

NAME

perlref - Perl Regular Expressions Reference

DESCRIPTION

This is a quick reference to Perl's regular expressions. For full information see *perlre* and *perllop*, as well as the *SEE ALSO* section in this document.

OPERATORS

`=~` determines to which variable the regex is applied. In its absence, `$_` is used.

```
$var =~ /foo/;
```

`!~` determines to which variable the regex is applied, and negates the result of the match; it returns false if the match succeeds, and true if it fails.

```
$var !~ /foo/;
```

`m/pattern/msixogc` searches a string for a pattern match, applying the given options.

```
m Multiline mode - ^ and $ match internal lines
s match as a Single line - . matches \n
i case-Insensitive
x eXtended legibility - free whitespace and comments
o compile pattern Once
g Global - all occurrences
c don't reset pos on failed matches when using /g
```

If 'pattern' is an empty string, the last *successfully* matched regex is used. Delimiters other than '/' may be used for both this operator and the following ones. The leading `m` can be omitted if the delimiter is '/'.

`qr/pattern/msixo` lets you store a regex in a variable, or pass one around. Modifiers as for `m//`, and are stored within the regex.

`s/pattern/replacement/msixogce` substitutes matches of 'pattern' with 'replacement'. Modifiers as for `m//`, with one addition:

```
e Evaluate 'replacement' as an expression
```

'e' may be specified multiple times. 'replacement' is interpreted as a double quoted string unless a single-quote (') is the delimiter.

`?pattern?` is like `m/pattern/` but matches only once. No alternate delimiters can be used. Must be reset with `reset()`.

SYNTAX

<code>\</code>	Escapes the character immediately following it
<code>.</code>	Matches any single character except a newline (unless <code>/s</code> is used)
<code>^</code>	Matches at the beginning of the string (or line, if <code>/m</code> is used)
<code>\$</code>	Matches at the end of the string (or line, if <code>/m</code> is used)
<code>*</code>	Matches the preceding element 0 or more times
<code>+</code>	Matches the preceding element 1 or more times
<code>?</code>	Matches the preceding element 0 or 1 times
<code>{...}</code>	Specifies a range of occurrences for the element preceding it
<code>[...]</code>	Matches any one of the characters contained within the brackets
<code>(...)</code>	Groups subexpressions for capturing to <code>\$1</code> , <code>\$2</code> ...

```
(?:...) Groups subexpressions without capturing (cluster)
|       Matches either the subexpression preceding or following it
\1, \2, \3 ...           Matches the text from the Nth group
```

ESCAPE SEQUENCES

These work as in normal strings.

```
\a      Alarm (beep)
\e      Escape
\f      Formfeed
\n      Newline
\r      Carriage return
\t      Tab
\037    Any octal ASCII value
\x7f    Any hexadecimal ASCII value
\x{263a} A wide hexadecimal value
\cx     Control-x
\N{name} A named character

\l      Lowercase next character
\u      Titlecase next character
\L      Lowercase until \E
\U      Uppercase until \E
\Q      Disable pattern metacharacters until \E
\E      End modification
```

For Titlecase, see *Titlecase*.

This one works differently from normal strings:

```
\b      An assertion, not backspace, except in a character class
```

CHARACTER CLASSES

```
[amy]    Match 'a', 'm' or 'y'
[f-j]    Dash specifies "range"
[f-j-]   Dash escaped or at start or end means 'dash'
[^f-j]   Caret indicates "match any character _except_ these"
```

The following sequences work within or without a character class. The first six are locale aware, all are Unicode aware. See *perllocale* and *perlunicode* for details.

```
\d      A digit
\D      A nondigit
\w      A word character
\W      A non-word character
\s      A whitespace character
\S      A non-whitespace character

\C      Match a byte (with Unicode, '.' matches a character)
\pP     Match P-named (Unicode) property
\p{...} Match Unicode property with long name
\pP     Match non-P
\p{...} Match lack of Unicode property with long name
\X      Match extended Unicode combining character sequence
```

POSIX character classes and their Unicode and Perl equivalents:

alnum	IsAlnum	Alphanumeric
alpha	IsAlpha	Alphabetic
ascii	IsASCII	Any ASCII char
blank	IsSpace [\t]	Horizontal whitespace (GNU extension)
cntrl	IsCntrl	Control characters
digit	IsDigit \d	Digits
graph	IsGraph	Alphanumeric and punctuation
lower	IsLower	Lowercase chars (locale and Unicode aware)
print	IsPrint	Alphanumeric, punct, and space
punct	IsPunct	Punctuation
space	IsSpace [\s\ck]	Whitespace
	IsSpacePerl \s	Perl's whitespace definition
upper	IsUpper	Uppercase chars (locale and Unicode aware)
word	IsWord \w	Alphanumeric plus _ (Perl extension)
xdigit	IsXDigit [0-9A-Fa-f]	Hexadecimal digit

Within a character class:

POSIX	traditional	Unicode
<code>[:digit:]</code>	<code>\d</code>	<code>\p{IsDigit}</code>
<code>[:^digit:]</code>	<code>\D</code>	<code>\P{IsDigit}</code>

ANCHORS

All are zero-width assertions.

<code>^</code>	Match string start (or line, if /m is used)
<code>\$</code>	Match string end (or line, if /m is used) or before newline
<code>\b</code>	Match word boundary (between \w and \W)
<code>\B</code>	Match except at word boundary (between \w and \w or \W and \W)
<code>\A</code>	Match string start (regardless of /m)
<code>\Z</code>	Match string end (before optional newline)
<code>\z</code>	Match absolute string end
<code>\G</code>	Match where previous m//g left off

QUANTIFIERS

Quantifiers are greedy by default -- match the **longest** leftmost.

Maximal	Minimal	Allowed range
-----	-----	-----
<code>{n,m}</code>	<code>{n,m}?</code>	Must occur at least n times but no more than m times
<code>{n,}</code>	<code>{n,}?</code>	Must occur at least n times
<code>{n}</code>	<code>{n}?</code>	Must occur exactly n times
<code>*</code>	<code>*?</code>	0 or more times (same as <code>{0,}</code>)
<code>+</code>	<code>+</code>	1 or more times (same as <code>{1,}</code>)
<code>?</code>	<code>??</code>	0 or 1 time (same as <code>{0,1}</code>)

There is no quantifier `{n}` -- that gets understood as a literal string.

EXTENDED CONSTRUCTS

<code>(?#text)</code>	A comment
<code>(?:...)</code>	Groups subexpressions without capturing (cluster)
<code>(?imsx-imsx:...)</code>	Enable/disable option (as per m// modifiers)
<code>(?=...)</code>	Zero-width positive lookahead assertion

<code>(?!...)</code>	Zero-width negative lookahead assertion
<code>(?<=...)</code>	Zero-width positive lookbehind assertion
<code>(?<!...)</code>	Zero-width negative lookbehind assertion
<code>(?>...)</code>	Grab what we can, prohibit backtracking
<code>(?{ code })</code>	Embedded code, return value becomes <code>^R</code>
<code>(??{ code })</code>	Dynamic regex, return value used as regex
<code>(?(cond)yes no)</code>	
<code>(?(cond)yes)</code>	Conditional expression, where "cond" can be: (N) subpattern N has matched something (? <code>{code}</code>) code condition

VARIABLES

<code>\$_</code>	Default variable for operators to use
<code>\$*</code>	Enable multiline matching (deprecated; not in 5.9.0 or later)
<code>\$`</code>	Everything prior to matched string
<code>\$&</code>	Entire matched string
<code>\$'</code>	Everything after to matched string

The use of `$``, `$&` or `$'` will slow down **all** regex use within your program. Consult *perlvar* for `@-` to see equivalent expressions that won't cause slow down. See also *Devel::SawAmpersand*. If you upgrade to Perl 5.10, you can also use the equivalent variables `${^PREMATCH}`, `${^MATCH}` and `${^POSTMATCH}`, but for them to be defined, you have to specify the `/p` (preserve) modifier on your regular expression.

<code>\$1</code> , <code>\$2</code> ...	hold the Xth captured expr
<code>\$+</code>	Last parenthesized pattern match
<code>^N</code>	Holds the most recently closed capture
<code>^R</code>	Holds the result of the last <code>(?{...})</code> expr
<code>@-</code>	Offsets of starts of groups. <code>\$-[0]</code> holds start of whole match
<code>@+</code>	Offsets of ends of groups. <code>\$+[0]</code> holds end of whole match

Captured groups are numbered according to their *opening* paren.

FUNCTIONS

<code>lc</code>	Lowercase a string
<code>lcfirst</code>	Lowercase first char of a string
<code>uc</code>	Uppercase a string
<code>ucfirst</code>	Titlecase first char of a string
<code>pos</code>	Return or set current match position
<code>quotemeta</code>	Quote metacharacters
<code>reset</code>	Reset <code>?pattern?</code> status
<code>study</code>	Analyze string for optimizing matching
<code>split</code>	Use a regex to split a string into parts

The first four of these are like the escape sequences `\L`, `\l`, `\U`, and `\u`. For Titlecase, see *Titlecase*.

TERMINOLOGY

Titlecase

Unicode concept which most often is equal to uppercase, but for certain characters like the German "sharp s" there is a difference.

AUTHOR

Iain Truskett. Updated by the Perl 5 Porters.

This document may be distributed under the same terms as Perl itself.

SEE ALSO

- *perlretut* for a tutorial on regular expressions.
- *perlrequick* for a rapid tutorial.
- *perlre* for more details.
- *perlvar* for details on the variables.
- *perlop* for details on the operators.
- *perlfunc* for details on the functions.
- *perlfaq6* for FAQs on regular expressions.
- *perlrebackslash* for a reference on backslash sequences.
- *perlrecharclass* for a reference on character classes.
- The *re* module to alter behaviour and aid debugging.
- "*Debugging regular expressions*" in *perldebug*
- *perluniintro*, *perlunicode*, *charnames* and *perllocale* for details on regexes and internationalisation.
- *Mastering Regular Expressions* by Jeffrey Friedl (<http://regex.info/>) for a thorough grounding and reference on the topic.

THANKS

David P.C. Wollmann, Richard Soderberg, Sean M. Burke, Tom Christiansen, Jim Cromie, and Jeffrey Goff for useful advice.