

## PROLOG

This manual page is part of the POSIX Programmer's Manual. The Cygwin implementation of this interface may differ (consult the corresponding Cygwin manual page for details of Cygwin behavior), or the interface may not be implemented on Cygwin.

## NAME

regcomp, regerror, regex, regfree — regular expression matching

## SYNOPSIS

```
#include <regex.h>

int regcomp(regex_t *restrict preg, const char *restrict pattern,
            int cflags);
size_t regerror(int errcode, const regex_t *restrict preg,
               char *restrict errbuf, size_t errbuf_size);
int regex(const regex_t *restrict preg, const char *restrict string,
          size_t nmatch, regmatch_t pmatch[restrict], int eflags);
void regfree(regex_t *preg);
```

## DESCRIPTION

These functions interpret *basic* and *extended* regular expressions as described in the Base Definitions volume of POSIX.1-2008, *Chapter 9, Regular Expressions*.

The **regex\_t** structure is defined in *<regex.h>* and contains at least the following member:

Member Type	Member Name	Description
<b>size_t</b>	<i>re_nsub</i>	Number of parenthesized subexpressions.

The **regmatch\_t** structure is defined in *<regex.h>* and contains at least the following members:

Member Type	Member Name	Description
<b>regoff_t</b>	<i>rm_so</i>	Byte offset from start of <i>string</i> to start of substring.
<b>regoff_t</b>	<i>rm_eo</i>	Byte offset from start of <i>string</i> of the first character after the end of substring.

The *regcomp()* function shall compile the regular expression contained in the string pointed to by the *pattern* argument and place the results in the structure pointed to by *preg*. The *cflags* argument is the bitwise-inclusive OR of zero or more of the following flags, which are defined in the *<regex.h>* header:

### REG\_EXTENDED

Use Extended Regular Expressions.

**REG\_ICASE** Ignore case in match (see the Base Definitions volume of POSIX.1-2008, *Chapter 9, Regular Expressions*).

**REG\_NOSUB** Report only success/fail in *regex()*.

### REG\_NEWLINE

Change the handling of <newline> characters, as described in the text.

The default regular expression type for *pattern* is a Basic Regular Expression. The application can specify Extended Regular Expressions using the **REG\_EXTENDED** *cflags* flag.

If the **REG\_NOSUB** flag was not set in *cflags*, then *regcomp()* shall set *re\_nsub* to the number of parenthesized subexpressions (delimited by "\(\)" in basic regular expressions or "(\)" in extended regular expressions) found in *pattern*.

The *regex()* function compares the null-terminated string specified by *string* with the compiled regular expression *preg* initialized by a previous call to *regcomp()*. If it finds a match, *regex()* shall return 0; otherwise, it shall return non-zero indicating either no match or an error. The *eflags* argument is the bitwise-inclusive OR of zero or more of the following flags, which are defined in the *<regex.h>* header:

**REG\_NOTBOL**

The first character of the string pointed to by *string* is not the beginning of the line. Therefore, the <circumflex> character (^), when taken as a special character, shall not match the beginning of *string*.

**REG\_NOTEOL** The last character of the string pointed to by *string* is not the end of the line. Therefore, the <dollar-sign> (\$), when taken as a special character, shall not match the end of *string*.

If *nmatch* is 0 or REG\_NOSUB was set in the *cflags* argument to *regcomp()*, then *regexexec()* shall ignore the *pmatch* argument. Otherwise, the application shall ensure that the *pmatch* argument points to an array with at least *nmatch* elements, and *regexexec()* shall fill in the elements of that array with offsets of the substrings of *string* that correspond to the parenthesized subexpressions of *pattern*: *pmatch[i].rm\_so* shall be the byte offset of the beginning and *pmatch[i].rm\_eo* shall be one greater than the byte offset of the end of substring *i*. (Subexpression *i* begins at the *i*th matched open parenthesis, counting from 1.) Offsets in *pmatch[0]* identify the substring that corresponds to the entire regular expression. Unused elements of *pmatch* up to *pmatch[nmatch-1]* shall be filled with -1. If there are more than *nmatch* subexpressions in *pattern* (*pattern* itself counts as a subexpression), then *regexexec()* shall still do the match, but shall record only the first *nmatch* substrings.

When matching a basic or extended regular expression, any given parenthesized subexpression of *pattern* might participate in the match of several different substrings of *string*, or it might not match any substring even though the pattern as a whole did match. The following rules shall be used to determine which substrings to report in *pmatch* when matching regular expressions:

1. If subexpression *i* in a regular expression is not contained within another subexpression, and it participated in the match several times, then the byte offsets in *pmatch[i]* shall delimit the last such match.
2. If subexpression *i* is not contained within another subexpression, and it did not participate in an otherwise successful match, the byte offsets in *pmatch[i]* shall be -1. A subexpression does not participate in the match when:

‘\*’ or “\{” appears immediately after the subexpression in a basic regular expression, or ‘\*’, ‘?’, or “{” appears immediately after the subexpression in an extended regular expression, and the subexpression did not match (matched 0 times)

or:

‘|’ is used in an extended regular expression to select this subexpression or another, and the other subexpression matched.

3. If subexpression *i* is contained within another subexpression *j*, and *i* is not contained within any other subexpression that is contained within *j*, and a match of subexpression *j* is reported in *pmatch[j]*, then the match or non-match of subexpression *i* reported in *pmatch[i]* shall be as described in 1. and 2. above, but within the substring reported in *pmatch[j]* rather than the whole string. The offsets in *pmatch[i]* are still relative to the start of *string*.
4. If subexpression *i* is contained in subexpression *j*, and the byte offsets in *pmatch[j]* are -1, then the pointers in *pmatch[i]* shall also be -1.
5. If subexpression *i* matched a zero-length string, then both byte offsets in *pmatch[i]* shall be the byte offset of the character or null terminator immediately following the zero-length string.

If, when *regexexec()* is called, the locale is different from when the regular expression was compiled, the result is undefined.

If REG\_NEWLINE is not set in *cflags*, then a <newline> in *pattern* or *string* shall be treated as an ordinary character. If REG\_NEWLINE is set, then <newline> shall be treated as an ordinary character except as follows:

1. A <newline> in *string* shall not be matched by a <period> outside a bracket expression or by any form of a non-matching list (see the Base Definitions volume of POSIX.1-2008, *Chapter 9, Regular*

*Expressions*).

2. A <circumflex> (^) in *pattern*, when used to specify expression anchoring (see the Base Definitions volume of POSIX.1-2008, *Section 9.3.8, BRE Expression Anchoring*), shall match the zero-length string immediately after a <newline> in *string*, regardless of the setting of REG\_NOTBOL.
3. A <dollar-sign> (\$) in *pattern*, when used to specify expression anchoring, shall match the zero-length string immediately before a <newline> in *string*, regardless of the setting of REG\_NOTEOL.

The *regfree()* function frees any memory allocated by *regcomp()* associated with *preg*.

The following constants are defined as the minimum set of error return values, although other errors listed as implementation extensions in <*regex.h*> are possible:

REG\_BADBR Content of "{\}" invalid: not a number, number too large, more than two numbers, first larger than second.

REG\_BADPAT Invalid regular expression.

REG\_BADRPT '?', '\*', or '+' not preceded by valid regular expression.

REG\_EBRACE "{\}" imbalance.

REG\_EBRACK "[\]" imbalance.

REG\_ECOLLATE Invalid collating element referenced.

REG\_ECTYPE Invalid character class type referenced.

REG\_EESCAPE Trailing <backslash> character in pattern.

REG\_EPAREN "\()" or "()" imbalance.

REG\_ERANGE Invalid endpoint in range expression.

REG\_ESPACE Out of memory.

REG\_ESUBREG Number in "\digit" invalid or in error.

REG\_NOMATCH *regexexec()* failed to match.

If more than one error occurs in processing a function call, any one of the possible constants may be returned, as the order of detection is unspecified.

The *regerror()* function provides a mapping from error codes returned by *regcomp()* and *regexexec()* to unspecified printable strings. It generates a string corresponding to the value of the *errcode* argument, which the application shall ensure is the last non-zero value returned by *regcomp()* or *regexexec()* with the given value of *preg*. If *errcode* is not such a value, the content of the generated string is unspecified.

If *preg* is a null pointer, but *errcode* is a value returned by a previous call to *regexexec()* or *regcomp()*, the *regerror()* still generates an error string corresponding to the value of *errcode*, but it might not be as detailed under some implementations.

If the *errbuf\_size* argument is not 0, *regerror()* shall place the generated string into the buffer of size *errbuf\_size* bytes pointed to by *errbuf*. If the string (including the terminating null) cannot fit in the buffer, *regerror()* shall truncate the string and null-terminate the result.

If *errbuf\_size* is 0, *regerror()* shall ignore the *errbuf* argument, and return the size of the buffer needed to hold the generated string.

If the *preg* argument to *regexexec()* or *regfree()* is not a compiled regular expression returned by *regcomp()*, the result is undefined. A *preg* is no longer treated as a compiled regular expression after it is given to

*regfree()*.

## RETURN VALUE

Upon successful completion, the *regcomp()* function shall return 0. Otherwise, it shall return an integer value indicating an error as described in *<regex.h>*, and the content of *preg* is undefined. If a code is returned, the interpretation shall be as given in *<regex.h>*.

If *regcomp()* detects an invalid RE, it may return REG\_BADPAT, or it may return one of the error codes that more precisely describes the error.

Upon successful completion, the *regexexec()* function shall return 0. Otherwise, it shall return REG\_NOMATCH to indicate no match.

Upon successful completion, the *regerror()* function shall return the number of bytes needed to hold the entire generated string, including the null termination. If the return value is greater than *errbuf\_size*, the string returned in the buffer pointed to by *errbuf* has been truncated.

The *regfree()* function shall not return a value.

## ERRORS

No errors are defined.

*The following sections are informative.*

## EXAMPLES

```
#include <regex.h>

/*
 * Match string against the extended regular expression in
 * pattern, treating errors as no match.
 *
 * Return 1 for match, 0 for no match.
 */

int
match(const char *string, char *pattern)
{
    int status;
    regex_t re;

    if (regcomp(&re, pattern, REG_EXTENDED|REG_NOSUB) != 0) {
        return(0);    /* Report error. */
    }
    status = regexexec(&re, string, (size_t) 0, NULL, 0);
    regfree(&re);
    if (status != 0) {
        return(0);    /* Report error. */
    }
    return(1);
}
```

The following demonstrates how the REG\_NOTBOL flag could be used with *regexexec()* to find all substrings in a line that match a pattern supplied by a user. (For simplicity of the example, very little error checking is done.)

```
(void) regcomp (&re, pattern, 0);
/* This call to regexexec() finds the first match on the line. */
error = regexexec (&re, &buffer[0], 1, &pm, 0);
while (error == 0) { /* While matches found. */
    /* Substring found between pm.rm_so and pm.rm_eo. */
    /* This call to regexexec() finds the next match. */
    error = regexexec (&re, buffer + pm.rm_eo, 1, &pm, REG_NOTBOL);
}
```

}

## APPLICATION USAGE

An application could use:

```
regerror(code,preg,(char *)NULL,(size_t)0)
```

to find out how big a buffer is needed for the generated string, *malloc()* a buffer to hold the string, and then call *regerror()* again to get the string. Alternatively, it could allocate a fixed, static buffer that is big enough to hold most strings, and then use *malloc()* to allocate a larger buffer if it finds that this is too small.

To match a pattern as described in the Shell and Utilities volume of POSIX.1-2008, *Section 2.13, Pattern Matching Notation*, use the *fnmatch()* function.

## RATIONALE

The *regexexec()* function must fill in all *nmatch* elements of *pmatch*, where *nmatch* and *pmatch* are supplied by the application, even if some elements of *pmatch* do not correspond to subexpressions in *pattern*. The application developer should note that there is probably no reason for using a value of *nmatch* that is larger than *preg->re\_nsub*+1.

The REG\_NEWLINE flag supports a use of RE matching that is needed in some applications like text editors. In such applications, the user supplies an RE asking the application to find a line that matches the given expression. An anchor in such an RE anchors at the beginning or end of any line. Such an application can pass a sequence of <newline>-separated lines to *regexexec()* as a single long string and specify REG\_NEWLINE to *regcomp()* to get the desired behavior. The application must ensure that there are no explicit <newline> characters in *pattern* if it wants to ensure that any match occurs entirely within a single line.

The REG\_NEWLINE flag affects the behavior of *regexexec()*, but it is in the *cflags* parameter to *regcomp()* to allow flexibility of implementation. Some implementations will want to generate the same compiled RE in *regcomp()* regardless of the setting of REG\_NEWLINE and have *regexexec()* handle anchors differently based on the setting of the flag. Other implementations will generate different compiled REs based on the REG\_NEWLINE.

The REG\_ICASE flag supports the operations taken by the *grep -i* option and the historical implementations of *ex* and *vi*. Including this flag will make it easier for application code to be written that does the same thing as these utilities.

The substrings reported in *pmatch[]* are defined using offsets from the start of the string rather than pointers. This allows type-safe access to both constant and non-constant strings.

The type **regoff\_t** is used for the elements of *pmatch[]* to ensure that the application can represent large arrays in memory (important for an application conforming to the Shell and Utilities volume of POSIX.1-2008).

The 1992 edition of this standard required **regoff\_t** to be at least as wide as **off\_t**, to facilitate future extensions in which the string to be searched is taken from a file. However, these future extensions have not appeared. The requirement rules out popular implementations with 32-bit **regoff\_t** and 64-bit **off\_t**, so it has been removed.

The standard developers rejected the inclusion of a *regsub()* function that would be used to do substitutions for a matched RE. While such a routine would be useful to some applications, its utility would be much more limited than the matching function described here. Both RE parsing and substitution are possible to implement without support other than that required by the ISO C standard, but matching is much more complex than substituting. The only difficult part of substitution, given the information supplied by *regexexec()*, is finding the next character in a string when there can be multi-byte characters. That is a much larger issue, and one that needs a more general solution.

The *errno* variable has not been used for error returns to avoid filling the *errno* name space for this feature.

The interface is defined so that the matched substrings *rm\_sp* and *rm\_ep* are in a separate **regmatch\_t** structure instead of in **regex\_t**. This allows a single compiled RE to be used simultaneously in several

contexts; in *main()* and a signal handler, perhaps, or in multiple threads of lightweight processes. (The *preg* argument to *regexexec()* is declared with type **const**, so the implementation is not permitted to use the structure to store intermediate results.) It also allows an application to request an arbitrary number of substrings from an RE. The number of subexpressions in the RE is reported in *re\_nsub* in *preg*. With this change to *regexexec()*, consideration was given to dropping the REG\_NOSUB flag since the user can now specify this with a zero *nmatch* argument to *regexexec()*. However, keeping REG\_NOSUB allows an implementation to use a different (perhaps more efficient) algorithm if it knows in *regcomp()* that no subexpressions need be reported. The implementation is only required to fill in *pmatch* if *nmatch* is not zero and if REG\_NOSUB is not specified. Note that the **size\_t** type, as defined in the ISO C standard, is unsigned, so the description of *regexexec()* does not need to address negative values of *nmatch*.

REG\_NOTBOL was added to allow an application to do repeated searches for the same pattern in a line. If the pattern contains a <circumflex> character that should match the beginning of a line, then the pattern should only match when matched against the beginning of the line. Without the REG\_NOTBOL flag, the application could rewrite the expression for subsequent matches, but in the general case this would require parsing the expression. The need for REG\_NOTEOL is not as clear; it was added for symmetry.

The addition of the *regerror()* function addresses the historical need for conforming application programs to have access to error information more than “Function failed to compile/match your RE for unknown reasons”.

This interface provides for two different methods of dealing with error conditions. The specific error codes (REG\_EBRACE, for example), defined in <*regex.h*>, allow an application to recover from an error if it is so able. Many applications, especially those that use patterns supplied by a user, will not try to deal with specific error cases, but will just use *regerror()* to obtain a human-readable error message to present to the user.

The *regerror()* function uses a scheme similar to *confstr()* to deal with the problem of allocating memory to hold the generated string. The scheme used by *strerror()* in the ISO C standard was considered unacceptable since it creates difficulties for multi-threaded applications.

The *preg* argument is provided to *regerror()* to allow an implementation to generate a more descriptive message than would be possible with *errcode* alone. An implementation might, for example, save the character offset of the offending character of the pattern in a field of *preg*, and then include that in the generated message string. The implementation may also ignore *preg*.

A REG\_FILENAME flag was considered, but omitted. This flag caused *regexexec()* to match patterns as described in the Shell and Utilities volume of POSIX.1-2008, *Section 2.13, Pattern Matching Notation* instead of REs. This service is now provided by the *fnmatch()* function.

Notice that there is a difference in philosophy between the ISO POSIX-2: 1993 standard and POSIX.1-2008 in how to handle a “bad” regular expression. The ISO POSIX-2: 1993 standard says that many bad constructs “produce undefined results”, or that “the interpretation is undefined”. POSIX.1-2008, however, says that the interpretation of such REs is unspecified. The term “undefined” means that the action by the application is an error, of similar severity to passing a bad pointer to a function.

The *regcomp()* and *regexexec()* functions are required to accept any null-terminated string as the *pattern* argument. If the meaning of the string is “undefined”, the behavior of the function is “unspecified”. POSIX.1-2008 does not specify how the functions will interpret the pattern; they might return error codes, or they might do pattern matching in some completely unexpected way, but they should not do something like abort the process.

## FUTURE DIRECTIONS

None.

## SEE ALSO

*fnmatch()*, *glob()*

The Base Definitions volume of POSIX.1-2008, *Chapter 9, Regular Expressions*, <**regex.h**>, <**sys\_types.h**>

The Shell and Utilities volume of POSIX.1-2008, *Section 2.13, Pattern Matching Notation*

**COPYRIGHT**

Portions of this text are reprinted and reproduced in electronic form from IEEE Std 1003.1, 2013 Edition, Standard for Information Technology -- Portable Operating System Interface (POSIX), The Open Group Base Specifications Issue 7, Copyright (C) 2013 by the Institute of Electrical and Electronics Engineers, Inc and The Open Group. (This is POSIX.1-2008 with the 2013 Technical Corrigendum 1 applied.) In the event of any discrepancy between this version and the original IEEE and The Open Group Standard, the original IEEE and The Open Group Standard is the referee document. The original Standard can be obtained online at <http://www.unix.org/online.html> .

Any typographical or formatting errors that appear in this page are most likely to have been introduced during the conversion of the source files to man page format. To report such errors, see [https://www.kernel.org/doc/man-pages/reporting\\_bugs.html](https://www.kernel.org/doc/man-pages/reporting_bugs.html) .