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

The I₄T_FX3 Project[†]

Released 2016/05/14

1 **I3regex** documentation

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

Let us give a few examples. After

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

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

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

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

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

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

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

^{*}This file describes v6492, last revised 2016/05/14.

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

1.1 Syntax of regular expressions

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

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

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

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

Characters.

 $x{hh...}$ Character with hex code hh...

- \xhh Character with hex code hh.
 - a Alarm (hex 07).
 - \e Escape (hex 1B).
 - f Form-feed (hex 0C).
 - n New line (hex 0A).
 - r Carriage return (hex 0D).
 - \t Horizontal tab (hex 09).

Character types.

- . A single period matches any token.
- d Any decimal digit.
- h Any horizontal space character, equivalent to $[\ \^I]$: space and tab.
- \s Any space character, equivalent to [\ $^1\M$].

- v Any vertical space character, equivalent to [J]. Note that K is a vertical space, but not a space, for compatibility with Perl.
- w Any word character, *i.e.*, alpha-numerics and underscore, equivalent to [A-Za-z0-9].
- D Any token not matched by d.
- H Any token not matched by h.
- N Any token other than the n character (hex 0A).
- S Any token not matched by s.
- \V Any token not matched by v.
- \W Any token not matched by $\w.$
- Of those, ., \D , \H , \N , \S , \V , and \W will match arbitrary control sequences. Character classes match exactly one token in the subject.
- [...] Positive character class. Matches any of the specified tokens.
- [^...] Negative character class. Matches any token other than the specified characters.
 - x-y Within a character class, this denotes a range (can be used with escaped characters).
- [:(name):] Within a character class (one more set of brackets), this denotes the POSIX character class (name), which can be alnum, alpha, ascii, blank, cntrl, digit, graph, lower, print, punct, space, upper, word, or xdigit.
- [: name :] Negative POSIX character class.

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

Quantifiers (repetition).

- ? 0 or 1, greedy.
- ?? 0 or 1, lazy.
- * 0 or more, greedy.
- *? 0 or more, lazy.
- + 1 or more, greedy.
- +? 1 or more, lazy.
- $\{n\}$ Exactly n.
- $\{n,\}$ n or more, greedy.
- $\{n,\}$? *n* or more, lazy.

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

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

Anchors and simple assertions.

- **\b** Word boundary: either the previous token is matched by \w and the next by \W , or the opposite. For this purpose, the ends of the token list are considered as \W .
- \B Not a word boundary: between two w tokens or two W tokens (including the boundary).
- r or A Start of the subject token list.
- , Zor z End of the subject token list.
 - \G Start of the current match. This is only different from ^ in the case of multiple matches: for instance \regex_count:nnN { \G a } { aaba } \l_tmpa_int yields 2, but replacing \G by ^ would result in \l_tmpa_int holding the value 1.

Alternation and capturing groups.

- A|B|C Either one of A, B, or C.
- (...) Capturing group.
- (?:...) Non-capturing group.
- (?|...) Non-capturing group which resets the group number for capturing groups in each alternative. The following group will be numbered with the first unused group number.

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

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

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

The c escape sequence is used as follows.

- - \cX Applies to the next object, which can be a character, character property, class, or group, and forces this object to only match tokens with category X (any of CBEMTPUDSLOA. For instance, \cL[A-Z\d] matches uppercase letters and digits of category code letter, \cC. matches any control sequence, and \cO(abc) matches abc where each character has category other.
 - \c[XYZ] Applies to the next object, and forces it to only match tokens with category X, Y, or Z (each being any of CBEMTPUDSLOA). For instance, \c[LSO](..) matches two tokens of category letter, space, or other.
 - $c[^XYZ]$ Applies to the next object and prevents it from matching any token with category X, Y, or Z (each being any of CBEMTPUDSLOA). For instance, $c[^0] d$ matches digits which have any category different from other.

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

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

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

In character classes, only [, ^, -,], $\$ and spaces are special, and should be escaped. Other non-alphanumeric characters can still be escaped without harm. Any escape sequence which matches a single character ($\d, \D, etc.$) is supported in character classes. If the first character is ^, then the meaning of the character class is inverted; ^ appearing anywhere else in the range is not special. If the first character (possibly following a leading ^) is] then it does not need to be escaped since ending the range there would make it empty. Ranges of characters can be expressed using -, for instance, [\D 0-5] and [^6-9] are equivalent.

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

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

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

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

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

results in 1_foo_seq containing the items {c3} and {bc}: the true match is {acbc3}, with first submatch {bc}, but K resets the beginning of the match to the last position where it appears.

1.2 Syntax of the replacement text

Most of the features described in regular expressions do not make sense within the replacement text. Escaped characters are supported as inside regular expressions. The whole match is accessed as $\0$, and the first 9 submatches are accessed as $1, \ldots, 9$. Further submatches are accessed through $g\{\langle number \rangle\}$ where $\langle number \rangle$ is any non-negative integer. If there are fewer than $\langle number \rangle$ capturing groups, the submatch is empty.

For instance,

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

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

Submatches keep the same category codes as in the original token list. The characters inserted by the replacement have category code 12 (other) by default, with the exception of space characters. Spaces inserted through \u have category code 10, while spaces inserted through $\x20$ or $\x{20}$ have category code 12. The escape sequence \c allows to insert characters with arbitrary category codes, as well as control sequences.

- \XY Produces the character Y (which can be given as an escape sequence such as t for tab, or (or) for a parenthesis) with category code X, which must be one of CBEMTPUDSLOA.
- $c{\langle text \rangle}$ Produces the control sequence with csname $\langle text \rangle$. The $\langle text \rangle$ may contain references to the submatches 0, 1, etc.

The escape sequence $u{\langle tl var name \rangle}$ allows to insert the contents of the token list with name $\langle tl var name \rangle$ directly into the replacement, avoiding the need to escape special characters. Within the construction $c{\langle text \rangle}$, the u escape sequence only expands its argument once, in effect performing $tl_to_str:v$. Submatches can be used within the argument of u. For instance,

```
\tl_set:Nn \l_my_one_tl { first }
\tl_set:Nn \l_my_two_tl { \emph{second} }
\tl_set:Nn \l_my_tl { one , two , one , one }
\regex_replace_all:nnN { [^,]+ } { \u{l_my_\0_tl} } \l_my_tl
```

results in \l_my_tl holding first, \emph{second}, first, first.

1.3 Pre-compiling regular expressions

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

\regex_new:N \regex_new:N (regex var)

Creates a new $\langle regex var \rangle$ or raises an error if the name is already taken. The declaration is global. The $\langle regex var \rangle$ will initially be such that it never matches.

\regex_set:Nn\regex_set:Nn $\langle regex var \rangle \{\langle regex \rangle\}$ \regex_gset:NnStores a compiled version of the $\langle regular \ expression \rangle$ in the $\langle regex \ var \rangle$. For instance, this function can be used as

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

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

\regex_show:n
\regex_show:N

 $regex_show:n {\langle regex \rangle}$

Shows how |3regex| interprets the $\langle regex \rangle$. For instance, $regex_show:n \{A X|Y\}$ shows

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

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

1.4 Matching

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

 \regex_match:nnTF
 \regex_match:nnTF
 \(token list)\) {(true code)} {(false code)}

 \regex_match:NnTF
 Tests whether the (regular expression) matches any part of the (token list). For instance,

 \regex_match:nnTF
 \b [cde]* } { abecdcx } { TRUE } { FALSE }

 \regex_match:nnTF { b [cde]* } { abecdcx } { TRUE } { FALSE }

 \regex_match:nnTF { b [cde]* } { example } { TRUE } { FALSE }

 leaves TRUE then FALSE in the input stream.

 \regex_count:nnN

 \regex_count:NNN

 Sets (int var) within the current TEX group level equal to the number of times (regular expression) appears in (token list). The search starts by finding the left-most longest match, respecting greedy and ungreedy operators. Then the search starts again from the character following the last character of the previous match, until reaching the end of

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

results in \l_foo_int taking the value 5.

1.5 Submatch extraction

\regex_extract_once:nnNTF
\regex_extract_once:NnNTF

 $\label{eq:linear} $$ \eqref{eq:linear} $$ \eqref{$

the token list. Infinite loops are prevented in the case where the regular expression can match an empty token list: then we count one match between each pair of characters.

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

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

Then the regular expression (anchored at the start with A and at the end with Z) will match the whole token list. The first capturing group, (La)?, matches La, and the second capturing group, (!*), matches !!!. Thus, 1_foo_seq will contain the items {LaTeX!!!}, {La}, and {!!!}, and the true branch is left in the input stream.

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

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

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

\regex_split:nnNTF
\regex_split:NnNTF

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

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

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

1.6 Replacement

\regex_replace_once:nnNTF
\regex_replace_once:NnNTF

 $\label{eq:lagrange} $$ \regex_replace_once:nnN {$ (regular expression) } {$ (replacement) } {tl var} $$ \regex_replace_once:nnNTF {$ (regular expression) } {$ (replacement) } {tl var} $$ (true code) } $$ {$ (false code) } $$$

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

```
      \regex_replace_all:nnNTF
      \regex_replace_all:nnN {<regular expression}} {</td>
      \treplacement} \treplace_all:nnNTF {

      \regex_replace_all:NnNTF
      \regex_replace_all:nnNTF {
      \treplacement} \treplacement} \treplacement} \treplacement} \treplacement
```

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

1.7 Bugs, misfeatures, future work, and other possibilities

The following need to be done now.

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

Additional error-checking to come.

- Currently, a{\x34} is recognized as a{4}.
- Cleaner error reporting in the replacement phase.
- Add tracing information.
- Detect attempts to use back-references.
- Test for the maximum register \c_max_register_int.
- Find out whether the fact that \W and friends match the end-marker leads to bugs. Possibly update __regex_item_reverse:n.
- Enforce that \cC can only be followed by a match-all dot.
- The empty cs should be matched by \c{}, not by \c{csname.?endcsname\s?}.

Code improvements to come.

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

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

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

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

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

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

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

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

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

Index

The italic numbers denote the pages where the corresponding entry is described, numbers underlined point to the definition, all others indicate the places where it is used.

В	$\ensuremath{\sc v}$
\begin 1, 5	$regex_extract_all:nnNTF \dots 9, 9$
	\regex_extract_once:nnN 8,9
\mathbf{C}	\regex_extract_once:NnNTF 8
cs commands:	\regex_extract_once:nnNTF 6, 8, 8
\cs_if_exist 11	\regex_gset:Nn
	$\sum_{regex_item_reverse:n} \dots 10$
\mathbf{F}	$\mbox{regex_match:NnTF} \dots \dots \dots \dots 8$
foo commands:	\regex_match:nnTF
\l_foo_int 8	\regex_new:N
\c_foo_regex $\dots \dots 1$	$\ensuremath{\columnwidth{regex_replace_all:nnN}\)$
\l_foo_seq 6, 6	\regex_replace_all:NnNTF 10
	\regex_replace_all:nnNTF 10, 10
I	$\ensuremath{\sc v}$
int commands:	$\ensuremath{\sc v}$
int_step 11	$\ensuremath{\sc v}$
	\regex_replacement_balance
M	$\texttt{one_match:n}$
max commands:	\regex_set:Nn
\c_max_register_int 10	$\sin \gamma$
my commands:	$regex_show:n \dots \gamma, 7, 7$
\1_my_tl 1, 6, 7	$\space{1.5} $
D	$\space{1.5} $
R	$\space{1.5} $
regex commands:	$\ensuremath{regular}_{\sqcup} expression$
\regex_(g)set:Nn 8	_
\regex_action_free:n 11	Τ
\regex_action_free_group:n 10	$T_{E}X$ and $L^{A}T_{E}X 2\varepsilon$ commands:
\regex_action_success: 11	\dimen 10, 11
\lregex_balance_tl 11	\skip 10, 10, 10
\regex_const:Nn 1, 7, 7	tl commands:
\regex_count:NnN 8	\tl_to_str:v 5, 7
\regex_count:nnN 8, 8	tmpa commands:
\regex_extract_all:nnN 9	\l_tmpa_int 4