

Integers and Floating-Point Numbers

Integers and floating-point values are the basic building blocks of arithmetic and computation. Built-in representations of such values are called numeric primitives, while representations of integers and floating-point numbers as immediate values in code are known as numeric literals. For example, `1` is an integer literal, while `1.0` is a floating-point literal; their binary in-memory representations as objects are numeric primitives.

Julia provides a broad range of primitive numeric types, and a full complement of arithmetic and bitwise operators as well as standard mathematical functions are defined over them. These map directly onto numeric types and operations that are natively supported on modern computers, thus allowing Julia to take full advantage of computational resources. Additionally, Julia provides software support for [Arbitrary Precision Arithmetic](#), which can handle operations on numeric values that cannot be represented effectively in native hardware representations, but at the cost of relatively slower performance.

The following are Julia's primitive numeric types:

- Integer types:

Type	Signed?	Number of bits	Smallest value	Largest value
Int8	✓	8	-2^7	$2^7 - 1$
UInt8		8	0	$2^8 - 1$
Int16	✓	16	-2^{15}	$2^{15} - 1$
UInt16		16	0	$2^{16} - 1$
Int32	✓	32	-2^{31}	$2^{31} - 1$
UInt32		32	0	$2^{32} - 1$
Int64	✓	64	-2^{63}	$2^{63} - 1$
UInt64		64	0	$2^{64} - 1$
Int128	✓	128	-2^{127}	$2^{127} - 1$

<code>UInt128</code>		128	0	$2^{128} - 1$
<code>Bool</code>	N/A	8	false (0)	true (1)

- Floating-point types:

Type	Precision	Number of bits
<code>Float16</code>	half	16
<code>Float32</code>	single	32
<code>Float64</code>	double	64

Additionally, full support for [Complex and Rational Numbers](#) is built on top of these primitive numeric types. All numeric types interoperate naturally without explicit casting, thanks to a flexible, user-extensible [type promotion system](#).

Integers

Literal integers are represented in the standard manner:

```
julia> 1
1

julia> 1234
1234
```

The default type for an integer literal depends on whether the target system has a 32-bit architecture or a 64-bit architecture:

```
# 32-bit system:
julia> typeof(1)
Int32

# 64-bit system:
julia> typeof(1)
Int64
```

The Julia internal variable `Sys.WORD_SIZE` indicates whether the target system is 32-bit or 64-bit:

```
# 32-bit system:
julia> Sys.WORD_SIZE
32

# 64-bit system:
julia> Sys.WORD_SIZE
64
```

Julia also defines the types `Int` and `UInt`, which are aliases for the system's signed and unsigned native integer types respectively:

```
# 32-bit system:
julia> Int
Int32
julia> UInt
UInt32

# 64-bit system:
julia> Int
Int64
julia> UInt
UInt64
```

Larger integer literals that cannot be represented using only 32 bits but can be represented in 64 bits always create 64-bit integers, regardless of the system type:

```
# 32-bit or 64-bit system:
julia> typeof(3000000000)
Int64
```

Unsigned integers are input and output using the `0x` prefix and hexadecimal (base 16) digits `0-9a-f` (the capitalized digits `A-F` also work for input). The size of the unsigned value is determined by the number of hex digits used:

```
julia> 0x1
0x01

julia> typeof(ans)
UInt8

julia> 0x123
0x0123
```

```
julia> typeof(ans)
UInt16

julia> 0x1234567
0x01234567

julia> typeof(ans)
UInt32

julia> 0x123456789abcdef
0x0123456789abcdef

julia> typeof(ans)
UInt64

julia> 0x11112222333344445555666677778888
0x11112222333344445555666677778888

julia> typeof(ans)
UInt128
```

This behavior is based on the observation that when one uses unsigned hex literals for integer values, one typically is using them to represent a fixed numeric byte sequence, rather than just an integer value.

Recall that the variable `ans` is set to the value of the last expression evaluated in an interactive session. This does not occur when Julia code is run in other ways.

Binary and octal literals are also supported:

```
julia> 0b10
0x02

julia> typeof(ans)
UInt8

julia> 0o010
0x08

julia> typeof(ans)
UInt8

julia> 0x00000000000000001111222233334444
0x00000000000000001111222233334444
```

```
julia> typeof(ans)
UInt128
```

As for hexadecimal literals, binary and octal literals produce unsigned integer types. The size of the binary data item is the minimal needed size, if the leading digit of the literal is not 0. In the case of leading zeros, the size is determined by the minimal needed size for a literal, which has the same length but leading digit 1. That allows the user to control the size. Values which cannot be stored in `UInt128` cannot be written as such literals.

Binary, octal, and hexadecimal literals may be signed by a `-` immediately preceding the unsigned literal. They produce an unsigned integer of the same size as the unsigned literal would do, with the two's complement of the value:

```
julia> -0x2
0xfe

julia> -0x0002
0xffffe
```

The minimum and maximum representable values of primitive numeric types such as integers are given by the `typemin` and `typemax` functions:

```
julia> (typemin{Int32}, typemax{Int32})
(-2147483648, 2147483647)

julia> for T in [Int8, Int16, Int32, Int64, Int128, UInt8, UInt16, UInt32, UInt64, UInt128]
    println("$(lpad(T, 7)): [$(typemin(T)), $(typemax(T))]" )
end
Int8: [-128, 127]
Int16: [-32768, 32767]
Int32: [-2147483648, 2147483647]
Int64: [-9223372036854775808, 9223372036854775807]
Int128: [-170141183460469231731687303715884105728, 170141183460469231731687303715884105727]
UInt8: [0, 255]
UInt16: [0, 65535]
UInt32: [0, 4294967295]
UInt64: [0, 18446744073709551615]
UInt128: [0, 340282366920938463463374607431768211455]
```

The values returned by `typemin` and `typemax` are always of the given argument type. (The above expression uses several features that have yet to be introduced, including [for loops](#), [Strings](#), and

[Interpolation](#), but should be easy enough to understand for users with some existing programming experience.)

Overflow behavior

In Julia, exceeding the maximum representable value of a given type results in a wraparound behavior:

```
julia> x = typemax{Int64}
9223372036854775807

julia> x + 1
-9223372036854775808

julia> x + 1 == typemin{Int64}
true
```

Thus, arithmetic with Julia integers is actually a form of [modular arithmetic](#). This reflects the characteristics of the underlying arithmetic of integers as implemented on modern computers. In applications where overflow is possible, explicit checking for wraparound produced by overflow is essential; otherwise, the [BigInt](#) type in [Arbitrary Precision Arithmetic](#) is recommended instead.

An example of overflow behavior and how to potentially resolve it is as follows:

```
julia> 10^19
-8446744073709551616

julia> big(10)^19
10000000000000000000
```

Division errors

Integer division (the `div` function) has two exceptional cases: dividing by zero, and dividing the lowest negative number (`typemin`) by `-1`. Both of these cases throw a [DivideError](#). The remainder and modulus functions (`rem` and `mod`) throw a [DivideError](#) when their second argument is zero.

Floating-Point Numbers

Literal floating-point numbers are represented in the standard formats, using [E-notation](#) when necessary:

```
julia> 1.0
1.0

julia> 1.
1.0

julia> 0.5
0.5

julia> .5
0.5

julia> -1.23
-1.23

julia> 1e10
1.0e10

julia> 2.5e-4
0.00025
```

The above results are all `Float64` values. Literal `Float32` values can be entered by writing an `f` in place of `e`:

```
julia> 0.5f0
0.5f0

julia> typeof(ans)
Float32

julia> 2.5f-4
0.00025f0
```

Values can be converted to `Float32` easily:

```
julia> Float32(-1.5)
-1.5f0

julia> typeof(ans)
Float32
```

Hexadecimal floating-point literals are also valid, but only as `Float64` values, with `p` preceding the base-2 exponent:

```
julia> 0x1p0
1.0

julia> 0x1.8p3
12.0

julia> 0x.4p-1
0.125

julia> typeof(ans)
Float64
```

Half-precision floating-point numbers are also supported (`Float16`), but they are implemented in software and use `Float32` for calculations.

```
julia> sizeof(Float16(4.))
2

julia> 2*Float16(4.)
Float16(8.0)
```

The underscore `_` can be used as digit separator:

```
julia> 10_000, 0.000_000_005, 0xdead_beef, 0b1011_0010
(10000, 5.0e-9, 0xdeadbeef, 0xb2)
```

Floating-point zero

Floating-point numbers have `two zeros`, positive zero and negative zero. They are equal to each other but have different binary representations, as can be seen using the `bitstring` function:

[illegible]

Special floating-point values

There are three specified standard floating-point values that do not correspond to any point on the real number line:

Float16	Float32	Float64	Name	Description
Inf16	Inf32	Inf	positive infinity	a value greater than all finite floating-point values
-Inf16	-Inf32	-Inf	negative infinity	a value less than all finite floating-point values
NaN16	NaN32	NaN	not a number	a value not == to any floating-point value (including itself)

For further discussion of how these non-finite floating-point values are ordered with respect to each other and other floats, see [Numeric Comparisons](#). By the [IEEE 754 standard](#), these floating-point values are the results of certain arithmetic operations:

```
julia> 1/Inf
0.0

julia> 1/0
Inf

julia> -5/0
-Inf

julia> 0.000001/0
Inf

julia> 0/0
NaN

julia> 500 + Inf
Inf

julia> 500 - Inf
-Inf

julia> Inf + Inf
Inf
```

```
julia> Inf - Inf
NaN

julia> Inf * Inf
Inf

julia> Inf / Inf
NaN

julia> 0 * Inf
NaN
```

The `typemin` and `typemax` functions also apply to floating-point types:

```
julia> (typemin(Float16), typemax(Float16))
(-Inf16, Inf16)

julia> (typemin(Float32), typemax(Float32))
(-Inf32, Inf32)

julia> (typemin(Float64), typemax(Float64))
(-Inf, Inf)
```

Machine epsilon

Most real numbers cannot be represented exactly with floating-point numbers, and so for many purposes it is important to know the distance between two adjacent representable floating-point numbers, which is often known as [machine epsilon](#).

Julia provides `eps`, which gives the distance between `1.0` and the next larger representable floating-point value:

```
julia> eps(Float32)
1.1920929f-7

julia> eps(Float64)
2.220446049250313e-16

julia> eps() # same as eps(Float64)
2.220446049250313e-16
```

These values are 2.0^{-23} and 2.0^{-52} as `Float32` and `Float64` values, respectively. The `eps` function can also take a floating-point value as an argument, and gives the absolute difference between that value and the next representable floating point value. That is, `eps(x)` yields a value of the same type as `x` such that `x + eps(x)` is the next representable floating-point value larger than `x`:

```
julia> eps(1.0)
2.220446049250313e-16

julia> eps(1000.)
1.1368683772161603e-13

julia> eps(1e-27)
1.793662034335766e-43

julia> eps(0.0)
5.0e-324
```

The distance between two adjacent representable floating-point numbers is not constant, but is smaller for smaller values and larger for larger values. In other words, the representable floating-point numbers are densest in the real number line near zero, and grow sparser exponentially as one moves farther away from zero. By definition, `eps(1.0)` is the same as `eps(Float64)` since `1.0` is a 64-bit floating-point value.

Julia also provides the `nextfloat` and `prevfloat` functions which return the next largest or smallest representable floating-point number to the argument respectively:

```
julia> x = 1.25f0
1.25f0

julia> nextfloat(x)
1.2500001f0

julia> prevfloat(x)
1.2499999f0

julia> bitstring(prevfloat(x))
"00111111100111111111111111111111"

julia> bitstring(x)
"00111111101000000000000000000000"

julia> bitstring(nextfloat(x))
"00111111101000000000000000000001"
```

This example highlights the general principle that the adjacent representable floating-point numbers also have adjacent binary integer representations.

Rounding modes

If a number doesn't have an exact floating-point representation, it must be rounded to an appropriate representable value. However, the manner in which this rounding is done can be changed if required according to the rounding modes presented in the [IEEE 754 standard](#).

The default mode used is always [RoundNearest](#), which rounds to the nearest representable value, with ties rounded towards the nearest value with an even least significant bit.

Background and References

Floating-point arithmetic entails many subtleties which can be surprising to users who are unfamiliar with the low-level implementation details. However, these subtleties are described in detail in most books on scientific computation, and also in the following references:

- The definitive guide to floating point arithmetic is the [IEEE 754-2008 Standard](#); however, it is not available for free online.
- For a brief but lucid presentation of how floating-point numbers are represented, see John D. Cook's [article](#) on the subject as well as his [introduction](#) to some of the issues arising from how this representation differs in behavior from the idealized abstraction of real numbers.

- Also recommended is Bruce Dawson's [series of blog posts on floating-point numbers](#).
- For an excellent, in-depth discussion of floating-point numbers and issues of numerical accuracy encountered when computing with them, see David Goldberg's paper [What Every Computer Scientist Should Know About Floating-Point Arithmetic](#).
- For even more extensive documentation of the history of, rationale for, and issues with floating-point numbers, as well as discussion of many other topics in numerical computing, see the [collected writings](#) of [William Kahan](#), commonly known as the "Father of Floating-Point". Of particular interest may be [An Interview with the Old Man of Floating-Point](#).

Arbitrary Precision Arithmetic

To allow computations with arbitrary-precision integers and floating point numbers, Julia wraps the [GNU Multiple Precision Arithmetic Library \(GMP\)](#) and the [GNU MPFR Library](#), respectively. The `BigInt` and `BigFloat` types are available in Julia for arbitrary precision integer and floating point numbers respectively.

Constructors exist to create these types from primitive numerical types, and the [string literal @big_str](#) or [parse](#) can be used to construct them from AbstractStrings. Once created, they participate in arithmetic with all other numeric types thanks to Julia's [type promotion and conversion mechanism](#):

[illegible]

However, type promotion between the primitive types above and `BigInt/BigFloat` is not automatic and must be explicitly stated.

```
julia> x = typemin(Int64)
-9223372036854775808

julia> x = x - 1
9223372036854775807

julia> typeof(x)
Int64

julia> y = BigInt(typemin(Int64))
-9223372036854775808

julia> y = y - 1
-9223372036854775809

julia> typeof(y)
BigInt
```

The default precision (in number of bits of the significand) and rounding mode of `BigFloat` operations can be changed globally by calling `setprecision` and `setrounding`, and all further calculations will take these changes in account. Alternatively, the precision or the rounding can be changed only within the execution of a particular block of code by using the same functions with a `do` block:

[illegible]

Numeric Literal Coefficients

To make common numeric formulae and expressions clearer, Julia allows variables to be immediately preceded by a numeric literal, implying multiplication. This makes writing polynomial expressions much cleaner:

```
julia> x = 3
3

julia> 2x^2 - 3x + 1
10

julia> 1.5x^2 - .5x + 1
13.0
```

It also makes writing exponential functions more elegant:

```
julia> 2^2x
64
```

The precedence of numeric literal coefficients is slightly lower than that of unary operators such as negation. So $-2x$ is parsed as $(-2) * x$ and $\sqrt{2}x$ is parsed as $(\sqrt{2}) * x$. However, numeric literal coefficients parse similarly to unary operators when combined with exponentiation. For example 2^3x is parsed as $2^{(3x)}$, and $2x^3$ is parsed as $2*(x^3)$.

Numeric literals also work as coefficients to parenthesized expressions:

```
julia> 2(x-1)^2 - 3(x-1) + 1
3
```

Note

The precedence of numeric literal coefficients used for implicit multiplication is higher than other binary operators such as multiplication (`*`), and division (`/`, `\`, and `//`). This means, for example, that `1 / 2im` equals `-0.5im` and `6 // 2(2 + 1)` equals `1 // 1`.

Additionally, parenthesized expressions can be used as coefficients to variables, implying multiplication of the expression by the variable:

```
julia> (x-1)x
6
```

Neither juxtaposition of two parenthesized expressions, nor placing a variable before a parenthesized expression, however, can be used to imply multiplication:

```
julia> (x-1)(x+1)
ERROR: MethodError: objects of type Int64 are not callable

julia> x(x+1)
ERROR: MethodError: objects of type Int64 are not callable
```

Both expressions are interpreted as function application: any expression that is not a numeric literal, when immediately followed by a parenthetical, is interpreted as a function applied to the values in parentheses (see [Functions](#) for more about functions). Thus, in both of these cases, an error occurs since the left-hand value is not a function.

The above syntactic enhancements significantly reduce the visual noise incurred when writing common mathematical formulae. Note that no whitespace may come between a numeric literal coefficient and the identifier or parenthesized expression which it multiplies.

Syntax Conflicts

Juxtaposed literal coefficient syntax may conflict with two numeric literal syntaxes: hexadecimal integer literals and engineering notation for floating-point literals. Here are some situations where syntactic conflicts arise:

- The hexadecimal integer literal expression `0xff` could be interpreted as the numeric literal `0` multiplied by the variable `xff`.
- The floating-point literal expression `1e10` could be interpreted as the numeric literal `1` multiplied by the variable `e10`, and similarly with the equivalent `E` form.
- The 32-bit floating-point literal expression `1.5f22` could be interpreted as the numeric literal `1.5` multiplied by the variable `f22`.

In all cases the ambiguity is resolved in favor of interpretation as numeric literals:

- Expressions starting with `0x` are always hexadecimal literals.
- Expressions starting with a numeric literal followed by `e` or `E` are always floating-point literals.
- Expressions starting with a numeric literal followed by `f` are always 32-bit floating-point literals.

Unlike `E`, which is equivalent to `e` in numeric literals for historical reasons, `F` is just another letter and does not behave like `f` in numeric literals. Hence, expressions starting with a numeric literal followed by `F` are interpreted as the numerical literal multiplied by a variable, which means that, for example, `1.5F22` is equal to `1.5 * F22`.

Literal zero and one

Julia provides functions which return literal 0 and 1 corresponding to a specified type or the type of a given variable.

Function	Description
<code>zero(x)</code>	Literal zero of type <code>x</code> or type of variable <code>x</code>
<code>one(x)</code>	Literal one of type <code>x</code> or type of variable <code>x</code>

These functions are useful in [Numeric Comparisons](#) to avoid overhead from unnecessary [type conversion](#).

Examples:

```
julia> zero(Float32)
0.0f0

julia> zero(1.0)
0.0

julia> one(Int32)
1

julia> one(BigFloat)
1.0
```

« [Variables](#)

[Mathematical Operations and Elementary Functions](#) »

Powered by [Documenter.jl](#) and the [Julia Programming Language](#).