

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

The L^AT_EX3 Project[†]

Released 2011/10/09

1 I3regex documentation

The I3regex package provides regular expression testing, extraction of submatches, splitting, and replacement, all acting on strings of characters. The syntax of regular expressions is mostly a subset of the PCRE syntax (and very close to POSIX). For performance reasons, only a limited set of features are implemented. Notably, back-references are not supported.

Let us give a few examples. After

```
\str_set:Nn \l_my_str { That~cat. }
\regex_replace_once:nnN { at } { is } \l_my_str
```

the string variable `\l_my_str` 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+) } { \1 , } \l_my_str
```

The `\w` sequence represents any “word” character, and `*` indicates that the `\w` sequence should be repeated as many times as possible, hence matching a word in the input string. The parentheses “capture” what their contents matched in the input string, and this can be used in the replacement text as `\1` (and higher numbers if several groups are used in the regular expression).

1.1 Syntax of regular expressions

Most characters match exactly themselves. Some characters are special and must be escaped with a backslash (*e.g.*, `*` matches an explicit 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, none of them should be escaped, because most of those escape sequences have special meanings;

^{*}This file describes v2900, last revised 2011/10/09.

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

- non-alphanumeric printable ascii characters can always be safely escaped, and for highest portability, they should always be escaped (this avoids problems if some currently “normal” characters is given a meaning in later releases);
- 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 literal string `$%\^\\abc#`.

TeXhackers note: When converting the regular expression to a string, the value of the escape character is set to be a backslash.

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

Characters.

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

`\x{hh}` 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 character, including newlines.¹

`\N` A character that is not the `\n` character (hex 0A).²

`\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 `[\^I-\^M]`. Note that `\^K` is a vertical space, but not a space, for compatibility with perl.

¹Should that be changed?

²Is that right?

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

\D Any character not matched by \d.

\H Any character not matched by \h.

\S Any character not matched by \s.

\V Any character not matched by \v.

\W Any character not matched by \w.

Character classes match exactly one character in the subject string.

[...] Positive character class.

[^...] Negative character class.

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

Quantifiers.

? 0 or 1, greedy.

?? 0 or 1, lazy.

* 0 or more, greedy.

*? 0 or more, lazy.

+ 1 or more, greedy.

+? 1 or more, lazy.

Not implemented yet:

{n} Exactly n .

{n,} n or more, greedy.

{n,}? n or more, lazy.

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

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

Anchors and simple assertions.

\b Word boundary.

\B Not a word boundary.

^\A Start of the subject string.³

³The multiline mode is not implemented yet, so those are currently identical.

`$\Z\z` End of the subject string.³

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

Alternation and capturing groups.

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

`(...)` Capturing group.

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

In character classes, only `^`, `-`, `]`, `\` and spaces are special, and should be escaped. Other non-alphanumeric characters can still be escaped without harm. The escape sequences `\d`, `\D`, `\w`, `\W` are also supported in character classes. If the first character is `^`, then the meaning of the character class is inverted. Ranges of characters can be expressed using `-`, for instance, `[\D 0-5]` is equivalent to `[^6-9]`.

Capturing groups are a means of extracting information about the match. Parenthesized groups are labelled in the order of their opening parenthesis, starting at 1. The contents of those groups corresponding to the “best” match (leftmost longest) can be extracted and stored in a sequence of strings using for instance `\regex_extract:nnNTF`.

1.2 Syntax of in the replacement text

Most of the features described in regular expressions do not make sense within the replacement text. Escaped characters are supported as inside regular expressions. The various submatches are accessed with `\1`, `\2`, etc., and the whole match is accessed using `\0`.

For instance,

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

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

Submatches with numbers higher than 10 are accessed in the same way, namely `\10`, `\11`, etc. To insert in the replacement text a submatch followed by a digit, the digit must be entered using the `\x` escape sequence: for instance, to get the first submatch followed by the digit 7, use `\1\x37`, because 7 has character code 37 (in hexadecimal).

1.3 Precompiling regular expressions

If a regular expression is to be used several times, it is better to compile it once rather than doing it each time the regular expression is used. The precompiled regular expression is stored as a token list variable. All of the `I3regex` module’s functions can be given their regular expression argument either as an explicit string or as a precompiled regular expression.

\regex_set:Nn \regex_set:Nn *tl var* {*regex*}

\regex_gset:Nn \regex_set:Nn *tl var* {*regex*}

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

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

The assignment is local for \regex_set:Nn and global for \regex_gset:Nn.

TeXhackers note: Precompiled regular expressions can be safely written to a file and read when the L^AT_EX3 syntax is active (as triggered by \ExplSyntaxOn).

1.4 String matching

\regex_match:nnTF \regex_match:nnTF {*regex*} {*string*} {*true code*} {*false code*}

\regex_match:NnTF \regex_match:nnTF {*regex*} {*string*} {*true code*} {*false code*}

Tests whether the *regular expression* matches any substring of *string*. For instance,

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

leaves TRUE FALSE in the input stream.

\regex_count:nnN \regex_count:nnN {*regex*} {*string*} *int var*

\regex_count:NnN \regex_count:nnN {*regex*} {*string*} *int var*

Sets *int var* equal to the number of times *regular expression* appears in *string*. 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 string. For instance,

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

results in \l_foo_int taking the value 5.

1.5 Submatch extraction

`\regex_extract:nnNTF` `\regex_extract:nnNTF {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩} {⟨true code⟩} {⟨false code⟩}`

`\regex_extract:NnNTF` `\regex_extract:nnNTF {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩} {⟨true code⟩} {⟨false code⟩}`

`\regex_extract:nnNTF` `\regex_extract:nnNTF {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩} {⟨true code⟩} {⟨false code⟩}`

`\regex_extract:NnNTF` `\regex_extract:nnNTF {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩} {⟨true code⟩} {⟨false code⟩}`

Finds the first match of the *⟨regular expression⟩* in the *⟨string⟩*. If it exists, the match is stored as the zeroeth item of the *⟨seq var⟩*, and further items are the contents of capturing groups, in the order of their opening left parenthesis. The *⟨seq var⟩* is assigned locally. If there is no match, the *⟨seq var⟩* is not altered. The testing versions return *⟨true⟩* if a match was found, and *⟨false⟩* otherwise. For instance, assume that you type

```
\regex_extract:nnNTF { ^(\La)?TeX(!*)$ } { LaTeX!!! }
\l_foo_seq { true } { false }
```

Then the regular expression (anchored at the start with `^` and at the end with `$`) will match the whole string. 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.

1.6 String splitting

`\regex_split:nnN` `\regex_split:nnN {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩}`

`\regex_split:NnN` `\regex_split:nnN {⟨regular expression⟩} {⟨string⟩} {⟨seq var⟩}`

Searches the *⟨string⟩* into a sequence of substrings, delimited by matches of the *⟨regular expression⟩*. If the *⟨regular expression⟩* has capturing groups, then the substrings 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 *⟨string⟩* as its sole item. If the *⟨regular expression⟩* matches the empty string, then the *⟨string⟩* is split into single character substrings.

1.7 String replacement

\regex_replace_once:nnTF
\regex_replace_once:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

\regex_replace_once:NnTF
\regex_replace_once:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

\regex_replace_once:nnNTF
\regex_replace_once:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

\regex_replace_once:NnNTF
\regex_replace_once:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

Searches for the *⟨regular expression⟩* in the *⟨string⟩* and replaces the matching part with the *⟨replacement⟩*. The result is assigned locally to *⟨str var⟩*. In the *⟨replacement⟩*, \0 represents the full match, \1 represent the contents of the first capturing group, \2 of the second, *etc.*

\regex_replace_all:nnN
\regex_replace_all:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

\regex_replace_all:NnN
\regex_replace_all:nnN {\langle regular expression \rangle} {\langle replacement \rangle} <str var>

Replaces all occurrences of the \regular expression in the *⟨string⟩* 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. The result is assigned locally to *⟨str var⟩*.

1.8 Bugs, misfeatures, future work, and other possibilities

The {} quantifiers are not implemented.

The following need discussion.

- Newline conventions are not done. In particular, . should not match newlines. Also, \A should differ from ^, and \Z, \z and \$ should differ.
- Caseless matching and more generally (*..) and (?.?) sequences to set some options.
- General look-ahead/behind assertions.
- Idea of # as a synonym of .*?.
- Do we need a facility for balanced groups? (That's non-regular.)

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

- Optimize simple strings: use less states (**abcade** should give two states, for **abc** and **ade**).
- Optimize groups with no alternative.

- Optimize the use of `\prg_stepwise_...` functions.
- Implement regex matching on external files.
- Conditional subpatterns with look ahead/behind: “if what follows is [...], then [...]”.
- `(?|...|...)` to reset the capturing group number at the start of each alternative.
- Todo: check space’s catcode in the result of a replacement.

The following features not present in PCRE nor perl (because they can be done in other ways in those languages) might be added.

- `#` as a synonym of `.*?`, in a way similar to macro parameters.
- Facility to match balanced groups: this is non-regular, but quite useful, and it can be implemented without too much performance loss. (*cf.* callout?)

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

- Other forms of escapes, and character classes. `\cx\simeq\^\^x`, `\ddd` (octal `ddd`). POSIX character classes. `\p{..}` and `\P{..}` for having/not having a Unicode property. `\X` for “extended” Unicode sequence.
- Callout with `(?C...)`.
- Conditional subpatterns (other than with a look-ahead or look-behind condition): this is non-regular, isn’t it?

The following features of PCRE or perl will not be implemented.

- Comments: `\TeX` already has its own system for comments.
- Named subpatterns: `\TeX` programmers have lived so far without any need for named macro parameters.
- `\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.
- 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.
- `\K` for resetting the beginning of the match: tied to backtracking.
- `\C` single byte in UTF-8 mode: well, UTF-8 mode is not implemented, and that seems like a very rarely useful feature anyways.

2 l3regex implementation

```
<*package>
  1 \ProvidesExplPackage
  2   {\ExplFileName}{\ExplFileVersion}{\ExplFileDescription}
  3 \RequirePackage{l3str}
```

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

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

2.1 Constants and variables

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

2.1.1 Variables used while building

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

\l_regex_max_state_int
\l_regex_left_state_int
\l_regex_right_state_int
\l_regex_left_state_seq
\l_regex_right_state_seq
  10 \int_new:N \l_regex_max_state_int
  11 \int_new:N \l_regex_left_state_int
  12 \int_new:N \l_regex_right_state_int
  13 \seq_new:N \l_regex_left_state_seq
  14 \seq_new:N \l_regex_right_state_seq
```

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

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

15 `\int_new:N \l_regex_capturing_group_int`
16 `\seq_new:N \l_regex_capturing_group_seq`

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

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

17 `\tl_new:N \l_regex_one_or_group_tl`

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

`\l_regex_repetition_int` Used when to distinguish threads at the same state but with different repetitions of some quantifiers.

18 `\int_new:N \l_regex_repetition_int`

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

`\l_regex_look_behind_bool` `\l_regex_look_behind_str` The boolean `\l_regex_look_behind_bool` indicates whether a look-behind assertion appears within the regular expression (currently only `\b` or `\B`). When matching, we will keep track of the part of the string before the current character in `\l_regex_look_behind_str`, stored backwards.

19 `\bool_new:N \l_regex_look_behind_bool`
20 `\tl_new:N \l_regex_look_behind_str`

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

2.1.2 Character classes

`\regex_build_tmp_class:n` Used when building character classes.

21 `\cs_new_eq:NN \regex_build_tmp_class:n \use_none:n`
22 `\bool_new:N \l_regex_class_bool`
23 `\tl_new:N \l_regex_class_tl`

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

`\c_regex_d_tl` These constant token lists encode which characters are recognized by `\d`, `\D`, `\w`, *etc.* in regular expressions. Namely, `\d=[0-9]`, `\w=[0-9A-Z_a-z]`, `\s=[\ \^\^I\^\^J\^\^L\^\^M]`, `\h=[\ \^\^I]`, `\v=[\^\^J-\^\^M]`, and the upper-case counterparts match anything that the lowercase does not match. The order in which the various ranges appear is optimized for usual mostly lowercase letter text.

24 `\tl_const:Nn \c_regex_d_tl`
25 {
26 `\regex_item_range:nn {48} {57} % 0--9`
27 }

`\c_regex_v_tl`
`\c_regex_V_tl`
`\c_regex_w_tl`
`\c_regex_W_tl`
`\c_regex_N_tl`

```

28 \tl_const:Nn \c_regex_D_tl
29 {
30   \regex_item_more:n {57} % 9
31   \regex_item_range:nn {0} {47} % 0
32 }
33 \tl_const:Nn \c_regex_h_tl
34 {
35   \regex_item_equal:n {32} % space
36   \regex_item_equal:n {9} % tab
37 }
38 \tl_const:Nn \c_regex_H_tl
39 {
40   \regex_item_neq:n {32} % space
41   \regex_item_neq:n {9} % tab
42   \regex_break_true:w
43 }
44 \tl_const:Nn \c_regex_s_tl
45 {
46   \regex_item_equal:n {32} % space
47   \regex_item_neq:n {11} % vtab
48   \regex_item_range:nn {9} {13} % tab, lf, vtab, ff, cr
49 }
50 \tl_const:Nn \c_regex_S_tl
51 {
52   \regex_item_more:n {32} % > space
53   \regex_item_range:nn {14} {31} % tab < ... < space
54   \regex_item_range:nn {0} {8} % < tab
55   \regex_item_equal:n {11} % vtab
56 }
57 \tl_const:Nn \c_regex_v_tl
58 {
59   \regex_item_range:nn {10} {13} % lf, vtab, ff, cr
60 }
61 \tl_const:Nn \c_regex_V_tl
62 {
63   \regex_item_more:n {13} % cr
64   \regex_item_range:nn {0} {9} % < lf
65 }
66 \tl_const:Nn \c_regex_w_tl
67 {
68   \regex_item_range:nn {97} {122} % a--z
69   \regex_item_range:nn {65} {90} % A--Z
70   \regex_item_range:nn {48} {57} % 0--9
71   \regex_item_equal:n {95} % _
72 }
73 \tl_const:Nn \c_regex_W_tl
74 {
75   \regex_item_range:nn {0} {47} % <'0
76   \regex_item_range:nn {58} {64} % ('9+1)--('A-1)
77   \regex_item_range:nn {91} {94} % ('Z+1)--('_-1)

```

```

78     \regex_item_equal:n {96}      % '
79     \regex_item_more:n {122}      % z
80   }
81 \tl_const:Nn \c_regex_N_tl
82 {
83   \regex_item_neq:n {10} % lf
84   \regex_break_true:w
85 }
(End definition for \c_regex_d_tl and \c_regex_D_tl. These functions are documented on page ??.)

```

2.1.3 Variables used when matching

\l_regex_query_str The string that is being matched.

```

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

```

\l_regex_start_step_int In the case of multiple matches, \l_regex_start_step_int is equal to the position where the current match attempt began.

```

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

```

\l_regex_current_char_int The character at the current position in the string, and the current position.

\l_regex_current_step_int

```

88 \int_new:N \l_regex_current_char_int
89 \int_new:N \l_regex_current_step_int
(End definition for \l_regex_current_char_int. This function is documented on page ??.)

```

\l_regex_unique_step_int This gives a unique identifier to every step in the loop over characters in the string. In the case of a single match, it is equal to \l_regex_current_step_int, but for multiple matches, it is not reset to the start of the next match, but incremented. A unique identifier is handy when considering whether a state of the automaton has been visited at that step or not, and when tracking submatch information.

```

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

```

\l_regex_current_state_int For every character in the string, 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.

```

91 \int_new:N \l_regex_current_state_int
92 \prop_new:N \l_regex_current_submatches_prop
(End definition for \l_regex_current_state_int. This function is documented on page ??.)

```

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

```

93 \int_new:N \l_regex_max_index_int

```

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

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

94 \tl_new:N \l_regex_every_match_t1

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

\regex_last_match_empty:F This function is most often \use:n, unless the current regular expression can match an empty string at the current position, and has already done so at the previous match attempt. The goal is to break infinite loops.

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

96 \cs_new_protected:Npn \regex_last_match_empty_yes:F

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

98 \cs_new_eq:NN \regex_last_match_empty:F \regex_last_match_empty_no:F

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

\l_regex_success_bool
\l_regex_success_step_int
\l_regex_success_submatches_prop
\l_regex_success_empty_bool The booleans \l_regex_success_bool and \l_regex_success_empty_bool are true if the current match attempt was successful, and if the match was empty, respectively. The other two variables hold the step number at which the match was successful and a property list of submatches.

99 \bool_new:N \l_regex_success_bool

100 \bool_new:N \l_regex_success_empty_bool

101 \int_new:N \l_regex_success_step_int

102 \prop_new:N \l_regex_success_submatches_prop

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

\l_regex_fresh_thread_bool This boolean marks when the current thread has started from the beginning of the regular expression at this character, in other words, it is true if the current thread has matched an empty string so far (and we only care about this boolean when a thread succeeds, hence matching an empty string overall).

103 \bool_new:N \l_regex_fresh_thread_bool

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

2.1.4 Variables used for user functions

\g_regex_submatches_seq This holds temporarily a sequence of submatches, global so that it exits the group.

104 \seq_new:N \g_regex_submatches_seq

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

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

105 \int_new:N \g_regex_match_count_int

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

`\g_regex_split_seq` The `\regex_split:nnN` function stores its result in that sequence variable before assigning it to the variable provided by the user.

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

`\l_regex_replacement_tl` The replacement code stores a processed version of the user's argument in `\l_regex_replacement_tl`. The result of the replacement is stored in `\g_regex_replaced_str`, global to exit the group.

```
107 \tl_new:N \l_regex_replacement_tl
108 \tl_new:N \g_regex_replaced_str
(End definition for \l_regex_replacement_tl. This function is documented on page ??.)
```

2.2 Helpers

2.2.1 Toks

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

```
109 \cs_new_eq:NN \s_regex_toks \scan_stop:
(End definition for \s_regex_toks. This function is documented on page ??.)
```

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

```
110 \cs_new_protected:Npn \regex_toks_put_left:Nx #1#2
111 {
112   \cs_set_nopar:Npx \regex_tmp:w \s_regex_toks { \s_regex_toks #2 }
113   \tex_toks:D #1 \exp_after:wN \exp_after:wN \exp_after:wN
114     { \exp_after:wN \regex_tmp:w \tex_the:D \tex_toks:D #1 }
115 }
116 \cs_new_protected:Npn \regex_toks_put_right:Nx #1#2
117 {
118   \cs_set_nopar:Npx \regex_tmp:w {#2}
119   \tex_toks:D #1 \exp_after:wN
120     { \tex_the:D \tex_toks:D \exp_after:wN #1 \regex_tmp:w }
121 }
```

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

2.2.2 Interrupting recursions

\regex_if_tail_stop:N	Test for the end of the <i><string></i> , and either stop or cause an error if it is reached.
\regex_if_tail_error:Nn	<pre> 122 \cs_new_eq:NN \regex_if_tail_stop:N \quark_if_recursion_tail_stop:N 123 \cs_new_protected_nopar:Npn \regex_if_tail_error:Nn #1#2 124 { \quark_if_recursion_tail_stop_do:Nn #1 { \regex_build_error:n {#2} } } (End definition for \regex_if_tail_stop:N. This function is documented on page ??.)</pre>
\regex_build_error:n	This macro is called if anything goes wrong when building the NFA corresponding to a given regular expression Negative codes are specific to the L ^A T _E X3 implementation. Other error codes match with the PCRE codes (not all of the PCRE errors can occur, since some constructions are not supported). <pre> 125 \cs_new_protected_nopar:Npn \regex_build_error:n #1 126 { 127 \msg_error:nnnx { regex } { build-error } {\int_eval:n{#1}} 128 { 129 \prg_case_int:nnn {#1} 130 { 131 {-999} {File-not-found} 132 {-998} {Unsupported-construct} 133 {-997} {The-regular-expression-is-too-large-(32768-states).} 134 {1} {\iow_char:N\at-end-of-pattern} 135 % {2} {\iow_char:N\c-at-end-of-pattern} 136 {4} {Numbers-out-of-order-in-\iow_char:N\{\iow_char\}-quantifier.} 137 {6} {Missing-terminating-\iow_char:N]-for-character-class} 138 {7} {Invalid-escape-sequence-in-character-class} 139 {8} {Range-out-of-order-in-character-class} 140 {22} {Mismatched-parentheses} 141 {34} {Character-value-in-\iow_char:N\x{...}-sequence-is-too-large} 142 % 143 {44} {Invalid-UTF-8-string} 144 % 145 {46} {Malformed-\iow_char:N\P-or\iow_char:N\p-sequence} 146 {47} {Unknown-property-after-\iow_char:N\P-or\iow_char:N\p} 147 {68} {\iow_char:N\c-must-be-followed-by-an-ASCII-character} 148 } 149 { Internal-bug. } 150 } 151 }</pre>

2.2.3 Testing characters

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

```

<test1> ... <testn>
\regex_break_point:TF {\<true code>} {\<false code>}
```

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

```

151 \cs_new_nopar:Npn \regex_break_true:w #1 \regex_break_point:TF #2 #3 {#2}
152 \cs_new_nopar:Npn \regex_break_false:w #1 \regex_break_point:TF #2 #3 {#3}
153 \cs_new_eq:NN \regex_break_point:TF \use_i:nn
(End definition for \regex_break_point:TF. This function is documented on page ??.)
```

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

```

\regex_item_range:nn
\regex_item_less:n
\regex_item_more:n
\regex_item_neq:n
154 \cs_new_nopar:Npn \regex_item_equal:n #1
155 {
156   \if_num:w #1 = \l_regex_current_char_int
157     \exp_after:wN \regex_break_true:w
158   \fi:
159 }
160 \cs_new_nopar:Npn \regex_item_range:nn #1 #2
161 {
162   \reverse_if:N \if_num:w #1 > \l_regex_current_char_int
163   \reverse_if:N \if_num:w #2 < \l_regex_current_char_int
164     \exp_after:wN \exp_after:wN \exp_after:wN \regex_break_true:w
165   \fi:
166   \fi:
167 }
168 \cs_new_nopar:Npn \regex_item_less:n #1
169 {
170   \if_num:w #1 > \l_regex_current_char_int
171     \exp_after:wN \regex_break_true:w
172   \fi:
173 }
174 \cs_new_nopar:Npn \regex_item_more:n #1
175 {
176   \if_num:w #1 < \l_regex_current_char_int
177     \exp_after:wN \regex_break_true:w
178   \fi:
179 }
180 \cs_new_nopar:Npn \regex_item_neq:n #1
181 {
182   \if_num:w #1 = \l_regex_current_char_int
183     \exp_after:wN \regex_break_false:w
184   \fi:
185 }
```

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

2.2.4 Grabbing digits

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

```

186 \cs_new_protected:Npn \regex_get_digits:nw #1
187 {
188     \tex_afterassignment:D \regex_tmp:w
189     \cs_set_nopar:Npx \regex_tmp:w
190     {
191         \exp_not:n {#1}
192         { \if_false: } } \fi:
193         \regex_get_digits_aux:N
194     }
195 \cs_new_nopar:Npn \regex_get_digits_aux:N #1
196 {
197     \if_num:w 9 < 1 \exp_not:N #1 \exp_stop_f:
198     \else:
199         \if_charcode:w \c_space_token \exp_not:N #1
200         \else:
201             \regex_get_digits_end:w
202             \fi:
203         \fi:
204         #1 \regex_get_digits_aux:N
205     }
206 \cs_new_nopar:Npn \regex_get_digits_end:w
207     \fi: \fi: #1 \regex_get_digits_aux:N
208 {
209     \fi: \fi:
210     \if_false: { { \fi: } }
211     #1
212 }

```

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

2.2.5 More char testing

```
\regex_aux_char_if_alphanumeric:NTF
\regex_aux_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).

```
213 \prg_new_conditional:Npnn \regex_aux_char_if_special:N #1 { TF }
```

```

214  {
215  \if_num:w '#1 < 97 \exp_stop_f:
216  \if_num:w '#1 < 58 \exp_stop_f:
217  \if_num:w \int_eval:w ('#1 - \c_eight)/\c_sixteen = \c_two
218  \prg_return_true:
219  \else:
220  \prg_return_false:
221  \fi:
222  \else:
223  \if_num:w \int_eval:w '#1 / 26 = \c_three
224  \prg_return_false:
225  \else:
226  \prg_return_true:
227  \fi:
228  \fi:
229  \else:
230  \if_num:w \int_eval:w '#1 / \c_five = 25 \exp_stop_f:
231  \prg_return_true:
232  \else:
233  \prg_return_false:
234  \fi:
235  \fi:
236  }
237 \prg_new_conditional:Npnn \regex_aux_char_if_alphanumeric:N #1 { TF }
238 {
239  \if_num:w '#1 < 91 \exp_stop_f:
240  \if_num:w '#1 < 65 \exp_stop_f:
241  \if_num:w \c_nine < 1 #1 \exp_stop_f:
242  \prg_return_true:
243  \else:
244  \prg_return_false:
245  \fi:
246  \else:
247  \prg_return_true:
248  \fi:
249  \else:
250  \if_num:w \int_eval:w ('#1-\c_six)/26 = \c_four
251  \prg_return_true:
252  \else:
253  \prg_return_false:
254  \fi:
255  \fi:
256 }

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

```

2.3 Building

2.3.1 Helpers for building an NFA

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

```
257 \cs_new_protected_nopar:Npn \regex_build_new_state:
258 {
259     \int_compare:nNnTF \l_regex_max_state_int > {32766}
260     { \regex_build_error:n {-997} }
261     {
262         \int_incr:N \l_regex_max_state_int
263         \tex_toks:D \l_regex_max_state_int { \s_regex_toks }
264     }
265     \int_set_eq:NN \l_regex_left_state_int \l_regex_right_state_int
266     \int_set_eq:NN \l_regex_right_state_int \l_regex_max_state_int
267 }
```

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

\regex_build_transition_aux:NN
\regex_build_transitions_aux:NNNN These functions create a new state, and put one or two transitions starting from the old current state.

```
268 \cs_new_protected_nopar:Npn \regex_build_transition_aux:NN #1#2
269 {
270     \regex_build_new_state:
271     \regex_toks_put_right:Nx \l_regex_left_state_int
272     { #1 { \int_use:N #2 } }
273 }
274 \cs_new_protected_nopar:Npn \regex_build_transitions_aux:NNNN #1#2#3#4
275 {
276     \regex_build_new_state:
277     \regex_toks_put_right:Nx \l_regex_left_state_int
278     {
279         #1 { \int_use:N #2 }
280         #3 { \int_use:N #4 }
281     }
282 }
```

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

2.3.2 From regex to NFA: framework

In order for the construction ab|cd to work, we enclose the whole pattern within parentheses. These have the added benefit to form a capturing group: hence we get the data of the whole match for free.

\regex_build:n First, reset a few variables. Then use the generic framework defined in l3str to parse the regular expression once, recognizing which characters are raw characters, and which have special meanings. The result is stored in \g_str_result_tl, and can be run directly.

The trailing `\prg_do_nothing:` ensure that the look-ahead done by some of the operations is harmless. Finally, `\regex_build_end:` adds the finishing code (checking that parentheses are properly nested, for instance).

```

283 \cs_new_protected:Npn \regex_build:n #1
284 {
285     \regex_build_setup:
286     \str_aux_escape>NNNn
287     \regex_build_i_unescaped:N
288     \regex_build_i_escaped:N
289     \regex_build_i_raw:N
290     { #1 }
291     \regex_build_open_aux:
292     \g_str_result_tl \prg_do_nothing: \prg_do_nothing:
293     \seq_push:Nn \l_regex_capturing_group_seq { 0 }
294     \regex_build_close_aux: \regex_build_group_:
295     \regex_build_end:
296 }
```

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

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

```

297 \cs_new_nopar:Npn \regex_build_i_unescaped:N #1
298 {
299     \regex_aux_char_if_special:NTF #1
300     { \exp_not:N \regex_build_control:N #1 }
301     { \exp_not:N \regex_build_raw:N #1 }
302 }
303 \cs_new_nopar:Npn \regex_build_i_escaped:N #1
304 {
305     \regex_aux_char_if_alphanumeric:NTF #1
306     { \exp_not:N \regex_build_control:N #1 }
307     { \exp_not:N \regex_build_raw:N #1 }
308 }
309 \cs_new_nopar:Npn \regex_build_i_raw:N #1
310 { \exp_not:N \regex_build_raw:N #1 }
```

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

`\regex_build_default_control:N` If the control character has a particular meaning in regular expressions, the corresponding function is used. Otherwise, it is interpreted as a raw character. The `\regex_build_default_raw:N` function is defined later.

```

311 \cs_new_protected_nopar:Npn \regex_build_default_control:N #1
312 {
313     \cs_if_exist_use:cF { \regex_build_#1: }
314     { \regex_build_default_raw:N #1 }
315 }
```

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

- \regex_build_setup: Hopefully, we didn't forget to initialize anything here. The search is not anchored: to achieve that, we insert state(s) responsible for repeating the match attempt on every character of the string: the first state has a free transition to the second state, where the regular expression really begins, and a costly transition to itself, to try again at the next character.

```
316 \cs_new_protected_nopar:Npn \regex_build_setup:
317 {
318     \cs_set_eq:NN \regex_build_control:N \regex_build_default_control:N
319     \cs_set_eq:NN \regex_build_raw:N \regex_build_default_raw:N
320     \int_set_eq:NN \l_regex_capturing_group_int \c_zero
321     \int_zero:N \l_regex_left_state_int
322     \int_zero:N \l_regex_right_state_int
323     \int_zero:N \l_regex_max_state_int
324     \regex_build_new_state:
325     \regex_build_new_state:
326     \regex_toks_put_right:Nx \l_regex_left_state_int
327     {
328         \regex_action_start_wildcard:nn
329         { \int_use:N \l_regex_left_state_int }
330         { \int_use:N \l_regex_right_state_int }
331     }
332 }
```

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

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

```
333 \cs_new_protected_nopar:Npn \regex_build_end:
334 {
335     \seq_if_empty:NF \l_regex_capturing_group_seq
336     {
337         \regex_build_error:n {22}
338         \prg_replicate:nn
339         { \seq_length:N \l_regex_capturing_group_seq } % (
340         { \regex_build_close_aux: \regex_build_group_:
341     }
342     \regex_toks_put_right:Nx \l_regex_right_state_int
343     { \regex_action_success: }
344 }
```

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

2.3.3 Anchoring and simple assertions

- \regex_build_A: Anchoring at the start corresponds to checking that the current character is the first in the string. Anchoring to the beginning of the match attempt uses \l_regex_start_step_int instead of \c_zero.

```
345 \cs_new_protected_nopar:cpn { regex_build_~: }
```

```

346   { \regex_build_anchor_start:N \c_zero }
347 \cs_new_protected_nopar:Npn \regex_build_A:
348   { \regex_build_anchor_start:N \c_zero }
349 \cs_new_protected_nopar:Npn \regex_build_G:
350   { \regex_build_anchor_start:N \l_regex_start_step_int }
351 \cs_new_protected_nopar:Npn \regex_build_anchor_start:N #1
352   {
353     \regex_build_new_state:
354     \regex_toks_put_right:Nx \l_regex_left_state_int
355     {
356       \exp_not:N \int_compare:nNnT {#1} = \l_regex_current_step_int
357       { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
358     }
359   }
(End definition for \regex_build_A:. This function is documented on page ??.)
```

\regex_build_Z: This matches the end of the string, marked by a character code of -1 .

\regex_build_z: `360 \cs_new_protected_nopar:cpn { regex_build_$: } % $`

\regex_build_\$: `361 {
362 \regex_build_new_state:
363 \regex_toks_put_right:Nx \l_regex_left_state_int
364 {
365 \exp_not:N \int_compare:nNnT
366 \c_minus_one = \l_regex_current_char_int
367 { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
368 }
369 }`

`370 \cs_new_eq:Nc \regex_build_Z: { regex_build_$: } %$`

`371 \cs_new_eq:Nc \regex_build_z: { regex_build_$: } %$`

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

\regex_build_b: Contrarily to \wedge and $\$, which could be implemented without really knowing what precedes$

\regex_build_B: in the string, this requires more information. We request it by setting `\l_regex_look_-` behind_bool. Then the matching code will keep store the characters that were already read, backwards, in `\l_regex_look_behind_str`, and we can analyse the first character of that string.

```

372 \cs_new_protected_nopar:Npn \regex_build_b:
373   {
374     \bool_set_true:N \l_regex_look_behind_bool
375     \regex_build_new_state:
376     \regex_toks_put_right:Nx \l_regex_left_state_int
377     {
378       \exp_not:N \regex_if_word_boundary:TF
379       { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
380       { }
381     }
382   }
383 \cs_new_protected_nopar:Npn \regex_build_B:
384   {
```

```

385 \bool_set_true:N \l_regex_look_behind_bool
386 \regex_build_new_state:
387 \regex_toks_put_right:Nx \l_regex_left_state_int
388 {
389     \exp_not:N \regex_if_word_boundary:TF
390     {
391         { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
392     }
393 }
394 \cs_new_protected_nopar:Npn \regex_if_word_boundary:TF
395 {
396     \tl_if_empty:NTF \l_regex_look_behind_str
397     { \c_regex_w_tl }
398     {
399         \group_begin:
400             \cs_set_nopar:Npx \regex_tmp:w
401             {
402                 \int_set:Nn \l_regex_current_char_int
403                 { ` \str_head:N \l_regex_look_behind_str }
404             }
405             \regex_tmp:w
406             \c_regex_w_tl
407             \regex_break_point:TF
408             { \group_end: \c_regex_W_tl }
409             { \group_end: \c_regex_w_tl }
410         }
411         \regex_break_point:TF
412     }

```

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

2.3.4 Normal character, and simple character classes

`\regex_build_default_raw:N` A normal alphanumeric or an escaped non-alphanumeric (actually, any unknown combination) will match itself and the thread will fail otherwise. We prepare `\regex_build_tmp_class:n` with the relevant test and commands. The program steps to be inserted in those commands will come as `##1` and `##2`: we don't know yet what those will be before checking for quantifiers.

```

413 \cs_new_protected_nopar:Npn \regex_build_default_raw:N #1
414 {
415     \cs_set:Npx \regex_build_tmp_class:n ##1
416     {
417         \exp_not:n { \exp_not:N \if_num:w }
418         \int_value:w '#1 = \l_regex_current_char_int
419         \regex_action_cost:n { ##1 }
420         \exp_not:n { \exp_not:N \fi: }
421     }
422     \regex_build_one_quantifier:
423 }

```

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

\regex_build_:: Similar to \regex_build_default_raw:N but accepts any character, and refuses -1, which marks the end of the string.

```

424 \cs_new_protected_nopar:cpn { regex_build_:: }
425   {
426     \cs_set:Npn \regex_build_tmp_class:n ##1
427     {
428       \exp_not:N \if_num:w \c_minus_one < \l_regex_current_char_int
429         \regex_action_cost:n {##1}
430       \exp_not:N \fi:
431     }
432     \regex_build_one_quantifier:
433   }
(End definition for \regex_build_:: This function is documented on page ??.)
```

\regex_build_d: The constants \c_regex_d_tl, etc. hold a list of tests which match the corresponding character class, and jump to the \regex_break_point:TF marker. As for a normal character, we check for quantifiers.

```

434 \cs_new_protected_nopar:Npn \regex_build_char_type:N #1
435   {
436     \cs_set:Npn \regex_build_tmp_class:n ##1
437     {
438       \exp_not:N #1
439       \exp_not:N \regex_break_point:TF
440         { \regex_action_cost:n {##1} }
441         { }
442     }
443     \regex_build_one_quantifier:
444   }
445 \tl_map_inline:nn { dDhHsSvVwWN }
446   {
447     \cs_new_protected_nopar:cpn { regex_build_#1: }
448     {
449       \exp_not:N \regex_build_char_type:N
450       \exp_not:c { c_regex_#1_tl }
451     }
452   }
(End definition for \regex_build_d: and \regex_build_D. These functions are documented on page ??.)
```

2.3.5 Character classes

\regex_build_[: This starts a class. The code for the class is collected in \l_regex_class_tl. The first character is special.

```

453 \cs_new_protected_nopar:cpn { regex_build_[: }
454   {
455     \tl_clear:N \l_regex_class_tl
456     \cs_set_eq:NN \regex_build_control:N \regex_class_control:N
457     \cs_set_eq:NN \regex_build_raw:N \regex_class_raw:N
458     \regex_class_first:NN
```

```

459     }
(End definition for \regex_build_[:. This function is documented on page ??.)
```

\regex_class_control:N This function is similar to \regex_build_control:N. If the control character has a meaning in character classes, call the corresponding function, otherwise, treat it as a raw character, with the \regex_class_raw:N function, defined later.

```

460 \cs_new_protected_nopar:Npn \regex_class_control:N #1
461 {
462     \cs_if_exist_use:cF { regex_class_#1: }
463     { \regex_class_raw:N #1 }
464 }
(End definition for \regex_class_control:N. This function is documented on page ??.)
```

\regex_class_]: If] appears as the first item of a class, then it doesn't end the class. Otherwise, it's the end, act just as for a single character, but with a more complicated test. And restore \regex_build_control:N and \regex_build_raw:N.

```

465 \cs_new_protected_nopar:cpn { regex_class_]: }
466 {
467     \tl_if_empty:NTF \l_regex_class_t1 %[
468     { \regex_class_raw:N } ]
469     {
470         \cs_set_eq:NN \regex_build_control:N \regex_build_default_control:N
471         \cs_set_eq:NN \regex_build_raw:N \regex_build_default_raw:N
472         \cs_set:Npn \regex_build_tmp_class:n ##1
473         {
474             \exp_not:o \l_regex_class_t1
475             \bool_if:NTF \l_regex_class_bool
476             {
477                 \exp_not:N \regex_break_point:TF
478                 { \regex_action_cost:n {##1} }
479                 { }
480             }
481             {
482                 \exp_not:N \regex_break_point:TF
483                 { }
484                 { \regex_action_cost:n {##1} }
485             }
486         }
487         \regex_build_one_quantifier:
488     }
489 }
```

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

\regex_class_first>NN If the first non-space character is ^, then the class is inverted. We keep track of this in \l_regex_class_bool.

```

490 \cs_new_protected_nopar:Npn \regex_class_first>NN #1#2
491 {
492     \str_if_eq:nnTF {#1#2} { \regex_build_control:N ^ }
```

```

493     { \bool_set_false:N \l_regex_class_bool }
494     {
495         \bool_set_true:N \l_regex_class_bool
496         #1 #2
497     }
498 }
```

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

\regex_class_raw:N Most characters are treated here. We look ahead for an unescaped dash. If there is none, then the character matches itself.

```

499 \cs_new_protected_nopar:Npn \regex_class_raw:N #1#2#3
500   {
501     \str_if_eq:nnTF {#2#3} { \regex_build_control:N - }
502     { \regex_class_range:Nw #1 }
503     {
504       \regex_class_single:N #1
505       #2 #3
506     }
507   }
508 \cs_new_protected_nopar:Npn \regex_class_single:N #1
509   {
510     \tl_put_right:Nx \l_regex_class_tl
511     { \exp_not:N \regex_item_equal:n { \int_value:w '#1 } }
512   }
```

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

\regex_class_range:Nw \regex_class_range_put:NN If the character is followed by a dash, we look for the end-point of the range. For “raw” characters, that’s simply #3. Most “control” characters also have no meaning, and can serve as an end-point, but those with a meaning interrupt the range. In the case of a true range, check whether the end-points are in the right order, and optimize in the case of equal end-points.

```

513 \cs_new_protected_nopar:Npn \regex_class_range:Nw #1#2#3
514   {
515     \token_if_eq_meaning:NNTF #2 \regex_build_control:N
516     {
517       \cs_if_exist:cTF { regex_class_#3: }
518       {
519         \regex_class_single:N #1
520         \regex_class_single:N -
521         #2#3
522       }
523       { \regex_class_range_put:NN #1#3 }
524     }
525     { \regex_class_range_put:NN #1#3 }
526   }
527 \cs_new_protected_nopar:Npn \regex_class_range_put:NN #1#2
528   {
529     \if_num:w '#1 > '#2 \exp_stop_f:
530     \regex_build_error:n {8}
```

```

531     \else:
532         \tl_put_right:Nx \l_regex_class_tl
533     {
534         \if_num:w '#1 = '#2 \exp_stop_f:
535             \exp_not:N \regex_item_equal:n
536         \else:
537             \exp_not:N \regex_item_range:nn { \int_value:w '#1 }
538         \fi:
539         { \int_value:w '#2 }
540     }
541     \fi:
542 }

```

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

`\regex_class_d:` Similar to `\regex_class_single:N`, adding the appropriate ranges of characters to the class. The token lists are not expanded because it is more memory efficient, with a tiny overhead on execution.
`\regex_class_D:` Similar to `\regex_class_d`, but the token lists are expanded.
`\regex_class_h:`
`\regex_class_H:` 543 \tl_map_inline:nn { dDhHsSvVwWN }
`\regex_class_s:` 544 {
`\regex_class_S:` 545 \cs_new_protected_nopar:cp { regex_class_#1: }
`\regex_class_v:` 546 {
`\regex_class_V:` 547 \tl_put_right:Nn \exp_not:N \l_regex_class_tl
`\regex_class_w:` 548 { \exp_not:c { c_regex_#1_tl } }
`\regex_class_W:` 549 }
550 }

(End definition for `\regex_class_d:` and `\regex_class_D:`. These functions are documented on page ??.)

2.3.6 Quantifiers

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

```

551 \cs_new_protected_nopar:Npn \regex_build_quantifier:w #1#2
552 {
553     \token_if_eq_meaning:NNTF #1 \regex_build_control:N
554     {
555         \cs_if_exist_use:cF { regex_build_quantifier_#2:w }
556         {
557             \regex_build_quantifier_end:n { }
558             #1 #2
559         }
560     }
561     {
562         \regex_build_quantifier_end:n { }
563         #1 #2
564     }
565 }

```

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

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

\regex_build_quantifier_+::w

```

566 \cs_new_protected_nopar:cpn { regex_build_quantifier_?:w }
567   { \regex_build_quantifier_aux:nnNN { } { ? } }
568 \cs_new_protected_nopar:cpn { regex_build_quantifier_*::w }
569   { \regex_build_quantifier_aux:nnNN { } { * } }
570 \cs_new_protected_nopar:cpn { regex_build_quantifier_+::w }
571   { \regex_build_quantifier_aux:nnNN { } { + } }

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

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

```

572 \cs_new_protected_nopar:Npn \regex_build_quantifier_aux:nnNN #1#2#3#4
573   {
574     \str_if_eq:nnTF { #3 #4 } { \regex_build_control:N ? }
575       { \regex_build_quantifier_end:n { #2 #4 } #1 }
576     {
577       \regex_build_quantifier_end:n { #2 } #1
578       #3 #4
579     }
580   }

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

\regex_build_quantifier_{::w}

Three possible syntaxes: {\langle int \rangle}, {\langle int \rangle,}, or {\langle int \rangle, \langle int \rangle}.

\regex_build_quantifier_lbrace:n
\regex_build_quantifier_lbrace:nw
\regex_build_quantifier_lbrace:nnw

```

581 \cs_new_protected_nopar:cpn { regex_build_quantifier_ \c_lbrace_str :w }
582   { \regex_get_digits:nw { \regex_build_quantifier_lbrace:n } }
583 \cs_new_protected_nopar:Npn \regex_build_quantifier_lbrace:n #1
584   {
585     \tl_if_empty:nTF {#1}
586     {
587       \regex_build_quantifier_end:n { }
588       \exp_after:wN \regex_build_raw:N \c_lbrace_str
589     }
590     { \regex_build_quantifier_lbrace:nw {#1} }
591   }
592 \cs_new_protected_nopar:Npx \regex_build_quantifier_lbrace:nw #1#2#3
593   {
594     \exp_not:N \prg_case_str:nnn { #2 #3 }
595     {
596       { \exp_not:N \regex_build_control:N , }
597       {
598         \exp_not:N \regex_get_digits:nw
599           { \exp_not:N \regex_build_quantifier_lbrace:nnw {#1} }
600       }
601       { \exp_not:N \regex_build_control:N \c_rbrace_str }
602         { \exp_not:N \regex_build_quantifier_end:n {n} {#1} }
603     }
604 }
```

```

604      {
605        \exp_not:N \regex_build_quantifier_end:n { }
606        \exp_not:N \regex_build_raw:N \c_lbrace_str #1#2
607      }
608    }
609  \cs_new_protected_nopar:Npn \regex_build_quantifier_lbrace:nnw #1#2#3
610  {
611    \str_if_eq:xxTF { \exp_not:N #3 } { \c_rbrace_str }
612    {
613      \tl_if_empty:nTF {#2}
614        { \regex_build_quantifier_aux:nnN {#1} {\c_max_int} } {nn}
615        {
616          \int_compare:nNnT {#1} > {#2}
617            { \regex_build_error:n {4} }
618            \regex_build_quantifier_aux:nnN {#1} {#2} } {nn}
619        }
620    }
621    {
622      \regex_build_quantifier_end:n { }
623      \use:x
624      {
625        \exp_not:n { \exp_args:No \tl_map_function:nN }
626        { \c_lbrace_str #1 #2 , }
627        \regex_build_raw:N
628      }
629    }
630  }

```

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

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

```

631  \cs_new_protected_nopar:Npn \regex_build_quantifier_end:n #1
632  { \use:c { regex_build_ \l_regex_one_or_group_tl _ #1 : } }

```

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

2.3.7 Quantifiers for one character or character class

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

```

633  \cs_new_protected_nopar:Npn \regex_build_one_quantifier:
634  {
635    \tl_set:Nx \l_regex_one_or_group_tl { one }
636    \regex_build_quantifier:w
637  }

```

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

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

```

638 \cs_new_protected_nopar:Npn \regex_build_one_:
639   {
640     \regex_build_transition_aux:NN
641       \regex_build_tmp_class:n \l_regex_right_state_int
642   }

```

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

\regex_build_one_? : The two transitions are a costly transition controlled by the character class, and a free

\regex_build_one_?? : transition, both going to a common new state. The only difference between the greedy and lazy operators is the order of transitions.

```

643 \cs_new_protected_nopar:cpn { regex_build_one_?: }
644   {
645     \regex_build_transitions_aux:NNNN
646       \regex_build_tmp_class:n \l_regex_right_state_int
647         \regex_action_free:n \l_regex_right_state_int
648   }
649 \cs_new_protected_nopar:cpn { regex_build_one_???: }
650   {
651     \regex_build_transitions_aux:NNNN
652       \regex_action_free:n \l_regex_right_state_int
653         \regex_build_tmp_class:n \l_regex_right_state_int
654   }

```

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

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

```

655 \cs_new_protected_nopar:cpn { regex_build_one_*: }
656   {
657     \regex_build_transitions_aux:NNNN
658       \regex_build_tmp_class:n \l_regex_left_state_int
659         \regex_action_free:n \l_regex_right_state_int
660   }
661 \cs_new_protected_nopar:cpn { regex_build_one_*?: }
662   {
663     \regex_build_transitions_aux:NNNN
664       \regex_action_free:n \l_regex_right_state_int
665         \regex_build_tmp_class:n \l_regex_left_state_int
666   }

```

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

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

```

667 \cs_new_protected_nopar:cpn { regex_build_one_+: }
668   {
669     \regex_build_one_:

```

```

670   \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
671   \regex_build_transitions_aux:NNNN
672     \regex_action_free:n \l_regex_tmpa_int
673     \regex_action_free:n \l_regex_right_state_int
674   }
675 \cs_new_protected_nopar:cpn { regex_build_one_+?: }
676 {
677   \regex_build_one_:
678   \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
679   \regex_build_transitions_aux:NNNN
680     \regex_action_free:n \l_regex_right_state_int
681     \regex_action_free:n \l_regex_tmpa_int
682 }
(End definition for \regex_build_one_+?: This function is documented on page ??.)
```

\regex_build_one_n: This function is called in case the syntax is $\{\langle int \rangle\}$. Greedy and lazy operators are identical, since the number of repetitions is fixed.

```

683 \cs_new_protected_nopar:Npn \regex_build_one_n: #1
684 {
685   \int_set_eq:NN \l_regex_tmpa_int \l_regex_right_state_int
686   \regex_build_new_state:
687   \regex_build_new_state:
688   \regex_toks_put_right:Nx \l_regex_tmpa_int
689   { %^A safe \if nesting?
690     \exp_not:N \if_num:w #1 > \l_regex_repetition_int
691       \regex_action_repeat_move:n
692         { \int_use:N \l_regex_left_state_int }
693     \exp_not:N \else:
694       \regex_action_no_repeat_move:n
695         { \int_use:N \l_regex_right_state_int }
696     \exp_not:N \fi:
697   }
698   \regex_toks_put_right:Nx \l_regex_left_state_int
699   {
700     \regex_build_tmp_class:n
701       { \int_use:N \l_regex_tmpa_int }
702   }
703 }
704 \cs_new_eq:cN { regex_build_one_n?: } \regex_build_one_n:
(End definition for \regex_build_one_n: This function is documented on page ??.)
```

\regex_build_one_nn: This function is called when the syntax is either $\{\langle int \rangle, \}$ or $\{\langle int \rangle, \langle int \rangle\}$.

\regex_build_one_nn?:

```

705 \cs_new_protected_nopar:Npn \regex_build_one_nn: #1#2
706 {
707   \regex_build_one_nn_aux:nn {#1}
708   {
709     \exp_not:N \if_num:w #2 > \l_regex_repetition_int
710       \regex_action_repeat_copy:n
711         { \int_use:N \l_regex_left_state_int }
```

```

712     \exp_not:N \fi:
713     \regex_action_no_repeat_move:n
714     { \int_use:N \l_regex_right_state_int }
715   }
716 }
717 \cs_new_protected_nopar:cpn { regex_build_one_nn?: } #1#2
718 {
719   \regex_build_one_nn_aux:nn {#1}
720   {
721     \regex_action_no_repeat_copy:n
722     { \int_use:N \l_regex_right_state_int }
723     \exp_not:N \if_num:w #2 > \l_regex_repetition_int
724     \regex_action_repeat_move:n
725     { \int_use:N \l_regex_left_state_int }
726     \exp_not:N \fi:
727   }
728 }
729 \cs_new_protected_nopar:Npn \regex_build_one_nn_aux:nn #1#2
730 {
731   \int_set_eq:NN \l_regex_tmpa_int \l_regex_right_state_int
732   \regex_build_new_state:
733   \regex_build_new_state:
734   \regex_toks_put_right:Nx \l_regex_tmpa_int
735   {
736     \exp_not:N \if_num:w #1 > \l_regex_repetition_int
737     \regex_action_repeat_move:n
738     { \int_use:N \l_regex_left_state_int }
739     \exp_not:N \else:
740     #2
741     \exp_not:N \fi:
742   }
743   \regex_toks_put_right:Nx \l_regex_left_state_int
744   {
745     \regex_build_tmp_class:n
746     { \int_use:N \l_regex_tmpa_int }
747   }
748 }

```

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

2.3.8 Groups and alternation

We support the syntax `(<expr1>...<exprn>)<quantifier>` for alternations.

- | | |
|---|---|
| <code>\regex_build_(:</code>
<code>\regex_build_):</code>
<code>\regex_build_open_aux:</code>
<code>\regex_build_ :</code>
<code>\regex_build_begin_alternation:</code>
<code>\regex_build_end_alternation:</code> | Grouping and alternation go together.
<ul style="list-style-type: none"> • Allocate the next available number for the end vertex of the alternation/group and store it on a stack (so that nested alternations work). • Put free transitions to separate all cases of the alternation. |
|---|---|

- Build each branch separately, and merge them to the common end-node.
- Test for a quantifier, and if needed, transfer the initial vertex to a new vertex.

```

749 \cs_new_protected_nopar:cpn { regex_build_(:) #1#2
750   {
751     \str_if_eq:nnTF { #1 #2 } { \regex_build_control:N ? }
752     { \regex_build_special_group:NN }
753     {
754       \int_incr:N \l_regex_capturing_group_int
755       \seq_push:Nx \l_regex_capturing_group_seq
756       { \int_use:N \l_regex_capturing_group_int }
757       \regex_build_open_aux:
758       #1 #2
759     }
760   }
761 \cs_new_protected_nopar:Npn \regex_build_open_aux:
762   {
763     \regex_build_new_state:
764     \seq_push:Nx \l_regex_left_state_seq
765     { \int_use:N \l_regex_left_state_int }
766     \seq_push:Nx \l_regex_right_state_seq
767     { \int_use:N \l_regex_right_state_int }
768     \regex_build_begin_alternation:
769   }
770 \cs_new_protected_nopar:cpn { regex_build_|: }
771   {
772     \regex_build_end_alternation:
773     \regex_build_begin_alternation:
774   }
775 \cs_new_protected_nopar:cpn { regex_build_): }
776   {
777     \seq_if_empty:NTF \l_regex_capturing_group_seq
778     { \regex_build_error:n {22} }
779     {
780       \regex_build_close_aux:
781       \regex_build_group_quantifier:
782     }
783   }
784 \cs_new_protected_nopar:Npn \regex_build_close_aux:
785   {
786     \regex_build_end_alternation:
787     \seq_pop:NN \l_regex_left_state_seq \l_regex_tmpa_tl
788     \int_set:Nn \l_regex_left_state_int \l_regex_tmpa_tl
789     \seq_pop:NN \l_regex_right_state_seq \l_regex_tmpa_tl
790     \int_set:Nn \l_regex_right_state_int \l_regex_tmpa_tl
791   }

```

Building each branch.

```

792 \cs_new_protected_nopar:Npn \regex_build_begin_alternation:
793   {

```

```

794     \regex_build_new_state:
795     \seq_get:NN \l_regex_left_state_seq \l_regex_tma_tl
796     \int_set:Nn \l_regex_left_state_int \l_regex_tma_tl
797     \regex_toks_put_right:Nx \l_regex_left_state_int
798         { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
799     }
800 \cs_new_protected_nopar:Npn \regex_build_end_alternation:
801 {
802     \seq_get:NN \l_regex_right_state_seq \l_regex_tma_tl
803     \regex_toks_put_right:Nx \l_regex_right_state_int
804         { \regex_action_free:n { \l_regex_tma_tl } }
805 }
(End definition for \regex_build_(: and \regex_build_):. These functions are documented on page ??.)

```

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

```

806 \cs_new_protected_nopar:Npn \regex_build_special_group:NN #1#2
807 {
808     \cs_if_exist_use:cF { regex_build_special_group_\token_to_str:N #2 : }
809     {
810         \regex_build_error:n { -998 }
811         \regex_build_control:N ( % )
812         \regex_build_raw:N ?
813         #1 #2
814     }
815 }
(End definition for \regex_build_special_group:NN. This function is documented on page ??.)

```

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

```

816 \cs_new_protected_nopar:cpn { regex_build_special_group:: }
817 {
818     \seq_push:Nx \l_regex_capturing_group_seq { * }
819     \regex_build_open_aux:
820 }
(End definition for \regex_build_special_group::: This function is documented on page ??.)

```

2.3.9 Quantifiers for groups

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

```

821 \cs_new_protected_nopar:Npn \regex_build_group_quantifier:
822 {
823     \tl_set:Nn \l_regex_one_or_group_tl { group }
824     \regex_build_quantifier:w
825 }

```

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

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

```
826 \cs_new_protected_nopar:Npn \regex_build_group_submatches:NN #1#2
827 {
828     \seq_pop:NN \l_regex_capturing_group_seq \l_regex_tmpa_tl
829     \str_if_eq:xxF { \l_regex_tmpa_tl } { * }
830     {
831         \regex_toks_put_left:Nx #1
832         { \regex_action_submatch:n { \l_regex_tmpa_tl < } }
833         \regex_toks_put_left:Nx #2
834         { \regex_action_submatch:n { \l_regex_tmpa_tl > } }
835     }
836 }
```

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

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

```
837 \cs_new_protected_nopar:Npn \regex_build_group_:
838 {
839     \regex_build_group_submatches:NN
840     \l_regex_left_state_int \l_regex_right_state_int
841     \regex_build_transition_aux:NN
842     \regex_action_free:n \l_regex_right_state_int
843 }
```

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

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

```
844 \cs_new_protected_nopar:Npn \regex_build_group_shift:N #1
845 {
846     \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
847     \regex_build_new_state:
848     \tex_toks:D \l_regex_right_state_int = \tex_toks:D \l_regex_tmpa_int
849     \use:x
850     {
851         \tex_toks:D \l_regex_tmpa_int
852         { \s_regex_toks #1 { \int_use:N \l_regex_right_state_int } }
853     }
854     \regex_build_group_submatches:NN
855     \l_regex_right_state_int \l_regex_left_state_int
856     \int_set_eq:NN \l_regex_right_state_int \l_regex_left_state_int
857 }
```

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

\regex_build_group_qs_aux:NN Shift the state at which the group begins using \regex_build_group_shift:N, then add two transitions. The first transition is taken once the group has been traversed: in the case of ? and ??, we should exit by going to \l_regex_right_state_int, while for * \regex_build_group_*: \regex_build_group_*?:

and `*?` we loop by going to `\l_regex_tmpa_int`. The second transition corresponds to skipping the group; it has lower priority (`put_right`) for greedy operators, and higher priority (`put_left`) for lazy operators.

```

858 \cs_new_protected_nopar:Npn \regex_build_group_qs_aux:NN #1#2
859   {
860     \regex_build_group_shift:N \regex_action_free:n
861     \regex_build_transition_aux:NN \regex_action_free:n #1
862     #2 \l_regex_tmpa_int
863     { \regex_action_free:n { \int_use:N \l_regex_right_state_int } }
864   }
865 \cs_new_protected_nopar:cpn { regex_build_group_?: }
866   {
867     \regex_build_group_qs_aux:NN
868     \l_regex_right_state_int \regex_toks_put_right:Nx
869   }
870 \cs_new_protected_nopar:cpn { regex_build_group_???: }
871   {
872     \regex_build_group_qs_aux:NN
873     \l_regex_right_state_int \regex_toks_put_left:Nx
874   }
875 \cs_new_protected_nopar:cpn { regex_build_group_*: }
876   {
877     \regex_build_group_qs_aux:NN
878     \l_regex_tmpa_int \regex_toks_put_right:Nx
879   }
880 \cs_new_protected_nopar:cpn { regex_build_group_*?: }
881   {
882     \regex_build_group_qs_aux:NN
883     \l_regex_tmpa_int \regex_toks_put_left:Nx
884   }

```

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

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

```

885 \cs_new_protected_nopar:cpn { regex_build_group_+: }
886   {
887     \regex_build_group_submatches:NN
888     \l_regex_left_state_int \l_regex_right_state_int
889     \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int
890     \regex_build_transitions_aux:NNNN
891     \regex_action_free:n \l_regex_tmpa_int
892     \regex_action_free:n \l_regex_right_state_int
893   }
894 \cs_new_protected_nopar:cpn { regex_build_group_+?: }
895   {
896     \regex_build_group_submatches:NN
897     \l_regex_left_state_int \l_regex_right_state_int
898     \int_set_eq:NN \l_regex_tmpa_int \l_regex_left_state_int

```

```

899   \regex_build_transitions_aux:NNNN
900     \regex_action_free:n \l_regex_right_state_int
901     \regex_action_free:n \l_regex_tmpa_int
902   }
(End definition for \regex_build_group_+:. This function is documented on page ??.)
```

\regex_build_group_n: These functions are called in case the syntax is $\{\langle int \rangle\}$. Greedy and lazy operators are identical, since the number of repetitions is fixed.

```

903 \cs_new_protected_nopar:Npn \regex_build_group_n: #1
904   {
905     \regex_build_group_shift:N \regex_action_repeat_move:n
906     \regex_build_transition_aux:NN
907       \regex_action_free:n \l_regex_tmpa_int
908     \use:x
909     {
910       \tex_toks:D \l_regex_tmpa_int
911       {
912         \s_regex_toks
913         \exp_not:N \if_num:w #1 > \l_regex_repetition_int
914           \tex_the:D \tex_toks:D \l_regex_tmpa_int
915         \exp_not:N \else:
916           \regex_action_no_repeat_move:n
917             { \int_use:N \l_regex_right_state_int }
918           \exp_not:N \fi:
919       }
920     }
921   }
922 \cs_new_eq:cN { regex_build_group_n?: } \regex_build_group_n:
(End definition for \regex_build_group_n:. This function is documented on page ??.)
```

\regex_build_group_nn: These functions are called when the syntax is either $\{\langle int \rangle, \}$ or $\{\langle int \rangle, \langle int \rangle\}$.

```

923 \cs_new_protected_nopar:Npn \regex_build_group_nn: #1#2
924   {
925     \regex_build_group_shift:N \regex_action_repeat_move:n
926     \regex_build_transition_aux:NN
927       \regex_action_free:n \l_regex_tmpa_int
928     \use:x
929     {
930       \tex_toks:D \l_regex_tmpa_int
931       {
932         \s_regex_toks
933         \exp_not:N \if_num:w #1 > \l_regex_repetition_int
934           \tex_the:D \tex_toks:D \l_regex_tmpa_int
935         \exp_not:N \else:
936           \exp_not:N \if_num:w #2 > \l_regex_repetition_int
937             \tex_the:D \tex_toks:D \l_regex_tmpa_int
938           \exp_not:N \fi:
939           \regex_action_no_repeat_copy:n
940             { \int_use:N \l_regex_right_state_int }
```

```

941           \exp_not:N \fi:
942       }
943   }
944 }
945 \cs_new_protected_nopar:cpn { regex_build_group_nn?: } #1#2
946 {
947   \regex_build_group_shift:N \regex_action_repeat_move:n
948   \regex_build_transition_aux:NN
949   \regex_action_free:n \l_regex_tmpa_int
950   \use:x
951   {
952     \tex_toks:D \l_regex_tmpa_int
953     {
954       \s_regex_toks
955       \exp_not:N \if_num:w #1 > \l_regex_repetition_int
956         \tex_the:D \tex_toks:D \l_regex_tmpa_int
957       \exp_not:N \else:
958         \regex_action_no_repeat_copy:n
959         { \int_use:N \l_regex_right_state_int }
960       \exp_not:N \if_num:w #2 > \l_regex_repetition_int
961         \tex_the:D \tex_toks:D \l_regex_tmpa_int
962       \exp_not:N \fi:
963     \exp_not:N \fi:
964   }
965 }
966

```

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

2.4 Matching

2.4.1 Use of `\TeX` registers when matching

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

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

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

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

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

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

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

2.4.2 Helpers for running the NFA

`\regex_if_state_free:nT` A state is free if it is not marker as taken, namely if the corresponding `\dimen` register is `0sp` rather than `1sp`.

```

967 \cs_new_protected:Npn \regex_if_state_free:nT #1
968 {
969     \if_num:w \tex_dimen:D #1 = \l_regex_unique_step_int
970         \exp_after:wN \use_none:n
971     \else:
972         \exp_after:wN \use:n
973     \fi:
974 }
```

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

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

```

975 \cs_new_protected:Npn \regex_store_state:n #1
976 {
977     \int_incr:N \l_regex_max_index_int
```

```

978     \tex_skip:D \l_regex_max_index_int #1 sp \scan_stop:
979     \regex_store_submatches:n {#1}
980 }

```

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

`\regex_state_use:` Use a given program instruction, unless it has already been executed at this step. The `\toks` registers begin with some submatch information, ignored by `\regex_state_use:`, but not by `\regex_state_use_submatches:`.

```

981 \cs_new_protected_nopar:Npn \regex_state_use_submatches:
982   { \regex_state_use_aux:n { } }
983 \cs_new_protected_nopar:Npn \regex_state_use:
984   { \regex_state_use_aux:n { \exp_after:wN \regex_state_use_aux_ii:w } }
985 \cs_new:Npn \regex_state_use_aux_ii:w #1 \s_regex_toks { }
986 \cs_new_protected:Npn \regex_state_use_aux:n #1
987   {
988     \regex_if_state_free:nT { \l_regex_current_state_int }
989   {
990     \tex_dimen:D \l_regex_current_state_int
991     = \l_regex_unique_step_int sp \scan_stop:
992     #1 \tex_the:D \tex_toks:D \l_regex_current_state_int \scan_stop:
993   }
994 }

```

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

2.4.3 Submatch tracking when running the NFA

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

```

995 \cs_new_protected_nopar:Npn \regex_disable_submatches:
996   {
997     \cs_set_eq:NN \regex_state_use_submatches: \regex_state_use:
998     \cs_set_eq:NN \regex_store_submatches:n
999     \regex_protected_use_none:n
1000    \cs_set_eq:NN \regex_action_submatches:n
1001    \regex_protected_use_none:n
1002  }
1003 \cs_new_protected:Npn \regex_protected_use_none:n #1 { }

```

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

`\regex_store_submatches:n` The submatch information pertaining to one given thread is moved from state to state as we execute the NFA.

```

1004 \cs_new_protected:Npn \regex_store_submatches:n #1
1005   {
1006     \exp_after:wN \regex_store_submatches_aux:Nw
1007     \exp_after:wN \l_regex_current_submatches_prop
1008     \tex_the:D \tex_toks:D #1 \scan_stop: \q_stop
1009   { #1 }

```

```

1010    }
1011 \cs_new_protected:Npn \regex_store_submatches_aux:Nw #1 #2 \s_regex_toks
1012 {
1013     \exp_args:No \regex_store_submatches_aux_ii:wnnwwn { #1 }
1014     #2
1015     \regex_state_submatches:nn { \c_minus_one } { \q_prop }
1016     \s_regex_toks
1017 }
1018 \cs_new_protected:Npn \regex_store_submatches_aux_ii:wnnwwn
1019     #1 \regex_state_submatches:nn #2#3 #4 \s_regex_toks #5 \q_stop #6
1020 {
1021     \tex_toks:D #6 \exp_after:wN
1022     {
1023         \exp_after:wN \regex_state_submatches:nn \exp_after:wN
1024         { \int_value:w \int_eval:w \l_regex_unique_step_int + \c_one }
1025         { #1 }
1026         \regex_state_submatches:nn {#2} {#3}
1027         \s_regex_toks
1028         #5
1029     }
1030 }

```

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

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

```

1031 \cs_new_protected:Npn \regex_state_submatches:nn #1#2
1032 {
1033     \if_num:w #1 = \l_regex_unique_step_int
1034         \tl_set:Nn \l_regex_current_submatches_prop {#2}
1035     \fi:
1036 }

```

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

2.4.4 Matching: framework

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

```

1037 \cs_new_protected:Npn \regex_match:n #1
1038 {
1039     \tl_set:Nx \l_regex_query_str { \tl_to_other_str:n {#1} }
1040     \regex_match_initial_setup:
1041     \regex_match_once:
1042 }

```

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

`\regex_match_once:` After setting up more variables in `\regex_match_setup:`, skip the `\l_regex_start_step_int` first characters of the query string, and loop over it. If there was a match,

use the token list `\l_regex_every_match_tl`, which may call `\regex_match_once:` to achieve multiple matches.

```

1043 \cs_new_protected_nopar:Npn \regex_match_once:
1044   {
1045     \regex_match_setup:
1046     \exp_after:wN \regex_match_once_aux: \l_regex_query_str
1047       \q_recursion_tail \q_recursion_stop
1048       \bool_if:NT \l_regex_success_bool { \l_regex_every_match_tl }
1049   }
1050 \cs_new_protected_nopar:Npn \regex_match_once_aux:
1051   { \str_skip_do:nn { \l_regex_start_step_int } { \regex_match_loop:N } }
(End definition for \regex_match_once:. This function is documented on page ??.)
```

`\regex_match_initial_setup:` This function holds the setup that should be done only once for one given pattern matching on a given string. It is called only once for the whole string. On the other hand, `\regex_match_setup:` is called for every match in the string in case of repeated matches, and `\regex_match_loop_setup:` is called at every step.

```

1052 \cs_new_protected_nopar:Npn \regex_match_initial_setup:
1053   {
1054     \tl_clear:N \l_regex_look_behind_str
1055     \prg_stepwise_inline:nnnn {1} {1} { \l_regex_max_state_int }
1056       { \tex_dimen:D ##1 \c_minus_one sp \scan_stop: }
1057     \int_set_eq:NN \l_regex_unique_step_int \c_minus_one
1058     \int_set_eq:NN \l_regex_start_step_int \c_minus_one
1059     \int_set_eq:NN \l_regex_current_step_int \c_zero
1060     \int_set_eq:NN \l_regex_success_step_int \c_zero
1061     \bool_set_false:N \l_regex_success_empty_bool
1062   }
(End definition for \regex_match_initial_setup:. This function is documented on page ??.)
```

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

```

1063 \cs_new_protected_nopar:Npn \regex_match_setup:
1064   {
1065     \prop_clear:N \l_regex_current_submatches_prop
1066     \bool_if:NTF \l_regex_success_empty_bool
1067       { \cs_set_eq:NN \regex_last_match_empty:F \regex_last_match_empty_yes:F }
1068       { \cs_set_eq:NN \regex_last_match_empty:F \regex_last_match_empty_no:F }
1069     \int_set_eq:NN \l_regex_start_step_int \l_regex_success_step_int
1070     \int_set_eq:NN \l_regex_current_step_int \l_regex_start_step_int
1071     \int_decr:N \l_regex_current_step_int
1072     \bool_set_false:N \l_regex_success_bool
1073     \int_zero:N \l_regex_max_index_int
1074     \regex_store_state:n {1}
1075   }
(End definition for \regex_match_setup:. This function is documented on page ??.)
```

`\regex_match_loop_setup:` This is called for every character in the string.

```
1076 \cs_new_protected_nopar:Npn \regex_match_loop_setup:
```

```

1077    {
1078      \int_incr:N \l_regex_current_step_int
1079      \int_incr:N \l_regex_unique_step_int
1080      \bool_set_false:N \l_regex_fresh_thread_bool
1081    }
(End definition for \regex_match_loop_setup:. This function is documented on page ??.)
```

\s_regex_step When a thread succeeds, all threads with lower precedence can be ignored, and the next character should be read. This is done by skipping to \s_regex_step. This marker does nothing if no thread succeeded at this step.

```

1082 \cs_new_eq:NN \s_regex_step \scan_stop:
1083 \cs_new:Npn \regex_break_step:w #1 \s_regex_step { }
(End definition for \s_regex_step. This function is documented on page ??.)
```

\regex_match_loop:N Setup what needs to be reset at every character, then set \l_regex_current_char_int to the character code of the character that is read (and -1 for the end of the string), and loop over the elements of the \skip array. Then repeat. There are a couple of tests to stop reading the string when no active state is left, or when the end is reached.

```

1084 \cs_new_protected_nopar:Npn \regex_match_loop:N #1
1085  {
1086    \regex_match_loop_setup:
1087    \token_if_eq_meaning:NNTF #1 \q_recursion_tail
1088      { \int_set_eq:NN \l_regex_current_char_int \c_minus_one }
1089      { \int_set:Nn \l_regex_current_char_int {'#1} }
1090    \cs_set_nopar:Npx \regex_tmp:w
1091    {
1092      \int_zero:N \l_regex_max_index_int
1093      \prg_stepwise_function:nnnN
1094        {1} {1} { \l_regex_max_index_int }
1095      \regex_match_one_index:n
1096    }
1097    \regex_tmp:w
1098    \s_regex_step
1099    \if_num:w \l_regex_max_index_int = \c_zero
1100      \exp_after:wN \use_none_delimit_by_q_recursion_stop:w
1101    \fi:
1102    \quark_if_recursion_tail_stop:N #1
1103    \bool_if:NT \l_regex_look_behind_bool
1104      { \tl_put_left:Nx \l_regex_look_behind_str {#1} }
1105    \regex_match_loop:N
1106  }
1107 \cs_new_nopar:Npn \regex_match_one_index:n #1
1108  {
1109    \regex_match_one_index_aux:n { \int_value:w \tex_skip:D #1 }
1110  }
1111 \cs_new_protected_nopar:Npn \regex_match_one_index_aux:n #1
1112  {
1113    \int_set:Nn \l_regex_current_state_int {#1}
1114    \prop_clear:N \l_regex_current_submatches_prop
```

```

1115     \regex_state_use_submatches:
1116 }
(End definition for \regex_match_loop:N. This function is documented on page ??.)
```

2.4.5 Actions when matching

\regex_action_start_wildcard:nn The search is made unanchored at the start by putting a free transition to the real start of the NFA, and a costly transition to the same state, waiting for the next character in the query string. This combination could be reused (with some changes). We sometimes need to know that the match for a given thread starts at this character. For that, we use the boolean \l_regex_fresh_thread_bool.

```

1117 \cs_new_protected_nopar:Npn \regex_action_start_wildcard:nn #1#2
1118 {
1119     \bool_set_true:N \l_regex_fresh_thread_bool
1120     \regex_action_free:n {#2}
1121     \bool_set_false:N \l_regex_fresh_thread_bool
1122     \regex_action_cost:n {#1}
1123 }
(End definition for \regex_action_start_wildcard:nn. This function is documented on page ??.)
```

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

```

1124 \cs_new_protected_nopar:Npn \regex_action_cost:n #1
1125 {
    \regex_store_state:n {#1}
(End definition for \regex_action_cost:n. This function is documented on page ??.)
```

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

```

1126 \cs_new_protected_nopar:Npn \regex_action_success:
1127 {
1128     \regex_last_match_empty:F
1129     {
1130         \bool_set_true:N \l_regex_success_bool
1131         \bool_set_eq:NN \l_regex_success_empty_bool
1132             \l_regex_fresh_thread_bool
1133             \int_set_eq:NN \l_regex_success_step_int
1134                 \l_regex_current_step_int
1135                 \prop_set_eq:NN \l_regex_success_submatches_prop
1136                     \l_regex_current_submatches_prop
1137                     \regex_break_step:w
1138     }
1139 }
```

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

\regex_action_free:n To copy a thread, check whether the program state has already been used at this character. If not, store submatches in the new state, and insert the instructions for that state in the input stream. Then restore the old value of \l_regex_current_state_int and of the current submatches.

```

1140 \cs_new_protected_nopar:Npn \regex_action_free:n #1
1141 {
1142     \regex_if_state_free:nT {#1}
1143     {
1144         \use:x
1145         {
1146             \int_set:Nn \l_regex_current_state_int {#1}
1147             \regex_store_submatches:n { \l_regex_current_state_int }
1148             \regex_state_use:
1149             \int_set:Nn \l_regex_current_state_int
1150                 { \int_use:N \l_regex_current_state_int }
1151             \tl_set:Nn \exp_not:N \l_regex_current_submatches_prop
1152                 { \exp_not:o \l_regex_current_submatches_prop }
1153         }
1154     }
1155 }
```

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

\regex_action_submatch:n

```

1156 \cs_new_protected_nopar:Npn \regex_action_submatch:n #1
1157 {
1158     \prop_put:Nno \l_regex_current_submatches_prop {#1}
1159         { \int_use:N \l_regex_current_step_int }
1160     \regex_store_submatches:n { \l_regex_current_state_int }
1161 }
```

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

```

1162 \cs_new_protected_nopar:Npn \regex_action_repeat_copy:n { \ERROR }
1163 \cs_new_protected_nopar:Npn \regex_action_repeat_move:n { \ERROR }
1164 \cs_new_protected_nopar:Npn \regex_action_no_repeat_copy:n { \ERROR }
1165 \cs_new_protected_nopar:Npn \regex_action_no_repeat_move:n { \ERROR }
```

2.5 Submatches, once the correct match is found

```

\regex_extract:
\regex_extract_aux:nTF
1166 \cs_new_protected_nopar:Npn \regex_extract:
1167 {
1168     \seq_gclear:N \g_regex_submatches_seq
1169     \prg_stepwise_inline:nnnn
1170         {0} {1} { \l_regex_capturing_group_int }
1171         {
1172             \regex_extract_aux:nTF { ##1 }
1173             {
1174                 \seq_gput_right:Nx \g_regex_submatches_seq
1175                 {
```

```

1176           \str_from_to:Nnn \l_regex_query_str
1177           { \l_regex_tmpa_tl }
1178           { \l_regex_tmpb_tl }
1179       }
1180   }
1181   { \seq_gput_right:Nn \g_regex_submatches_seq { } }
1182 }
1183 }
1184 \cs_new_protected_nopar:Npn \regex_extract_aux:nTF #1#2#3
1185 {
1186   \prop_get:NnNTF \l_regex_success_submatches_prop
1187   { #1 < } \l_regex_tmpa_tl
1188   {
1189     \prop_get:NnNTF \l_regex_success_submatches_prop
1190     { #1 > } \l_regex_tmpb_tl
1191     {#2}
1192     {#3}
1193   }
1194   {#3}
1195 }
(End definition for \regex_extract:. This function is documented on page ??.)
```

2.6 User commands

2.6.1 Precompiled pattern

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

```

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

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

<code>\regex_set:Nn</code> <code>\regex_gset:Nn</code> <code>\regex_set_aux:NNn</code> <code>\regex_set_aux:n</code> <code> \regex_nfa:Nw</code> <code>\regex_nfa_end:</code>	Within a group, build the NFA corresponding to the given regular expression, with submatch tracking. Then save the contents of all relevant <code>\toks</code> registers into <code>\g_regex_tmpa_tl</code> , then transferred to the user's <code>tl</code> variable. The auxiliary <code>\regex_nfa:Nw</code> is not protected: this ensures that the NFA will properly be replaced by an error message in expansion contexts.
---	--

```
1196 \cs_new_protected_nopar:Npn \regex_set:Nn
```

```

1197 { \regex_set_aux:NNn \tl_set_eq:NN }
1198 \cs_new_protected_nopar:Npn \regex_gset:Nn
1199 { \regex_set_aux:NNn \tl_gset_eq:NN }
1200 \cs_new_protected:Npn \regex_set_aux:NNn #1#2#3
1201 {
1202   \group_begin:
1203     \regex_build:n {#3}
1204     \tl_gset:Nx \g_regex_tmpa_tl
1205     {
1206       \exp_not:N \regex_nfa:Nw \exp_not:N #2
1207       \l_regex_max_state_int
1208       = \int_use:N \l_regex_max_state_int
1209       \l_regex_capturing_group_int
1210       = \int_use:N \l_regex_capturing_group_int
1211       \bool_if:NTF \l_regex_look_behind_bool
1212       { \bool_set_true:N \l_regex_look_behind_bool }
1213       { \bool_set_false:N \l_regex_look_behind_bool }
1214       \prg_stepwise_function:nnnN
1215       {1} {1} {\l_regex_max_state_int}
1216       \regex_set_aux:n
1217       \regex_nfa_end:
1218     }
1219   \group_end:
1220   #1 #2 \g_regex_tmpa_tl
1221 }
1222 \cs_new_nopar:Npn \regex_set_aux:n #1
1223 { \tex_toks:D #1 { \tex_the:D \tex_toks:D #1 } }
1224 \cs_new:Npn \regex_nfa:Nw #1 #2 \regex_nfa_end:
1225 { \msg_expandable_error:n { Automaton~#1 used~incorrectly. } }
1226 \cs_new_eq:NN \regex_nfa_end: \scan_stop:
(End definition for \regex_set:Nn. This function is documented on page ??.)
```

\regex_check_nfa:N If a token list variable starts with \regex_nfa:Nw, then it most likely holds the data for a precompiled pattern.

```

1227 \prg_new_protected_conditional:Npnn \regex_check_nfa:N #1 { TF }
1228 { \exp_after:wN \regex_check_nfa_aux:Nw #1 \q_stop }
1229 \cs_new:Npn \regex_check_nfa_aux:Nw #1 #2 \q_stop
1230 {
1231   \if_meaning:w \regex_nfa:Nw #1
1232     \prg_return_true:
1233   \else:
1234     \msg_error:nnx { regex } { not-nfa } { \token_to_str:N #1 }
1235     \prg_return_false:
1236   \fi:
1237 }
1238 \msg_new:nnn { regex } { not-nfa }
1239 {
1240   I~was~expecting~a~regular~expression~variable.\\
1241   Instead,~I~got~#1.
1242 }
```

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

\regex_use:N No error-checking.

```
1243 \cs_new_protected_nopar:Npn \regex_use:N #1
1244   { \exp_after:wN \use_none:nn #1 }
```

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

2.6.2 Generic auxiliary functions

\regex_user_aux:n This is an auxiliary used by most user functions. The first and second arguments control whether we should track submatches, and whether we track one or multiple submatches. Everything is done within a group, so that #1 can perform “unsafe” assignments. Most user functions return a result using \group_insert_after:N.

```
1245 \cs_new_protected:Npn \regex_user_aux:n #1
1246   {
1247     \group_begin:
1248       \tl_clear:N \l_regex_every_match_tl
1249       #1
1250     \group_end:
1251   }
```

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

\regex_return_after_group: Most of l3regex’s work is done within a group. This function triggers either \prg_return_false: or \prg_return_true: as appropriate to whether a match was found or not.

```
1252 \cs_new_protected_nopar:Npn \regex_return_after_group:
1253   {
1254     \if_bool:N \l_regex_success_bool
1255       \group_insert_after:N \prg_return_true:
1256     \else:
1257       \group_insert_after:N \prg_return_false:
1258     \fi:
1259   }
```

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

\regex_extract_after_group:N Extract submatches, and store them in the user-given variable after the group has ended.

```
1260 \cs_new_protected_nopar:Npn \regex_extract_after_group:N #1
1261   {
1262     \if_bool:N \l_regex_success_bool
1263       \regex_extract:
1264       \group_insert_after:N \seq_set_eq:NN
1265       \group_insert_after:N #1
1266       \group_insert_after:N \g_regex_submatches_seq
1267     \fi:
1268   }
```

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

\regex_count_after_group:N Same procedure as \regex_extract_after_group:N, but simpler since the match counting has already taken place.

```

1269 \cs_new_protected_nopar:Npn \regex_count_after_group:N #1
1270   {
1271     \group_insert_after:N \int_set_eq:NN
1272     \group_insert_after:N #1
1273     \group_insert_after:N \g_regex_match_count_int
1274   }
(End definition for \regex_count_after_group:N. This function is documented on page ??.)
```

2.6.3 Matching

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

```

1275 \prg_new_protected_conditional:Npnn \regex_match:nn #1#2 { T , F , TF }
1276   {
1277     \regex_user_aux:n
1278     {
1279       \regex_disable_submatches:
1280       \regex_build:n {#1}
1281       \regex_match:n {#2}
1282       \regex_return_after_group:
1283     }
1284   }
1285 \prg_new_protected_conditional:Npnn \regex_match:Nn #1#2 { T , F , TF }
1286   {
1287     \regex_check_nfa:NTF #1
1288     {
1289       \regex_user_aux:n
1290       {
1291         \regex_disable_submatches:
1292         \regex_use:N #1
1293         \regex_match:n {#2}
1294         \regex_return_after_group:
1295       }
1296     }
1297   { \prg_return_false: }
1298 }
```

(End definition for \regex_match:nn. This function is documented on page 5.)

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

```

1299 \cs_new_protected:Npn \regex_count:nnN #1#2#3
1300   {
1301     \regex_user_aux:n
1302     {
1303       \regex_disable_submatches:
```

```

1304     \int_gzero:N \g_regex_match_count_int
1305     \tl_set:Nn \l_regex_every_match_tl
1306     {
1307         \int_gincr:N \g_regex_match_count_int
1308         \regex_match_once:
1309     }
1310     \regex_build:n {#1}
1311     \regex_match:n {#2}
1312     \regex_count_after_group:N #3
1313 }
1314 }
1315 \cs_new_protected:Npn \regex_count:NnN #1#2#3
1316 {
1317     \regex_check_nfa:NTF #1
1318     {
1319         \regex_user_aux:n
1320         {
1321             \regex_disable_submatches:
1322             \int_gzero:N \g_regex_match_count_int
1323             \tl_set:Nn \l_regex_every_match_tl
1324             {
1325                 \int_gincr:N \g_regex_match_count_int
1326                 \regex_match_once:
1327             }
1328             \regex_use:N #1
1329             \regex_match:n {#2}
1330             \regex_count_after_group:N #3
1331         }
1332     }
1333     { }
1334 }
```

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

2.6.4 Submatch extraction

```

\regex_extract:nnN
\regex_extract:NnN
\regex_extract:nnN
\regex_extract:NnN
1335 \cs_new_protected:Npn \regex_extract:nnN #1#2#3
1336 {
1337     \regex_user_aux:n
1338     {
1339         \regex_build:n {#1}
1340         \regex_match:n {#2}
1341         \regex_extract_after_group:N #3
1342     }
1343 }
1344 \prg_new_protected_conditional:Npnn \regex_extract:nnN #1#2#3 { T , F , TF }
1345 {
1346     \regex_user_aux:n
1347     {
```

```

1348         \regex_build:n {#1}
1349         \regex_match:n {#2}
1350         \regex_extract_after_group:N #3
1351         \regex_return_after_group:
1352     }
1353 }
1354 \cs_new_protected:Npn \regex_extract:NnN #1#2#3
1355 {
1356     \regex_check_nfa:NTF #1
1357     {
1358         \regex_user_aux:n
1359         {
1360             \regex_use:N #1
1361             \regex_match:n {#2}
1362             \regex_extract_after_group:N #3
1363         }
1364     }
1365     { }
1366 }
1367 \prg_new_protected_conditional:Npnn \regex_extract:NnN #1#2#3 { T , F , TF }
1368 {
1369     \regex_check_nfa:NTF #1
1370     {
1371         \regex_user_aux:n
1372         {
1373             \regex_use:N #1
1374             \regex_match:n {#2}
1375             \regex_extract_after_group:N #3
1376             \regex_return_after_group:
1377         }
1378     }
1379     { \prg_return_false: }
1380 }

```

(End definition for `\regex_extract:nnN`. This function is documented on page 6.)

2.6.5 Splitting a string by matches of a regex

`\regex_split:nnN` Similarly to `\regex_count:nnN` functions, recurse through matches of the pattern. Then
`\regex_split:NnN` we do something slightly different, extracting submatches. Submatches are not extracted
`\regex_split_aux:` if the pattern matched an empty string at the start of the match attempt (to avoid adding
`\regex_split_after_group:N` spurious empty items to the resulting sequence).

```

1381 \cs_new_protected:Npn \regex_split:nnN #1#2#3
1382 {
1383     \regex_user_aux:n
1384     {
1385         \seq_gclear:N \g_regex_split_seq
1386         \tl_set:Nn \l_regex_every_match_tl { \regex_split_aux: }
1387         \regex_build:n {#1}
1388         \regex_match:n {#2}

```

```

1389         \regex_split_after_group:N #3
1390     }
1391 }
1392 \cs_new_protected:Npn \regex_split:NnN #1#2#3
1393 {
1394     \regex_check_nfa:NTF #1
1395     {
1396         \regex_user_aux:n
1397         {
1398             \seq_gclear:N \g_regex_split_seq
1399             \tl_set:Nn \l_regex_every_match_tl { \regex_split_aux: }
1400             \regex_use:N #1
1401             \regex_match:n {#2}
1402             \regex_split_after_group:N #3
1403         }
1404     }
1405     { }
1406 }
1407 \cs_new_protected_nopar:Npn \regex_split_aux:
1408 {
1409     \int_compare:nNnF \l_regex_start_step_int = \l_regex_success_step_int
1410     {
1411         \regex_extract:
1412         \seq_pop:NN \g_regex_submatches_seq \l_regex_tmpa_tl
1413         \regex_extract_aux:nTF {0}
1414         {
1415             \seq_gput_left:Nx \g_regex_submatches_seq
1416             {
1417                 \str_from_to:Nnn \l_regex_query_str
1418                 { \int_use:N \l_regex_start_step_int }
1419                 { \l_regex_tmpa_tl }
1420             }
1421         }
1422         { \msg_error:nn { regex } { internal } }
1423         \seq_gconcat:NNN \g_regex_split_seq
1424             \g_regex_split_seq \g_regex_submatches_seq
1425         }
1426         \regex_match_once:
1427     }
1428 \cs_new_protected_nopar:Npn \regex_split_after_group:N #1
1429 {
1430     \int_compare:nNnTF \l_regex_start_step_int = \l_regex_current_step_int
1431     {
1432         \bool_if:NF \l_regex_success_empty_bool
1433             { \seq_gput_right:Nn \g_regex_split_seq { } }
1434     }
1435     {
1436         \seq_gput_right:Nx \g_regex_split_seq
1437             {
1438                 \str_from_to:Nnn \l_regex_query_str

```

```

1439         { \int_use:N \l_regex_start_step_int }
1440         { \int_use:N \l_regex_current_step_int }
1441     }
1442     }
1443     \group_insert_after:N \seq_set_eq:NN
1444     \group_insert_after:N #1
1445     \group_insert_after:N \g_regex_split_seq
1446   }
(End definition for \regex_split:nnN. This function is documented on page ??.)
```

2.6.6 String replacement

```

\regex_replacement:n
\regex_replacement_loop:N
\regex_replacement_normal:N
1447 \cs_new_protected:Npn \regex_replacement:n #1
1448   {
1449     \str_aux_escape>NNNn
1450     \prg_do_nothing:
1451     \regex_replacement_escaped:N
1452     \regex_replacement_raw:N
1453     {#1}
1454     \tl_set_eq:NN \l_regex_replacement_tl \g_str_result_tl
1455     \tl_set:Nx \l_regex_replacement_tl
1456     { \l_regex_replacement_tl \prg_do_nothing: }
1457   }
1458 \cs_new_nopar:Npn \regex_replacement_raw:N { \exp_not:N \prg_do_nothing: }
1459 \cs_new_nopar:Npn \regex_replacement_escaped:N #1
1460   {
1461     \if_num:w 9 < 1 #1 \exp_stop_f:
1462     \exp_not:N \regex_replacement_escaped_ii:nN {#1}
1463     \else:
1464     #1
1465     \fi:
1466   }
1467 \cs_new_nopar:Npn \regex_replacement_escaped_ii:nN #1#2
1468   {
1469     \regex_token_if_other_digit:NTF #2
1470     { \regex_replacement_escaped_ii:nN {#1#2} }
1471     { \exp_not:N \seq_item:Nn \exp_not:N \g_regex_submatches_seq {#1} #2 }
1472   }
1473 \prg_new_conditional:Npnn \regex_token_if_other_digit:N #1 { TF }
1474   {
1475     \if_num:w 9 < 1 \exp_not:N #1 \exp_stop_f:
1476     \prg_return_true:
1477     \else:
1478     \prg_return_false:
1479     \fi:
1480   }
(End definition for \regex_replacement:n. This function is documented on page ??.)
```

```

\regex_replace_after_group:N
1481 \cs_new_protected_nopar:Npn \regex_replace_after_group:N #1
1482 {
1483   \group_insert_after:N \tl_set_eq:NN
1484   \group_insert_after:N #1
1485   \group_insert_after:N \g_regex_replaced_str
1486 }
(End definition for \regex_replace_after_group:N. This function is documented on page ??..)

\regex_replace_once:nnN
\regex_replace_once:NnN
\regex_replace_once:nnN
\regex_replace_once:NnN
\regex_replace_once_aux:Nn
1487 \cs_new_protected:Npn \regex_replace_once_aux:Nn #1#2
1488 {
1489   \group_begin:
1490     \regex_replace_after_group:N #1
1491     \tl_clear:N \l_regex_every_match_tl
1492     #2
1493     \exp_args:No \regex_match:n {#1}
1494     \regex_extract:
1495     \regex_extract_aux:nTF {0}
1496     {
1497       \tl_gset:Nx \g_regex_replaced_str
1498       {
1499         \str_from_to:Nnn \l_regex_query_str {0} { \l_regex_tma_tl }
1500         \l_regex_replacement_tl
1501         \str_from_to:Nnn \l_regex_query_str
1502           { \l_regex_tmpb_tl } { \c_max_int }
1503       }
1504     }
1505     { \tl_gset_eq:NN \g_regex_replaced_str \l_regex_query_str }
1506   \group_end:
1507 }
1508 \cs_new_protected:Npn \regex_replace_once:nnN #1#2#3
1509 {
1510   \regex_replace_once_aux:Nn #3
1511   {
1512     \regex_build:n {#1}
1513     \regex_replacement:n {#2}
1514   }
1515 }
1516 \cs_new_protected:Npn \regex_replace_once:NnN #1#2#3
1517 {
1518   \regex_replace_once_aux:Nn #3
1519   {
1520     \regex_use:N #1
1521     \regex_replacement:n {#2}
1522   }
1523 }
1524 \prg_new_protected_conditional:Npnn \regex_replace_once:nnN #1#2#3 {T,F,TF}
1525 {

```

```

1526     \regex_replace_once_aux:Nn #3
1527     {
1528         \regex_build:n {#1}
1529         \regex_replacement:n {#2}
1530         \regex_return_after_group:
1531     }
1532 }
1533 \prg_new_protected_conditional:Npnn \regex_replace_once:NnN #1#2#3 {T,F,TF}
1534 {
1535     \regex_replace_once_aux:Nn #3
1536     {
1537         \regex_use:N #1
1538         \regex_replacement:n {#2}
1539         \regex_return_after_group:
1540     }
1541 }
1542 (End definition for \regex_replace_once:nnN. This function is documented on page ??..)

\regex_replace_all:nnN
\regex_replace_all:NnN
\regex_replace_all_aux:Nn
\regex_replace_all_aux:
1543 {
1544     \group_begin:
1545         \regex_replace_after_group:N #1
1546         \tl_set:Nn \l_regex_every_match_tl { \regex_replace_all_aux: }
1547         \tl_gclear:N \g_regex_replaced_str
1548         #2
1549         \exp_args:No \regex_match:n {#1}
1550         \tl_gput_right:Nx \g_regex_replaced_str
1551         {
1552             \str_from_to:Nnn \l_regex_query_str
1553                 { \int_use:N \l_regex_start_step_int }
1554                 { \c_max_int }
1555         }
1556     \group_end:
1557 }
1558 \cs_new_protected_nopar:Npn \regex_replace_all_aux:
1559 {
1560     \regex_extract:
1561     \regex_extract_aux:nTF {0}
1562     {
1563         \tl_gput_right:Nx \g_regex_replaced_str
1564         {
1565             \str_from_to:Nnn \l_regex_query_str
1566                 { \int_use:N \l_regex_start_step_int }
1567                 { \l_regex_tmpa_tl }
1568             \l_regex_replacement_tl
1569         }
1570     }
1571     { \msg_error:nn { regex } { internal } }
1572 \regex_match_once:

```

```

1573    }
1574 \cs_new_protected:Npn \regex_replace_all:nnN #1#2#3
1575 {
1576   \regex_replace_all_aux:Nn #3
1577   {
1578     \regex_build:n {#1}
1579     \regex_replacement:n {#2}
1580   }
1581 }
1582 \cs_new_protected:Npn \regex_replace_all:NnN #1#2#3
1583 {
1584   \regex_replace_all_aux:Nn #3
1585   {
1586     \regex_use:N #1
1587     \regex_replacement:n {#2}
1588   }
1589 }
(End definition for \regex_replace_all:nnN. This function is documented on page ??.)
1590 </package>

```