn }
521     \tl_build_one:n { \regex_branch:n { \if_false: } \fi: }
522 }
523 \cs_new_protected_nopar:Npn \regex_compile_end:
524 {
525     \regex_if_in_class:TF
526     {
527         \msg_kernel_error:nn { regex } { missing-rbrack }
528         \use:c { regex_compile_}: }
529         \prg_do_nothing: \prg_do_nothing:
530     }
531     { }
532     \if_num:w \l_regex_group_level_int > \c_zero
533         \msg_kernel_error:nnx { regex } { missing-rparen }
534         { \int_use:N \l_regex_group_level_int }
535         \prg_replicate:nn
536         { \l_regex_group_level_int }
537         {
538             \tl_build_one:n
539             {
540                 \if_false: { \fi: }
541                 \if_false: { \fi: } { 1 } { 0 } \c_true_bool
542             }
543             \tl_build_end:
544             \tl_build_one:o \l_regex_internal_regex
545         }
546         \fi:
547         \tl_build_one:n { \if_false: { \fi: } }
548         \tl_build_end:
549     }

```

(End definition for \regex\_compile:w and \regex\_compile\_end:. These functions are documented on page ??.)

\regex\_compile:n The compilation is done between \regex\_compile:w and \regex\_compile\_end:, starting in mode 0. Then \regex\_escape\_use:nnnn distinguishes special characters, escaped

alphanumerics, and raw characters, interpreting \a, \x and other sequences. The 4 trailing \prg\_do\_nothing: are needed because some functions defined later look up to 4 tokens ahead. Before ending, make sure that any \c{...} is properly closed.

```

550 \cs_new_protected:Npn \regex_compile:n #1
551   {
552     \regex_compile:w
553       \int_set:Nn \tex_escapechar:D { 92 }
554       \int_set_eq:NN \l_regex_mode_int \c_zero
555       \regex_escape_use:nnnn
556       {
557         \regex_char_if_special:NTF ##1
558           \regex_compile_special:N \regex_compile_raw:N ##1
559       }
560       {
561         \regex_char_if_alphanumeric:NTF ##1
562           \regex_compile_escaped:N \regex_compile_raw:N ##1
563         }
564         { \regex_compile_raw:N ##1 }
565         { #1 }
566       \prg_do_nothing: \prg_do_nothing:
567       \prg_do_nothing: \prg_do_nothing:
568       \int_compare:nNnT \l_regex_mode_int < \c_zero
569       {
570         \msg_kernel_error:nn { regex } { c-missing-rbrace }
571         \regex_compile_end:
572         \regex_compile_one:x
573           { \regex_item_cs:n { \exp_not:o \l_regex_internal_regex } }
574           \prg_do_nothing: \prg_do_nothing:
575           \prg_do_nothing: \prg_do_nothing:
576         }
577       \regex_compile_end:
578     }

```

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

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

```

579 \cs_new_protected:Npn \regex_compile_special:N #1
580   {
581     \cs_if_exist_use:cF { \regex_compile_#1: }
582       { \regex_compile_raw:N #1 }
583     }
584 \cs_new_protected:Npn \regex_compile_escaped:N #1
585   {
586     \cs_if_exist_use:cF { \regex_compile_/#1: }
587       { \regex_compile_raw:N #1 }
588     }

```

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

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

```

589 \cs_new_protected:Npn \regex_compile_one:x #1
590   {
591     \regex_mode_quit_c:
592     \regex_if_in_class:TF { }
593     {
594       \tl_build_one:n
595       { \regex_class:NnnnN \c_true_bool { \if_false: } \fi: }
596     }
597     \tl_build_one:x
598     {
599       \if_num:w \l_regex_catcodes_int < \c_regex_all_catcodes_int
600         \regex_item_catcode:nT { \int_use:N \l_regex_catcodes_int }
601         { \exp_not:N \exp_not:n {#1} }
602       \else:
603         \exp_not:N \exp_not:n {#1}
604       \fi:
605     }
606     \int_set_eq:NN \l_regex_catcodes_int \l_regex_default_catcodes_int
607     \regex_if_in_class:TF { } { \regex_compile_quantifier:w }
608   }

```

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

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

```

609 \cs_new_protected:Npn \regex_compile_abort_tokens:n #1
610   {
611     \use:x
612     {
613       \exp_args:No \tl_map_function:nN { \tl_to_str:n {#1} }
614       \regex_compile_raw:N
615     }
616   }
617 \cs_generate_variant:Nn \regex_compile_abort_tokens:n { x }

```

(End definition for `\regex_compile_abort_tokens:n` and `\regex_compile_abort_tokens:x`. These functions are documented on page ??.)

### 2.3.5 Quantifiers

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

```

618 \cs_new_protected:Npn \regex_compile_quantifier:w #1#2
619   {
620     \token_if_eq_meaning:NNTF #1 \regex_compile_special:N

```

```

621     {
622         \cs_if_exist_use:cF { regex_compile_quantifier_#2:w }
623             { \regex_compile_quantifier_none: #1 #2 }
624     }
625     { \regex_compile_quantifier_none: #1 #2 }
626 }
```

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

```
\regex_compile_quantifier_none:
\regex_compile_quantifier_abort:xNN
```

Those functions are called whenever there is no quantifier, or a braced construction is invalid (equivalent to no quantifier, and whatever characters were grabbed are left raw).

```

627 \cs_new_protected:Npn \regex_compile_quantifier_none:
628     { \tl_build_one:n { \if_false: { \fi: } { 1 } { 0 } \c_false_bool } }
629 \cs_new_protected:Npn \regex_compile_quantifier_abort:xNN #1#2#3
630     {
631         \regex_compile_quantifier_none:
632         \msg_kernel_warning:nxxx { regex } { invalid-quantifier } {#1} {#3}
633         \regex_compile_abort_tokens:x {#1}
634         #2 #3
635     }
```

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

```
\regex_compile_quantifier_laziness:nnNN
```

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

```

636 \cs_new_protected:Npn \regex_compile_quantifier_laziness:nnNN #1#2#3#4
637     {
638         \str_if_eq:nnTF { #3 #4 } { \regex_compile_special:N ? }
639             { \tl_build_one:n { \if_false: { \fi: } {#1} {#2} \c_true_bool } }
640             {
641                 \tl_build_one:n { \if_false: { \fi: } {#1} {#2} \c_false_bool }
642                 #3 #4
643             }
644     }
```

(End definition for `\regex_compile_quantifier_laziness:nnNN`.)

```
\regex_compile_quantifier_?:w
\regex_compile_quantifier_*:w
\regex_compile_quantifier_+:w
```

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

```

645 \cs_new_protected_nopar:cpn { regex_compile_quantifier_?:w }
646     { \regex_compile_quantifier_laziness:nnNN { 0 } { 1 } }
647 \cs_new_protected_nopar:cpn { regex_compile_quantifier_*:w }
648     { \regex_compile_quantifier_laziness:nnNN { 0 } { -1 } }
649 \cs_new_protected_nopar:cpn { regex_compile_quantifier_+:w }
650     { \regex_compile_quantifier_laziness:nnNN { 1 } { -1 } }
```

(End definition for `\regex_compile_quantifier_?:w`, `\regex_compile_quantifier_*:w`, and `\regex_compile_quantifier_+:w`.)

```
\regex_compile_quantifier_{:w
\regex_compile_quantifier_braced_i:w
\regex_compile_quantifier_braced_ii:w
\regex_compile_quantifier_braced_iii:w
```

Three possible syntaxes: `{<int>}`, `{<int>,}`, or `{<int>,<int>}`. Any other syntax causes us to abort and put whatever we collected back in the input stream, as raw characters, including the opening brace. Grab a number into `\l_regex_internal_a_int`. If the number is followed by a right brace, the range is  $[a, a]$ . If followed by a comma, grab one more number, and call the `_ii` or `_iii` auxiliary. Those auxiliaries check for a closing brace, leading to the range  $[a, \infty)$  or  $[a, b]$ , encoded as  $\{a\}\{-1\}$  and  $\{a\}\{b - a\}$ .

```
651 \cs_new_protected:cpn { regex_compile_quantifier_ \c_lbrace_str :w }
652 {
653     \regex_get_digits:NTFw \l_regex_internal_a_int
654     { \regex_compile_quantifier_braced_i:w }
655     { \regex_compile_quantifier_abort:xNN { \c_lbrace_str } }
656 }
657 \cs_new_protected:Npn \regex_compile_quantifier_braced_i:w #1#2
658 {
659     \prg_case_str:xxn { #1 #2 }
660     {
661         { \regex_compile_special:N \c_rbrace_str }
662         {
663             \exp_args:No \regex_compile_quantifier_lazyness:nnNN
664             { \int_use:N \l_regex_internal_a_int } { 0 }
665         }
666         { \regex_compile_special:N , }
667         {
668             \regex_get_digits:NTFw \l_regex_internal_b_int
669             { \regex_compile_quantifier_braced_iii:w }
670             { \regex_compile_quantifier_braced_ii:w }
671         }
672     }
673     {
674         \regex_compile_quantifier_abort:xNN
675         { \c_lbrace_str \int_use:N \l_regex_internal_a_int }
676         #1 #2
677     }
678 }
679 \cs_new_protected:Npn \regex_compile_quantifier_braced_ii:w #1#2
680 {
681     \str_if_eq:xxTF
682     { #1 #2 } { \regex_compile_special:N \c_rbrace_str }
683     {
684         \exp_args:No \regex_compile_quantifier_lazyness:nnNN
685         { \int_use:N \l_regex_internal_a_int } { -1 }
686     }
687     {
688         \regex_compile_quantifier_abort:xNN
689         { \c_lbrace_str \int_use:N \l_regex_internal_a_int , }
690         #1 #2
691     }
692 }
693 \cs_new_protected:Npn \regex_compile_quantifier_braced_iii:w #1#2
```

```

694  {
695  \str_if_eq:xxTF
696  { #1 #2 } { \regex_compile_special:N \c_rbrace_str }
697  {
698  \if_num:w \l_regex_internal_a_int > \l_regex_internal_b_int
699  \msg_kernel_error:nxxx { regex } { backwards-quantifier }
700  { \int_use:N \l_regex_internal_a_int }
701  { \int_use:N \l_regex_internal_b_int }
702  \int_zero:N \l_regex_internal_b_int
703  \else:
704  \int_sub:Nn \l_regex_internal_b_int \l_regex_internal_a_int
705  \fi:
706  \exp_args:Noo \regex_compile_quantifier_lazyness:nnNN
707  { \int_use:N \l_regex_internal_a_int }
708  { \int_use:N \l_regex_internal_b_int }
709  }
710  {
711  \regex_compile_quantifier_abort:xNN
712  {
713  \c_lbrace_str
714  \int_use:N \l_regex_internal_a_int ,
715  \int_use:N \l_regex_internal_b_int
716  }
717  #1 #2
718  }
719  }

```

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

### 2.3.6 Raw characters

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

```

720 \cs_new_protected:Npn \regex_compile_raw_error:N #1
721  {
722  \msg_kernel_error:nnx { regex } { bad-escape } {#1}
723  \regex_compile_raw:N #1
724  }

```

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

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

```

725 \cs_new_protected:Npn \regex_compile_raw:N #1#2#3
726  {
727  \regex_if_in_class:TF
728  {
729  \str_if_eq:nnTF {#2#3} { \regex_compile_special:N - }
730  { \regex_compile_range:Nw #1 }
731  {

```

```

732     \regex_compile_one:x
733         { \regex_item_equal:n { \int_value:w '#1 ~ } }
734         #2 #3
735     }
736 }
737 {
738     \regex_compile_one:x
739         { \regex_item_equal:n { \int_value:w '#1 ~ } }
740         #2 #3
741     }
742 }
```

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

`\regex_compile_range:Nw`  
`\regex_if_end_range:NNTF`

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

```

743 \prg_new_protected_conditional:Npnn \regex_if_end_range>NN #1#2 { TF }
744   {
745     \if_meaning:w \regex_compile_raw:N #1
746       \prg_return_true:
747     \else:
748       \if_meaning:w \regex_compile_special:N #1
749         \if_charcode:w ] #2
750           \prg_return_false:
751         \else:
752           \prg_return_true:
753         \fi:
754       \else:
755         \prg_return_false:
756       \fi:
757     \fi:
758   }
759 \cs_new_protected:Npn \regex_compile_range:Nw #1#2#3
760   {
761     \regex_if_end_range:NNTF #2 #3
762     {
763       \if_num:w '#1 > '#3 \exp_stop_f:
764         \msg_kernel_error:nnnx { regex } { range-backwards } {#1} {#3}
765       \else:
766         \tl_build_one:x
767         {
768           \if_num:w '#1 = '#3 \exp_stop_f:
769             \regex_item_equal:n
770           \else:
771             \regex_item_range:nn { \int_value:w '#1 ~ }
772           \fi:
773             { \int_value:w '#3 ~ }
774         }
775     \fi:
```

```

776     }
777 {
778     \msg_kernel_warning:nnxx { regex } { range-missing-end }
779     {#1} { \c_backslash_str #3 }
780     \tl_build_one:x
781     {
782         \regex_item_equal:n { \int_value:w '#1 ~ }
783         \regex_item_equal:n { \int_value:w '- ~ }
784     }
785     #2#3
786 }
787 }
```

(End definition for \regex\_compile\_range:Nw and \regex\_if\_end\_range:NNTF.)

### 2.3.7 Character properties

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

```

788 \cs_new_protected_nopar:cp{ \regex_compile_:: }
789 {
790     \exp_not:N \regex_if_in_class:TF
791     { \regex_compile_raw:N . }
792     { \regex_compile_one:x \exp_not:c { \regex_prop_:: } }
793 }
794 \cs_new_protected_nopar:cpn { \regex_prop_:: }
795 {
796     \if_num:w \l_regex_current_char_int > - \c_two
797     \exp_after:wN \regex_break_true:w
798     \fi:
799 }
```

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

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

```

800 \cs_set_protected:Npn \regex_tmp:w #1#2
801 {
802     \cs_new_protected_nopar:cp{ \regex_compile_/#1: }
803     { \regex_compile_one:x \exp_not:c { \regex_prop_#/1: } }
804     \cs_new_protected_nopar:cp{ \regex_compile_/#2: }
805     {
806         \regex_compile_one:x
807         { \regex_item_reverse:n \exp_not:c { \regex_prop_#/1: } }
808     }
809 }
810 \regex_tmp:w d D
811 \regex_tmp:w h H
812 \regex_tmp:w s S
813 \regex_tmp:w v V
```

```

814 \regex_tmp:w w W
815 \cs_new_protected_nopar:cpn { regex_compile_/N: }
816   { \regex_compile_one:x \regex_prop_N: }
(End definition for \regex_compile_/d: and others.)

```

### 2.3.8 Anchoring and simple assertions

\regex\_compile\_anchor:NF  
 \regex\_compile\_~:  
 \regex\_compile/\_A:  
 \regex\_compile/\_G:  
 \regex\_compile\_\$:  
 \regex\_compile/\_Z:  
 \regex\_compile/\_z:

```

817 \cs_new_protected:Npn \regex_compile_anchor:NF #1#2
818   {
819     \regex_if_in_class_or_catcode:TF {#2}
820     {
821       \tl_build_one:n
822         { \regex_assertion:Nn \c_true_bool { \regex_anchor:N #1 } }
823     }
824   }
825 \cs_set_protected:Npn \regex_tmp:w #1#2
826   {
827     \cs_new_protected_nopar:cpn { regex_compile_/#1: }
828       { \regex_compile_anchor:NF #2 { \regex_compile_raw_error:N #1 } }
829   }
830 \regex_tmp:w A \l_regex_min_pos_int
831 \regex_tmp:w G \l_regex_start_pos_int
832 \regex_tmp:w Z \l_regex_max_pos_int
833 \regex_tmp:w z \l_regex_max_pos_int
834 \cs_set_protected:Npn \regex_tmp:w #1#2
835   {
836     \cs_new_protected_nopar:cpn { regex_compile_#1: }
837       { \regex_compile_anchor:NF #2 { \regex_compile_raw:N #1 } }
838   }
839 \exp_args:Nx \regex_tmp:w { \iow_char:N \^ } \l_regex_min_pos_int
840 \exp_args:Nx \regex_tmp:w { \iow_char:N \$ } \l_regex_max_pos_int
(End definition for \regex_compile_anchor:NF. This function is documented on page ??.)
```

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

```

841 \cs_new_protected_nopar:cpn { regex_compile/_b: }
842   {
843     \regex_if_in_class_or_catcode:TF
844       { \regex_compile_raw_error:N b }
845     {
846       \tl_build_one:n
847         { \regex_assertion:Nn \c_true_bool { \regex_b_test: } }
848     }
849 }
```

```

850 \cs_new_protected_nopar:cpn { regex_compile_/B: }
851   {
852     \regex_if_in_class_or_catcode:TF
853     { \regex_compile_raw_error:N B }
854     {
855       \tl_build_one:n
856       { \regex_assertion:Nn \c_false_bool { \regex_b_test: } }
857     }
858   }
(End definition for \regex_compile_/B: and \regex_compile_/B:.)
```

### 2.3.9 Character classes

\regex\_compile\_[:]: 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. Look for quantifiers, unless we are still in a class after leaving one (the case of [...] \cL [...] ...]). quantifiers.

```

859 \cs_new_protected:cpn { regex_compile_]: }
860   {
861     \regex_if_in_class:TF
862     {
863       \if_num:w \l_regex_mode_int > \c_sixteen
864         \tl_build_one:n { \if_false: { \fi: } }
865       \fi:
866       \tex_advance:D \l_regex_mode_int - \c_fifteen
867       \tex_divide:D \l_regex_mode_int \c_thirteen
868       \if_int_odd:w \l_regex_mode_int \else:
869         \exp_after:wN \regex_compile_quantifier:w
870       \fi:
871     }
872     { \regex_compile_raw:N ] }
873   }
(End definition for \regex_compile_]:.)
```

\regex\_compile\_[[:]:]: In a class, left brackets might introduce a POSIX character class, or mean nothing. Immediately following \c(category), we must insert the appropriate catcode test, then parse the class; we pre-expand the catcode as an optimization. Otherwise (modes 0, -2 and -6) just parse the class. The mode is updated later.

```

874 \cs_new_protected_nopar:cpn { regex_compile_[:]}
875   {
876     \regex_if_in_class:TF
877     { \regex_compile_class_posix_test:w }
878     {
879       \regex_if_within_catcode:TF
880       {
881         \exp_after:wN \regex_compile_class_catcode:w
882         \int_use:N \l_regex_catcodes_int ;
883       }
884     { \regex_compile_class_normal:w }
885   }
```

```

886    }
(End definition for \regex_compile_[:])
```

\regex\_compile\_class\_normal:w In the “normal” case, we will insert \regex\_class:NnnnN *<boolean>* in the compiled code. The *<boolean>* is true for positive classes, and false for negative classes, characterized by a leading ^ . The auxiliary \regex\_compile\_class:TFNN also checks for a leading ] which has a special meaning.

```

887 \cs_new_protected:Npn \regex_compile_class_normal:w
888 {
889     \regex_compile_class:TFNN
890     { \regex_class:NnnnN \c_true_bool }
891     { \regex_class:NnnnN \c_false_bool }
892 }
(End definition for \regex_compile_class_normal:w.)
```

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

```

893 \cs_new_protected:Npn \regex_compile_class_catcode:w #1;
894 {
895     \if_int_compare:w \l_regex_mode_int = \c_two
896         \tl_build_one:n
897         { \regex_class:NnnnN \c_true_bool { \if_false: } \fi: }
898     \fi:
899     \int_set_eq:NN \l_regex_catcodes_int \l_regex_default_catcodes_int
900     \regex_compile_class:TFNN
901     { \regex_item_catcode:nT {#1} }
902     { \regex_item_catcode_reverse:nT {#1} }
903 }
(End definition for \regex_compile_class_catcode:w.)
```

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

```

904 \cs_new_protected:Npn \regex_compile_class:TFNN #1#2#3#4
905 {
906     \l_regex_mode_int = \int_value:w \l_regex_mode_int 3 \exp_stop_f:
907     \str_if_eq:nnTF { #3 #4 } { \regex_compile_special:N ^ }
908     {
909         \tl_build_one:n { #2 { \if_false: } \fi: }
910         \regex_compile_class_i:NN
911     }
912     {
913         \tl_build_one:n { #1 { \if_false: } \fi: }
914         \regex_compile_class_i:NN #3 #4
915     }
916 }
```

```

917 \cs_new_protected:Npn \regex_compile_class_ii:NN #1#2
918   {
919     \token_if_eq_charcode:NNTF #2 ]
920     { \regex_compile_raw:N #2 }
921     { #1 #2 }
922   }
(End definition for \regex_compile_class:TFNN and \regex_compile_class_ii:NN.)
```

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

```

923 \cs_new_protected:Npn \regex_compile_class_posix_test:w #1#2
924   {
925     \token_if_eq_meaning:NNT \regex_compile_special:N #1
926     {
927       \prg_case_str:nnn { #2 }
928       {
929         : { \regex_compile_class_posix:NNNNw }
930         = { \msg_kernel_warning:nnx { regex } { posix-unsupported } { = } }
931         . { \msg_kernel_warning:nnx { regex } { posix-unsupported } { . } }
932       }
933     }
934   }
935   \regex_compile_raw:N [ #1 #2
936   }
937 \cs_new_protected:Npn \regex_compile_class_posix:NNNNw #1#2#3#4#5#6
938   {
939     \str_if_eq:nnTF { #5 #6 } { \regex_compile_special:N ^ }
940     {
941       \bool_set_false:N \l_regex_internal_bool
942       \tl_set:Nx \l_regex_internal_a_tl { \if_false: } \fi:
943       \regex_compile_class_posix_loop:w
944     }
945   {
946     \bool_set_true:N \l_regex_internal_bool
947     \tl_set:Nx \l_regex_internal_a_tl { \if_false: } \fi:
948     \regex_compile_class_posix_loop:w #5 #6
949   }
950   }
951 \cs_new:Npn \regex_compile_class_posix_loop:w #1#2
952   {
953     \token_if_eq_meaning:NNTF \regex_compile_raw:N #1
954     { #2 \regex_compile_class_posix_loop:w }
955     { \if_false: { \fi: } \regex_compile_class_posix_end:w #1 #2 }
956   }
957 \cs_new_protected:Npn \regex_compile_class_posix_end:w #1#2#3#4
958   {
```

```

959   \str_if_eq:nntF { #1 #2 #3 #4 }
960   { \regex_compile_special:N : \regex_compile_special:N ] }
961   {
962     \cs_if_exist:cTF { regex_posix_ \l_regex_internal_a_t1 : }
963     {
964       \regex_compile_one:x
965       {
966         \bool_if:NF \l_regex_internal_bool \regex_item_reverse:n
967         \exp_not:c { regex_posix_ \l_regex_internal_a_t1 : }
968       }
969     }
970   {
971     \msg_kernel_warning:nnx { regex } { posix-unknown }
972     { \l_regex_internal_a_t1 }
973     \regex_compile_abort_tokens:x
974     {
975       [: \bool_if:NF \l_regex_internal_bool { ^ } ]
976       \l_regex_internal_a_t1 :]
977     }
978   }
979 }
980 {
981   \msg_kernel_error:nnxx { regex } { posix-missing-close }
982   { [: \l_regex_internal_a_t1 ] { #2 #4 }
983   \regex_compile_abort_tokens:x { [: \l_regex_internal_a_t1 ] }
984   #1 #2 #3 #4
985 }
986 }

```

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

### 2.3.10 Groups and alternations

```
\regex_compile_group_begin:N
\regex_compile_group_end:
```

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

```

987 \cs_new_protected:Npn \regex_compile_group_begin:N #1
988 {
989   \tl_build_one:n { #1 { \if_false: } \fi: }
990   \regex_mode_quit_c:
991   \tl_set_build:Nw \l_regex_internal_regex
992   \int_set_eq:NN \l_regex_default_catcodes_int \l_regex_catcodes_int
993   \int_incr:N \l_regex_group_level_int
994   \tl_build_one:n { \regex_branch:n { \if_false: } \fi: }
995 }
```

```

996 \cs_new_protected:Npn \regex_compile_group_end:
997 {
998     \if_num:w \l_regex_group_level_int > \c_zero
999         \tl_build_one:n { \if_false: { \fi: } }
1000     \tl_build_end:
1001     \int_set_eq:NN \l_regex_catcodes_int \l_regex_default_catcodes_int
1002     \tl_build_one:o \l_regex_internal_regex
1003     \exp_after:wN \regex_compile_quantifier:w
1004 \else:
1005     \msg_kernel_warning:nn { regex } { extra-rparen }
1006     \exp_after:wN \regex_compile_raw:N \exp_after:wN )
1007     \fi:
1008 }
(End definition for \regex_compile_group_begin:N and \regex_compile_group_end:.)
```

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

```

1009 \cs_new_protected_nopar:cpn { regex_compile_(: }
1010 {
1011     \regex_if_in_class:TF { \regex_compile_raw:N ( }
1012     { \regex_compile_lparen:w }
1013 }
1014 \cs_new_protected:Npn \regex_compile_lparen:w #1#2#3#4
1015 {
1016     \str_if_eq:nnTF { #1 #2 } { \regex_compile_special:N ? }
1017     {
1018         \cs_if_exist_use:cF
1019         { \regex_compile_special_group_\token_to_str:N #4 :w }
1020         {
1021             \msg_kernel_warning:nnx { regex } { special-group-unknown }
1022             { (? \token_to_str:N #4 ) }
1023             \regex_compile_group_begin:N \regex_group:nnnN
1024             \regex_compile_raw:N ? #3 #4
1025         }
1026     }
1027     {
1028         \regex_compile_group_begin:N \regex_group:nnnN
1029         #1 #2 #3 #4
1030     }
1031 }
```

(End definition for \regex\_compile\_(:.)

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

```

1032 \cs_new_protected_nopar:cpn { regex_compile_|: }
1033 {
1034     \regex_if_in_class:TF { \regex_compile_raw:N | }
1035     {
1036         \tl_build_one:n
```

```

1037           { \if_false: { \fi: } \regex_branch:n { \if_false: } \fi: }
1038       }
1039   }
(End definition for \regex_compile_/:.)
```

\regex\_compile\_(): Within a class, parentheses are not special. Outside, close a group.

```

1040 \cs_new_protected_nopar:cpn { regex_compile_(): }
1041   {
1042     \regex_if_in_class:TF { \regex_compile_raw:N }
1043     { \regex_compile_group_end: }
1044   }
(End definition for \regex_compile_():)
```

\regex\_compile\_special\_group\_::w Non-capturing, and resetting groups are easy to take care of during compilation; for those groups, the harder parts will come when building.

```

1045 \cs_new_protected_nopar:cpn { regex_compile_special_group_::w }
1046   { \regex_compile_group_begin:N \regex_group_no_capture:nnN }
1047 \cs_new_protected_nopar:cpn { regex_compile_special_group_|:w }
1048   { \regex_compile_group_begin:N \regex_group_resetting:nnN }
(End definition for \regex_compile_special_group_::w. This function is documented on page ??.)
```

\regex\_compile\_special\_group\_i:w The match can be made case-insensitive by setting the option with (?i); the original behaviour is restored by (?-i). This is the only supported option.

```

1049 \cs_new_protected:Npn \regex_compile_special_group_i:w #1#2
1050   {
1051     \str_if_eq:nnTF { #1 #2 } { \regex_compile_special:N }
1052     {
1053       \cs_set_nopar:Npn \regex_item_equal:n { \regex_item_caseless_equal:n }
1054       \cs_set_nopar:Npn \regex_item_range:nn { \regex_item_caseless_range:nn }
1055     }
1056     {
1057       \msg_kernel_warning:nnx { regex } { unknown-option } { (?i #2 }
1058       \regex_compile_raw:N (
1059       \regex_compile_raw:N ?
1060       \regex_compile_raw:N i
1061       #1 #2
1062     }
1063   }
1064 \cs_new_protected_nopar:cpn { regex_compile_special_group_-:w } #1#2#3#4
1065   {
1066     \str_if_eq:nnTF { #1 #2 #3 #4 }
1067     { \regex_compile_raw:N i \regex_compile_special:N }
1068     {
1069       \cs_set_nopar:Npn \regex_item_equal:n { \regex_item_caseful_equal:n }
1070       \cs_set_nopar:Npn \regex_item_range:nn { \regex_item_caseful_range:nn }
1071     }
1072     {
1073       \msg_kernel_warning:nnx { regex } { unknown-option } { (?-#2#4 }
1074       \regex_compile_raw:N (
```

```

1075     \regex_compile_raw:N ?
1076     \regex_compile_raw:N -
1077     #1 #2 #3 #4
1078   }
1079 }
(End definition for \regex_compile_special_group_i:w and \regex_compile_special_group_-:w.)
```

### 2.3.11 Catcodes and csnames

\regex\_compile/\_c:  
\regex\_compile\_c\_test:NN

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.

```

1080 \cs_new_protected:cpn { regex_compile/_c: }
1081   { \regex_chk_c_allowed:T { \regex_compile_c_test:NN } }
1082 \cs_new_protected:Npn \regex_compile_c_test:NN #1#2
1083   {
1084     \token_if_eq_meaning:NNTF #1 \regex_compile_raw:N
1085     {
1086       \int_if_exist:cTF { c_regex_catcode_#2_int }
1087       {
1088         \int_set_eq:Nc \l_regex_catcodes_int { c_regex_catcode_#2_int }
1089         \l_regex_mode_int
1090         = \if_case:w \l_regex_mode_int \c_two \else: \c_six \fi:
1091       }
1092     }
1093   { \cs_if_exist_use:cF { regex_compile_c_#2:w } }
1094   {
1095     \msg_kernel_error:nnx { regex } { c-missing-category } {#2}
1096     #1 #2
1097   }
1098 }
```

(End definition for \regex\_compile/\_c: and \regex\_compile\_c\_test:NN.)

\regex\_compile\_c\_[:w]  
\regex\_compile\_c\_lbrack\_loop:NN  
\regex\_compile\_c\_lbrack\_add:N  
\regex\_compile\_c\_lbrack\_end:

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

```

1099 \cs_new_protected:cpn { regex_compile_c_[:w] } #1#2
1100   {
1101     \l_regex_mode_int
1102     = \if_case:w \l_regex_mode_int \c_two \else: \c_six \fi:
1103     \int_zero:N \l_regex_catcodes_int
1104     \str_if_eq:nnTF { #1 #2 } { \regex_compile_special:N ^ }
1105     {
1106       \bool_set_false:N \l_regex_catcodes_bool
1107       \regex_compile_c_lbrack_loop:NN
1108     }
1109   {
1110     \bool_set_true:N \l_regex_catcodes_bool
1111     \regex_compile_c_lbrack_loop:NN
1112     #1 #2
1113 }
```

```

1113     }
1114   }
1115 \cs_new_protected:Npn \regex_compile_c_lbrack_loop:NN #1#2
1116   {
1117     \token_if_eq_meaning:NNTF #1 \regex_compile_raw:N
1118     {
1119       \int_if_exist:cTF { c_regex_catcode_#2_int }
1120       {
1121         \exp_args:Nc \regex_compile_c_lbrack_add:N
1122           { c_regex_catcode_#2_int }
1123         \regex_compile_c_lbrack_loop:NN
1124       }
1125     }
1126   {
1127     \token_if_eq_charcode:NNTF #2 ]
1128       { \regex_compile_c_lbrack_end: }
1129   }
1130   {
1131     \msg_kernel_error:nnx { regex } { c-missing-rbrack } {#2}
1132     \regex_compile_c_lbrack_end:
1133     #1 #2
1134   }
1135 }
1136 \cs_new_protected:Npn \regex_compile_c_lbrack_add:N #1
1137   {
1138     \if_int_odd:w \int_eval:w \l_regex_catcodes_int / #1 \int_eval_end:
1139     \else:
1140       \tex_advance:D \l_regex_catcodes_int #1
1141     \fi:
1142   }
1143 \cs_new_protected_nopar:Npn \regex_compile_c_lbrack_end:
1144   {
1145     \if_meaning:w \c_false_bool \l_regex_catcodes_bool
1146       \int_set:Nn \l_regex_catcodes_int
1147         { \c_regex_all_catcodes_int - \l_regex_catcodes_int }
1148     \fi:
1149   }

```

(End definition for \regex\_compile\_c\_[:w and others.)

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

```

1150 \cs_new_protected_nopar:cpx { regex_compile_c_ \c_lbrace_str :w }
1151   {
1152     \regex_compile:w
1153       \regex_disable_submatches:
1154     \l_regex_mode_int
1155       = - \if_case:w \l_regex_mode_int \c_two \else: \c_six \fi:
1156   }

```

(End definition for \regex\_compile\_c\_{{}. This function is documented on page ??.)

- \regex\_compile\_{{}: Non-escaped right braces are only special if they appear when compiling the regular expression for a csname, but not within a class: \c{{}} matches the control sequences \} and \{... Admittedly, that would be better done as \c{{}}. So, end compiling the inner regex (this closes any dangling class or group). Then insert the corresponding test in the outer regex.

```
1157 \cs_new_protected:cpn { regex_compile_ \c_rbrace_str : }
1158   {
1159     \regex_if_in_cs:TF
1160     {
1161       \regex_compile_end:
1162       \regex_compile_one:x
1163       { \regex_item_cs:n { \exp_not:o \l_regex_internal_regex } }
1164     }
1165     { \exp_after:wN \regex_compile_raw:N \c_rbrace_str }
1166   }
```

(End definition for \regex\_compile\_{{}. This function is documented on page ??.)

### 2.3.12 Raw token lists with \u

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

```
1167 \cs_new_protected:cpn { regex_compile/_u: } #1#2
1168   {
1169     \regex_if_in_class_or_catcode:TF
1170     { \regex_compile_error:N u #1 #2 }
1171     {
1172       \str_if_eq:xxTF {#1#2} { \regex_compile_special:N \c_lbrace_str }
1173       {
1174         \tl_set:Nx \l_regex_internal_a_tl { \if_false: } \fi:
1175         \regex_compile_u_loop:NN
1176       }
1177       {
1178         \msg_kernel_error:nn { regex } { u-missing-lbrace }
1179         \regex_compile_error:N u #1 #2
1180       }
1181     }
1182   }
1183 \cs_new:Npn \regex_compile_u_loop:NN #1#2
1184   {
1185     \token_if_eq_meaning:NNTF #1 \regex_compile_raw:N
1186     { #2 \regex_compile_u_loop:NN }
```

```

1187    {
1188        \token_if_eq_meaning:NNTF #1 \regex_compile_special:N
1189        {
1190            \exp_after:wN \token_if_eq_charcode:NNTF \c_rbrace_str #2
1191            { \if_false: { \fi: } \regex_compile_u_end: }
1192            { #2 \regex_compile_u_loop>NN }
1193        }
1194        {
1195            \if_false: { \fi: }
1196            \msg_kernel_error:nnx { regex } { u-missing-rbrace } {#2}
1197            \regex_compile_u_end:
1198            #1 #2
1199        }
1200    }
1201 }

```

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

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

```

1202 \cs_new_protected:Npn \regex_compile_u_end:
1203 {
1204     \tl_set:Nv \l_regex_internal_a_tl { \l_regex_internal_a_tl }
1205     \if_num:w \l_regex_mode_int = \c_zero
1206         \regex_compile_u_not_cs:
1207     \else:
1208         \regex_compile_u_in_cs:
1209     \fi:
1210 }

```

(End definition for `\regex_compile_u_end::`)

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

```

1211 \cs_new_protected:Npn \regex_compile_u_in_cs:
1212 {
1213     \exp_args:NNo \str_gset_other:Nn \g_regex_internal_tl
1214     { \l_regex_internal_a_tl }
1215     \tl_build_one:x
1216     {
1217         \tl_map_function:NN \g_regex_internal_tl
1218         \regex_compile_u_in_cs_aux:n
1219     }
1220 }
1221 \cs_new:Npn \regex_compile_u_in_cs_aux:n #1
1222 {
1223     \regex_class:NnnnN \c_true_bool
1224     { \regex_item_caseful_equal:n { \int_value:w '#1 } }
1225     { 1 } { 0 } \c_false_bool

```

```

1226    }
(End definition for \regex_compile_u_in_cs::)
```

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

```

1227 \cs_new_protected:Npn \regex_compile_u_not_cs:
1228 {
1229     \exp_args:No \tl_analysis_map_inline:nn { \l_regex_internal_a_tl }
1230     {
1231         \tl_build_one:n
1232         {
1233             \regex_class:NnnnN \c_true_bool
1234             {
1235                 \if_num:w ##2 = \c_zero
1236                     \regex_item_exact_cs:c { \exp_after:wN \cs_to_str:N ##1 }
1237                 \else:
1238                     \regex_item_exact:nn { \int_value:w "##2 } { ##3 }
1239                 \fi:
1240             }
1241             { 1 } { 0 } \c_false_bool
1242         }
1243     }
1244 }
```

(End definition for \regex\_compile\_u\_not\_cs::)

### 2.3.13 Other

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

```

1245 \cs_new_protected_nopar:cpn { regex_compile_/_K: }
1246 {
1247     \int_compare:nNnTF \l_regex_mode_int = \c_zero
1248     {
1249         \tl_build_one:n { \regex_command_K: }
1250     }
(End definition for \regex_compile_/_K::)
```

### 2.3.14 Showing regexes

\regex\_show\_aux:Nx Within a \tl\_set\_build:Nw ... \tl\_build\_end: group, we redefine all the function that can appear in a compiled regex, then run the regex. The result is then shown.

```

1251 \cs_new_protected:Npn \regex_show_aux:Nx #1#2
1252 {
1253     \tl_set_build:Nw \l_regex_internal_a_tl
1254     \cs_set_protected_nopar:Npn \regex_branch:n
1255     {
```

```

1256     \seq_pop_right:NN \l_regex_show_prefix_seq \l_regex_internal_a_t1
1257     \regex_show_one:n { +-branch }
1258     \seq_put_right:No \l_regex_show_prefix_seq \l_regex_internal_a_t1
1259     \use:n
1260   }
1261   \cs_set_protected_nopar:Npn \regex_group:nnnN
1262     { \regex_show_group_aux:nnnnN { } }
1263   \cs_set_protected_nopar:Npn \regex_group_no_capture:nnnN
1264     { \regex_show_group_aux:nnnnN { ~(no~capture) } }
1265   \cs_set_protected_nopar:Npn \regex_group_resetting:nnnN
1266     { \regex_show_group_aux:nnnnN { ~(resetting) } }
1267   \cs_set_eq:NN \regex_class:NnnnN \regex_show_class:NnnnN
1268   \cs_set_protected_nopar:Npn \regex_command_K:
1269     { \regex_show_one:n { reset~match~start~(\iow_char:N\K) } }
1270   \cs_set_protected:Npn \regex_assertion:Nn ##1##2
1271     { \regex_show_one:n { \bool_if:NF ##1 { negative~ } assertion:~##2 } }
1272   \cs_set_nopar:Npn \regex_b_test: { word~boundary }
1273   \cs_set_eq:NN \regex_anchor:N \regex_show_anchor_to_str:N
1274   \cs_set_protected:Npn \regex_item_caseful_equal:n ##1
1275     { \regex_show_one:n { char~code~\int_eval:n{##1} } }
1276   \cs_set_protected:Npn \regex_item_caseful_range:nn ##1##2
1277     { \regex_show_one:n { range~[\int_eval:n{##1}, \int_eval:n{##2}] } }
1278   \cs_set_protected:Npn \regex_item_caseless_equal:n ##1
1279     { \regex_show_one:n { char~code~\int_eval:n{##1}~(caseless) } }
1280   \cs_set_protected:Npn \regex_item_caseless_range:nn ##1##2
1281     {
1282       \regex_show_one:n
1283         { Range~[\int_eval:n{##1}, \int_eval:n{##2}]~(caseless) }
1284     }
1285   \cs_set_protected:Npn \regex_item_catcode:nT
1286     { \regex_show_item_catcode:NnT \c_true_bool }
1287   \cs_set_protected:Npn \regex_item_catcode_reverse:nT
1288     { \regex_show_item_catcode:NnT \c_false_bool }
1289   \cs_set_protected:Npn \regex_item_reverse:n
1290     { \regex_show_scope:nn { Reversed~match } }
1291   \cs_set_protected:Npn \regex_item_exact:nn ##1##2
1292     { \regex_show_one:n { char~##2,~catcode~##1 } }
1293   \cs_set_protected:Npn \regex_item_exact_cs:c ##1
1294     { \regex_show_one:n { control~sequence~\iow_char:N\##1 } }
1295   \cs_set_protected:Npn \regex_item_cs:n
1296     { \regex_show_scope:nn { control~sequence } }
1297   \cs_set_cpn { regex_prop_..: } { \regex_show_one:n { any~token } }
1298   \seq_clear:N \l_regex_show_prefix_seq
1299   \regex_show_push:n { ~ }
1300   #1
1301 \tl_build_end:
1302   \msg_aux_show:x { > Compiled~regex~#2: \l_regex_internal_a_t1 }
1303 }

```

(End definition for \regex\_show\_aux:Nx.)

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

```

1304 \cs_new_protected:Npn \regex_show_one:n #1
1305   {
1306     \int_incr:N \l_regex_show_lines_int
1307     \tl_build_one:x
1308     { \iow_newline: \seq_use:N \l_regex_show_prefix_seq #1 }
1309   }
(End definition for \regex_show_one:n.)
```

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

```

1310 \cs_new_protected:Npn \regex_show_push:n #1
1311   { \seq_put_right:Nx \l_regex_show_prefix_seq { #1 ~ } }
1312 \cs_new_protected:Npn \regex_show_pop:
1313   { \seq_pop_right:NN \l_regex_show_prefix_seq \l_regex_internal_a_tl }
1314 \cs_new_protected:Npn \regex_show_scope:nn #1#2
1315   {
1316     \regex_show_one:n {#1}
1317     \regex_show_push:n { ~ }
1318     #2
1319     \regex_show_pop:
1320   }
(End definition for \regex_show_push:n, \regex_show_pop:, and \regex_show_scope:nn.)
```

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

```

1321 \cs_new_protected:Npn \regex_show_group_aux:nnnnN #1#2#3#4#5
1322   {
1323     \regex_show_one:n { ,-group~begin #1 }
1324     \regex_show_push:n { | }
1325     \use_ii:nn #2
1326     \regex_show_pop:
1327     \regex_show_one:n
1328     { '-group~end \regex_msg_repeated:nnN {#3} {#4} #5 }
1329   }
(End definition for \regex_show_group_aux:nnnnN.)
```

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

```

1330 \cs_set:Npn \regex_show_class:NnnnN #1#2#3#4#5
1331   {
1332     \tl_set_build:Nw \l_regex_internal_a_tl
```

```

1333 \int_zero:N \l_regex_show_lines_int
1334 \regex_show_push:n {~}
1335 #2
1336 \exp_last_unbraced:Nf
1337 \prg_case_int:nnn { \l_regex_show_lines_int }
1338 {
1339   {0}
1340   {
1341     \tl_build_end:
1342     \regex_show_one:n { \bool_if:NTF #1 { Fail } { Pass } }
1343   }
1344   {1}
1345   {
1346     \tl_build_end:
1347     \bool_if:NTF #1
1348     {
1349       #2
1350       \tl_build_one:n { \regex_msg_repeated:nnN {#3} {#4} #5 }
1351     }
1352     {
1353       \regex_show_one:n
1354       { Don't~match~\regex_msg_repeated:nnN {#3} {#4} #5 }
1355       \tl_build_one:o \l_regex_internal_a_tl
1356     }
1357   }
1358 }
1359 {
1360   \tl_build_end:
1361   \regex_show_one:n
1362   {
1363     \bool_if:NTF #1 { M } { Don't~m } attach
1364     \regex_msg_repeated:nnN {#3} {#4} #5
1365   }
1366   \tl_build_one:o \l_regex_internal_a_tl
1367 }
1368 }

(End definition for \regex_show_class:NnnnN.)

```

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

```

1369 \cs_new:Npn \regex_show_anchor_to_str:N #1
1370 {
1371   anchor-at~
1372   \prg_case_str:nnn { #1 }
1373   {
1374     { \l_regex_min_pos_int } { start~(\iow_char:N\\A) }
1375     { \l_regex_start_pos_int } { start-of-match~(\iow_char:N\\G) }
1376     { \l_regex_max_pos_int } { end~(\iow_char:N\\Z) }
1377   }

```

```

1378     { <error:~'#1'~not~recognized> }
1379   }

```

(End definition for \regex\_show\_anchor\_to\_str:N.)

\regex\_show\_item\_catcode:NnT

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

```

1380 \cs_new_protected:Npn \regex_show_item_catcode:NnT #1#2
1381   {
1382     \seq_set_split:Nnn \l_regex_internal_seq { } { CBEMTPUDSLOA }
1383     \seq_set_filter:NNn \l_regex_internal_seq \l_regex_internal_seq
1384       { \int_if_odd_p:n { #2 / \int_use:c { c_regex_catcode_##1_int } } }
1385     \regex_show_scope:nn
1386     {
1387       categories~\seq_use:N \l_regex_internal_seq, ~
1388       \bool_if:NF #1 { negative~ } class
1389     }
1390   }

```

(End definition for \regex\_show\_item\_catcode:NnT.)

## 2.4 Building

### 2.4.1 Variables used while building

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

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

```

1391 \int_new:N \l_regex_min_state_int
1392 \int_new:N \l_regex_max_state_int

```

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

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

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

```

1393 \int_new:N \l_regex_left_state_int
1394 \int_new:N \l_regex_right_state_int
1395 \seq_new:N \l_regex_left_state_seq
1396 \seq_new:N \l_regex_right_state_seq

```

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

\l\_regex\_capturing\_group\_int

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

```
1397 \int_new:N \l_regex_capturing_group_int
```

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

### 2.4.2 Framework

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

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

We strive to preserve the following properties while building.

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

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

```

1398 \cs_new_protected:Npn \regex_build:n #1
1399  {
1400    \regex_compile:n {#1}
1401    \regex_build:N \l_regex_internal_regex
1402  }
1403 \cs_new_protected:Npn \regex_build:N #1
1404  {

```

```

1405 <trace>    \trace_push:nnn { regex } { 1 } { regex_build }
1406     \int_set:Nn \tex_escapechar:D { 92 }
1407     \int_zero:N \l_regex_capturing_group_int
1408     \int_set_eq:NN \l_regex_max_state_int \l_regex_min_state_int
1409     \regex_build_new_state:
1410     \regex_build_new_state:
1411     \regex_toks_put_right:Nn \l_regex_left_state_int
1412         { \regex_action_start_wildcard: }
1413     \regex_group:nnN {#1} { 1 } { 0 } \c_false_bool
1414     \regex_toks_put_right:Nn \l_regex_right_state_int
1415         { \regex_action_success: }
1416 <trace>    \regex_trace_states:n { 2 }
1417 <trace>    \trace_pop:nnn { regex } { 1 } { regex_build }
1418 }
```

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

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

```

1419 \cs_new_protected:Npn \regex_build_for_cs:n #1
1420 {
1421 <trace>    \trace_push:nnn { regex } { 1 } { regex_build_for_cs }
1422     \int_set_eq:NN \l_regex_max_state_int \l_regex_min_state_int
1423     \regex_build_new_state:
1424     \regex_build_new_state:
1425     \regex_push_lr_states:
1426         #1
1427     \regex_pop_lr_states:
1428     \regex_toks_put_right:Nn \l_regex_right_state_int
1429         {
1430             \if_num:w \l_regex_current_pos_int = \l_regex_max_pos_int
1431                 \exp_after:wN \regex_action_success:
1432             \fi:
1433         }
1434 <trace>    \regex_trace_states:n { 2 }
1435 <trace>    \trace_pop:nnn { regex } { 1 } { regex_build_for_cs }
1436 }
```

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

### 2.4.3 Helpers for building an nfa

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

```

1437 \cs_new_protected_nopar:Npn \regex_push_lr_states:
1438 {
1439     \seq_push:No \l_regex_left_state_seq
1440         { \int_use:N \l_regex_left_state_int }
1441     \seq_push:No \l_regex_right_state_seq
```

```

1442           { \int_use:N \l_regex_right_state_int }
1443     }
1444 \cs_new_protected_nopar:Npn \regex_pop_lr_states:
1445   {
1446     \seq_pop>NN \l_regex_left_state_seq \l_regex_internal_a_t1
1447     \int_set:Nn \l_regex_left_state_int \l_regex_internal_a_t1
1448     \seq_pop>NN \l_regex_right_state_seq \l_regex_internal_a_t1
1449     \int_set:Nn \l_regex_right_state_int \l_regex_internal_a_t1
1450   }
(End definition for \regex_push_lr_states: and \regex_pop_lr_states:. These functions are documented on page ??.)
```

\regex\_toks\_put\_left:Nx  
\regex\_toks\_put\_right:Nx  
\regex\_toks\_put\_right:Nn

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

```

1451 \cs_new_protected:Npn \regex_toks_put_left:Nx #1#2
1452   {
1453     \cs_set_nopar:Npx \regex_tmp:w { #2 }
1454     \tex_toks:D #1 \exp_after:wN \exp_after:wN \exp_after:wN
1455       { \exp_after:wN \regex_tmp:w \tex_the:D \tex_toks:D #1 }
1456   }
1457 \cs_new_protected:Npn \regex_toks_put_right:Nx #1#2
1458   {
1459     \cs_set_nopar:Npx \regex_tmp:w { #2 }
1460     \tex_toks:D #1 \exp_after:wN
1461       { \tex_the:D \tex_toks:D \exp_after:wN #1 \regex_tmp:w }
1462   }
1463 \cs_new_protected:Npn \regex_toks_put_right:Nn #1#2
1464   { \tex_toks:D #1 \exp_after:wN { \tex_the:D \tex_toks:D #1 #2 } }
(End definition for \regex_toks_put_left:Nx. This function is documented on page ??.)
```

\regex\_build\_transition\_left:NNN  
\regex\_build\_transition\_right:nNn

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

```

1465 \cs_new_protected:Npn \regex_build_transition_left:NNN #1#2#3
1466   { \regex_toks_put_left:Nx #2 { #1 { \int_eval:n { #3 - #2 } } } }
1467 \cs_new_protected:Npn \regex_build_transition_right:nNn #1#2#3
1468   { \regex_toks_put_right:Nx #2 { #1 { \int_eval:n { #3 - #2 } } } }
(End definition for \regex_build_transition_left:NNN and \regex_build_transition_right:nNn. These functions are documented on page ??.)
```

\regex\_build\_new\_state:

Add a new empty state to the NFA. Then update the `left`, `right`, and `max` states, so that the `right` state is the new empty state, and the `left` state points to the previously “current” state.

```

1469 \cs_new_protected_nopar:Npn \regex_build_new_state:
1470   {
```

```

1471 /*trace>
1472     \trace:nnx { regex } { 2 }
1473     {
1474         regex~new~state~
1475         L=\int_use:N \l_regex_left_state_int ~ -> ~
1476         R=\int_use:N \l_regex_right_state_int ~ -> ~
1477         M=\int_use:N \l_regex_max_state_int ~ -> ~
1478         \int_eval:n { \l_regex_max_state_int + \c_one }
1479     }
1480 
```

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

`\regex_build_transitions_lazyness:NNNNN` This function creates a new state, and puts two transitions starting from the old current state. The order of the transitions is controlled by #1, true for lazy quantifiers, and false for greedy quantifiers.

```

1486 \cs_new_protected:Npn \regex_build_transitions_lazyness:NNNNN #1#2#3#4#5
1487 {
1488     \regex_build_new_state:
1489     \regex_toks_put_right:Nx \l_regex_left_state_int
1490     {
1491         \if_meaning:w \c_true_bool #1
1492             #2 { \int_eval:n { #3 - \l_regex_left_state_int } }
1493             #4 { \int_eval:n { #5 - \l_regex_left_state_int } }
1494         \else:
1495             #4 { \int_eval:n { #5 - \l_regex_left_state_int } }
1496             #2 { \int_eval:n { #3 - \l_regex_left_state_int } }
1497         \fi:
1498     }
1499 }
```

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

#### 2.4.4 Building classes

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

```

1500 \cs_new_protected:Npn \regex_class:NnnnN #1#2#3#4#5
1501 {
1502     \cs_set_nopar:Npx \regex_tests_action_cost:n ##1
1503     {
1504         \exp_not:n { \exp_not:n {\#2} }
```

```

1505           \bool_if:NTF #1
1506             { \regex_break_point:TF { \regex_action_cost:n {##1} } { } }
1507             { \regex_break_point:TF { } { \regex_action_cost:n {##1} } }
1508         }
1509     \if_case:w - #4 \exp_stop_f:
1510       \regex_class_repeat:n {#3}
1511     \or: \regex_class_repeat:nN {#3}      #5
1512     \else: \regex_class_repeat:nnN {#3} {#4} #5
1513     \fi:
1514   }
1515 \cs_new:Npn \regex_tests_action_cost:n { \regex_action_cost:n }
(End definition for \regex_class:NnnN. This function is documented on page ??.)
```

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

```

1516 \cs_new_protected:Npn \regex_class_repeat:n #1
1517   {
1518     \prg_replicate:nn {#1}
1519     {
1520       \regex_build_new_state:
1521       \regex_build_transition_right:nNn \regex_tests_action_cost:n
1522         \l_regex_left_state_int \l_regex_right_state_int
1523     }
1524   }
(End definition for \regex_class_repeat:n.)
```

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

```

1525 \cs_new_protected:Npn \regex_class_repeat:nN #1#2
1526   {
1527     \if_num:w #1 = \c_zero
1528       \regex_build_transitions_lazyness:NNNNN #2
1529         \regex_action_free:n      \l_regex_right_state_int
1530         \regex_tests_action_cost:n \l_regex_left_state_int
1531     \else:
1532       \regex_class_repeat:n {#1}
1533       \int_set_eq:NN \l_regex_internal_a_int \l_regex_left_state_int
1534       \regex_build_transitions_lazyness:NNNNN #2
1535         \regex_action_free:n \l_regex_right_state_int
1536         \regex_action_free:n \l_regex_internal_a_int
1537     \fi:
1538   }
(End definition for \regex_class_repeat:nN.)
```

```
\regex_class_repeat:nnN
```

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

```
1539 \cs_new_protected:Npn \regex_class_repeat:nnN #1#2#3
1540 {
1541     \regex_class_repeat:n {#1}
1542     \int_set:Nn \l_regex_internal_a_int
1543         { \l_regex_max_state_int + #2 - \c_one }
1544     \prg_replicate:nn {#2}
1545     {
1546         \regex_build_transitions_laziness:NNNNN #3
1547             \regex_action_free:n      \l_regex_internal_a_int
1548             \regex_tests_action_cost:n \l_regex_right_state_int
1549     }
1550 }
```

(End definition for `\regex_class_repeat:nnN`.)

#### 2.4.5 Building groups

```
\regex_group_aux:nnnnN
```

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

```
1551 \cs_new_protected:Npn \regex_group_aux:nnnnN #1#2#3#4#5
1552 {
1553     \trace_push:nnn { regex } { 1 } { regex_group }
1554     \if_num:w #3 = \c_zero
1555         \regex_build_new_state:
1556     \assert\assert_int:n { \l_regex_max_state_int = \l_regex_right_state_int + 1 }
1557         \regex_build_transition_right:nNn \regex_action_free_group:n
1558             \l_regex_left_state_int \l_regex_right_state_int
1559     \fi:
1560     \regex_build_new_state:
1561     \regex_push_lr_states:
1562     #2
1563     \regex_pop_lr_states:
1564     \if_case:w - #4 \exp_stop_f:
1565         \regex_group_repeat:nn {#1} {#3}
1566     \or: \regex_group_repeat:nnN {#1} {#3}      #5
1567     \else: \regex_group_repeat:nnnN {#1} {#3} {#4} #5
```

```

1568           \fi:
1569   <trace>      \trace_pop:nnn { regex } { 1 } { regex_group }
1570   }
(End definition for \regex_group_aux:nnnnN.)
```

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

```

1571 \cs_new_protected:Npn \regex_group:nnnN #1
1572   {
1573     \exp_args:No \regex_group_aux:nnnnN
1574       { \int_use:N \l_regex_capturing_group_int }
1575       {
1576         \int_incr:N \l_regex_capturing_group_int
1577         #1
1578       }
1579   }
1580 \cs_new_protected_nopar:Npn \regex_group_no_capture:nnnN
1581   { \regex_group_aux:nnnnN { -1 } }
```

(End definition for \regex\_group:nnnN and \regex\_group\_no\_capture:nnnN. These functions are documented on page ??.)

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

```

1582 \cs_new_protected:Npn \regex_group_resetting:nnnN #1
1583   {
1584     \regex_group_aux:nnnnN { -1 }
1585     {
1586       \exp_args:Noo \regex_group_resetting_loop:nnNn
1587         { \int_use:N \l_regex_capturing_group_int }
1588         { \int_use:N \l_regex_capturing_group_int }
1589         #1
1590         { ?? \prg_map_break:n } { }
1591         \prg_break_point:n { }
1592     }
1593   }
1594 \cs_new_protected:Npn \regex_group_resetting_loop:nnNn #1#2#3#4
1595   {
1596     \use_none:nn #3 { \int_set:Nn \l_regex_capturing_group_int {#1} }
1597     \int_set:Nn \l_regex_capturing_group_int {#2}
1598     #3 {#4}
1599     \exp_args:Nf \regex_group_resetting_loop:nnNn
1600       { \int_max:nn {#1} { \l_regex_capturing_group_int } }
1601       {#2}
1602   }
```

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

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

```

1603 \cs_new_protected:Npn \regex_branch:n #1
1604 {
1605   \trace_push:nnn { regex } { 1 } { regex_branch }
1606   \regex_build_new_state:
1607   \seq_get:NN \l_regex_left_state_seq \l_regex_internal_a_tl
1608   \int_set:Nn \l_regex_left_state_int \l_regex_internal_a_tl
1609   \regex_build_transition_right:nNn \regex_action_free:n
1610     \l_regex_left_state_int \l_regex_right_state_int
1611   #1
1612   \seq_get:NN \l_regex_right_state_seq \l_regex_internal_a_tl
1613   \regex_build_transition_right:nNn \regex_action_free:n
1614     \l_regex_right_state_int \l_regex_internal_a_tl
1615   \trace_pop:nnn { regex } { 1 } { regex_branch }
1616 }
```

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

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

```

1617 \cs_new_protected:Npn \regex_group_repeat:nn #1#2
1618 {
1619   \if_num:w #2 = \c_zero
1620     \int_set:Nn \l_regex_max_state_int
1621       { \l_regex_left_state_int - \c_one }
1622     \regex_build_new_state:
1623   \else:
1624     \regex_group_repeat_aux:n {#2}
1625     \regex_group_submatches:nNN {#1}
1626       \l_regex_internal_a_int \l_regex_right_state_int
1627     \regex_build_new_state:
1628   \fi:
1629 }
```

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

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

```

1630 \cs_new_protected:Npn \regex_group_submatches:nNN #1#2#3
1631 {
1632   \if_num:w #1 > \c_minus_one
1633     \regex_toks_put_left:Nx #2 { \regex_action_submatch:n { #1 < } }
1634     \regex_toks_put_left:Nx #3 { \regex_action_submatch:n { #1 > } }
```

```

1635     \fi:
1636 }
(End definition for \regex_group_submatches:nNN.)
```

### \regex\_group\_repeat\_aux:n

Here we repeat `\toks` ranging from `left_state` to `max_state`,  $\#1 > 0$  times. First add a transition so that the copies will “chain” properly. Compute the shift `c` between the original copy and the last copy we want. Shift the `right_state` and `max_state` to their final values. We then want to perform `c` copy operations. At the end, `b` is equal to the `max_state`, and `a` points to the left of the last copy of the group.

```

1637 \cs_new_protected:Npn \regex_group_repeat_aux:n #1
1638 {
1639     \regex_build_transition_right:nNn \regex_action_free:n
1640         \l_regex_right_state_int \l_regex_max_state_int
1641     \int_set_eq:NN \l_regex_internal_a_int \l_regex_left_state_int
1642     \int_set_eq:NN \l_regex_internal_b_int \l_regex_max_state_int
1643     \if_num:w \int_eval:w #1 > \c_one
1644         \int_set:Nn \l_regex_internal_c_int
1645         {
1646             ( #1 - \c_one )
1647             * ( \l_regex_internal_b_int - \l_regex_internal_a_int )
1648         }
1649         \tex_advance:D \l_regex_right_state_int \l_regex_internal_c_int
1650         \tex_advance:D \l_regex_max_state_int \l_regex_internal_c_int
1651         \prg_replicate:nn \l_regex_internal_c_int
1652         {
1653             \tex_toks:D \l_regex_internal_b_int
1654                 = \tex_toks:D \l_regex_internal_a_int
1655             \tex_advance:D \l_regex_internal_a_int \c_one
1656             \tex_advance:D \l_regex_internal_b_int \c_one
1657         }
1658     \fi:
1659 }
```

(End definition for \regex\_group\_repeat\_aux:n.)

### \regex\_group\_repeat:nnN

This function is called to repeat a group at least  $n$  times; the case  $n = 0$  is very different from  $n > 0$ . Assume first that  $n = 0$ . Insert submatch tracking information at the start and end of the group, add a free transition from the right end to the “true” left state `a` (remember: in this case we had added an extra state before the left state). This forms the loop, which we break away from by adding a free transition from `a` to a new state.

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

```

1660 \cs_new_protected:Npn \regex_group_repeat:nnN #1#2#3
1661 {
```

```

1662 \if_num:w #2 = \c_zero
1663   \regex_group_submatches:nNN {#1}
1664     \l_regex_left_state_int \l_regex_right_state_int
1665   \int_set:Nn \l_regex_internal_a_int
1666     { \l_regex_left_state_int - \c_one }
1667   \regex_build_transition_right:nNn \regex_action_free:n
1668     \l_regex_right_state_int \l_regex_internal_a_int
1669   \regex_build_new_state:
1670   \if_meaning:w \c_true_bool #3
1671     \regex_build_transition_left:NNN \regex_action_free:n
1672       \l_regex_internal_a_int \l_regex_right_state_int
1673   \else:
1674     \regex_build_transition_right:nNn \regex_action_free:n
1675       \l_regex_internal_a_int \l_regex_right_state_int
1676   \fi:
1677 \else:
1678   \regex_group_repeat_aux:n {#2}
1679   \regex_group_submatches:nNN {#1}
1680     \l_regex_internal_a_int \l_regex_right_state_int
1681   \if_meaning:w \c_true_bool #3
1682     \regex_build_transition_right:nNn \regex_action_free_group:n
1683       \l_regex_right_state_int \l_regex_internal_a_int
1684   \else:
1685     \regex_build_transition_left:NNN \regex_action_free_group:n
1686       \l_regex_right_state_int \l_regex_internal_a_int
1687   \fi:
1688   \regex_build_new_state:
1689   \fi:
1690 }
(End definition for \regex_group_repeat:nnN.)
```

\regex\_group\_repeat:nnN

We wish to repeat the group between #2 and #2 + #3 times, with a laziness controlled by #4. We insert submatch tracking up front: in principle, we could avoid recording submatches for the first #2 copies of the group, but that forces us to treat specially the case #2 = 0. Repeat that group with submatch tracking #2 + #3 times (the maximum number of repetitions). Then our goal is to add #3 transitions from the end of the #2-th group, and each subsequent groups, to the end. For a lazy quantifier, we add those transitions to the left states, before submatch tracking. For the greedy case, we add the transitions to the right states, after submatch tracking and the transitions which go on with more repetitions. In the greedy case with #2 = 0, the transition which skips over all copies of the group must be added separately, because its starting state does not follow the normal pattern: we had to add it “by hand” earlier.

```

1691 \cs_new_protected:Npn \regex_group_repeat:nnN #1#2#3#4
1692 {
1693   \regex_group_submatches:nNN {#1}
1694     \l_regex_left_state_int \l_regex_right_state_int
1695   \regex_group_repeat_aux:n { #2 + #3 }
1696   \if_meaning:w \c_true_bool #4
1697     \int_set_eq:NN \l_regex_left_state_int \l_regex_max_state_int
```

```

1698 \prg_replicate:nn { #3 }
1699 {
1700     \int_sub:Nn \l_regex_left_state_int
1701     { \l_regex_internal_b_int - \l_regex_internal_a_int }
1702     \regex_build_transition_left:NNN \regex_action_free:n
1703     \l_regex_left_state_int \l_regex_max_state_int
1704 }
1705 \else:
1706     \prg_replicate:nn { #3 - \c_one }
1707     {
1708         \int_sub:Nn \l_regex_right_state_int
1709         { \l_regex_internal_b_int - \l_regex_internal_a_int }
1710         \regex_build_transition_right:nNn \regex_action_free:n
1711         \l_regex_right_state_int \l_regex_max_state_int
1712     }
1713 \if_num:w #2 = \c_zero
1714     \int_set:Nn \l_regex_right_state_int
1715     { \l_regex_left_state_int - \c_one }
1716 \else:
1717     \int_sub:Nn \l_regex_right_state_int
1718     { \l_regex_internal_b_int - \l_regex_internal_a_int }
1719 \fi:
1720     \regex_build_transition_right:nNn \regex_action_free:n
1721     \l_regex_right_state_int \l_regex_max_state_int
1722 \fi:
1723     \regex_build_new_state:
1724 }
(End definition for \regex_group_repeat:nnnN.)

```

#### 2.4.6 Others

\regex\_assertion:Nn Usage: \regex\_assertion:Nn *<boolean>* {*<test>*}, where the *<test>* is either of the two other functions. Add a free transition to a new state, conditionally to the assertion test.  
\regex\_b\_test: The \regex\_b\_test: test is used by the \b and \B escape: check if the last character was a word character or not, and do the same to the current character. The boundary-markers of the string are non-word characters for this purpose. Anchors at the start or end of match use \regex\_anchor:N, with a position controlled by the integer #1.

```

1725 \cs_new_protected:Npn \regex_assertion:Nn #1#2
1726 {
1727     \regex_build_new_state:
1728     \regex_toks_put_right:Nx \l_regex_left_state_int
1729     {
1730         \exp_not:n {#2}
1731         \regex_break_point:TF
1732         \bool_if:NF #1 { { } }
1733         {
1734             \regex_action_free:n
1735             {
1736                 \int_eval:n

```

```

1737             { \l_regex_right_state_int - \l_regex_left_state_int }
1738         }
1739     }
1740   \bool_if:NT #1 { { } }
1741 }
1742 }
1743 \cs_new_protected:Npn \regex_anchor:N #1
1744 {
1745   \if_num:w #1 = \l_regex_current_pos_int
1746     \exp_after:wN \regex_break_true:w
1747   \fi:
1748 }
1749 \cs_new_protected_nopar:Npn \regex_b_test:
1750 {
1751   \group_begin:
1752     \int_set_eq:NN \l_regex_current_char_int \l_regex_last_char_int
1753     \regex_prop_w:
1754     \regex_break_point:TF
1755     { \group_end: \regex_item_reverse:n \regex_prop_w: }
1756     { \group_end: \regex_prop_w: }
1757 }
(End definition for \regex_assertion:Nn, \regex_b_test:, and \regex_anchor:N. These functions are documented on page ??.)
```

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

```

1758 \cs_new_protected_nopar:Npn \regex_command_K:
1759 {
1760   \regex_build_new_state:
1761   \regex_toks_put_right:Nx \l_regex_left_state_int
1762   {
1763     \regex_action_submatch:n { 0< }
1764     \bool_set_true:N \l_regex_fresh_thread_bool
1765     \regex_action_free:n
1766     { \int_eval:n { \l_regex_right_state_int - \l_regex_left_state_int } }
1767     \bool_set_false:N \l_regex_fresh_thread_bool
1768   }
1769 }
```

(End definition for \regex\_command\_K:. This function is documented on page ??.)

## 2.5 Matching

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

query. At every step (for each token in the query), we unpack that list of active states and the corresponding submatch props, and empty those.

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

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

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

### 2.5.1 Variables used when matching

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

```
1770 \int_new:N \l_regex_min_pos_int
1771 \int_new:N \l_regex_max_pos_int
1772 \int_new:N \l_regex_current_pos_int
1773 \int_new:N \l_regex_start_pos_int
1774 \int_new:N \l_regex_success_pos_int
```

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

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

```
1775 \int_new:N \l_regex_current_char_int
1776 \int_new:N \l_regex_current_catcode_int
1777 \int_new:N \l_regex_last_char_int
1778 \int_new:N \l_regex_case_changed_char_int
```

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

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

1779 `\int_new:N \l_regex_current_state_int`

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

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

1780 `\prop_new:N \l_regex_current_submatches_prop`

1781 `\prop_new:N \l_regex_success_submatches_prop`

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

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

1782 `\int_new:N \l_regex_step_int`

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

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

1783 `\int_new:N \l_regex_min_active_int`

1784 `\int_new:N \l_regex_max_active_int`

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

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

1785 `\tl_new:N \l_regex_every_match_tl`

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

\l\_regex\_fresh\_thread\_bool  
\l\_regex\_empty\_success\_bool  
\regex\_if\_two\_empty\_matches:F

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

```
1786 \bool_new:N \l_regex_fresh_thread_bool
1787 \bool_new:N \l_regex_empty_success_bool
1788 \cs_new_eq:NN \regex_if_two_empty_matches:F \use:n
(End definition for \l_regex_fresh_thread_bool and \l_regex_empty_success_bool. These functions
are documented on page ??.)
```

\g\_regex\_success\_bool  
\l\_regex\_saved\_success\_bool  
\l\_regex\_match\_success\_bool

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

```
1789 \bool_new:N \g_regex_success_bool
1790 \bool_new:N \l_regex_saved_success_bool
1791 \bool_new:N \l_regex_match_success_bool
(End definition for \g_regex_success_bool, \l_regex_saved_success_bool, and \l_regex_match_success_bool.
These variables are documented on page ??.)
```

## 2.5.2 Matching: framework

\regex\_match:n

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

```
1792 \cs_new_protected:Npn \regex_match:n #1
1793 {
1794 \trace \trace_push:nnx { regex } { 1 } { regex_match }
1795 \trace \trace:nnx { regex } { 1 } { analyzing~query~token~list }
1796 \int_zero:N \l_regex_balance_int
1797 \int_set:Nn \l_regex_current_pos_int { \c_two * \l_regex_max_state_int }
1798 \regex_query_set:nnn { } { -1 } { -2 }
1799 \int_set_eq:NN \l_regex_min_pos_int \l_regex_current_pos_int
1800 \tl_analysis_map_inline:nn {#1}
1801 { \regex_query_set:nnn {##1} {##2} {##3} }
```

```

1802   \int_set_eq:NN \l_regex_max_pos_int \l_regex_current_pos_int
1803   \regex_query_set:nnn { } { -1 } { -2 }
1804 <trace>   \trace:nnx { regex } { 1 } { initializing }
1805   \bool_gset_false:N \g_regex_success_bool
1806   \prg_stepwise_inline:nnnn
1807     \l_regex_min_state_int \c_one { \l_regex_max_state_int - \c_one }
1808     { \tex_dimen:D ##1 \c_one sp \scan_stop: }
1809   \int_set_eq:NN \l_regex_min_active_int \l_regex_max_state_int
1810   \int_set_eq:NN \l_regex_step_int \c_zero
1811   \int_set_eq:NN \l_regex_success_pos_int \l_regex_min_pos_int
1812   \int_set:Nn \l_regex_submatch_int
1813     { \c_two * \l_regex_max_state_int }
1814   \bool_set_false:N \l_regex_empty_success_bool
1815   \regex_match_once:
1816 <trace>   \trace_pop:nnx { regex } { 1 } { regex_match }
1817 }

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

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

```

1818 \cs_new_protected_nopar:Npn \regex_match_once:
1819 {
1820   \if_meaning:w \c_true_bool \l_regex_empty_success_bool
1821   \cs_set_nopar:Npn \regex_if_two_empty_matches:F
1822     { \int_compare:nNnF \l_regex_start_pos_int = \l_regex_current_pos_int }
1823 \else:
1824   \cs_set_eq:NN \regex_if_two_empty_matches:F \use:n
1825 \fi:
1826   \int_set_eq:NN \l_regex_start_pos_int \l_regex_success_pos_int
1827   \bool_set_false:N \l_regex_match_success_bool
1828   \prop_clear:N \l_regex_current_submatches_prop
1829   \int_set_eq:NN \l_regex_max_active_int \l_regex_min_active_int
1830   \regex_store_state:n { \l_regex_min_state_int }
1831   \int_set:Nn \l_regex_current_pos_int
1832     { \l_regex_start_pos_int - \c_one }
1833   \regex_query_get:
1834   \regex_match_loop:
1835   \l_regex_every_match_tl
1836 }
```

(End definition for \regex\_match\_once:. This function is documented on page ??.)

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

```

1837 \cs_new_protected_nopar:Npn \regex_single_match:
1838   {
1839     \tl_set:Nn \l_regex_every_match_tl
1840     { \bool_gset_eq:NN \g_regex_success_bool \l_regex_match_success_bool }
1841   }
1842 \cs_new_protected:Npn \regex_multi_match:n #1
1843   {
1844     \tl_set:Nn \l_regex_every_match_tl
1845     {
1846       \if_meaning:w \c_true_bool \l_regex_match_success_bool
1847         \bool_gset_true:N \g_regex_success_bool
1848         #1
1849         \exp_after:wN \regex_match_once:
1850       \fi:
1851     }
1852   }
(End definition for \regex_single_match: and \regex_multi_match:n. These functions are documented on page ??.)
```

\regex\_match\_loop: At each new position, set some variables and get the new character and category from the query. Then unpack the array of active threads, and clear it by resetting its length (`max_active`). This results in a sequence of `\regex_use_state_and_submatches:nn` `{<state>}` `{<prop>}`, and we consider those states one by one in order. As soon as a thread succeeds, exit the step, and, if there are threads to consider at the next position, and we have not reached the end of the string, repeat the loop. Otherwise, the last thread that succeeded is what `\regex_match_once:` matches. We explain the `fresh_thread` business when describing `\regex_action_wildcard::`

```

1853 \cs_new_protected_nopar:Npn \regex_match_loop:
1854   {
1855     \tex_advance:D \l_regex_step_int \c_two
1856     \int_incr:N \l_regex_current_pos_int
1857     \int_set_eq:NN \l_regex_last_char_int \l_regex_current_char_int
1858     \int_set_eq:NN \l_regex_case_changed_char_int \c_max_int
1859     \regex_query_get:
1860     \use:x
1861     {
1862       \int_set_eq:NN \l_regex_max_active_int \l_regex_min_active_int
1863       \exp_after:wN \regex_match_one_active:w
1864       \int_use:N \l_regex_min_active_int ;
1865     }
1866     \prg_break_point:n { \bool_set_false:N \l_regex_fresh_thread_bool }
1867     \if_num:w \l_regex_max_active_int > \l_regex_min_active_int
1868       \if_num:w \l_regex_current_pos_int < \l_regex_max_pos_int
1869         \exp_after:wN \exp_after:wN \exp_after:wN \regex_match_loop:
1870       \fi:
1871     \fi:
```

```

1872    }
1873 \cs_new:Npn \regex_match_one_active:w #1;
1874 {
1875     \if_num:w #1 < \l_regex_max_active_int
1876         \regex_use_state_and_submatches:nn
1877             { \int_value:w \tex_skip:D #1 }
1878             { \tex_the:D \tex_toks:D #1 }
1879         \exp_after:wN \regex_match_one_active:w
1880             \int_use:N \int_eval:w #1 + \c_one \exp_after:wN ;
1881     \fi:
1882 }
(End definition for \regex_match_loop:. This function is documented on page ??.)
```

\regex\_query\_set:nnn

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

```

1883 \cs_new_protected:Npn \regex_query_set:nnn #1#2#3
1884 {
1885     \tex_muskip:D \l_regex_current_pos_int
1886         = \etex_gluetomu:D
1887         #3 sp
1888         plus #2 sp
1889         minus \l_regex_balance_int sp
1890         \scan_stop:
1891     \tex_toks:D \l_regex_current_pos_int {#1}
1892     \int_incr:N \l_regex_current_pos_int
1893     \if_case:w #2 \exp_stop_f:
1894     \or: \int_incr:N \l_regex_balance_int
1895     \or: \int_decr:N \l_regex_balance_int
1896     \fi:
1897 }
(End definition for \regex_query_set:nnn.)
```

\regex\_query\_get:

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

```

1898 \cs_new_protected_nopar:Npn \regex_query_get:
1899 {
1900     \l_regex_current_char_int
1901         = \etex_mutoglu:D \tex_muskip:D \l_regex_current_pos_int
1902     \l_regex_current_catcode_int = \etex_gluestretch:D
1903         \etex_mutoglu:D \tex_muskip:D \l_regex_current_pos_int
1904 }
(End definition for \regex_query_get:.)
```

### 2.5.3 Using states of the nfa

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

```

1905 \cs_new_protected_nopar:Npn \regex_use_state:
1906  {
1907  (*trace)
1908    \trace:n { regex } { 2 } { state-\int_use:N \l_regex_current_state_int }
1909  
```

```

1910    \tex_dimen:D \l_regex_current_state_int
1911      = \l_regex_step_int sp \scan_stop:
1912    \tex_the:D \tex_toks:D \l_regex_current_state_int
1913    \tex_dimen:D \l_regex_current_state_int
1914      = \int_eval:w \l_regex_step_int + \c_one \int_eval_end: sp \scan_stop:
1915  }

```

(End definition for \regex\_use\_state:. This function is documented on page ??.)

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

```

1916 \cs_new_protected:Npn \regex_use_state_and_submatches:nn #1 #2
1917  {
1918    \int_set:Nn \l_regex_current_state_int {#1}
1919    \if_num:w \tex_dimen:D \l_regex_current_state_int < \l_regex_step_int
1920      \tl_set:Nn \l_regex_current_submatches_prop {#2}
1921      \exp_after:wN \regex_use_state:
1922    \fi:
1923    \scan_stop:
1924  }

```

(End definition for \regex\_use\_state\_and\_submatches:nn. This function is documented on page ??.)

### 2.5.4 Actions when matching

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

```

1925 \cs_new_protected_nopar:Npn \regex_action_start_wildcard:
1926  {
1927    \bool_set_true:N \l_regex_fresh_thread_bool
1928    \regex_action_free:n {1}
1929    \bool_set_false:N \l_regex_fresh_thread_bool
1930    \regex_action_cost:n {0}
1931  }

```

(End definition for \regex\_action\_start\_wildcard:. This function is documented on page ??.)

```
\regex_action_free:n
\regex_action_free_group:n
\regex_action_free_aux:nn
```

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

```
1932 \cs_new_protected_nopar:Npn \regex_action_free:n
1933   { \regex_action_free_aux:nn { > \l_regex_step_int \else: } }
1934 \cs_new_protected_nopar:Npn \regex_action_free_group:n
1935   { \regex_action_free_aux:nn { < \l_regex_step_int } }
1936 \cs_new_protected:Npn \regex_action_free_aux:nn #1#2
1937   {
1938     \use:x
1939     {
1940       \int_add:Nn \l_regex_current_state_int {#2}
1941       \exp_not:n
1942       {
1943         \if_num:w \tex_dimen:D \l_regex_current_state_int #1
1944           \exp_after:wN \regex_use_state:
1945         \fi:
1946       }
1947       \int_set:Nn \l_regex_current_state_int
1948         { \int_use:N \l_regex_current_state_int }
1949       \tl_set:Nn \exp_not:N \l_regex_current_submatches_prop
1950         { \exp_not:o \l_regex_current_submatches_prop }
1951     }
1952 }
```

(End definition for `\regex_action_free:n` and `\regex_action_free_group:n`. These functions are documented on page ??.)

```
\regex_action_cost:n
```

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

```
1953 \cs_new_protected:Npn \regex_action_cost:n #1
1954   {
1955     \exp_args:No \regex_store_state:n
1956       { \int_use:N \int_eval:w \l_regex_current_state_int + #1 }
1957   }
```

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

```
\regex_store_state:n
\regex_store_submatches:
```

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

```
1958 \cs_new_protected:Npn \regex_store_state:n #1
1959   {
1960     \regex_store_submatches:
```

```

1961   \tex_skip:D \l_regex_max_active_int = #1 sp \scan_stop:
1962   \int_incr:N \l_regex_max_active_int
1963 }
1964 \cs_new_protected_nopar:Npn \regex_store_submatches:
1965 {
1966   \tex_toks:D \l_regex_max_active_int \exp_after:wN
1967   { \l_regex_current_submatches_prop }
1968 }
(End definition for \regex_store_state:n. This function is documented on page ??.)
```

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

```

1969 \cs_new_protected_nopar:Npn \regex_disable_submatches:
1970 {
1971   \cs_set_protected_nopar:Npn \regex_store_submatches: { }
1972   \cs_set_protected:Npn \regex_action_submatch:n ##1 { }
1973 }
```

(End definition for \regex\_disable\_submatches:. This function is documented on page ??.)

**\regex\_action\_submatch:n** Update the current submatches with the information from the current position. Maybe a bottleneck.

```

1974 \cs_new_protected:Npn \regex_action_submatch:n #1
1975 {
1976   \prop_put:Nno \l_regex_current_submatches_prop {#1}
1977   { \int_use:N \l_regex_current_pos_int }
1978 }
```

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

**\regex\_action\_success:** There is a successful match when an execution path reaches the last state in the NFA, unless this marks a second identical empty match. Then mark that there was a successful match; it is empty if it is “fresh”; and we store the current position and submatches. The current step is then interrupted with \prg\_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.

```

1979 \cs_new_protected_nopar:Npn \regex_action_success:
1980 {
1981   \regex_if_two_empty_matches:F
1982   {
1983     \bool_set_true:N \l_regex_match_success_bool
1984     \bool_set_eq:NN \l_regex_empty_success_bool
1985     \l_regex_fresh_thread_bool
1986     \int_set_eq:NN \l_regex_success_pos_int \l_regex_current_pos_int
1987     \prop_set_eq:NN \l_regex_success_submatches_prop
1988     \l_regex_current_submatches_prop
1989     \prg_map_break:
1990   }
1991 }
```

(End definition for \regex\_action\_success:. This function is documented on page ??.)

## 2.6 Replacement

### 2.6.1 Variables and helpers used in replacement

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

```
1992 \int_new:N \l_regex_replacement_csnames_int  
(End definition for \l_regex_replacement_csnames_int. This variable is documented on page ??.)
```

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

```
1993 \tl_new:N \l_regex_balance_tl  
(End definition for \l_regex_balance_tl. This variable is documented on page ??.)
```

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

```
1994 \cs_new:Npn \regex_replacement_balance_one_match:n #1  
1995 { - \regex_submatch_balance:n {#1} }  
(End definition for \regex_replacement_balance_one_match:n.)
```

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

```
1996 \cs_new:Npn \regex_replacement_do_one_match:n #1  
1997 {  
1998   \regex_query_range:nn  
1999   { \etex_glueshrink:D \tex_skip:D #1 }  
2000   { \tex_skip:D #1 }  
2001 }  
(End definition for \regex_replacement_do_one_match:n.)
```

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

```
2002 \cs_new:Npn \regex_replacement_exp_not:N #1 { \exp_not:n {#1} }  
(End definition for \regex_replacement_exp_not:N.)
```

### 2.6.2 Query and brace balance

```
\regex_query_range:nn
\regex_query_range_loop:ww
```

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

```
2003 \cs_new:Npn \regex_query_range:nn #1#2
2004 {
2005   \exp_after:wN \regex_query_range_loop:ww
2006   \int_use:N \int_eval:w #1 \exp_after:wN ;
2007   \int_use:N \int_eval:w #2 ;
2008   \prg_break_point:n { }
2009 }
2010 \cs_new:Npn \regex_query_range_loop:ww #1 ; #2 ;
2011 {
2012   \if_num:w #1 < #2 \exp_stop_f:
2013   \else:
2014     \exp_after:wN \prg_map_break:
2015   \fi:
2016   \tex_the:D \tex_toks:D #1 \exp_stop_f:
2017   \exp_after:wN \regex_query_range_loop:ww
2018     \int_use:N \int_eval:w #1 + \c_one ; #2 ;
2019 }
```

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

```
\regex_query_submatch:n
```

When this function is called, `\skipi` holds the start and end positions for the *i*-th overall submatch as its main and stretch components. In the case of repeated matches, submatches from all the matches are put one after the other in blocks of `\l_regex_capturing_group_int` `\skip` registers.

```
2020 \cs_new:Npn \regex_query_submatch:n #1
2021 {
2022   \regex_query_range:nn
2023   { \tex_skip:D \int_eval:w #1 }
2024   { \etex_gluestretch:D \tex_skip:D \int_eval:w #1 }
2025 }
```

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

```
\regex_submatch_balance:n
```

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

```
2026 \cs_new_protected:Npn \regex_submatch_balance:n #1
```

```

2027   {
2028     \etex_glueshrink:D \etex_mutoglu:D \etex_muexpr:D
2029       \tex_muskip:D \etex_gluestretch:D \tex_skip:D #1
2030       - \tex_muskip:D \tex_skip:D #1
2031     \scan_stop:
2032   }

```

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

### 2.6.3 Framework

`\regex_replacement:n` The replacement text is built incrementally by abusing `\toks` within a group (see `\l3tl-build`). We keep track in `\l1_regex_balance_int` of the balance of explicit begin- and end-group tokens and `\l1_regex_balance_tl` will consist of some code to compute the brace balance from submatches (see its description). Detect unescaped right braces, and escaped characters, with trailing `\prg_do_nothing`: because some of the later function look-ahead. Once the whole replacement text has been parsed, make sure that there is no open csname. Finally, define the `balance_one_match` and `do_one_match` functions.

```

2033 \cs_new_protected:Npn \regex_replacement:n #1
2034   {
2035     \trace_push:nnn { regex } { 1 } { regex_replacement:n }
2036     \tl_set_build:Nw \l1_regex_internal_a_tl
2037       \int_zero:N \l1_regex_balance_int
2038       \tl_clear:N \l1_regex_balance_tl
2039       \regex_escape_use:nnnn
2040         {
2041           \if_charcode:w \c_rbrace_str ##1
2042             \regex_replacement_rbrace:N \else: \tl_build_one:n \fi: ##1
2043           }
2044           { \regex_replacement_escaped:N ##1 }
2045           { \tl_build_one:n ##1 }
2046           {##1}
2047         \prg_do_nothing: \prg_do_nothing:
2048         \if_int_compare:w \l1_regex_replacement_csnames_int > \c_zero
2049           \msg_kernel_error:nnx { regex } { replacement-missing-rbrace }
2050             { \int_use:N \l1_regex_replacement_csnames_int }
2051             \tl_build_one:x
2052               { \prg_replicate:nn \l1_regex_replacement_csnames_int \cs_end: }
2053             \fi:
2054             \cs_gset:Npx \regex_replacement_balance_one_match:n ##1
2055               {
2056                 + \int_use:N \l1_regex_balance_int
2057                 \l1_regex_balance_tl
2058                 - \regex_submatch_balance:n {##1}
2059               }
2060             \tl_build_end:
2061             \exp_args:No \regex_replacement_aux:n \l1_regex_internal_a_tl
2062     \trace_pop:nnn { regex } { 1 } { regex_replacement:n }
2063   }
2064 \cs_new_protected:Npn \regex_replacement_aux:n #1

```

```

2065   {
2066     \cs_set:Npn \regex_replacement_do_one_match:n ##1
2067     {
2068       \regex_query_range:nn
2069         { \etex_glueshrink:D \tex_skip:D ##1 }
2070         { \tex_skip:D ##1 }
2071       #1
2072     }
2073   }

```

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

### `\regex_replacement_escaped:N`

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

```

2074 \cs_new_protected:Npn \regex_replacement_escaped:N #1
2075   {
2076     \cs_if_exist_use:cF { regex_replacement_#1:w }
2077     {
2078       \if_num:w \c_one < 1#1 \exp_stop_f:
2079         \regex_replacement_put_submatch:n {#1}
2080       \else:
2081         \tl_build_one:n #1
2082       \fi:
2083     }
2084   }

```

(End definition for `\regex_replacement_escaped:N`.)

#### 2.6.4 Submatches

##### `\regex_replacement_put_submatch:n`

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

```

2085 \cs_new_protected:Npn \regex_replacement_put_submatch:n #1
2086   {
2087     \if_num:w #1 < \l_regex_capturing_group_int
2088       \tl_build_one:n { \regex_query_submatch:n { #1 + ##1 } }
2089       \if_num:w \l_regex_replacement_csnames_int = \c_zero
2090         \tl_put_right:Nn \l_regex_balance_tl
2091           { + \regex_submatch_balance:n { \int_eval:w #1+##1 \int_eval_end: } }
2092       \fi:
2093     \fi:
2094   }

```

(End definition for `\regex_replacement_put_submatch:n`.)

```
\regex_replacement_g:w
\regex_replacement_g_digits>NN
```

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

```
2095 \cs_new_protected:Npn \regex_replacement_g:w #1#2
2096  {
2097    \str_if_eq:xxTF { #1#2 } { \tl_build_one:n \c_lbrace_str }
2098    {
2099      \int_zero:N \l_regex_internal_a_int
2100      \regex_replacement_g_digits>NN
2101    }
2102    { \regex_replacement_error>NNN g #1 #2 }
2103  }
2104 \cs_new_protected:Npn \regex_replacement_g_digits>NN #1#2
2105  {
2106    \token_if_eq_meaning:NNTF #1 \tl_build_one:n
2107    {
2108      \if_num:w \c_one < 1#2 \exp_stop_f:
2109      \int_set:Nn \l_regex_internal_a_int
2110      { \c_ten * \l_regex_internal_a_int + #2 }
2111      \exp_after:wN \use_i:nnn
2112      \exp_after:wN \regex_replacement_g_digits>NN
2113    \else:
2114      \exp_after:wN \regex_replacement_error>NNN
2115      \exp_after:wN g
2116    \fi:
2117  }
2118  {
2119    \if_meaning:w \regex_replacement_rbrace:N #1
2120    \exp_args:No \regex_replacement_put_submatch:n
2121    { \int_use:N \l_regex_internal_a_int }
2122    \exp_after:wN \use_none:nn
2123  \else:
2124    \exp_after:wN \regex_replacement_error>NNN
2125    \exp_after:wN g
2126  \fi:
2127 }
2128 #1 #2
2129 }
```

(End definition for \regex\_replacement\_g:w and \regex\_replacement\_g\_digits>NN.)

### 2.6.5 Csnames in replacement

```
\regex_replacement_c:w
\regex_replacement_c_{{:w}}
```

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

```
2130 \cs_new_protected:Npn \regex_replacement_c:w #1#2
2131  {
2132    \token_if_eq_meaning:NNTF #1 \tl_build_one:n
```

```

2133 {
2134   \cs_if_exist_use:cF { regex_replacement_c_#2:w }
2135   { \regex_error:NNN c #1#2 }
2136 }
2137 { \regex_error:NNN c #1#2 }
2138 }
2139 \cs_new_protected_nopar:cpn { regex_replacement_c_ \c_lbrace_str :w }
2140 {
2141   \if_case:w \l_regex_replacement_csnames_int
2142   \tl_build_one:n
2143   { \exp_not:n { \exp_after:wN \regex_error:N \cs:w } }
2144 \else:
2145   \tl_build_one:n { \exp_not:n { \exp_after:wN \tl_to_str:N \cs:w } }
2146 \fi:
2147 \int_incr:N \l_regex_replacement_csnames_int
2148 }

```

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

`\regex_replacement_u:w` Check that `\u` is followed by a left brace. If so, start a control sequence with `\cs:w`, which is then unpacked either with `\exp_not:V` or `\tl_to_str:V` depending on the current context.

```

2149 \cs_new_protected:Npn \regex_replacement_u:w #1#2
2150 {
2151   \str_if_eq:xxTF { #1#2 } { \tl_build_one:n \c_lbrace_str }
2152   {
2153     \if_case:w \l_regex_replacement_csnames_int
2154     \tl_build_one:n { \exp_not:n { \exp_after:wN \exp_not:V \cs:w } }
2155     \else:
2156     \tl_build_one:n { \exp_not:n { \exp_after:wN \tl_to_str:V \cs:w } }
2157     \fi:
2158     \int_incr:N \l_regex_replacement_csnames_int
2159   }
2160   { \regex_error:NNN u #1#2 }
2161 }

```

(End definition for `\regex_replacement_u:w`.)

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

```

2162 \cs_new_protected:Npn \regex_replacement_rbrace:N #1
2163 {
2164   \if_int_compare:w \l_regex_replacement_csnames_int > \c_zero
2165   \tl_build_one:n \cs_end:
2166   \int_decr:N \l_regex_replacement_csnames_int
2167   \else:
2168   \tl_build_one:n #1
2169   \fi:
2170 }

```

(End definition for `\regex_replacement_rbrace:N`.)

### 2.6.6 Characters in replacement

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

```
2171 \group_begin:
```

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

```
2172 \cs_new_protected:Npn \regex_replacement_char:nNN #1#2#3
2173 {
2174     \if_meaning:w \prg_do_nothing: #3
2175     \msg_kernel_error:nn { regex } { replacement-catcode-end }
2176     \else:
2177         \tex_lccode:D \c_zero = '#3 \scan_stop:
2178         \tl_to_lowercase:n { \tl_build_one:n {#1} }
2179         \fi:
2180     }
(End definition for \regex_replacement_char:nNN.)
```

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

```
2181 \char_set_catcode_active:N \^^@
2182 \cs_new_protected_nopar:Npn \regex_replacement_c_A:w
2183     { \regex_replacement_char:nNN { \exp_not:n { \exp_not:N \^^@ } } }
(End definition for \regex_replacement_c_A:w.)
```

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

```
2184 \char_set_catcode_group_begin:N \^^@
2185 \cs_new_protected_nopar:Npn \regex_replacement_c_B:w
2186     {
2187         \if_num:w \l_regex_replacement_csnames_int = \c_zero
2188             \int_incr:N \l_regex_balance_int
2189         \fi:
2190         \regex_replacement_char:nNN
2191             { \exp_not:n { \exp_after:wN \^^@ \if_false: } \fi: } }
2192     }
(End definition for \regex_replacement_c_B:w.)
```

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

```
2193 \cs_new_protected:Npn \regex_replacement_c_C:w #1#2
2194   { \tl_build_one:n { \exp_not:N \exp_not:N \exp_not:c {#2} } }
(End definition for \regex_replacement_c_C:w.)
```

\regex\_replacement\_c\_D:w Subscripts fit the mould: \lowercase the null byte with the correct category.

```
2195 \char_set_catcode_math_subscript:N \^^@
2196 \cs_new_protected_nopar:Npn \regex_replacement_c_D:w
2197   { \regex_replacement_char:nNN { \^^@ } }
(End definition for \regex_replacement_c_D:w.)
```

\regex\_replacement\_c\_E:w Similar to the begin-group case, the second x-expansion produces the bare end-group token.

```
2198 \char_set_catcode_group_end:N \^^@
2199 \cs_new_protected_nopar:Npn \regex_replacement_c_E:w
2200   {
2201     \if_num:w \l_regex_replacement_csnames_int = \c_zero
2202       \int_decr:N \l_regex_balance_int
2203     \fi:
2204     \regex_replacement_char:nNN
2205       { \exp_not:n { \if_false: { \fi: \^^@ } } }
2206   }
(End definition for \regex_replacement_c_E:w.)
```

\regex\_replacement\_c\_L:w Simply \lowercase a letter null byte to produce an arbitrary letter.

```
2207 \char_set_catcode_letter:N \^^@
2208 \cs_new_protected_nopar:Npn \regex_replacement_c_L:w
2209   { \regex_replacement_char:nNN { \^^@ } }
(End definition for \regex_replacement_c_L:w.)
```

\regex\_replacement\_c\_M:w No surprise here, we lowercase the null math toggle.

```
2210 \char_set_catcode_math_toggle:N \^^@
2211 \cs_new_protected_nopar:Npn \regex_replacement_c_M:w
2212   { \regex_replacement_char:nNN { \^^@ } }
(End definition for \regex_replacement_c_M:w.)
```

\regex\_replacement\_c\_0:w Lowercase an other null byte.

```
2213 \char_set_catcode_other:N \^^@
2214 \cs_new_protected_nopar:Npn \regex_replacement_c_0:w
2215   { \regex_replacement_char:nNN { \^^@ } }
(End definition for \regex_replacement_c_0:w.)
```

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

```

2216  \char_set_catcode_parameter:N \^@\^@^@
2217  \cs_new_protected_nopar:Npn \regex_replacement_c_P:w
2218  {
2219    \regex_replacement_char:nNN
2220    { \exp_not:n { \exp_not:n { ^@\^@\^@\^@ } } }
2221  }
(End definition for \regex_replacement_c_P:w.)
```

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

```

2222  \cs_new_protected:Npn \regex_replacement_c_S:w #1#2
2223  {
2224    \if_meaning:w \prg_do_nothing: #2
2225      \msg_kernel_error:nn { regex } { replacement-catcode-end }
2226    \else:
2227      \if_num:w '#2 = \c_zero
2228        \msg_kernel_error:nn { regex } { replacement-null-space }
2229      \fi:
2230      \tex_lccode:D 32 = '#2 \scan_stop:
2231      \tl_to_lowercase:n { \tl_build_one:n {~} }
2232    \fi:
2233  }
(End definition for \regex_replacement_c_S:w.)
```

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

```

2234  \char_set_catcode_alignment:N \^@\^@^@
2235  \cs_new_protected_nopar:Npn \regex_replacement_c_T:w
2236  { \regex_replacement_char:nNN { ^@\^@ } }
(End definition for \regex_replacement_c_T:w.)
```

\regex\_replacement\_c\_U:w Simple call to \regex\_replacement\_char:nNN which lowers the math superscript ^@\^@.

```

2237  \char_set_catcode_math_superscript:N \^@\^@^@
2238  \cs_new_protected_nopar:Npn \regex_replacement_c_U:w
2239  { \regex_replacement_char:nNN { ^@\^@ } }
(End definition for \regex_replacement_c_U:w.)
```

Restore the catcode of the null byte.

```
2240  \group_end:
```

### 2.6.7 An error

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

```
2241 \cs_new_protected:Npn \regex_replacement_error:NNN #1#2#3
2242 {
2243     \msg_kernel_error:nnx { regex } { replacement-#1 } {#3}
2244     #2 #3
2245 }
(End definition for \regex_replacement_error:NNN.)
```

## 2.7 User functions

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

```
2246 \cs_new_protected:Npn \regex_new:N #1
2247 { \cs_new_eq:NN #1 \c_regex_no_match_regex }
(End definition for \regex_new:N. This function is documented on page 6.)
```

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

```
2248 \cs_new_protected_nopar:Npn \regex_set:Nn #1#2
2249 {
2250     \regex_compile:n {#2}
2251     \tl_set_eq:NN #1 \l_regex_internal_regex
2252 }
2253 \cs_new_protected_nopar:Npn \regex_gset:Nn #1#2
2254 {
2255     \regex_compile:n {#2}
2256     \tl_gset_eq:NN #1 \l_regex_internal_regex
2257 }
2258 \cs_new_protected_nopar:Npn \regex_const:Nn #1#2
2259 {
2260     \regex_compile:n {#2}
2261     \tl_const:Nx #1 { \exp_not:o \l_regex_internal_regex }
2262 }
(End definition for \regex_set:Nn, \regex_gset:Nn, and \regex_const:Nn. These functions are documented on page 7.)
```

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

```
2263 \cs_new_protected:Npn \regex_show:n #1
2264 {
2265     \regex_compile:n {#1}
2266     \regex_show_aux:Nx \l_regex_internal_regex
2267     { { \tl_to_str:n {#1} } }
2268 }
2269 \cs_new_protected:Npn \regex_show:N #1
2270 { \regex_show_aux:Nx #1 { variable~\token_to_str:N #1 } }
```

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

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

```

2271 \prg_new_protected_conditional:Npnn \regex_match:nn #1#2 { T , F , TF }
2272   {
2273     \regex_match_aux:nn { \regex_build:n {#1} } {#2}
2274     \regex_aux_return:
2275   }
2276 \prg_new_protected_conditional:Npnn \regex_match:Nn #1#2 { T , F , TF }
2277   {
2278     \regex_match_aux:nn { \regex_build:N #1 } {#2}
2279     \regex_aux_return:
2280   }

```

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

`\regex_count:nnN` Again, use an auxiliary whose first argument builds the NFA.

`\regex_count:NnN`

```

2281 \cs_new_protected:Npn \regex_count:nnN #1
2282   { \regex_count_aux:nnN { \regex_build:n {#1} } }
2283 \cs_new_protected:Npn \regex_count:NnN #1
2284   { \regex_count_aux:nnN { \regex_build:N #1 } }

```

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

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

```

2285 \cs_set_protected:Npn \regex_tmp:w #1#2#3
2286   {
2287     \cs_new_protected:Npn #2 ##1 { #1 { \regex_build:n {##1} } }
2288     \cs_new_protected:Npn #3 ##1 { #1 { \regex_build:N ##1 } }
2289     \prg_new_protected_conditional:Npnn #2 ##1##2##3 { T , F , TF }
2290       { #1 { \regex_build:n {##1} } {##2} ##3 \regex_aux_return: }
2291     \prg_new_protected_conditional:Npnn #3 ##1##2##3 { T , F , TF }
2292       { #1 { \regex_build:N ##1 } {##2} ##3 \regex_aux_return: }
2293   }
2294 \regex_tmp:w \regex_extract_once_aux:nnN
2295   \regex_extract_once:nnN \regex_extract_once:NnN
2296 \regex_tmp:w \regex_extract_all_aux:nnN
2297   \regex_extract_all:nnN \regex_extract_all:NnN
2298 \regex_tmp:w \regex_replace_once_aux:nnN
2299   \regex_replace_once:nnN \regex_replace_once:NnN
2300 \regex_tmp:w \regex_replace_all_aux:nnN
2301   \regex_replace_all:nnN \regex_replace_all:NnN

```

2302 \regex\_tmp:w \regex\_split\_aux:nnN \regex\_split:nnN \regex\_split:NnN  
*(End definition for \regex\_extract\_once:nnN and others. These functions are documented on page ??.)*

### 2.7.1 Variables and helpers for user functions

\l\_regex\_match\_count\_int The number of matches found so far is stored in \l\_regex\_match\_count\_int. This is only used in the \regex\_count:nnN functions.

2303 \int\_new:N \l\_regex\_match\_count\_int  
*(End definition for \l\_regex\_match\_count\_int. This variable is documented on page ??.)*

regex\_begin Those flags are raised to indicate extra begin-group or end-group tokens when extracting submatches.

2304 \flag\_new:n { regex\_begin }  
 2305 \flag\_new:n { regex\_end }  
*(End definition for regex\_begin and regex\_end. These variables are documented on page ??.)*

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

2306 \int\_new:N \l\_regex\_submatch\_int  
 2307 \int\_new:N \l\_regex\_zeroth\_submatch\_int  
*(End definition for \l\_regex\_submatch\_int and \l\_regex\_zeroth\_submatch\_int. These variables are documented on page ??.)*

\regex\_aux\_return: This function triggers either \prg\_return\_false: or \prg\_return\_true: as appropriate to whether a match was found or not. It is used by all user conditionals.

2308 \cs\_new\_protected\_nopar:Npn \regex\_aux\_return:  
 2309 {  
 2310 \if\_meaning:w \c\_true\_bool \g\_regex\_success\_bool  
 2311 \prg\_return\_true:  
 2312 \else:  
 2313 \prg\_return\_false:  
 2314 \fi:  
 2315 }

*(End definition for \regex\_aux\_return:.)*

### 2.7.2 Matching

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

```

2316 \cs_new_protected:Npn \regex_match_aux:nn #1#2
2317 {
2318   \group_begin:
2319   \regex_disable_submatches:
2320   \regex_single_match:
2321   #1
2322   \regex_match:n {#2}
2323   \group_end:
2324 }
```

(End definition for \regex\_match\_aux:nn.)

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

```

2325 \cs_new_protected:Npn \regex_count_aux:nnN #1#2#3
2326 {
2327   \group_begin:
2328   \regex_disable_submatches:
2329   \int_zero:N \l_regex_match_count_int
2330   \regex_multi_match:n { \int_incr:N \l_regex_match_count_int }
2331   #1
2332   \regex_match:n {#2}
2333   \exp_args:NNNo
2334   \group_end:
2335   \int_set:Nn #3 { \int_use:N \l_regex_match_count_int }
2336 }
```

(End definition for \regex\_count\_aux:nnN.)

### 2.7.3 Extracting submatches

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

```

2337 \cs_new_protected:Npn \regex_extract_once_aux:nnN #1#2#3
2338 {
2339   \group_begin:
2340   \regex_single_match:
2341   #1
2342   \regex_match:n {#2}
2343   \regex_extract:
2344   \regex_group_end_extract_seq:N #3
2345 }
2346 \cs_new_protected:Npn \regex_extract_all_aux:nnN #1#2#3
2347 {
```

```

2348   \group_begin:
2349     \regex_multi_match:n { \regex_extract: }
2350     #1
2351     \regex_match:n {#2}
2352   \regex_group_end_extract_seq:N #3
2353 }

(End definition for \regex_extract_once_aux:nnN and \regex_extract_all_aux:nnN.)
```

\regex\_split\_aux:nnN

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

```

2354 \cs_new_protected:Npn \regex_split_aux:nnN #1#2#3
2355   {
2356     \group_begin:
2357       \regex_multi_match:n
2358       {
2359         \if_num:w \l_regex_start_pos_int < \l_regex_success_pos_int
2360           \regex_extract:
2361           \tex_skip:D \l_regex_zeroth_submatch_int
2362             = \l_regex_start_pos_int sp
2363               plus \tex_skip:D \l_regex_zeroth_submatch_int \scan_stop:
2364             \fi:
2365       }
2366     #1
2367     \regex_match:n {#2}
2368   \assert\assert_int:n { \l_regex_current_pos_int = \l_regex_max_pos_int }
2369   \tex_skip:D \l_regex_submatch_int
2370     = \l_regex_start_pos_int sp plus \l_regex_max_pos_int sp \scan_stop:
2371   \int_incr:N \l_regex_submatch_int
2372   \if_meaning:w \c_true_bool \l_regex_empty_success_bool
2373     \if_num:w \l_regex_start_pos_int = \l_regex_max_pos_int
2374       \int_decr:N \l_regex_submatch_int
2375     \fi:
2376   \fi:
2377   \regex_group_end_extract_seq:N #3
2378 }
```

(End definition for \regex\_split\_aux:nnN.)

\regex\_group\_end\_extract\_seq:N

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

```

2379 \cs_new_protected:Npn \regex_group_end_extract_seq:N #1
2380 {
2381     \cs_set_eq:NN \seq_item:n \scan_stop:
2382     \flag_clear:n { regex_begin }
2383     \flag_clear:n { regex_end }
2384     \tl_set:Nx \l_regex_internal_a_tl
2385     {
2386         \prg_stepwise_function:nnnN
2387             { \c_two * \l_regex_max_state_int }
2388             \c_one
2389             { \l_regex_submatch_int - \c_one }
2390             \regex_extract_seq_aux:n
2391     }
2392     \int_compare:nNnF
2393         { \flag_height:n { regex_begin } + \flag_height:n { regex_end } }
2394     = \c_zero
2395     {
2396         \msg_kernel_error:nnxxx { regex } { result-unbalanced }
2397             { splitting-or-extracting-submatches }
2398             { \flag_height:n { regex_end } }
2399             { \flag_height:n { regex_begin } }
2400     }
2401     \use:x
2402     {
2403         \group_end:
2404         \tl_set:Nn \exp_not:N #1 { \l_regex_internal_a_tl }
2405     }
2406 }
(End definition for \regex_group_end_extract_seq:N.)
```

\regex\_extract\_seq\_aux:n  
\regex\_extract\_seq\_aux:ww

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

```

2407 \cs_new:Npn \regex_extract_seq_aux:n #1
2408 {
2409     \seq_item:n
2410     {
2411         \exp_after:wN \regex_extract_seq_aux:ww
2412         \int_value:w \regex_submatch_balance:n {#1} ; #1;
2413     }
2414 }
2415 \cs_new:Npn \regex_extract_seq_aux:ww #1; #2;
2416 {
2417     \if_num:w #1 < \c_zero
2418         \flag_raise:n { regex_end }
2419         \prg_replicate:nn {-#1} { \exp_not:n { \if_false: } \fi: }
2420     \fi:
2421     \regex_query_submatch:n {#2}
2422     \if_num:w #1 > \c_zero
```

```

2423     \flag_raise:n { regex_begin }
2424     \prg_replicate:nn {#1} { \exp_not:n { \if_false: { \fi: } } }
2425     \fi:
2426   }
(End definition for \regex_extract_seq_aux:n and \regex_extract_seq_aux:ww.)
```

```
\regex_extract:
\regex_extract_aux_b:wn
\regex_extract_aux_e:wn
```

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

```

2427 \cs_new_protected:Npn \regex_extract:
2428   {
2429     \if_meaning:w \c_true_bool \g_regex_success_bool
2430       \int_set_eq:NN \l_regex_zeroth_submatch_int \l_regex_submatch_int
2431       \prg_replicate:nn \l_regex_capturing_group_int
2432         {
2433           \tex_skip:D \l_regex_submatch_int \c_zero sp \scan_stop:
2434           \int_incr:N \l_regex_submatch_int
2435         }
2436       \prop_map_inline:Nn \l_regex_success_submatches_prop
2437         {
2438           \if_num:w ##1 \c_minus_one
2439             \exp_after:wN \regex_extract_aux_e:wn \int_value:w
2440           \else:
2441             \exp_after:wN \regex_extract_aux_b:wn \int_value:w
2442           \fi:
2443             \int_eval:w \l_regex_zeroth_submatch_int + ##1 {##2}
2444           }
2445           \tex_skip:D \l_regex_zeroth_submatch_int
2446             = \tex_the:D \tex_skip:D \l_regex_zeroth_submatch_int
2447               minus \l_regex_start_pos_int sp \scan_stop:
2448           \fi:
2449         }
2450       \cs_new_protected:Npn \regex_extract_aux_b:wn #1 < #2
2451         {
2452           \tex_skip:D #1 = #2 sp
2453             plus \etex_gluestretch:D \tex_skip:D #1 \scan_stop:
2454         }
2455       \cs_new_protected:Npn \regex_extract_aux_e:wn #1 > #2
2456         {
2457           \tex_skip:D #1
2458             = 1 \tex_skip:D #1 plus #2 sp \scan_stop:
2459         }
(End definition for \regex_extract:, \regex_extract_aux_b:wn, and \regex_extract_aux_e:wn.)
```

### 2.7.4 Replacement

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

```

2460 \cs_new_protected:Npn \regex_replace_once_aux:nnN #1#2#3
2461 {
2462     \group_begin:
2463         \regex_single_match:
2464             #1
2465             \regex_replacement:n {#2}
2466             \exp_args:No \regex_match:n { #3 }
2467             \if_meaning:w \c_false_bool \g_regex_success_bool
2468                 \group_end:
2469             \else:
2470                 \regex_extract:
2471                 \int_set:Nn \l_regex_balance_int
2472                 {
2473                     \regex_replacement_balance_one_match:n
2474                         { \l_regex_zeroth_submatch_int }
2475                 }
2476                 \tl_set:Nx \l_regex_internal_a_tl
2477                 {
2478                     \regex_replacement_do_one_match:n { \l_regex_zeroth_submatch_int }
2479                     \regex_query_range:nn
2480                         { \etex_gluestretch:D \tex_skip:D \l_regex_zeroth_submatch_int }
2481                         { \l_regex_max_pos_int }
2482                 }
2483                 \regex_group_end_replace:N #3
2484             \fi:
2485 }
```

(End definition for \regex\_replace\_once\_aux:nnN.)

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

```

2486 \cs_new_protected:Npn \regex_replace_all_aux:nnN #1#2#3
2487 {
```

```

2488 \group_begin:
2489   \regex_multi_match:n { \regex_extract: }
2490   #1
2491   \regex_replacement:n {#2}
2492   \exp_args:No \regex_match:n {#3}
2493   \int_set:Nn \l_regex_balance_int
2494   {
2495     0
2496     \prg_stepwise_function:nnnN
2497       { \c_two * \l_regex_max_state_int }
2498       \l_regex_capturing_group_int
2499       { \l_regex_submatch_int - \c_one }
2500       \regex_replacement_balance_one_match:n
2501   }
2502   \tl_set:Nx \l_regex_internal_a_tl
2503   {
2504     \prg_stepwise_function:nnnN
2505       { \c_two * \l_regex_max_state_int }
2506       \l_regex_capturing_group_int
2507       { \l_regex_submatch_int - \c_one }
2508       \regex_replacement_do_one_match:n
2509       \regex_query_range:nn
2510       \l_regex_start_pos_int \l_regex_max_pos_int
2511   }
2512   \regex_group_end_replace:N #3
2513 }
(End definition for \regex_replace_all_aux:nnN.)

```

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

```

2514 \cs_new_protected_nopar:Npn \regex_group_end_replace:N #1
2515   {
2516     \if_num:w \l_regex_balance_int = \c_zero
2517     \else:
2518       \msg_kernel_error:nnxxx { regex } { result-unbalanced }
2519         { replacing }
2520         { \int_max:nn { - \l_regex_balance_int } { \c_zero } }
2521         { \int_max:nn { \l_regex_balance_int } { \c_zero } }
2522     \fi:
2523     \use:x
2524     {
2525       \group_end:
2526       \tl_set:Nn \exp_not:N #1
2527       {
2528         \if_int_compare:w \l_regex_balance_int < \c_zero
2529           \prg_replicate:nn { - \l_regex_balance_int }
2530             { { \if_false: } \fi: }
2531       \fi:

```

```

2532         \l_regex_internal_a_t1
2533         \if_int_compare:w \l_regex_balance_int > \c_zero
2534             \prg_replicate:nn { \l_regex_balance_int }
2535                 { \if_false: { \fi: } }
2536             \fi:
2537         }
2538     }
2539 }
(End definition for \regex_group_end_replace:N.)
```

## 2.7.5 Storing and showing compiled patterns

## 2.8 Messages

Messages for the preparsing phase.

```

2540 \msg_kernel_new:nnnn { regex } { trailing-backslash }
2541   { Trailing~escape~character~\iow_char:N\\. }
2542   {
2543     A~regular~expression~or~its~replacement~text~ends~with~
2544     the~escape~character~\iow_char:N\\. It~will~be~ignored.
2545   }
2546 \msg_kernel_new:nnnn { regex } { x-missing-rbrace }
2547   { Missing~closing~brace~in~\iow_char:N\x~hexadecimal~sequence. }
2548   {
2549     You~wrote~something~like~
2550     '\iow_char:N\x{\int_to_hexadecimal:n{#1}}'.~
2551     The~closing~brace~is~missing.
2552   }
2553 \msg_kernel_new:nnnn { regex } { x-overflow }
2554   { Character~code~'#1'~too~large~in~\iow_char:N\x~hexadecimal~sequence. }
2555   {
2556     You~wrote~something~like~
2557     '\iow_char:N\x{\int_to_hexadecimal:n{#1}\}'.~
2558     The~character~code~'#1'~is~larger~than~\int_use:N \c_max_char_int.
2559 }
```

Invalid quantifier.

```

2560 \msg_kernel_new:nnnn { regex } { invalid-quantifier }
2561   { Braced-quantifier~'#1'~may~not~be~followed~by~'#2'. }
2562   {
2563     The~character~'#2'~is~invalid~in~the~braced~quantifier~'#1'.~
2564     The~only~valid~quantifiers~are~'*',~'?',~'+',~'{<int>}',~
2565     '{<min>},~and~'{<min>,<max>}',~followed~or~not~by~'?'.
2566 }
```

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

```

2567 \msg_kernel_new:nnnn { regex } { missing-rbrack }
2568   { Missing~right~bracket~inserted~in~regular~expression. }
2569   {
```

```

2570   LaTeX-was-given-a-regular-expression-where-a-character-class-
2571   was-started-with-'[,~but-the-matching']'~is~missing.
2572 }
2573 \msg_kernel_new:nnnn { regex } { missing-rparen }
2574 {
2575   Missing-right-parenthes\int_compare:nTF{#1=1}{i}{e}s-
2576   inserted-in-regular-expression.
2577 }
2578 {
2579   LaTeX-was-given-a-regular-expression-with-\int_eval:n{#1}-
2580   more-left-parenthes\int_compare:nTF{#1=1}{i}{e}s-than-right-
2581   parenthes\int_compare:nTF{#1=1}{i}{e}s.
2582 }

```

```

2583 \msg_kernel_new:nnnn { regex } { extra-rparen }
2584 {
2585   Extra-right-parenthesis-ignored-in-regular-expression. }
2586 {
2587   LaTeX-came-across-a-closing-parenthesis-when-no-submatch-group-
2588   was-open.-The-parenthesis-will-be-ignored.
2589 }

```

Some escaped alphanumerics are not allowed everywhere.

```

2589 \msg_kernel_new:nnnn { regex } { bad-escape }
2590 {
2591   Invalid-escape-\c_backslash_str #1-
2592   \regex_if_in_cs:TF { within-a-control-sequence. }
2593   {
2594     \regex_if_in_class:TF
2595     { in-a-character-class. }
2596     { following-a-category-test. }
2597   }
2598 }
2599 {
2600   The-escape-sequence-\iow_char:N\#1-may-not-appear-
2601   \regex_if_in_cs:TF
2602   {
2603     within-a-control-sequence-test-introduced-by-
2604     \iow_char:N\c\iow_char:N\{.
2605   }
2606   {
2607     \regex_if_in_class:TF
2608     { within-a-character-class~ }
2609     { following-a-category-test-such-as-\iow_char:N\cL~ }
2610     because-it-does-not-match-exactly-one-character.
2611   }
2612 }

```

Range errors.

```

2613 \msg_kernel_new:nnnn { regex } { range-missing-end }
2614   { Invalid-end-point-for-range-#1-#2-in-character-class. }
2615   {
2616     The-end-point-#2-of-the-range-#1-#2-may-not-serve-as-an-

```

```

2617     end-point-for-a-range:-alphanumeric-characters-should-not-be-
2618     escaped,-and-non-alphanumeric-characters-should-be-escaped.
2619   }
2620 \msg_kernel_new:nnnn { regex } { range-backwards }
2621   { Range-[#1-#2]-out-of-order-in-character-class. }
2622   {
2623     In-ranges-of-characters-[x-y]-appearing-in-character-classes,-
2624     the-first-character-code-must-not-be-larger-than-the-second.-.
2625     Here,-#1-has-character-code-\int_eval:n {'#1},-while-#2-has-
2626     character-code-\int_eval:n {'#2}.
2627   }
2628 Errors related to \c and \u.
2629
2630 \msg_kernel_new:nnnn { regex } { c-bad-mode }
2631   { Invalid-nested-\iow_char:N\c-escape-in-regular-expression. }
2632   {
2633     The-\iow_char:N\c-escape-cannot-be-used-within-
2634     a-control-sequence-test-\iow_char:N\c{...}.-
2635     To-combine-several-category-tests,-use-\iow_char:N\c[...].
2636   }
2637 \msg_kernel_new:nnnn { regex } { c-missing-rbrace }
2638   { Missing-right-brace-inserted-for-\iow_char:N\c-escape. }
2639   {
2640     LaTeX-was-given-a-regular-expression-where-a-
2641     '\iow_char:N\c\iow_char:N\{\...'-construction-was-not-ended-
2642     with-a-closing-brace-\iow_char:N\}'.
2643   }
2644 \msg_kernel_new:nnnn { regex } { c-missing-rbrack }
2645   { Missing-right-bracket-inserted-for-\iow_char:N\c-escape. }
2646   {
2647     A-construction-\iow_char:N\c[...]-appears-in-a-
2648     regular-expression,-but-the-closing-]-is-not-present.
2649   }
2650 \msg_kernel_new:nnnn { regex } { c-missing-category }
2651   { Invalid-character-'#1'-following-\iow_char:N\c-escape. }
2652   {
2653     In-regular-expressions,-the-\iow_char:N\c-escape-sequence-
2654     may-only-be-followed-by-a-left-brace,-a-left-bracket,-or-a-
2655     capital-letter-representing-a-character-category,-namely-
2656     one-of-ABCDELMOPSTU.
2657   }
2658 \msg_kernel_new:nnnn { regex } { u-missing-lbrace }
2659   { Missing-left-brace-following-\iow_char:N\u-escape. }
2660   {
2661     The-\iow_char:N\u-escape-sequence-must-be-followed-by-
2662     a-brace-group-with-the-name-of-the-variable-to-use.
2663   }
2664 \msg_kernel_new:nnnn { regex } { u-missing-rbrace }
2665   { Missing-right-brace-inserted-for-\iow_char:N\u-escape. }
2666   {

```

```

2665     LaTeX-
2666     \tl_if_empty:xTF {#2}
2667         { reached~the~end~of~the~string~ }
2668         { encountered~an~escaped~alphanumeric~character '\iow_char:N\\#2'~ }
2669         when~parsing~the~argument~of~an~'\iow_char:N\\u\iow_char:N\{...\}'~escape.
2670     }

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

2671 \msg_kernel_new:nnnn { regex } { posix-unsupported }
2672     { POSIX~collating~element~'[#1 ~ #1]'~not~supported. }
2673     {
2674         The~[.foo.]~and~[=bar=]~syntaxes~have~a~special~meaning~in~POSIX~
2675         regular~expressions.~This~is~not~supported~by~LaTeX.~Maybe~you~
2676         forgot~to~escape~a~left~bracket~in~a~character~class?
2677     }
2678 \msg_kernel_new:nnnn { regex } { posix-unknown }
2679     { POSIX-class~[:#1:]~unknown. }
2680     {
2681         [:#1:]~is~not~among~the~known~POSIX~classes~
2682         [:alnum:],~[:alpha:],~[:ascii:],~[:blank:],~
2683         [:cntrl:],~[:digit:],~[:graph:],~[:lower:],~
2684         [:print:],~[:punct:],~[:space:],~[:upper:],~
2685         [:word:],~and~[:xdigit:].
2686     }
2687 \msg_kernel_new:nnnn { regex } { posix-missing-close }
2688     { Missing~closing~':]'~for~POSIX~class. }
2689     { The~POSIX~syntax~'#1'~must~be~followed~by~':]',~not~'#2'. }


```

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

```

2690 \msg_kernel_new:nnnn { regex } { result-unbalanced }
2691     { Missing~brace~inserted~when~#1. }
2692     {
2693         LaTeX~was~asked~to~do~some~regular~expression~operation,~
2694         and~the~resulting~token~list~would~not~have~the~same~number~
2695         of~begin-group~and~end-group~tokens.~Braces~were~inserted:~
2696         #2~left,~#3~right.
2697     }


```

Error message for unknown options.

```

2698 \msg_kernel_new:nnnn { regex } { unknown-option }
2699     { Unknown~option~'#1'~for~regular~expressions. }
2700     {
2701         The~only~available~option~is~'case-insensitive',~toggled~by~
2702         '(?i)'~and~'(?-i)'.
2703     }


```

Errors in the replacement text.

```

2704 \msg_kernel_new:nnnn { regex } { replacement-c }
2705     { Misused~\iow_char:N\\c~command~in~a~replacement~text. }


```

```

2706  {
2707    In~a~replacement~text,~the~\iow_char:N\\c~escape~sequence~
2708    can~be~followed~by~one~of~the~letters~ABCDELMOPSTU~
2709    or~a~brace~group,~not~by~'#1'.
2710  }
2711 \msg_kernel_new:nnn { regex } { replacement-u }
2712  { Misused~\iow_char:N\\u~command~in~a~replacement~text. }
2713  {
2714    In~a~replacement~text,~the~\iow_char:N\\u~escape~sequence~
2715    must~be~followed~by~a~brace~group~holding~the~name~of~the~
2716    variable~to~use.
2717  }
2718 \msg_kernel_new:nnn { regex } { replacement-g }
2719  { Missing~brace~for~the~\iow_char:N\\g~construction~in~a~replacement~text. }
2720  {
2721    In~the~replacement~text~for~a~regular~expression~search,~
2722    submatches~are~represented~either~as~\iow_char:N \\g{dd..d},~
2723    or~\d,~where~'d'~are~single~digits.~Here,~a~brace~is~missing.
2724  }
2725 \msg_kernel_new:nnn { regex } { replacement-catcode-end }
2726  {
2727    Missing~character~for~the~\iow_char:N\\c<category><character>~
2728    construction~in~a~replacement~text.
2729  }
2730  {
2731    In~a~replacement~text,~the~\iow_char:N\\c~escape~sequence~
2732    can~be~followed~by~one~of~the~letters~ABCDELMOPSTU~representing~
2733    the~character~category.~Then,~a~character~must~follow.~LaTeX~
2734    reached~the~end~of~the~replacement~when~looking~for~that.
2735  }
2736 \msg_kernel_new:nnn { regex } { replacement-null-space }
2737  { TeX~cannot~build~a~space~token~with~character~code~0. }
2738  {
2739    You~asked~for~a~character~token~with~category~'space',~
2740    and~character~code~0,~for~instance~through~
2741    '\iow_char:N\\cS\iow_char:N\\x00'.~
2742    This~specific~case~is~impossible~and~will~be~replaced~
2743    by~a~normal~space.
2744  }
2745 \msg_kernel_new:nnn { regex } { replacement-missing-rbrace }
2746  { Missing~right~brace~inserted~in~replacement~text. }
2747  {
2748    There~were~\int_use:N \\l_regex_replacement_csnames_int \\
2749    missing~right~braces.
2750  }

```

\regex\_msg\_repeated:nnN

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

```

2751 \cs_new:Npn \regex_msg_repeated:nnN #1#2#3
2752 {
2753     \str_if_eq:xxF { #1 #2 } { 1 0 }
2754     {
2755         , ~ repeated ~
2756         \prg_case_int:nnn {#2}
2757         {
2758             { -1 } { #1~or~more~times,~\bool_if:NTF #3 { lazy } { greedy } }
2759             { 0 } { #1~times }
2760         }
2761         {
2762             between~#1~and~\int_eval:n {#1+#2}~times,~
2763             \bool_if:NTF #3 { lazy } { greedy }
2764         }
2765     }
2766 }
```

(End definition for \regex\_msg\_repeated:nnN.)

## 2.9 Code for tracing

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

\regex\_trace\_states:n This function lists the contents of all states of the NFA, stored in \toks from 0 to \l\_regex\_max\_state\_int (excluded).

```

2767 {*trace}
2768 \cs_new_protected:Npn \regex_trace_states:n #1
2769 {
2770     \prg_stepwise_inline:nnnn
2771         \l_regex_min_state_int
2772         \c_one
2773         { \l_regex_max_state_int - 1 }
2774         {
2775             \trace:nnx { regex } { #1 }
2776             { \iow_char:N \\toks ##1 = { \tex_the:D \tex_toks:D ##1 } }
2777         }
2778     }
2779 
```

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

```
2780 
```