Manual  /  Strings                                               ⓞ Edit on GitHub  ⚙  ☰

# Strings

Strings are finite sequences of characters. Of course, the real trouble comes when one asks what a character is. The characters that English speakers are familiar with are the letters A, B, C, etc., together with numerals and common punctuation symbols. These characters are standardized together with a mapping to integer values between 0 and 127 by the ASCII standard. There are, of course, many other characters used in non-English languages, including variants of the ASCII characters with accents and other modifications, related scripts such as Cyrillic and Greek, and scripts completely unrelated to ASCII and English, including Arabic, Chinese, Hebrew, Hindi, Japanese, and Korean. The Unicode standard tackles the complexities of what exactly a character is, and is generally accepted as the definitive standard addressing this problem. Depending on your needs, you can either ignore these complexities entirely and just pretend that only ASCII characters exist, or you can write code that can handle any of the characters or encodings that one may encounter when handling non-ASCII text. Julia makes dealing with plain ASCII text simple and efficient, and handling Unicode is as simple and efficient as possible. In particular, you can write C-style string code to process ASCII strings, and they will work as expected, both in terms of performance and semantics. If such code encounters non-ASCII text, it will gracefully fail with a clear error message, rather than silently introducing corrupt results. When this happens, modifying the code to handle non-ASCII data is straightforward.

There are a few noteworthy high-level features about Julia's strings:

- The built-in concrete type used for strings (and string literals) in Julia is `String`. This supports the full range of Unicode characters via the UTF-8 encoding. (A `transcode` function is provided to convert to/from other Unicode encodings.)
- All string types are subtypes of the abstract type `AbstractString`, and external packages define additional `AbstractString` subtypes (e.g. for other encodings). If you define a function expecting a string argument, you should declare the type as `AbstractString` in order to accept any string type.
- Like C and Java, but unlike most dynamic languages, Julia has a first-class type for representing a single character, called `AbstractChar`. The built-in `Char` subtype of `AbstractChar` is a 32-bit primitive type that can represent any Unicode character (and which is based on the UTF-8 encoding).
- As in Java, strings are immutable: the value of an `AbstractString` object cannot be changed. To construct a different string value, you construct a new string from parts of other strings.
- Conceptually, a string is a *partial function* from indices to characters: for some index values, no character value is returned, and instead an exception is thrown. This allows for efficient indexing into strings by the byte index of an encoded representation rather than by a character index, which

cannot be implemented both efficiently and simply for variable-width encodings of Unicode strings.

# Characters

A `Char` value represents a single character: it is just a 32-bit primitive type with a special literal representation and appropriate arithmetic behaviors, and which can be converted to a numeric value representing a Unicode code point. (Julia packages may define other subtypes of `AbstractChar`, e.g. to optimize operations for other text encodings.) Here is how `Char` values are input and shown:

```julia
julia> 'x'
'x': ASCII/Unicode U+0078 (category Ll: Letter, lowercase)

julia> typeof(ans)
Char
```

You can easily convert a `Char` to its integer value, i.e. code point:

```julia
julia> Int('x')
120

julia> typeof(ans)
Int64
```

On 32-bit architectures, `typeof(ans)` will be `Int32`. You can convert an integer value back to a `Char` just as easily:

```julia
julia> Char(120)
'x': ASCII/Unicode U+0078 (category Ll: Letter, lowercase)
```

Not all integer values are valid Unicode code points, but for performance, the `Char` conversion does not check that every character value is valid. If you want to check that each converted value is a valid code point, use the `isvalid` function:

```julia
julia> Char(0x110000)
'\U110000': Unicode U+110000 (category In: Invalid, too high)

julia> isvalid(Char, 0x110000)
false
```

As of this writing, the valid Unicode code points are `U+0000` through `U+D7FF` and `U+E000` through

U+10FFFF. These have not all been assigned intelligible meanings yet, nor are they necessarily interpretable by applications, but all of these values are considered to be valid Unicode characters.

You can input any Unicode character in single quotes using `\u` followed by up to four hexadecimal digits or `\U` followed by up to eight hexadecimal digits (the longest valid value only requires six):

```julia
julia> '\u0'
'\0': ASCII/Unicode U+0000 (category Cc: Other, control)

julia> '\u78'
'x': ASCII/Unicode U+0078 (category Ll: Letter, lowercase)

julia> '\u2200'
'∀': Unicode U+2200 (category Sm: Symbol, math)

julia> '\U10ffff'
'\U10ffff': Unicode U+10FFFF (category Cn: Other, not assigned)
```

Julia uses your system's locale and language settings to determine which characters can be printed as-is and which must be output using the generic, escaped `\u` or `\U` input forms. In addition to these Unicode escape forms, all of [C's traditional escaped input forms](#) can also be used:

```julia
julia> Int('\0')
0

julia> Int('\t')
9

julia> Int('\n')
10

julia> Int('\e')
27

julia> Int('\x7f')
127

julia> Int('\177')
127
```

You can do comparisons and a limited amount of arithmetic with `Char` values:

```julia
julia> 'A' < 'a'
```

```
true

julia> 'A' <= 'a' <= 'Z'
false

julia> 'A' <= 'X' <= 'Z'
true

julia> 'x' - 'a'
23

julia> 'A' + 1
'B': ASCII/Unicode U+0042 (category Lu: Letter, uppercase)
```

# String Basics

String literals are delimited by double quotes or triple double quotes:

```
julia> str = "Hello, world.\n"
"Hello, world.\n"

julia> """Contains "quote" characters"""
"Contains \"quote\" characters"
```

If you want to extract a character from a string, you index into it:

```
julia> str[begin]
'H': ASCII/Unicode U+0048 (category Lu: Letter, uppercase)

julia> str[1]
'H': ASCII/Unicode U+0048 (category Lu: Letter, uppercase)

julia> str[6]
',': ASCII/Unicode U+002C (category Po: Punctuation, other)

julia> str[end]
'\n': ASCII/Unicode U+000A (category Cc: Other, control)
```

Many Julia objects, including strings, can be indexed with integers. The index of the first element (the first character of a string) is returned by `firstindex(str)`, and the index of the last element (character) with `lastindex(str)`. The keywords `begin` and `end` can be used inside an indexing operation as shorthand for the first and last indices, respectively, along the given dimension. String indexing, like

most indexing in Julia, is 1-based: `firstindex` always returns `1` for any `AbstractString`. As we will see below, however, `lastindex(str)` is *not* in general the same as `length(str)` for a string, because some Unicode characters can occupy multiple "code units".

You can perform arithmetic and other operations with `end`, just like a normal value:

```
julia> str[end-1]
'.': ASCII/Unicode U+002E (category Po: Punctuation, other)

julia> str[end÷2]
' ': ASCII/Unicode U+0020 (category Zs: Separator, space)
```

Using an index less than `begin` (`1`) or greater than `end` raises an error:

```
julia> str[begin-1]
ERROR: BoundsError: attempt to access String
  at index [0]
[...]

julia> str[end+1]
ERROR: BoundsError: attempt to access String
  at index [15]
[...]
```

You can also extract a substring using range indexing:

```
julia> str[4:9]
"lo, wo"
```

Notice that the expressions `str[k]` and `str[k:k]` do not give the same result:

```
julia> str[6]
',': ASCII/Unicode U+002C (category Po: Punctuation, other)

julia> str[6:6]
","
```

The former is a single character value of type `Char`, while the latter is a string value that happens to contain only a single character. In Julia these are very different things.

Range indexing makes a copy of the selected part of the original string. Alternatively, it is possible to create a view into a string using the type `SubString`, for example:

```
julia> str = "long string"
"long string"

julia> substr = SubString(str, 1, 4)
"long"

julia> typeof(substr)
SubString{String}
```

Several standard functions like `chop`, `chomp` or `strip` return a `SubString`.

# Unicode and UTF-8

Julia fully supports Unicode characters and strings. As discussed above, in character literals, Unicode code points can be represented using Unicode `\u` and `\U` escape sequences, as well as all the standard C escape sequences. These can likewise be used to write string literals:

```
julia> s = "\u2200 x \u2203 y"
"∀ x ∃ y"
```

Whether these Unicode characters are displayed as escapes or shown as special characters depends on your terminal's locale settings and its support for Unicode. String literals are encoded using the UTF-8 encoding. UTF-8 is a variable-width encoding, meaning that not all characters are encoded in the same number of bytes ("code units"). In UTF-8, ASCII characters — i.e. those with code points less than 0x80 (128) – are encoded as they are in ASCII, using a single byte, while code points 0x80 and above are encoded using multiple bytes — up to four per character.

String indices in Julia refer to code units (= bytes for UTF-8), the fixed-width building blocks that are used to encode arbitrary characters (code points). This means that not every index into a `String` is necessarily a valid index for a character. If you index into a string at such an invalid byte index, an error is thrown:

```
julia> s[1]
'∀': Unicode U+2200 (category Sm: Symbol, math)

julia> s[2]
ERROR: StringIndexError("∀ x ∃ y", 2)
[...]

julia> s[3]
ERROR: StringIndexError("∀ x ∃ y", 3)
```

```
Stacktrace:
[...]

julia> s[4]
' ': ASCII/Unicode U+0020 (category Zs: Separator, space)
```

In this case, the character ∀ is a three-byte character, so the indices 2 and 3 are invalid and the next character's index is 4; this next valid index can be computed by `nextind(s,1)`, and the next index after that by `nextind(s,4)` and so on.

Since `end` is always the last valid index into a collection, `end-1` references an invalid byte index if the second-to-last character is multibyte.

```
julia> s[end-1]
' ': ASCII/Unicode U+0020 (category Zs: Separator, space)

julia> s[end-2]
ERROR: StringIndexError("∀ x ∃ y", 9)
Stacktrace:
[...]

julia> s[prevind(s, end, 2)]
'∃': Unicode U+2203 (category Sm: Symbol, math)
```

The first case works, because the last character `y` and the space are one-byte characters, whereas `end-2` indexes into the middle of the ∃ multibyte representation. The correct way for this case is using `prevind(s, lastindex(s), 2)` or, if you're using that value to index into `s` you can write `s[prevind(s, end, 2)]` and `end` expands to `lastindex(s)`.

Extraction of a substring using range indexing also expects valid byte indices or an error is thrown:

```
julia> s[1:1]
"∀"

julia> s[1:2]
ERROR: StringIndexError("∀ x ∃ y", 2)
Stacktrace:
[...]

julia> s[1:4]
"∀ "
```

Because of variable-length encodings, the number of characters in a string (given by `length(s)`) is not

always the same as the last index. If you iterate through the indices 1 through `lastindex(s)` and index into `s`, the sequence of characters returned when errors aren't thrown is the sequence of characters comprising the string `s`. Thus we have the identity that `length(s) <= lastindex(s)`, since each character in a string must have its own index. The following is an inefficient and verbose way to iterate through the characters of `s`:

```julia
julia> for i = firstindex(s):lastindex(s)
           try
               println(s[i])
           catch
               # ignore the index error
           end
       end
∀

x

∃

y
```

The blank lines actually have spaces on them. Fortunately, the above awkward idiom is unnecessary for iterating through the characters in a string, since you can just use the string as an iterable object, no exception handling required:

```julia
julia> for c in s
           println(c)
       end
∀

x

∃

y
```

If you need to obtain valid indices for a string, you can use the `nextind` and `prevind` functions to increment/decrement to the next/previous valid index, as mentioned above. You can also use the `eachindex` function to iterate over the valid character indices:

```julia
julia> collect(eachindex(s))
7-element Array{Int64,1}:
```

```
 1
 4
 5
 6
 7
10
11
```

To access the raw code units (bytes for UTF-8) of the encoding, you can use the `codeunit(s,i)` function, where the index `i` runs consecutively from 1 to `ncodeunits(s)`. The `codeunits(s)` function returns an `AbstractVector{UInt8}` wrapper that lets you access these raw codeunits (bytes) as an array.

Strings in Julia can contain invalid UTF-8 code unit sequences. This convention allows to treat any byte sequence as a `String`. In such situations a rule is that when parsing a sequence of code units from left to right characters are formed by the longest sequence of 8-bit code units that matches the start of one of the following bit patterns (each `x` can be `0` or `1`):

- `0xxxxxxx`;
- `110xxxxx 10xxxxxx`;
- `1110xxxx 10xxxxxx 10xxxxxx`;
- `11110xxx 10xxxxxx 10xxxxxx 10xxxxxx`;
- `10xxxxxx`;
- `11111xxx`.

In particular this means that overlong and too-high code unit sequences and prefixes thereof are treated as a single invalid character rather than multiple invalid characters. This rule may be best explained with an example:

```
julia> s = "\xc0\xa0\xe2\x88\xe2|"
"\xc0\xa0\xe2\x88\xe2|"

julia> foreach(display, s)
'\xc0\xa0': [overlong] ASCII/Unicode U+0020 (category Zs: Separator, space)
'\xe2\x88': Malformed UTF-8 (category Ma: Malformed, bad data)
'\xe2': Malformed UTF-8 (category Ma: Malformed, bad data)
'|': ASCII/Unicode U+007C (category Sm: Symbol, math)

julia> isvalid.(collect(s))
4-element BitArray{1}:
 0
 0
```

```
 0
 1

julia> s2 = "\xf7\xbf\xbf\xbf"
"\U1fffff"

julia> foreach(display, s2)
'\U1fffff': Unicode U+1FFFFF (category In: Invalid, too high)
```

We can see that the first two code units in the string `s` form an overlong encoding of space character. It is invalid, but is accepted in a string as a single character. The next two code units form a valid start of a three-byte UTF-8 sequence. However, the fifth code unit `\xe2` is not its valid continuation. Therefore code units 3 and 4 are also interpreted as malformed characters in this string. Similarly code unit 5 forms a malformed character because `|` is not a valid continuation to it. Finally the string `s2` contains one too high code point.

Julia uses the UTF-8 encoding by default, and support for new encodings can be added by packages. For example, the LegacyStrings.jl package implements `UTF16String` and `UTF32String` types. Additional discussion of other encodings and how to implement support for them is beyond the scope of this document for the time being. For further discussion of UTF-8 encoding issues, see the section below on byte array literals. The `transcode` function is provided to convert data between the various UTF-xx encodings, primarily for working with external data and libraries.

# Concatenation

One of the most common and useful string operations is concatenation:

```
julia> greet = "Hello"
"Hello"

julia> whom = "world"
"world"

julia> string(greet, ", ", whom, ".\n")
"Hello, world.\n"
```

It's important to be aware of potentially dangerous situations such as concatenation of invalid UTF-8 strings. The resulting string may contain different characters than the input strings, and its number of characters may be lower than sum of numbers of characters of the concatenated strings, e.g.:

```
julia> a, b = "\xe2\x88", "\x80"
```

```
("\xe2\x88", "\x80")

julia> c = a*b
"∀"

julia> collect.([a, b, c])
3-element Array{Array{Char,1},1}:
 ['\xe2\x88']
 ['\x80']
 ['∀']

julia> length.([a, b, c])
3-element Array{Int64,1}:
 1
 1
 1
```

This situation can happen only for invalid UTF-8 strings. For valid UTF-8 strings concatenation preserves all characters in strings and additivity of string lengths.

Julia also provides `*` for string concatenation:

```
julia> greet * ", " * whom * ".\n"
"Hello, world.\n"
```

While `*` may seem like a surprising choice to users of languages that provide + for string concatenation, this use of `*` has precedent in mathematics, particularly in abstract algebra.

In mathematics, + usually denotes a *commutative* operation, where the order of the operands does not matter. An example of this is matrix addition, where `A + B == B + A` for any matrices `A` and `B` that have the same shape. In contrast, `*` typically denotes a *noncommutative* operation, where the order of the operands *does* matter. An example of this is matrix multiplication, where in general `A * B != B * A`. As with matrix multiplication, string concatenation is noncommutative: `greet * whom != whom * greet`. As such, `*` is a more natural choice for an infix string concatenation operator, consistent with common mathematical use.

More precisely, the set of all finite-length strings *S* together with the string concatenation operator `*` forms a [free monoid](free monoid) (*S*, `*`). The identity element of this set is the empty string, `""`. Whenever a free monoid is not commutative, the operation is typically represented as `\cdot`, `*`, or a similar symbol, rather than +, which as stated usually implies commutativity.

# Interpolation

Constructing strings using concatenation can become a bit cumbersome, however. To reduce the need for these verbose calls to `string` or repeated multiplications, Julia allows interpolation into string literals using $, as in Perl:

```
julia> "$greet, $whom.\n"
"Hello, world.\n"
```

This is more readable and convenient and equivalent to the above string concatenation – the system rewrites this apparent single string literal into the call `string(greet, ", ", whom, ".\n")`.

The shortest complete expression after the $ is taken as the expression whose value is to be interpolated into the string. Thus, you can interpolate any expression into a string using parentheses:

```
julia> "1 + 2 = $(1 + 2)"
"1 + 2 = 3"
```

Both concatenation and string interpolation call `string` to convert objects into string form. However, `string` actually just returns the output of `print`, so new types should add methods to `print` or `show` instead of `string`.

Most non-`AbstractString` objects are converted to strings closely corresponding to how they are entered as literal expressions:

```
julia> v = [1,2,3]
3-element Array{Int64,1}:
 1
 2
 3

julia> "v: $v"
"v: [1, 2, 3]"
```

`string` is the identity for `AbstractString` and `AbstractChar` values, so these are interpolated into strings as themselves, unquoted and unescaped:

```
julia> c = 'x'
'x': ASCII/Unicode U+0078 (category Ll: Letter, lowercase)

julia> "hi, $c"
"hi, x"
```

To include a literal $ in a string literal, escape it with a backslash:

```julia
julia> print("I have \$100 in my account.\n")
I have $100 in my account.
```

# Triple-Quoted String Literals

When strings are created using triple-quotes ("""..."""") they have some special behavior that can be useful for creating longer blocks of text.

First, triple-quoted strings are also dedented to the level of the least-indented line. This is useful for defining strings within code that is indented. For example:

```julia
julia> str = """
           Hello,
           world.
         """
"  Hello,\n  world.\n"
```

In this case the final (empty) line before the closing """ sets the indentation level.

The dedentation level is determined as the longest common starting sequence of spaces or tabs in all lines, excluding the line following the opening """ and lines containing only spaces or tabs (the line containing the closing """ is always included). Then for all lines, excluding the text following the opening """, the common starting sequence is removed (including lines containing only spaces and tabs if they start with this sequence), e.g.:

```julia
julia> """    This
         is
           a test"""
"    This\nis\n  a test"
```

Next, if the opening """ is followed by a newline, the newline is stripped from the resulting string.

```
"""hello"""
```

is equivalent to

```
"""
hello"""
```

but

```
"""

hello"""
```

will contain a literal newline at the beginning.

Stripping of the newline is performed after the dedentation. For example:

```
julia> """
         Hello,
         world."""
"Hello,\nworld."
```

Trailing whitespace is left unaltered.

Triple-quoted string literals can contain `"` characters without escaping.

Note that line breaks in literal strings, whether single- or triple-quoted, result in a newline (LF) character `\n` in the string, even if your editor uses a carriage return `\r` (CR) or CRLF combination to end lines. To include a CR in a string, use an explicit escape `\r`; for example, you can enter the literal string `"a CRLF line ending\r\n"`.

# Common Operations

You can lexicographically compare strings using the standard comparison operators:

```
julia> "abracadabra" < "xylophone"
true

julia> "abracadabra" == "xylophone"
false

julia> "Hello, world." != "Goodbye, world."
true

julia> "1 + 2 = 3" == "1 + 2 = $(1 + 2)"
true
```

You can search for the index of a particular character using the `findfirst` and `findlast` functions:

```julia
julia> findfirst(isequal('o'), "xylophone")
4

julia> findlast(isequal('o'), "xylophone")
7

julia> findfirst(isequal('z'), "xylophone")
```

You can start the search for a character at a given offset by using the functions `findnext` and `findprev`:

```julia
julia> findnext(isequal('o'), "xylophone", 1)
4

julia> findnext(isequal('o'), "xylophone", 5)
7

julia> findprev(isequal('o'), "xylophone", 5)
4

julia> findnext(isequal('o'), "xylophone", 8)
```

You can use the `occursin` function to check if a substring is found within a string:

```julia
julia> occursin("world", "Hello, world.")
true

julia> occursin("o", "Xylophon")
true

julia> occursin("a", "Xylophon")
false

julia> occursin('o', "Xylophon")
true
```

The last example shows that `occursin` can also look for a character literal.

Two other handy string functions are `repeat` and `join`:

```julia
julia> repeat(".:Z:.", 10)
".:Z:...:Z:...:Z:...:Z:...:Z:...:Z:...:Z:...:Z:...:Z:...:Z:."

julia> join(["apples", "bananas", "pineapples"], ", ", " and ")
```

```
"apples, bananas and pineapples"
```

Some other useful functions include:

- `firstindex(str)` gives the minimal (byte) index that can be used to index into `str` (always 1 for strings, not necessarily true for other containers).
- `lastindex(str)` gives the maximal (byte) index that can be used to index into `str`.
- `length(str)` the number of characters in `str`.
- `length(str, i, j)` the number of valid character indices in `str` from `i` to `j`.
- `ncodeunits(str)` number of code units in a string.
- `codeunit(str, i)` gives the code unit value in the string `str` at index `i`.
- `thisind(str, i)` given an arbitrary index into a string find the first index of the character into which the index points.
- `nextind(str, i, n=1)` find the start of the `n`th character starting after index `i`.
- `prevind(str, i, n=1)` find the start of the `n`th character starting before index `i`.

# Non-Standard String Literals

There are situations when you want to construct a string or use string semantics, but the behavior of the standard string construct is not quite what is needed. For these kinds of situations, Julia provides non-standard string literals. A non-standard string literal looks like a regular double-quoted string literal, but is immediately prefixed by an identifier, and doesn't behave quite like a normal string literal. Regular expressions, byte array literals and version number literals, as described below, are some examples of non-standard string literals. Other examples are given in the Metaprogramming section.

# Regular Expressions

Julia has Perl-compatible regular expressions (regexes), as provided by the PCRE library (a description of the syntax can be found here). Regular expressions are related to strings in two ways: the obvious connection is that regular expressions are used to find regular patterns in strings; the other connection is that regular expressions are themselves input as strings, which are parsed into a state machine that can be used to efficiently search for patterns in strings. In Julia, regular expressions are input using non-standard string literals prefixed with various identifiers beginning with `r`. The most basic regular expression literal without any options turned on just uses `r"..."`:

```
julia> r"^\s*(?:#|$)"
r"^\s*(?:#|$)"
```

```
julia> typeof(ans)
Regex
```

To check if a regex matches a string, use `occursin`:

```
julia> occursin(r"^\s*(?:#|$)", "not a comment")
false

julia> occursin(r"^\s*(?:#|$)", "# a comment")
true
```

As one can see here, `occursin` simply returns true or false, indicating whether a match for the given regex occurs in the string. Commonly, however, one wants to know not just whether a string matched, but also *how* it matched. To capture this information about a match, use the `match` function instead:

```
julia> match(r"^\s*(?:#|$)", "not a comment")

julia> match(r"^\s*(?:#|$)", "# a comment")
RegexMatch("#")
```

If the regular expression does not match the given string, `match` returns `nothing` – a special value that does not print anything at the interactive prompt. Other than not printing, it is a completely normal value and you can test for it programmatically:

```
m = match(r"^\s*(?:#|$)", line)
if m === nothing
    println("not a comment")
else
    println("blank or comment")
end
```

If a regular expression does match, the value returned by `match` is a `RegexMatch` object. These objects record how the expression matches, including the substring that the pattern matches and any captured substrings, if there are any. This example only captures the portion of the substring that matches, but perhaps we want to capture any non-blank text after the comment character. We could do the following:

```
julia> m = match(r"^\s*(?:#\s*(.*?)\s*$|$)", "# a comment ")
RegexMatch("# a comment ", 1="a comment")
```

When calling `match`, you have the option to specify an index at which to start the search. For example:

```julia
julia> m = match(r"[0-9]","aaaa1aaaa2aaaa3",1)
RegexMatch("1")

julia> m = match(r"[0-9]","aaaa1aaaa2aaaa3",6)
RegexMatch("2")

julia> m = match(r"[0-9]","aaaa1aaaa2aaaa3",11)
RegexMatch("3")
```

You can extract the following info from a `RegexMatch` object:

- the entire substring matched: `m.match`
- the captured substrings as an array of strings: `m.captures`
- the offset at which the whole match begins: `m.offset`
- the offsets of the captured substrings as a vector: `m.offsets`

For when a capture doesn't match, instead of a substring, `m.captures` contains `nothing` in that position, and `m.offsets` has a zero offset (recall that indices in Julia are 1-based, so a zero offset into a string is invalid). Here is a pair of somewhat contrived examples:

```julia
julia> m = match(r"(a|b)(c)?(d)", "acd")
RegexMatch("acd", 1="a", 2="c", 3="d")

julia> m.match
"acd"

julia> m.captures
3-element Array{Union{Nothing, SubString{String}},1}:
 "a"
 "c"
 "d"

julia> m.offset
1

julia> m.offsets
3-element Array{Int64,1}:
 1
 2
 3

julia> m = match(r"(a|b)(c)?(d)", "ad")
RegexMatch("ad", 1="a", 2=nothing, 3="d")
```

```
julia> m.match
"ad"

julia> m.captures
3-element Array{Union{Nothing, SubString{String}},1}:
 "a"
 nothing
 "d"

julia> m.offset
1

julia> m.offsets
3-element Array{Int64,1}:
 1
 0
 2
```

It is convenient to have captures returned as an array so that one can use destructuring syntax to bind them to local variables:

```
julia> first, second, third = m.captures; first
"a"
```

Captures can also be accessed by indexing the `RegexMatch` object with the number or name of the capture group:

```
julia> m=match(r"(?<hour>\d+):(?<minute>\d+)","12:45")
RegexMatch("12:45", hour="12", minute="45")

julia> m[:minute]
"45"

julia> m[2]
"45"
```

Captures can be referenced in a substitution string when using `replace` by using `\n` to refer to the nth capture group and prefixing the substitution string with `s`. Capture group 0 refers to the entire match object. Named capture groups can be referenced in the substitution with `\g<groupname>`. For example:

```julia
julia> replace("first second", r"(\w+) (?<agroup>\w+)" => s"\g<agroup> \1")
"second first"
```

Numbered capture groups can also be referenced as `\g<n>` for disambiguation, as in:

```julia
julia> replace("a", r"." => s"\g<0>1")
"a1"
```

You can modify the behavior of regular expressions by some combination of the flags `i`, `m`, `s`, and `x` after the closing double quote mark. These flags have the same meaning as they do in Perl, as explained in this excerpt from the perlre manpage:

```
i   Do case-insensitive pattern matching.

    If locale matching rules are in effect, the case map is taken
    from the current locale for code points less than 255, and
    from Unicode rules for larger code points. However, matches
    that would cross the Unicode rules/non-Unicode rules boundary
    (ords 255/256) will not succeed.

m   Treat string as multiple lines.  That is, change "^" and "$"
    from matching the start or end of the string to matching the
    start or end of any line anywhere within the string.

s   Treat string as single line.  That is, change "." to match any
    character whatsoever, even a newline, which normally it would
    not match.

    Used together, as r""ms, they let the "." match any character
    whatsoever, while still allowing "^" and "$" to match,
    respectively, just after and just before newlines within the
    string.

x   Tells the regular expression parser to ignore most whitespace
    that is neither backslashed nor within a character class. You
    can use this to break up your regular expression into
    (slightly) more readable parts. The '#' character is also
    treated as a metacharacter introducing a comment, just as in
    ordinary code.
```

For example, the following regex has all three flags turned on:

```
julia> r"a+.*b+.*?d$"ism
r"a+.*b+.*?d$"ims

julia> match(r"a+.*b+.*?d$"ism, "Goodbye,\nOh, angry,\nBad world\n")
RegexMatch("angry,\nBad world")
```

The `r"..."` literal is constructed without interpolation and unescaping (except for quotation mark `"` which still has to be escaped). Here is an example showing the difference from standard string literals:

```
julia> x = 10
10

julia> r"$x"
r"$x"

julia> "$x"
"10"

julia> r"\x"
r"\x"

julia> "\x"
ERROR: syntax: invalid escape sequence
```

Triple-quoted regex strings, of the form `r"""..."""`, are also supported (and may be convenient for regular expressions containing quotation marks or newlines).

The `Regex()` constructor may be used to create a valid regex string programmatically. This permits using the contents of string variables and other string operations when constructing the regex string. Any of the regex codes above can be used within the single string argument to `Regex()`. Here are some examples:

```
julia> using Dates

julia> d = Date(1962,7,10)
1962-07-10

julia> regex_d = Regex("Day " * string(day(d)))
r"Day 10"

julia> match(regex_d, "It happened on Day 10")
RegexMatch("Day 10")
```

```
julia> name = "Jon"
"Jon"

julia> regex_name = Regex("[\"( ]$name[\") ]")  # interpolate value of name
r"[\"( ]Jon[\") ]"

julia> match(regex_name," Jon ")
RegexMatch(" Jon ")

julia> match(regex_name,"[Jon]") === nothing
true
```

# Byte Array Literals

Another useful non-standard string literal is the byte-array string literal: `b"..."`. This form lets you use string notation to express read only literal byte arrays – i.e. arrays of `UInt8` values. The type of those objects is `CodeUnits{UInt8, String}`. The rules for byte array literals are the following:

- ASCII characters and ASCII escapes produce a single byte.
- `\x` and octal escape sequences produce the *byte* corresponding to the escape value.
- Unicode escape sequences produce a sequence of bytes encoding that code point in UTF-8.

There is some overlap between these rules since the behavior of `\x` and octal escapes less than 0x80 (128) are covered by both of the first two rules, but here these rules agree. Together, these rules allow one to easily use ASCII characters, arbitrary byte values, and UTF-8 sequences to produce arrays of bytes. Here is an example using all three:

```
julia> b"DATA\xff\u2200"
8-element Base.CodeUnits{UInt8,String}:
 0x44
 0x41
 0x54
 0x41
 0xff
 0xe2
 0x88
 0x80
```

The ASCII string "DATA" corresponds to the bytes 68, 65, 84, 65. `\xff` produces the single byte 255. The Unicode escape `\u2200` is encoded in UTF-8 as the three bytes 226, 136, 128. Note that the resulting byte array does not correspond to a valid UTF-8 string:

```julia
julia> isvalid("DATA\xff\u2200")
false
```

As it was mentioned `CodeUnits{UInt8,String}` type behaves like read only array of `UInt8` and if you need a standard vector you can convert it using `Vector{UInt8}`:

```julia
julia> x = b"123"
3-element Base.CodeUnits{UInt8,String}:
 0x31
 0x32
 0x33

julia> x[1]
0x31

julia> x[1] = 0x32
ERROR: setindex! not defined for Base.CodeUnits{UInt8,String}
[...]

julia> Vector{UInt8}(x)
3-element Array{UInt8,1}:
 0x31
 0x32
 0x33
```

Also observe the significant distinction between `\xff` and `\uff`: the former escape sequence encodes the *byte 255*, whereas the latter escape sequence represents the *code point 255*, which is encoded as two bytes in UTF-8:

```julia
julia> b"\xff"
1-element Base.CodeUnits{UInt8,String}:
 0xff

julia> b"\uff"
2-element Base.CodeUnits{UInt8,String}:
 0xc3
 0xbf
```

Character literals use the same behavior.

For code points less than `\u80`, it happens that the UTF-8 encoding of each code point is just the single byte produced by the corresponding `\x` escape, so the distinction can safely be ignored. For the escapes

`\x80` through `\xff` as compared to `\u80` through `\uff`, however, there is a major difference: the former escapes all encode single bytes, which – unless followed by very specific continuation bytes – do not form valid UTF-8 data, whereas the latter escapes all represent Unicode code points with two-byte encodings.

If this is all extremely confusing, try reading "The Absolute Minimum Every Software Developer Absolutely, Positively Must Know About Unicode and Character Sets". It's an excellent introduction to Unicode and UTF-8, and may help alleviate some confusion regarding the matter.

# Version Number Literals

Version numbers can easily be expressed with non-standard string literals of the form `v"..."`. Version number literals create `VersionNumber` objects which follow the specifications of semantic versioning, and therefore are composed of major, minor and patch numeric values, followed by pre-release and build alpha-numeric annotations. For example, `v"0.2.1-rc1+win64"` is broken into major version 0, minor version 2, patch version 1, pre-release `rc1` and build `win64`. When entering a version literal, everything except the major version number is optional, therefore e.g. `v"0.2"` is equivalent to `v"0.2.0"` (with empty pre-release/build annotations), `v"2"` is equivalent to `v"2.0.0"`, and so on.

`VersionNumber` objects are mostly useful to easily and correctly compare two (or more) versions. For example, the constant `VERSION` holds Julia version number as a `VersionNumber` object, and therefore one can define some version-specific behavior using simple statements as:

```
if v"0.2" <= VERSION < v"0.3-"
    # do something specific to 0.2 release series
end
```

Note that in the above example the non-standard version number `v"0.3-"` is used, with a trailing `-`: this notation is a Julia extension of the standard, and it's used to indicate a version which is lower than any 0.3 release, including all of its pre-releases. So in the above example the code would only run with stable 0.2 versions, and exclude such versions as `v"0.3.0-rc1"`. In order to also allow for unstable (i.e. pre-release) 0.2 versions, the lower bound check should be modified like this: `v"0.2-" <= VERSION`.

Another non-standard version specification extension allows one to use a trailing `+` to express an upper limit on build versions, e.g. `VERSION > v"0.2-rc1+"` can be used to mean any version above 0.2-rc1 and any of its builds: it will return `false` for version `v"0.2-rc1+win64"` and `true` for `v"0.2-rc2"`.

It is good practice to use such special versions in comparisons (particularly, the trailing `-` should always be used on upper bounds unless there's a good reason not to), but they must not be used as the actual version number of anything, as they are invalid in the semantic versioning scheme.

Besides being used for the `VERSION` constant, `VersionNumber` objects are widely used in the `Pkg` module, to specify packages versions and their dependencies.

# Raw String Literals

Raw strings without interpolation or unescaping can be expressed with non-standard string literals of the form `raw"..."`. Raw string literals create ordinary `String` objects which contain the enclosed contents exactly as entered with no interpolation or unescaping. This is useful for strings which contain code or markup in other languages which use `$` or `\` as special characters.

The exception is that quotation marks still must be escaped, e.g. `raw"\""` is equivalent to `"\""`. To make it possible to express all strings, backslashes then also must be escaped, but only when appearing right before a quote character:

```
julia> println(raw"\\ \\\"")
\\ \"
```

Notice that the first two backslashes appear verbatim in the output, since they do not precede a quote character. However, the next backslash character escapes the backslash that follows it, and the last backslash escapes a quote, since these backslashes appear before a quote.

---