

# Performance Tips

In the following sections, we briefly go through a few techniques that can help make your Julia code run as fast as possible.

## Avoid global variables

A global variable might have its value, and therefore its type, change at any point. This makes it difficult for the compiler to optimize code using global variables. Variables should be local, or passed as arguments to functions, whenever possible.

Any code that is performance critical or being benchmarked should be inside a function.

We find that global names are frequently constants, and declaring them as such greatly improves performance:

```
const DEFAULT_VAL = 0
```

Uses of non-constant globals can be optimized by annotating their types at the point of use:

```
global x = rand(1000)

function loop_over_global()
    s = 0.0
    for i in x::Vector{Float64}
        s += i
    end
    return s
end
```

Passing arguments to functions is better style. It leads to more reusable code and clarifies what the inputs and outputs are.

### Note

All code in the REPL is evaluated in global scope, so a variable defined and assigned at top level will be a global variable. Variables defined at top level scope inside modules are also global.

In the following REPL session:

```
julia> x = 1.0
```

is equivalent to:

```
julia> global x = 1.0
```

so all the performance issues discussed previously apply.

## Measure performance with `@time` and pay attention to memory allocation

A useful tool for measuring performance is the `@time` macro. We here repeat the example with the global variable above, but this time with the type annotation removed:

```
julia> x = rand(1000);

julia> function sum_global()
    s = 0.0
    for i in x
        s += i
    end
    return s
end;

julia> @time sum_global()
0.017705 seconds (15.28 k allocations: 694.484 KiB)
496.84883432553846

julia> @time sum_global()
0.000140 seconds (3.49 k allocations: 70.313 KiB)
496.84883432553846
```

On the first call (`@time sum_global()`) the function gets compiled. (If you've not yet used `@time` in this session, it will also compile functions needed for timing.) You should not take the results of this run

seriously. For the second run, note that in addition to reporting the time, it also indicated that a significant amount of memory was allocated. We are here just computing a sum over all elements in a vector of 64-bit floats so there should be no need to allocate memory (at least not on the heap which is what `@time` reports).

Unexpected memory allocation is almost always a sign of some problem with your code, usually a problem with type-stability or creating many small temporary arrays. Consequently, in addition to the allocation itself, it's very likely that the code generated for your function is far from optimal. Take such indications seriously and follow the advice below.

If we instead pass `x` as an argument to the function it no longer allocates memory (the allocation reported below is due to running the `@time` macro in global scope) and is significantly faster after the first call:

```
julia> x = rand(1000);

julia> function sum_arg(x)
    s = 0.0
    for i in x
        s += i
    end
    return s
end;

julia> @time sum_arg(x)
0.007701 seconds (821 allocations: 43.059 KiB)
496.84883432553846

julia> @time sum_arg(x)
0.000006 seconds (5 allocations: 176 bytes)
496.84883432553846
```

The 5 allocations seen are from running the `@time` macro itself in global scope. If we instead run the timing in a function, we can see that indeed no allocations are performed:

```
julia> time_sum(x) = @time sum_arg(x);

julia> time_sum(x)
0.000001 seconds
496.84883432553846
```

In some situations, your function may need to allocate memory as part of its operation, and this can

complicate the simple picture above. In such cases, consider using one of the [tools](#) below to diagnose problems, or write a version of your function that separates allocation from its algorithmic aspects (see [Pre-allocating outputs](#)).

### Note

For more serious benchmarking, consider the [BenchmarkTools.jl](#) package which among other things evaluates the function multiple times in order to reduce noise.

## Tools

Julia and its package ecosystem includes tools that may help you diagnose problems and improve the performance of your code:

- [Profiling](#) allows you to measure the performance of your running code and identify lines that serve as bottlenecks. For complex projects, the [ProfileView](#) package can help you visualize your profiling results.
- The [Traceur](#) package can help you find common performance problems in your code.
- Unexpectedly-large memory allocations—as reported by `@time`, `@allocated`, or the profiler (through calls to the garbage-collection routines)—hint that there might be issues with your code. If you don't see another reason for the allocations, suspect a type problem. You can also start Julia with the `--track-allocation=user` option and examine the resulting `*.mem` files to see information about where those allocations occur. See [Memory allocation analysis](#).
- `@code_warntype` generates a representation of your code that can be helpful in finding expressions that result in type uncertainty. See [@code\\_warntype](#) below.

## Avoid containers with abstract type parameters

When working with parameterized types, including arrays, it is best to avoid parameterizing with abstract types where possible.

Consider the following:

```
julia> a = Real[]
Real[]

julia> push!(a, 1); push!(a, 2.0); push!(a, π)
3-element Array{Real,1}:
 1
```

```
2.0
π = 3.1415926535897...
```

Because `a` is an array of abstract type `Real`, it must be able to hold any `Real` value. Since `Real` objects can be of arbitrary size and structure, `a` must be represented as an array of pointers to individually allocated `Real` objects. However, if we instead only allow numbers of the same type, e.g. `Float64`, to be stored in `a` these can be stored more efficiently:

```
julia> a = Float64[]
Float64[]

julia> push!(a, 1); push!(a, 2.0); push!(a, π)
3-element Array{Float64,1}:
 1.0
 2.0
 3.141592653589793
```

Assigning numbers into `a` will now convert them to `Float64` and `a` will be stored as a contiguous block of 64-bit floating-point values that can be manipulated efficiently.

See also the discussion under [Parametric Types](#).

## Type declarations

In many languages with optional type declarations, adding declarations is the principal way to make code run faster. This is *not* the case in Julia. In Julia, the compiler generally knows the types of all function arguments, local variables, and expressions. However, there are a few specific instances where declarations are helpful.

### Avoid fields with abstract type

Types can be declared without specifying the types of their fields:

```
julia> struct MyAmbiguousType
    a
end
```

This allows `a` to be of any type. This can often be useful, but it does have a downside: for objects of type `MyAmbiguousType`, the compiler will not be able to generate high-performance code. The reason is that the compiler uses the types of objects, not their values, to determine how to build code. Unfortunately,

very little can be inferred about an object of type `MyAmbiguousType`:

```
julia> b = MyAmbiguousType("Hello")
MyAmbiguousType("Hello")

julia> c = MyAmbiguousType(17)
MyAmbiguousType(17)

julia> typeof(b)
MyAmbiguousType

julia> typeof(c)
MyAmbiguousType
```

The values of `b` and `c` have the same type, yet their underlying representation of data in memory is very different. Even if you stored just numeric values in field `a`, the fact that the memory representation of a `UInt8` differs from a `Float64` also means that the CPU needs to handle them using two different kinds of instructions. Since the required information is not available in the type, such decisions have to be made at run-time. This slows performance.

You can do better by declaring the type of `a`. Here, we are focused on the case where `a` might be any one of several types, in which case the natural solution is to use parameters. For example:

```
julia> mutable struct MyType{T<:AbstractFloat}
    a::T
end
```

This is a better choice than

```
julia> mutable struct MyStillAmbiguousType
    a::AbstractFloat
end
```

because the first version specifies the type of `a` from the type of the wrapper object. For example:

```
julia> m = MyType(3.2)
MyType{Float64}(3.2)

julia> t = MyStillAmbiguousType(3.2)
MyStillAmbiguousType(3.2)

julia> typeof(m)
```

```
MyType{Float64}

julia> typeof(t)
MyStillAmbiguousType
```

The type of field `a` can be readily determined from the type of `m`, but not from the type of `t`. Indeed, in `t` it's possible to change the type of the field `a`:

```
julia> typeof(t.a)
Float64

julia> t.a = 4.5f0
4.5f0

julia> typeof(t.a)
Float32
```

In contrast, once `m` is constructed, the type of `m.a` cannot change:

```
julia> m.a = 4.5f0
4.5f0

julia> typeof(m.a)
Float64
```

The fact that the type of `m.a` is known from `m`'s type—coupled with the fact that its type cannot change mid-function—allows the compiler to generate highly-optimized code for objects like `m` but not for objects like `t`.

Of course, all of this is true only if we construct `m` with a concrete type. We can break this by explicitly constructing it with an abstract type:

```
julia> m = MyType{AbstractFloat}(3.2)
MyType{AbstractFloat}(3.2)

julia> typeof(m.a)
Float64

julia> m.a = 4.5f0
4.5f0

julia> typeof(m.a)
Float32
```

For all practical purposes, such objects behave identically to those of `MyStillAmbiguousType`.

It's quite instructive to compare the sheer amount code generated for a simple function

```
func(m::MyType) = m.a+1
```

using

```
code_llvm(func, Tuple{MyType{Float64}})
code_llvm(func, Tuple{MyType{AbstractFloat}})
```

For reasons of length the results are not shown here, but you may wish to try this yourself. Because the type is fully-specified in the first case, the compiler doesn't need to generate any code to resolve the type at run-time. This results in shorter and faster code.

## Avoid fields with abstract containers

The same best practices also work for container types:

```
julia> struct MySimpleContainer{A<:AbstractVector}
    a::A
end

julia> struct MyAmbiguousContainer{T}
    a::AbstractVector{T}
end
```

For example:



```
julia> c = MySimpleContainer(1:3);

julia> typeof(c)
MySimpleContainer{UnitRange{Int64}}

julia> c = MySimpleContainer([1:3;]);

julia> typeof(c)
MySimpleContainer{Array{Int64,1}}

julia> b = MyAmbiguousContainer(1:3);

julia> typeof(b)
MyAmbiguousContainer{Int64}

julia> b = MyAmbiguousContainer([1:3;]);

julia> typeof(b)
MyAmbiguousContainer{Int64}
```

For `MySimpleContainer`, the object is fully-specified by its type and parameters, so the compiler can generate optimized functions. In most instances, this will probably suffice.

While the compiler can now do its job perfectly well, there are cases where *you* might wish that your code could do different things depending on the *element type* of `a`. Usually the best way to achieve this is to wrap your specific operation (here, `foo`) in a separate function:

```
julia> function sumfoo(c::MySimpleContainer)
    s = 0
    for x in c.a
        s += foo(x)
    end
    s
end
sumfoo (generic function with 1 method)

julia> foo(x::Integer) = x
foo (generic function with 1 method)

julia> foo(x::AbstractFloat) = round(x)
foo (generic function with 2 methods)
```

This keeps things simple, while allowing the compiler to generate optimized code in all cases.

However, there are cases where you may need to declare different versions of the outer function for different element types or types of the `AbstractVector` of the field `a` in `MySimpleContainer`. You could do it like this:

```
julia> function myfunc(c::MySimpleContainer{<:AbstractArray{<:Integer}})
    return c.a[1]+1
end
myfunc (generic function with 1 method)

julia> function myfunc(c::MySimpleContainer{<:AbstractArray{<:AbstractFloat}})
    return c.a[1]+2
end
myfunc (generic function with 2 methods)

julia> function myfunc(c::MySimpleContainer{Vector{T}}) where T <: Integer
    return c.a[1]+3
end
myfunc (generic function with 3 methods)
```

```
julia> myfunc(MySimpleContainer(1:3))
2

julia> myfunc(MySimpleContainer(1.0:3))
3.0

julia> myfunc(MySimpleContainer([1:3;]))
4
```

## Annotate values taken from untyped locations

It is often convenient to work with data structures that may contain values of any type (arrays of type `Array{Any}`). But, if you're using one of these structures and happen to know the type of an element, it helps to share this knowledge with the compiler:

```
function foo(a::Array{Any,1})
    x = a[1]::Int32
    b = x+1
    ...
end
```

Here, we happened to know that the first element of `a` would be an `Int32`. Making an annotation like

this has the added benefit that it will raise a run-time error if the value is not of the expected type, potentially catching certain bugs earlier.

In the case that the type of `a[1]` is not known precisely, `x` can be declared via `x = convert{Int32, a[1]}::Int32`. The use of the `convert` function allows `a[1]` to be any object convertible to an `Int32` (such as `UInt8`), thus increasing the genericity of the code by loosening the type requirement. Notice that `convert` itself needs a type annotation in this context in order to achieve type stability. This is because the compiler cannot deduce the type of the return value of a function, even `convert`, unless the types of all the function's arguments are known.

Type annotation will not enhance (and can actually hinder) performance if the type is constructed at run-time. This is because the compiler cannot use the annotation to specialize the subsequent code, and the type-check itself takes time. For example, in the code:

```
function nr(a, prec)
    ctype = prec == 32 ? Float32 : Float64
    b = Complex{ctype}(a)
    c = (b + 1.0f0)::Complex{ctype}
    abs(c)
end
```

the annotation of `c` harms performance. To write performant code involving types constructed at run-time, use the [function-barrier technique](#) discussed below, and ensure that the constructed type appears among the argument types of the kernel function so that the kernel operations are properly specialized by the compiler. For example, in the above snippet, as soon as `b` is constructed, it can be passed to another function `k`, the kernel. If, for example, function `k` declares `b` as an argument of type `Complex{T}`, where `T` is a type parameter, then a type annotation appearing in an assignment statement within `k` of the form:

```
c = (b + 1.0f0)::Complex{T}
```

does not hinder performance (but does not help either) since the compiler can determine the type of `c` at the time `k` is compiled.

## Be aware of when Julia avoids specializing

As a heuristic, Julia avoids automatically specializing on argument type parameters in three specific cases: `Type`, `Function`, and `Vararg`. Julia will always specialize when the argument is used within the method, but not if the argument is just passed through to another function. This usually has no performance impact at runtime and [improves compiler performance](#). If you find it does have a

performance impact at runtime in your case, you can trigger specialization by adding a type parameter to the method declaration. Here are some examples:

This will not specialize:

```
function f_type(t) # or t::Type
    x = ones(t, 10)
    return sum(map(sin, x))
end
```

but this will:

```
function g_type(t::Type{T}) where T
    x = ones(T, 10)
    return sum(map(sin, x))
end
```

These will not specialize:

```
f_func(f, num) = ntuple(f, div(num, 2))
g_func(g::Function, num) = ntuple(g, div(num, 2))
```

but this will:

```
h_func(h::H, num) where {H} = ntuple(h, div(num, 2))
```

This will not specialize:

```
f_vararg(x::Int...) = tuple(x...)
```

but this will:

```
g_vararg(x::Vararg{Int, N}) where {N} = tuple(x...)
```

One only needs to introduce a single type parameter to force specialization, even if the other types are unconstrained. For example, this will also specialize, and is useful when the arguments are not all of the same type:

```
h_vararg(x::Vararg{Any, N}) where {N} = tuple(x...)
```

Note that `@code_typed` and friends will always show you specialized code, even if Julia would not normally specialize that method call. You need to check the [method internals](#) if you want to see whether specializations are generated when argument types are changed, i.e., if `(@which f(...)).specializations` contains specializations for the argument in question.

## Break functions into multiple definitions

Writing a function as many small definitions allows the compiler to directly call the most applicable code, or even inline it.

Here is an example of a "compound function" that should really be written as multiple definitions:

```
using LinearAlgebra

function mynorm(A)
    if isa(A, Vector)
        return sqrt(real(dot(A,A)))
    elseif isa(A, Matrix)
        return maximum(svdvals(A))
    else
        error("mynorm: invalid argument")
    end
end
```

This can be written more concisely and efficiently as:

```
norm(x::Vector) = sqrt(real(dot(x, x)))
norm(A::Matrix) = maximum(svdvals(A))
```

It should however be noted that the compiler is quite efficient at optimizing away the dead branches in code written as the `mynorm` example.

## Write "type-stable" functions

When possible, it helps to ensure that a function always returns a value of the same type. Consider the following definition:

```
pos(x) = x < 0 ? 0 : x
```

Although this seems innocent enough, the problem is that `0` is an integer (of type `Int`) and `x` might be of

any type. Thus, depending on the value of `x`, this function might return a value of either of two types. This behavior is allowed, and may be desirable in some cases. But it can easily be fixed as follows:

```
pos(x) = x < 0 ? zero(x) : x
```

There is also a `oneunit` function, and a more general `oftype(x, y)` function, which returns `y` converted to the type of `x`.

## Avoid changing the type of a variable

An analogous "type-stability" problem exists for variables used repeatedly within a function:

```
function foo()
    x = 1
    for i = 1:10
        x /= rand()
    end
    return x
end
```

Local variable `x` starts as an integer, and after one loop iteration becomes a floating-point number (the result of `/` operator). This makes it more difficult for the compiler to optimize the body of the loop.

There are several possible fixes:

- Initialize `x` with `x = 1.0`
- Declare the type of `x` explicitly as `x::Float64 = 1`
- Use an explicit conversion by `x = oneunit{Float64}()`
- Initialize with the first loop iteration, to `x = 1 / rand()`, then loop for `i = 2:10`

## Separate kernel functions (aka, function barriers)

Many functions follow a pattern of performing some set-up work, and then running many iterations to perform a core computation. Where possible, it is a good idea to put these core computations in separate functions. For example, the following contrived function returns an array of a randomly-chosen type:

```
julia> function strange_twos(n)
    a = Vector{rand{Bool} ? Int64 : Float64}(undef, n)
    for i = 1:n
        a[i] = 2
    end
end
```

```
        end
        return a
    end;

julia> strange_twos(3)
3-element Array{Float64,1}:
 2.0
 2.0
 2.0
```

This should be written as:

```
julia> function fill_twos!(a)
    for i = eachindex(a)
        a[i] = 2
    end
end;

julia> function strange_twos(n)
    a = Vector{rand{Bool} ? Int64 : Float64}(undef, n)
    fill_twos!(a)
    return a
end;

julia> strange_twos(3)
3-element Array{Float64,1}:
 2.0
 2.0
 2.0
```

Julia's compiler specializes code for argument types at function boundaries, so in the original implementation it does not know the type of `a` during the loop (since it is chosen randomly). Therefore the second version is generally faster since the inner loop can be recompiled as part of `fill_twos!` for different types of `a`.

The second form is also often better style and can lead to more code reuse.

This pattern is used in several places in Julia Base. For example, see `vcut` and `hcat` in [abstractarray.jl](#), or the `fill!` function, which we could have used instead of writing our own `fill_twos!`.

Functions like `strange_twos` occur when dealing with data of uncertain type, for example data loaded from an input file that might contain either integers, floats, strings, or something else.

# Types with values-as-parameters

Let's say you want to create an N-dimensional array that has size 3 along each axis. Such arrays can be created like this:

```
julia> A = fill(5.0, (3, 3))
3×3 Array{Float64,2}:
 5.0  5.0  5.0
 5.0  5.0  5.0
 5.0  5.0  5.0
```

This approach works very well: the compiler can figure out that A is an `Array{Float64,2}` because it knows the type of the fill value (`5.0::Float64`) and the dimensionality (`(3, 3)::NTuple{2,Int}`). This implies that the compiler can generate very efficient code for any future usage of A in the same function.

But now let's say you want to write a function that creates a `3×3×...` array in arbitrary dimensions; you might be tempted to write a function

```
julia> function array3(fillval, N)
    fill(fillval, ntuple(d->3, N))
end
array3 (generic function with 1 method)

julia> array3(5.0, 2)
3×3 Array{Float64,2}:
 5.0  5.0  5.0
 5.0  5.0  5.0
 5.0  5.0  5.0
```

This works, but (as you can verify for yourself using `@code_warntype array3(5.0, 2)`) the problem is that the output type cannot be inferred: the argument N is a *value* of type `Int`, and type-inference does not (and cannot) predict its value in advance. This means that code using the output of this function has to be conservative, checking the type on each access of A; such code will be very slow.

Now, one very good way to solve such problems is by using the [function-barrier technique](#). However, in some cases you might want to eliminate the type-instability altogether. In such cases, one approach is to pass the dimensionality as a parameter, for example through `Val{T}()` (see "[Value types](#)"):

```
julia> function array3(fillval, ::Val{N}) where N
    fill(fillval, ntuple(d->3, Val(N)))
end
array3 (generic function with 1 method)
```



```
julia> array3(5.0, Val(2))
3×3 Array{Float64,2}:
 5.0  5.0  5.0
 5.0  5.0  5.0
 5.0  5.0  5.0
```

Julia has a specialized version of `ntuple` that accepts a `Val{::Int}` instance as the second parameter; by passing `N` as a type-parameter, you make its "value" known to the compiler. Consequently, this version of `array3` allows the compiler to predict the return type.

However, making use of such techniques can be surprisingly subtle. For example, it would be of no help if you called `array3` from a function like this:

```
function call_array3(fillval, n)
    A = array3(fillval, Val(n))
end
```

Here, you've created the same problem all over again: the compiler can't guess what `n` is, so it doesn't know the *type* of `Val(n)`. Attempting to use `Val`, but doing so incorrectly, can easily make performance *worse* in many situations. (Only in situations where you're effectively combining `Val` with the function-barrier trick, to make the kernel function more efficient, should code like the above be used.)

An example of correct usage of `Val` would be:

```
function filter3(A::AbstractArray{T,N}) where {T,N}
    kernel = array3(1, Val(N))
    filter(A, kernel)
end
```

In this example, `N` is passed as a parameter, so its "value" is known to the compiler. Essentially, `Val(T)` works only when `T` is either hard-coded/literal (`Val(3)`) or already specified in the type-domain.

## The dangers of abusing multiple dispatch (aka, more on types with values-as-parameters)

Once one learns to appreciate multiple dispatch, there's an understandable tendency to go overboard and try to use it for everything. For example, you might imagine using it to store information, e.g.

```
struct Car{Make, Model}
    year::Int
```

```
...more fields...  
end
```

and then dispatch on objects like `Car{ :Honda, :Accord }(year, args...)`.

This might be worthwhile when either of the following are true:

- You require CPU-intensive processing on each `Car`, and it becomes vastly more efficient if you know the `Make` and `Model` at compile time and the total number of different `Make` or `Model` that will be used is not too large.
- You have homogenous lists of the same type of `Car` to process, so that you can store them all in an `Array{Car{ :Honda, :Accord }, N}`.

When the latter holds, a function processing such a homogenous array can be productively specialized: Julia knows the type of each element in advance (all objects in the container have the same concrete type), so Julia can "look up" the correct method calls when the function is being compiled (obviating the need to check at run-time) and thereby emit efficient code for processing the whole list.

When these do not hold, then it's likely that you'll get no benefit; worse, the resulting "combinatorial explosion of types" will be counterproductive. If `items[i+1]` has a different type than `item[i]`, Julia has to look up the type at run-time, search for the appropriate method in method tables, decide (via type intersection) which one matches, determine whether it has been JIT-compiled yet (and do so if not), and then make the call. In essence, you're asking the full type- system and JIT-compilation machinery to basically execute the equivalent of a switch statement or dictionary lookup in your own code.

Some run-time benchmarks comparing (1) type dispatch, (2) dictionary lookup, and (3) a "switch" statement can be found [on the mailing list](#).

Perhaps even worse than the run-time impact is the compile-time impact: Julia will compile specialized functions for each different `Car{Make, Model}`; if you have hundreds or thousands of such types, then every function that accepts such an object as a parameter (from a custom `get_year` function you might write yourself, to the generic `push!` function in Julia Base) will have hundreds or thousands of variants compiled for it. Each of these increases the size of the cache of compiled code, the length of internal lists of methods, etc. Excess enthusiasm for values-as-parameters can easily waste enormous resources.

## Access arrays in memory order, along columns

Multidimensional arrays in Julia are stored in column-major order. This means that arrays are stacked one column at a time. This can be verified using the `vec` function or the syntax `[ : ]` as shown below (notice that the array is ordered `[ 1 3 2 4 ]`, not `[ 1 2 3 4 ]`):

```
julia> x = [1 2; 3 4]
2×2 Array{Int64,2}:
 1  2
 3  4

julia> x[:]
4-element Array{Int64,1}:
 1
 3
 2
 4
```

This convention for ordering arrays is common in many languages like Fortran, Matlab, and R (to name a few). The alternative to column-major ordering is row-major ordering, which is the convention adopted by C and Python (numpy) among other languages. Remembering the ordering of arrays can have significant performance effects when looping over arrays. A rule of thumb to keep in mind is that with column-major arrays, the first index changes most rapidly. Essentially this means that looping will be faster if the inner-most loop index is the first to appear in a slice expression. Keep in mind that indexing an array with `:` is an implicit loop that iteratively accesses all elements within a particular dimension; it can be faster to extract columns than rows, for example.

Consider the following contrived example. Imagine we wanted to write a function that accepts a [Vector](#) and returns a square [Matrix](#) with either the rows or the columns filled with copies of the input vector. Assume that it is not important whether rows or columns are filled with these copies (perhaps the rest of the code can be easily adapted accordingly). We could conceivably do this in at least four ways (in addition to the recommended call to the built-in [repeat](#)):

```
function copy_cols(x::Vector{T}) where T
    inds = axes(x, 1)
    out = similar{Array{T}, inds, inds}
    for i = inds
        out[:, i] = x
    end
    return out
end

function copy_rows(x::Vector{T}) where T
    inds = axes(x, 1)
    out = similar{Array{T}, inds, inds}
    for i = inds
        out[i, :] = x
    end
end
```

```

        return out
    end

    function copy_col_row(x::Vector{T}) where T
        inds = axes(x, 1)
        out = similar(Array{T}, inds, inds)
        for col = inds, row = inds
            out[row, col] = x[row]
        end
        return out
    end

    function copy_row_col(x::Vector{T}) where T
        inds = axes(x, 1)
        out = similar(Array{T}, inds, inds)
        for row = inds, col = inds
            out[row, col] = x[col]
        end
        return out
    end
end

```

Now we will time each of these functions using the same random 10000 by 1 input vector:

```

julia> x = randn(10000);

julia> fmt(f) = println(rpad(string(f)*": ", 14, ' '), @elapsed f(x))

julia> map(fmt, [copy_cols, copy_rows, copy_col_row, copy_row_col]);
copy_cols:      0.331706323
copy_rows:      1.799009911
copy_col_row:   0.415630047
copy_row_col:   1.721531501

```

Notice that `copy_cols` is much faster than `copy_rows`. This is expected because `copy_cols` respects the column-based memory layout of the `Matrix` and fills it one column at a time. Additionally, `copy_col_row` is much faster than `copy_row_col` because it follows our rule of thumb that the first element to appear in a slice expression should be coupled with the inner-most loop.

## Pre-allocating outputs

If your function returns an `Array` or some other complex type, it may have to allocate memory.

Unfortunately, oftentimes allocation and its converse, garbage collection, are substantial bottlenecks.

Sometimes you can circumvent the need to allocate memory on each function call by preallocating the output. As a trivial example, compare

```
julia> function xinc(x)
    return [x, x+1, x+2]
end;

julia> function loopinc()
    y = 0
    for i = 1:10^7
        ret = xinc(i)
        y += ret[2]
    end
    return y
end;
```

with

```
julia> function xinc!(ret::AbstractVector{T}, x::T) where T
    ret[1] = x
    ret[2] = x+1
    ret[3] = x+2
    nothing
end;

julia> function loopinc_prealloc()
    ret = Vector{Int}(undef, 3)
    y = 0
    for i = 1:10^7
        xinc!(ret, i)
        y += ret[2]
    end
    return y
end;
```

Timing results:

```
julia> @time loopinc()
0.529894 seconds (40.00 M allocations: 1.490 GiB, 12.14% gc time)
50000015000000

julia> @time loopinc_prealloc()
0.030850 seconds (6 allocations: 288 bytes)
```

```
50000015000000
```

Preallocation has other advantages, for example by allowing the caller to control the "output" type from an algorithm. In the example above, we could have passed a `SubArray` rather than an `Array`, had we so desired.

Taken to its extreme, pre-allocation can make your code uglier, so performance measurements and some judgment may be required. However, for "vectorized" (element-wise) functions, the convenient syntax `x .= f.(y)` can be used for in-place operations with fused loops and no temporary arrays (see the [dot syntax for vectorizing functions](#)).

## More dots: Fuse vectorized operations

Julia has a special [dot syntax](#) that converts any scalar function into a "vectorized" function call, and any operator into a "vectorized" operator, with the special property that nested "dot calls" are *fusing*: they are combined at the syntax level into a single loop, without allocating temporary arrays. If you use `.=` and similar assignment operators, the result can also be stored in-place in a pre-allocated array (see above).

In a linear-algebra context, this means that even though operations like `vector + vector` and `vector * scalar` are defined, it can be advantageous to instead use `vector .+ vector` and `vector .* scalar` because the resulting loops can be fused with surrounding computations. For example, consider the two functions:

```
julia> f(x) = 3x.^2 + 4x + 7x.^3;

julia> fdot(x) = @. 3x^2 + 4x + 7x^3 # equivalent to 3 .* x.^2 .+ 4 .* x .+ 7 .* x.^
```

Both `f` and `fdot` compute the same thing. However, `fdot` (defined with the help of the `@.` macro) is significantly faster when applied to an array:

```
julia> x = rand(10^6);

julia> @time f(x);
0.019049 seconds (16 allocations: 45.777 MiB, 18.59% gc time)

julia> @time fdot(x);
0.002790 seconds (6 allocations: 7.630 MiB)

julia> @time f.(x);
0.002626 seconds (8 allocations: 7.630 MiB)
```

That is, `fdot(x)` is ten times faster and allocates 1/6 the memory of `f(x)`, because each `*` and `+` operation in `f(x)` allocates a new temporary array and executes in a separate loop. (Of course, if you just do `f.(x)` then it is as fast as `fdot(x)` in this example, but in many contexts it is more convenient to just sprinkle some dots in your expressions rather than defining a separate function for each vectorized operation.)

## Consider using views for slices

In Julia, an array "slice" expression like `array[1:5, :]` creates a copy of that data (except on the left-hand side of an assignment, where `array[1:5, :] = ...` assigns in-place to that portion of `array`). If you are doing many operations on the slice, this can be good for performance because it is more efficient to work with a smaller contiguous copy than it would be to index into the original array. On the other hand, if you are just doing a few simple operations on the slice, the cost of the allocation and copy operations can be substantial.

An alternative is to create a "view" of the array, which is an array object (a `SubArray`) that actually references the data of the original array in-place, without making a copy. (If you write to a view, it modifies the original array's data as well.) This can be done for individual slices by calling `view`, or more simply for a whole expression or block of code by putting `@views` in front of that expression. For example:

```
julia> fcopy(x) = sum(x[2:end-1]);

julia> @views fview(x) = sum(x[2:end-1]);

julia> x = rand(10^6);

julia> @time fcopy(x);
0.003051 seconds (7 allocations: 7.630 MB)

julia> @time fview(x);
0.001020 seconds (6 allocations: 224 bytes)
```

Notice both the 3× speedup and the decreased memory allocation of the `fview` version of the function.

## Copying data is not always bad

Arrays are stored contiguously in memory, lending themselves to CPU vectorization and fewer memory accesses due to caching. These are the same reasons that it is recommended to access arrays in column-major order (see above). Irregular access patterns and non-contiguous views can drastically slow down

computations on arrays because of non-sequential memory access.

Copying irregularly-accessed data into a contiguous array before operating on it can result in a large speedup, such as in the example below. Here, a matrix and a vector are being accessed at 800,000 of their randomly-shuffled indices before being multiplied. Copying the views into plain arrays speeds up the multiplication even with the cost of the copying operation.

```
julia> using Random

julia> x = randn(1_000_000);

julia> inds = shuffle(1:1_000_000)[1:800000];

julia> A = randn(50, 1_000_000);

julia> xtmp = zeros(800_000);

julia> Atmp = zeros(50, 800_000);

julia> @time sum(view(A, :, inds) * view(x, inds))
0.412156 seconds (14 allocations: 960 bytes)
-4256.759568345458

julia> @time begin
        copyto!(xtmp, view(x, inds))
        copyto!(Atmp, view(A, :, inds))
        sum(Atmp * xtmp)
    end
0.285923 seconds (14 allocations: 960 bytes)
-4256.759568345134
```

Provided there is enough memory for the copies, the cost of copying the view to an array is far outweighed by the speed boost from doing the matrix multiplication on a contiguous array.

## Avoid string interpolation for I/O

When writing data to a file (or other I/O device), forming extra intermediate strings is a source of overhead. Instead of:

```
println(file, "$a $b")
```

use:



```
println(file, a, " ", b)
```

The first version of the code forms a string, then writes it to the file, while the second version writes values directly to the file. Also notice that in some cases string interpolation can be harder to read. Consider:

```
println(file, "$(f(a))$(f(b))")
```

versus:

```
println(file, f(a), f(b))
```

## Optimize network I/O during parallel execution

When executing a remote function in parallel:

```
using Distributed

responses = Vector{Any}(undef, nworkers())
@sync begin
    for (idx, pid) in enumerate(workers())
        @async responses[idx] = remotecall_fetch(foo, pid, args...)
    end
end
```

is faster than:

```
using Distributed

refs = Vector{Any}(undef, nworkers())
for (idx, pid) in enumerate(workers())
    refs[idx] = @spawnat pid foo(args...)
end
responses = [fetch(r) for r in refs]
```

The former results in a single network round-trip to every worker, while the latter results in two network calls - first by the `@spawnat` and the second due to the `fetch` (or even a `wait`). The `fetch/wait` is also being executed serially resulting in an overall poorer performance.

# Fix deprecation warnings

A deprecated function internally performs a lookup in order to print a relevant warning only once. This extra lookup can cause a significant slowdown, so all uses of deprecated functions should be modified as suggested by the warnings.

## Tweaks

These are some minor points that might help in tight inner loops.

- Avoid unnecessary arrays. For example, instead of `sum([x, y, z])` use `x+y+z`.
- Use `abs2(z)` instead of `abs(z)^2` for complex `z`. In general, try to rewrite code to use `abs2` instead of `abs` for complex arguments.
- Use `div(x, y)` for truncating division of integers instead of `trunc(x/y)`, `fld(x, y)` instead of `floor(x/y)`, and `cld(x, y)` instead of `ceil(x/y)`.

## Performance Annotations

Sometimes you can enable better optimization by promising certain program properties.

- Use `@inbounds` to eliminate array bounds checking within expressions. Be certain before doing this. If the subscripts are ever out of bounds, you may suffer crashes or silent corruption.
- Use `@fastmath` to allow floating point optimizations that are correct for real numbers, but lead to differences for IEEE numbers. Be careful when doing this, as this may change numerical results. This corresponds to the `-ffast-math` option of clang.
- Write `@simd` in front of `for` loops to promise that the iterations are independent and may be reordered. Note that in many cases, Julia can automatically vectorize code without the `@simd` macro; it is only beneficial in cases where such a transformation would otherwise be illegal, including cases like allowing floating-point re-associativity and ignoring dependent memory accesses (`@simd ivdep`). Again, be very careful when asserting `@simd` as erroneously annotating a loop with dependent iterations may result in unexpected results. In particular, note that `setindex!` on some `AbstractArray` subtypes is inherently dependent upon iteration order. This feature is experimental and could change or disappear in future versions of Julia.

The common idiom of using `1:n` to index into an `AbstractArray` is not safe if the `Array` uses unconventional indexing, and may cause a segmentation fault if bounds checking is turned off. Use `LinearIndices(x)` or `eachindex(x)` instead (see also [Arrays with custom indices](#)).

### Note

While `@simd` needs to be placed directly in front of an innermost `for` loop, both `@inbounds` and `@fastmath` can be applied to either single expressions or all the expressions that appear within nested blocks of code, e.g., using `@inbounds begin` or `@inbounds for ....`

Here is an example with both `@inbounds` and `@simd` markup (we here use `@noinline` to prevent the optimizer from trying to be too clever and defeat our benchmark):

```
@noinline function inner(x, y)
    s = zero(eltype(x))
    for i=eachindex(x)
        @inbounds s += x[i]*y[i]
    end
    return s
end

@noinline function innersimd(x, y)
    s = zero(eltype(x))
    @simd for i = eachindex(x)
        @inbounds s += x[i] * y[i]
    end
    return s
end

function timeit(n, reps)
    x = rand{Float32}(n)
    y = rand{Float32}(n)
    s = zero{Float64}
    time = @elapsed for j in 1:reps
        s += inner(x, y)
    end
    println("GFlop/sec          = ", 2n*reps / time*1E-9)
    time = @elapsed for j in 1:reps
        s += innersimd(x, y)
    end
    println("GFlop/sec (SIMD) = ", 2n*reps / time*1E-9)
end

timeit(1000, 1000)
```

On a computer with a 2.4GHz Intel Core i5 processor, this produces:

```
GFlop/sec          = 1.9467069505224963
```

```
GFlop/sec (SIMD) = 17.578554163920018
```

(GFlop/sec measures the performance, and larger numbers are better.)

Here is an example with all three kinds of markup. This program first calculates the finite difference of a one-dimensional array, and then evaluates the L2-norm of the result:

```
function init!(u::Vector)
    n = length(u)
    dx = 1.0 / (n-1)
    @fastmath @inbounds @simd for i in 1:n #by asserting that `u` is a `Vector` we c
        u[i] = sin(2pi*dx*i)
    end
end

function deriv!(u::Vector, du)
    n = length(u)
    dx = 1.0 / (n-1)
    @fastmath @inbounds du[1] = (u[2] - u[1]) / dx
    @fastmath @inbounds @simd for i in 2:n-1
        du[i] = (u[i+1] - u[i-1]) / (2*dx)
    end
    @fastmath @inbounds du[n] = (u[n] - u[n-1]) / dx
end

function mynorm(u::Vector)
    n = length(u)
    T = eltype(u)
    s = zero(T)
    @fastmath @inbounds @simd for i in 1:n
        s += u[i]^2
    end
    @fastmath @inbounds return sqrt(s)
end

function main()
    n = 2000
    u = Vector{Float64}(undef, n)
    init!(u)
    du = similar(u)

    deriv!(u, du)
    nu = mynorm(du)
end
```

```
@time for i in 1:10^6
    deriv!(u, du)
    nu = mynorm(du)
end

println(nu)
end

main()
```

On a computer with a 2.7 GHz Intel Core i7 processor, this produces:

```
$ julia wave.jl;
 1.207814709 seconds
4.443986180758249

$ julia --math-mode=ieee wave.jl;
 4.487083643 seconds
4.443986180758249
```

Here, the option `--math-mode=ieee` disables the `@fastmath` macro, so that we can compare results.

In this case, the speedup due to `@fastmath` is a factor of about 3.7. This is unusually large – in general, the speedup will be smaller. (In this particular example, the working set of the benchmark is small enough to fit into the L1 cache of the processor, so that memory access latency does not play a role, and computing time is dominated by CPU usage. In many real world programs this is not the case.) Also, in this case this optimization does not change the result – in general, the result will be slightly different. In some cases, especially for numerically unstable algorithms, the result can be very different.

The annotation `@fastmath` re-arranges floating point expressions, e.g. changing the order of evaluation, or assuming that certain special cases (inf, nan) cannot occur. In this case (and on this particular computer), the main difference is that the expression  $1 / (2 * dx)$  in the function `deriv` is hoisted out of the loop (i.e. calculated outside the loop), as if one had written `idx = 1 / (2 * dx)`. In the loop, the expression  $\dots / (2 * dx)$  then becomes  $\dots * idx$ , which is much faster to evaluate. Of course, both the actual optimization that is applied by the compiler as well as the resulting speedup depend very much on the hardware. You can examine the change in generated code by using Julia's `code_native` function.

Note that `@fastmath` also assumes that NaNs will not occur during the computation, which can lead to surprising behavior:

```
julia> f(x) = isnan(x);
```

```
julia> f(NaN)
true

julia> f_fast(x) = @fastmath isnan(x);

julia> f_fast(NaN)
false
```

## Treat Subnormal Numbers as Zeros

Subnormal numbers, formerly called [denormal numbers](#), are useful in many contexts, but incur a performance penalty on some hardware. A call `set_zero_subnormals(true)` grants permission for floating-point operations to treat subnormal inputs or outputs as zeros, which may improve performance on some hardware. A call `set_zero_subnormals(false)` enforces strict IEEE behavior for subnormal numbers.

Below is an example where subnormals noticeably impact performance on some hardware:

```

function timestep(b::Vector{T}, a::Vector{T}, Δt::T) where T
    @assert length(a)==length(b)
    n = length(b)
    b[1] = 1                                # Boundary condition
    for i=2:n-1
        b[i] = a[i] + (a[i-1] - T(2)*a[i] + a[i+1]) * Δt
    end
    b[n] = 0                                # Boundary condition
end

function heatflow(a::Vector{T}, nstep::Integer) where T
    b = similar(a)
    for t=1:div(nstep,2)                    # Assume nstep is even
        timestep(b,a,T(0.1))
        timestep(a,b,T(0.1))
    end
end

heatflow(zeros(Float32,10),2)              # Force compilation
for trial=1:6
    a = zeros(Float32,1000)
    set_zero_subnormals(iseven(trial))     # Odd trials use strict IEEE arithmetic
    @time heatflow(a,1000)
end

```

This gives an output similar to

```

0.002202 seconds (1 allocation: 4.063 KiB)
0.001502 seconds (1 allocation: 4.063 KiB)
0.002139 seconds (1 allocation: 4.063 KiB)
0.001454 seconds (1 allocation: 4.063 KiB)
0.002115 seconds (1 allocation: 4.063 KiB)
0.001455 seconds (1 allocation: 4.063 KiB)

```

Note how each even iteration is significantly faster.

This example generates many subnormal numbers because the values in `a` become an exponentially decreasing curve, which slowly flattens out over time.

Treating subnormals as zeros should be used with caution, because doing so breaks some identities, such as `x-y == 0` implies `x == y`:

```
julia> x = 3f-38; y = 2f-38;
```

```
julia> set_zero_subnormals(true); (x - y, x == y)
(0.0f0, false)

julia> set_zero_subnormals(false); (x - y, x == y)
(1.0000001f-38, false)
```

In some applications, an alternative to zeroing subnormal numbers is to inject a tiny bit of noise. For example, instead of initializing `a` with zeros, initialize it with:

```
a = rand(Float32, 1000) * 1.f-9
```

## @code\_warntype

The macro `@code_warntype` (or its function variant `code_warntype`) can sometimes be helpful in diagnosing type-related problems. Here's an example:

```
julia> @noinline pos(x) = x < 0 ? 0 : x;

julia> function f(x)
    y = pos(x)
    return sin(y*x + 1)
end;

julia> @code_warntype f(3.2)
Variables
  #self#::Core.Compiler.Const(f, false)
  x::Float64
  y::UNION{Float64, Int64}

Body::Float64
1 -      (y = Main.pos(x))
|   %2 = (y * x)::Float64
|   %3 = (%2 + 1)::Float64
|   %4 = Main.sin(%3)::Float64
└─      return %4
```

Interpreting the output of `@code_warntype`, like that of its cousins `@code_lowered`, `@code_typed`, `@code_llvm`, and `@code_native`, takes a little practice. Your code is being presented in form that has been heavily digested on its way to generating compiled machine code. Most of the expressions are annotated by a type, indicated by the `::T` (where `T` might be `Float64`, for example). The most important



characteristic of `@code_warntype` is that non-concrete types are displayed in red; since this document is written in Markdown, which has no color, in this document, red text is denoted by uppercase.

At the top, the inferred return type of the function is shown as `Body::Float64`. The next lines represent the body of `f` in Julia's SSA IR form. The numbered boxes are labels and represent targets for jumps (via `goto`) in your code. Looking at the body, you can see that the first thing that happens is that `pos` is called and the return value has been inferred as the Union type `UNION{FLOAT64, INT64}` shown in uppercase since it is a non-concrete type. This means that we cannot know the exact return type of `pos` based on the input types. However, the result of `y*x` is a `Float64` no matter if `y` is a `Float64` or `Int64`. The net result is that `f(x::Float64)` will not be type-unstable in its output, even if some of the intermediate computations are type-unstable.

How you use this information is up to you. Obviously, it would be far and away best to fix `pos` to be type-stable: if you did so, all of the variables in `f` would be concrete, and its performance would be optimal. However, there are circumstances where this kind of *ephemeral* type instability might not matter too much: for example, if `pos` is never used in isolation, the fact that `f`'s output is type-stable (for `Float64` inputs) will shield later code from the propagating effects of type instability. This is particularly relevant in cases where fixing the type instability is difficult or impossible. In such cases, the tips above (e.g., adding type annotations and/or breaking up functions) are your best tools to contain the "damage" from type instability. Also, note that even Julia Base has functions that are type unstable. For example, the function `findfirst` returns the index into an array where a key is found, or `nothing` if it is not found, a clear type instability. In order to make it easier to find the type instabilities that are likely to be important, Unions containing either `missing` or `nothing` are color highlighted in yellow, instead of red.

The following examples may help you interpret expressions marked as containing non-leaf types:

- Function body starting with `Body::UNION{T1, T2}`
  - Interpretation: function with unstable return type
  - Suggestion: make the return value type-stable, even if you have to annotate it
- `invoke Main.g(%x::Int64)::UNION{FLOAT64, INT64}`
  - Interpretation: call to a type-unstable function `g`.
  - Suggestion: fix the function, or if necessary annotate the return value
- `invoke Base.getindex(%x::Array{Any, 1}, 1::Int64)::ANY`
  - Interpretation: accessing elements of poorly-typed arrays
  - Suggestion: use arrays with better-defined types, or if necessary annotate the type of individual element accesses
- `Base.getfield(%x, (:data))::ARRAY{FLOAT64, N} WHERE N`

- Interpretation: getting a field that is of non-leaf type. In this case, `ArrayContainer` had a field `data::Array{T}`. But `Array` needs the dimension `N`, too, to be a concrete type.
- Suggestion: use concrete types like `Array{T, 3}` or `Array{T, N}`, where `N` is now a parameter of `ArrayContainer`

## Performance of captured variable

Consider the following example that defines an inner function:

```
function abmult(r::Int)
    if r < 0
        r = -r
    end
    f = x -> x * r
    return f
end
```

Function `abmult` returns a function `f` that multiplies its argument by the absolute value of `r`. The inner function assigned to `f` is called a "closure". Inner functions are also used by the language for `do`-blocks and for generator expressions.

This style of code presents performance challenges for the language. The parser, when translating it into lower-level instructions, substantially reorganizes the above code by extracting the inner function to a separate code block. "Captured" variables such as `r` that are shared by inner functions and their enclosing scope are also extracted into a heap-allocated "box" accessible to both inner and outer functions because the language specifies that `r` in the inner scope must be identical to `r` in the outer scope even after the outer scope (or another inner function) modifies `r`.

The discussion in the preceding paragraph referred to the "parser", that is, the phase of compilation that takes place when the module containing `abmult` is first loaded, as opposed to the later phase when it is first invoked. The parser does not "know" that `Int` is a fixed type, or that the statement `r = -r` transforms an `Int` to another `Int`. The magic of type inference takes place in the later phase of compilation.

Thus, the parser does not know that `r` has a fixed type (`Int`), nor that `r` does not change value once the inner function is created (so that the box is unneeded). Therefore, the parser emits code for box that holds an object with an abstract type such as `Any`, which requires run-time type dispatch for each occurrence of `r`. This can be verified by applying `@code_warntype` to the above function. Both the boxing and the run-time type dispatch can cause loss of performance.

If captured variables are used in a performance-critical section of the code, then the following tips help

ensure that their use is performant. First, if it is known that a captured variable does not change its type, then this can be declared explicitly with a type annotation (on the variable, not the right-hand side):

```
function abmult2(r0::Int)
    r::Int = r0
    if r < 0
        r = -r
    end
    f = x -> x * r
    return f
end
```

The type annotation partially recovers lost performance due to capturing because the parser can associate a concrete type to the object in the box. Going further, if the captured variable does not need to be boxed at all (because it will not be reassigned after the closure is created), this can be indicated with `let` blocks as follows.

```
function abmult3(r::Int)
    if r < 0
        r = -r
    end
    f = let r = r
        x -> x * r
    end
    return f
end
```

The `let` block creates a new variable `r` whose scope is only the inner function. The second technique recovers full language performance in the presence of captured variables. Note that this is a rapidly evolving aspect of the compiler, and it is likely that future releases will not require this degree of programmer annotation to attain performance. In the mean time, some user-contributed packages like [FastClosures](#) automate the insertion of `let` statements as in `abmult3`.

## Checking for equality with a singleton

When checking if a value is equal to some singleton it can be better for performance to check for identity (`===`) instead of equality (`==`). The same advice applies to using `!==` over `!=`. These type of checks frequently occur e.g. when implementing the iteration protocol and checking if nothing is returned from [iterate](#).

---

[« Stack Traces](#)[Workflow Tips »](#)

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