

Optimizing DiffEq Code

Chris Rackauckas

October 1, 2020

In this notebook we will walk through some of the main tools for optimizing your code in order to efficiently solve DifferentialEquations.jl. User-side optimizations are important because, for sufficiently difficult problems, most of the time will be spent inside of your `f` function, the function you are trying to solve. "Efficient" integrators are those that reduce the required number of `f` calls to hit the error tolerance. The main ideas for optimizing your DiffEq code, or any Julia function, are the following:

- Make it non-allocating
- Use StaticArrays for small arrays
- Use broadcast fusion
- Make it type-stable
- Reduce redundant calculations
- Make use of BLAS calls
- Optimize algorithm choice

We'll discuss these strategies in the context of small and large systems. Let's start with small systems.

0.1 Optimizing Small Systems (<100 DEs)

Let's take the classic Lorenz system from before. Let's start by naively writing the system in its out-of-place form:

```
function lorenz(u,p,t)
    dx = 10.0*(u[2]-u[1])
    dy = u[1]*(28.0-u[3]) - u[2]
    dz = u[1]*u[2] - (8/3)*u[3]
    [dx,dy,dz]
end
```

```
lorenz (generic function with 1 method)
```

Here, `lorenz` returns an object, `[dx,dy,dz]`, which is created within the body of `lorenz`.

This is a common code pattern from high-level languages like MATLAB, SciPy, or R's `deSolve`. However, the issue with this form is that it allocates a vector, `[dx,dy,dz]`, at each step. Let's benchmark the solution process with this choice of function:

```
using DifferentialEquations, BenchmarkTools
u0 = [1.0;0.0;0.0]
tspan = (0.0,100.0)
prob = ODEProblem(lorenz,u0,tspan)
@benchmark solve(prob,Tsit5())
```

```
BenchmarkTools.Trial:
  memory estimate: 10.83 MiB
  allocs estimate: 101434
  -----
  minimum time:      3.327 ms (0.00% GC)
  median time:       3.389 ms (0.00% GC)
  mean time:         4.063 ms (16.11% GC)
  maximum time:      11.319 ms (43.71% GC)
  -----
  samples:           1230
  evals/sample:      1
```

The BenchmarkTools package's `@benchmark` runs the code multiple times to get an accurate measurement. The minimum time is the time it takes when your OS and other background processes aren't getting in the way. Notice that in this case it takes about 5ms to solve and allocates around 11.11 MiB. However, if we were to use this inside of a real user code we'd see a lot of time spent doing garbage collection (GC) to clean up all of the arrays we made. Even if we turn off saving we have these allocations.

```
@benchmark solve(prob,Tsit5(),save_everystep=false)
```

```
BenchmarkTools.Trial:
  memory estimate: 9.47 MiB
  allocs estimate: 88652
  -----
  minimum time:      2.911 ms (0.00% GC)
  median time:       2.920 ms (0.00% GC)
  mean time:         3.450 ms (14.23% GC)
  maximum time:      9.042 ms (41.86% GC)
  -----
  samples:           1448
  evals/sample:      1
```

The problem of course is that arrays are created every time our derivative function is called. This function is called multiple times per step and is thus the main source of memory usage. To fix this, we can use the in-place form to *****make our code non-allocating*****:

```
function lorenz!(du,u,p,t)
  du[1] = 10.0*(u[2]-u[1])
  du[2] = u[1]*(28.0-u[3]) - u[2]
  du[3] = u[1]*u[2] - (8/3)*u[3]
end

lorenz! (generic function with 1 method)
```

Here, instead of creating an array each time, we utilized the cache array `du`. When the inplace form is used, `DifferentialEquations.jl` takes a different internal route that minimizes

the internal allocations as well. When we benchmark this function, we will see quite a difference.

```
u0 = [1.0;0.0;0.0]
tspan = (0.0,100.0)
prob = ODEProblem(lorenz!,u0,tspan)
@benchmark solve(prob,Tsit5())

BenchmarkTools.Trial:
  memory estimate:  1.38 MiB
  allocs estimate:  12996
  -----
  minimum time:     780.576 μs (0.00% GC)
  median time:      786.674 μs (0.00% GC)
  mean time:        874.915 μs (9.35% GC)
  maximum time:     5.344 ms (80.74% GC)
  -----
  samples:          5702
  evals/sample:     1

@benchmark solve(prob,Tsit5(),save_everystep=false)
```

```
BenchmarkTools.Trial:
  memory estimate:  5.22 KiB
  allocs estimate:  54
  -----
  minimum time:     346.009 μs (0.00% GC)
  median time:      348.332 μs (0.00% GC)
  mean time:        348.743 μs (0.00% GC)
  maximum time:     371.806 μs (0.00% GC)
  -----
  samples:          10000
  evals/sample:     1
```

There is a 4x time difference just from that change! Notice there are still some allocations and this is due to the construction of the integration cache. But this doesn't scale with the problem size:

```
tspan = (0.0,500.0) # 5x longer than before
prob = ODEProblem(lorenz!,u0,tspan)
@benchmark solve(prob,Tsit5(),save_everystep=false)
```

```
BenchmarkTools.Trial:
  memory estimate:  5.22 KiB
  allocs estimate:  54
  -----
  minimum time:     1.737 ms (0.00% GC)
  median time:      1.744 ms (0.00% GC)
  mean time:        1.746 ms (0.00% GC)
  maximum time:     1.974 ms (0.00% GC)
  -----
  samples:          2863
  evals/sample:     1
```

since that's all just setup allocations.

But if the system is small we can optimize even more. Allocations are only expensive if they are "heap allocations". For a more in-depth definition of heap allocations, [there are](#)

a lot of sources online. But a good working definition is that heap allocations are variable-sized slabs of memory which have to be pointed to, and this pointer indirection costs time. Additionally, the heap has to be managed and the garbage controllers has to actively keep track of what's on the heap.

However, there's an alternative to heap allocations, known as stack allocations. The stack is statically-sized (known at compile time) and thus its accesses are quick. Additionally, the exact block of memory is known in advance by the compiler, and thus re-using the memory is cheap. This means that allocating on the stack has essentially no cost!

Arrays have to be heap allocated because their size (and thus the amount of memory they take up) is determined at runtime. But there are structures in Julia which are stack-allocated. `structs` for example are stack-allocated "value-type"s. `Tuples` are a stack-allocated collection. The most useful data structure for `DiffEq` though is the `StaticArray` from the package [StaticArrays.jl](#). These arrays have their length determined at compile-time. They are created using macros attached to normal array expressions, for example:

```
using StaticArrays
A = @SVector [2.0,3.0,5.0]

3-element StaticArrays.SArray{Tuple{3},Float64,1,3} with indices SOneTo(3):
 2.0
 3.0
 5.0
```

Notice that the 3 after `SVector` gives the size of the `SVector`. It cannot be changed. Additionally, `SVectors` are immutable, so we have to create a new `SVector` to change values. But remember, we don't have to worry about allocations because this data structure is stack-allocated. `SArrays` have a lot of extra optimizations as well: they have fast matrix multiplication, fast QR factorizations, etc. which directly make use of the information about the size of the array. Thus, when possible they should be used.

Unfortunately static arrays can only be used for sufficiently small arrays. After a certain size, they are forced to heap allocate after some instructions and their compile time balloons. Thus static arrays shouldn't be used if your system has more than 100 variables. Additionally, only the native Julia algorithms can fully utilize static arrays.

Let's `***optimize lorenz using static arrays***`. Note that in this case, we want to use the out-of-place allocating form, but this time we want to output a static array:

```
function lorenz_static(u,p,t)
    dx = 10.0*(u[2]-u[1])
    dy = u[1]*(28.0-u[3]) - u[2]
    dz = u[1]*u[2] - (8/3)*u[3]
    @SVector [dx,dy,dz]
end

lorenz_static (generic function with 1 method)
```

To make the solver internally use static arrays, we simply give it a static array as the initial condition:

```
u0 = @SVector [1.0,0.0,0.0]
tspan = (0.0,100.0)
prob = ODEProblem(lorenz_static,u0,tspan)
@benchmark solve(prob,Tsit5())
```

```

BenchmarkTools.Trial:
  memory estimate:  466.28 KiB
  allocs estimate:  2595
  -----
  minimum time:     338.382 μs (0.00% GC)
  median time:      342.823 μs (0.00% GC)
  mean time:        364.389 μs (4.33% GC)
  maximum time:     2.950 ms (85.26% GC)
  -----
  samples:          10000
  evals/sample:     1

@benchmark solve(prob,Tsit5(),save_everystep=false)

```

```

BenchmarkTools.Trial:
  memory estimate:  3.33 KiB
  allocs estimate:  28
  -----
  minimum time:     236.592 μs (0.00% GC)
  median time:      240.215 μs (0.00% GC)
  mean time:        240.409 μs (0.00% GC)
  maximum time:     254.428 μs (0.00% GC)
  -----
  samples:          10000
  evals/sample:     1

```

And that's pretty much all there is to it. With static arrays you don't have to worry about allocating, so use operations like `*` and don't worry about fusing operations (discussed in the next section). Do "the vectorized code" of R/MATLAB/Python and your code in this case will be fast, or directly use the numbers/values.

Exercise 1 Implement the out-of-place array, in-place array, and out-of-place static array forms for the [Henon-Heiles System](#) and time the results.

0.2 Optimizing Large Systems

0.2.1 Interlude: Managing Allocations with Broadcast Fusion

When your system is sufficiently large, or you have to make use of a non-native Julia algorithm, you have to make use of `Arrays`. In order to use arrays in the most efficient manner, you need to be careful about temporary allocations. Vectorized calculations naturally have plenty of temporary array allocations. This is because a vectorized calculation outputs a vector. Thus:

```

A = rand(1000,1000); B = rand(1000,1000); C = rand(1000,1000)
test(A,B,C) = A + B + C
@benchmark test(A,B,C)

```

```

BenchmarkTools.Trial:
  memory estimate:  7.63 MiB
  allocs estimate:  2
  -----
  minimum time:     1.135 ms (0.00% GC)
  median time:      1.160 ms (0.00% GC)
  mean time:        1.295 ms (10.24% GC)

```

```

maximum time:      2.824 ms (56.78% GC)
-----
samples:           3843
evals/sample:      1

```

That expression `A + B + C` creates 2 arrays. It first creates one for the output of `A + B`, then uses that result array to `+ C` to get the final result. 2 arrays! We don't want that! The first thing to do to fix this is to use broadcast fusion. [Broadcast fusion](#) puts expressions together. For example, instead of doing the `+` operations separately, if we were to add them all at the same time, then we would only have a single array that's created. For example:

```

test2(A,B,C) = map((a,b,c)->a+b+c,A,B,C)
@benchmark test2(A,B,C)

```

```

BenchmarkTools.Trial:
 memory estimate:  7.63 MiB
 allocs estimate:  8
-----
 minimum time:     2.210 ms (0.00% GC)
 median time:     2.226 ms (0.00% GC)
 mean time:       2.352 ms (5.29% GC)
 maximum time:     3.802 ms (40.63% GC)
-----
 samples:         2123
 evals/sample:    1

```

Puts the whole expression into a single function call, and thus only one array is required to store output. This is the same as writing the loop:

```

function test3(A,B,C)
    D = similar(A)
    @inbounds for i in eachindex(A)
        D[i] = A[i] + B[i] + C[i]
    end
    D
end
@benchmark test3(A,B,C)

```

```

BenchmarkTools.Trial:
 memory estimate:  7.63 MiB
 allocs estimate:  2
-----
 minimum time:     1.126 ms (0.00% GC)
 median time:     1.151 ms (0.00% GC)
 mean time:       1.287 ms (10.32% GC)
 maximum time:     2.795 ms (57.40% GC)
-----
 samples:         3867
 evals/sample:    1

```

However, Julia's broadcast is syntactic sugar for this. If multiple expressions have a `.`, then it will put those vectorized operations together. Thus:

```

test4(A,B,C) = A .+ B .+ C
@benchmark test4(A,B,C)

```

```

BenchmarkTools.Trial:
 memory estimate:  7.63 MiB
 allocs estimate:  2

```

```

-----
minimum time:      1.130 ms (0.00% GC)
median time:       1.156 ms (0.00% GC)
mean time:         1.292 ms (10.33% GC)
maximum time:      2.837 ms (56.65% GC)
-----

samples:           3852
evals/sample:      1

```

is a version with only 1 array created (the output). Note that `.s` can be used with function calls as well:

```
sin.(A) .+ sin.(B)
```

```

1000x0*(1000 Array{Float64,2}:
 1.3318      1.03918  0.887918  1.00151  ...0*( 1.03369 0.348665 1.18041 1.31843 0.180731
1.61303 0.432489 0.743137 0.947268 0.882330 0.837925 0.476861 0.425674 0.202355 0.764645
0.861877 0.175019 0.951978 0.190512 1.57814 0.341034 0.796808 0.801442 1.213020 0.959764
1.4384 0.651705 0.862134 0.54159 0.48242 1.294691 0.27836 1.06678 1.61041 0.759755 (*@...@*(
0.774674 1.05993 0.747471 0.909648 1.23279 0.239356 0.947454 0.689775 1.59072
0.617020 0.0566593 0.143266 0.782544 0.363328 0.734336 1.39576 0.310102 0.650645 1.21329
0.547903 1.13496 0.240357 1.3828 0.486759 1.28042 1.20366 0.745039 0.782717 1.28894
0.611047 1.10041(*@:*( (*@`'.*(0.475611 0.455008 0.863672 0.370343 0.878027 0.434347
1.486551 1.32962 0.462008 0.561171 1.31755 0.721338 0.746222 0.766714 0.472661 0.688161
1.15282 0.283092 0.801634 0.796207 0.546280 0.679418 1.17586 0.900804 0.536497 1.33048
0.988912 0.823366 1.55493 1.24402 1.15303 1.50803 (*@...@*( 0.800089 0.102445
0.866568 1.15568 0.34324 0.364801 1.05421 0.826784 0.558172 0.510031 0.29874 1.15521
0.847367 0.835584 0.421068 0.272068 0.649116 1.14135 0.860285 1.04829 0.982288 1.28133
1.36746 0.872593 1.13586 0.874325 1.1923 1.23324 1.05135 0.857278 0.432765

```

Also, the `@.` macro applies a dot to every operator:

```
test5(A,B,C) = @. A + B + C #only one array allocated
@benchmark test5(A,B,C)
```

```

BenchmarkTools.Trial:
 memory estimate:  7.63 MiB
 allocs estimate:  2
-----
minimum time:      1.133 ms (0.00% GC)
median time:       1.160 ms (0.00% GC)
mean time:         1.296 ms (10.29% GC)
maximum time:      2.823 ms (56.66% GC)
-----

samples:           3840
evals/sample:      1

```

Using these tools we can get rid of our intermediate array allocations for many vectorized function calls. But we are still allocating the output array. To get rid of that allocation, we can instead use mutation. Mutating broadcast is done via `.=`. For example, if we pre-allocate the output:

```
D = zeros(1000,1000);
```

Then we can keep re-using this cache for subsequent calculations. The mutating broadcasting form is:

```
test6!(D,A,B,C) = D .= A .+ B .+ C #only one array allocated
@benchmark test6!(D,A,B,C)
```

```
BenchmarkTools.Trial:
  memory estimate: 0 bytes
  allocs estimate: 0
  -----
  minimum time:      1.059 ms (0.00% GC)
  median time:       1.067 ms (0.00% GC)
  mean time:         1.069 ms (0.00% GC)
  maximum time:      1.323 ms (0.00% GC)
  -----
  samples:           4655
  evals/sample:      1
```

If we use `@.` before the `=`, then it will turn it into `.=`:

```
test7!(D,A,B,C) = @. D = A + B + C #only one array allocated
@benchmark test7!(D,A,B,C)
```

```
BenchmarkTools.Trial:
  memory estimate: 0 bytes
  allocs estimate: 0
  -----
  minimum time:      1.056 ms (0.00% GC)
  median time:       1.063 ms (0.00% GC)
  mean time:         1.065 ms (0.00% GC)
  maximum time:      1.320 ms (0.00% GC)
  -----
  samples:           4670
  evals/sample:      1
```

Notice that in this case, there is no "output", and instead the values inside of `D` are what are changed (like with the `DiffEq` inplace function). Many Julia functions have a mutating form which is denoted with a `!`. For example, the mutating form of the `map` is `map!`:

```
test8!(D,A,B,C) = map!((a,b,c)->a+b+c,D,A,B,C)
@benchmark test8!(D,A,B,C)
```

```
BenchmarkTools.Trial:
  memory estimate: 32 bytes
  allocs estimate: 1
  -----
  minimum time:      2.374 ms (0.00% GC)
  median time:       2.385 ms (0.00% GC)
  mean time:         2.391 ms (0.00% GC)
  maximum time:      3.235 ms (0.00% GC)
  -----
  samples:           2088
  evals/sample:      1
```

Some operations require using an alternate mutating form in order to be fast. For example, matrix multiplication via `*` allocates a temporary:

```
@benchmark A*B
```

```
BenchmarkTools.Trial:
  memory estimate: 7.63 MiB
  allocs estimate: 2
  -----
  minimum time:      6.343 ms (0.00% GC)
  median time:       6.433 ms (0.00% GC)
  mean time:         6.599 ms (2.07% GC)
```



```

maximum time:      8.515 ms (19.17% GC)
-----
samples:          757
evals/sample:      1

```

Instead, we can use the mutating form `mul!` into a cache array to avoid allocating the output:

```

using LinearAlgebra
@benchmark mul!(D,A,B) # same as D = A * B

```

```

BenchmarkTools.Trial:
 memory estimate:  0 bytes
 allocs estimate:  0
-----
 minimum time:     5.943 ms (0.00% GC)
 median time:     5.985 ms (0.00% GC)
 mean time:       5.988 ms (0.00% GC)
 maximum time:     6.542 ms (0.00% GC)
-----
 samples:         834
 evals/sample:    1

```

For repeated calculations this reduced allocation can stop GC cycles and thus lead to more efficient code. Additionally, ***we can fuse together higher level linear algebra operations using BLAS***. The package [SugarBLAS.jl](#) makes it easy to write higher level operations like `alpha*B*A + beta*C` as mutating BLAS calls.

0.2.2 Example Optimization: Gierer-Meinhardt Reaction-Diffusion PDE Discretization

Let's optimize the solution of a Reaction-Diffusion PDE's discretization. In its discretized form, this is the ODE:

$$du = D_1(A_y u + u A_x) + \frac{au^2}{v} + \bar{u} - \alpha u \quad (1)$$

$$dv = D_2(A_y v + v A_x) + au^2 + \beta v \quad (2)$$

where u , v , and A are matrices. Here, we will use the simplified version where A is the tridiagonal stencil $[1, -2, 1]$, i.e. it's the 2D discretization of the Laplacian. The native code would be something along the lines of:

```

# Generate the constants
p = (1.0,1.0,1.0,10.0,0.001,100.0) # a,α,ubar,β,D1,D2
N = 100
Ax = Array{Tridiagonal{Float64}}{1,N,N}([1.0 for i in 1:N-1],[-2.0 for i in 1:N],[1.0 for i in 1:N-1]))
Ay = copy(Ax)
Ax[2,1] = 2.0
Ax[end-1,end] = 2.0
Ay[1,2] = 2.0
Ay[end,end-1] = 2.0

function basic_version!(dr,r,p,t)
    a,α,ubar,β,D1,D2 = p
    u = r[:, :, 1]

```

```

v = r[:, :, 2]
Du = D1*(Ay*u + u*Ax)
Dv = D2*(Ay*v + v*Ax)
dr[:, :, 1] = Du ./ a.*u.*u./v ./ ubar ./ alpha*u
dr[:, :, 2] = Dv ./ a.*u.*u ./ beta*v
end

a, alpha, ubar, beta, D1, D2 = p
uss = (ubar+beta)/alpha
vss = (a/beta)*uss^2
r0 = zeros(100,100,2)
r0[:, :, 1] = uss.+0.1.*rand.()
r0[:, :, 2] = vss

prob = ODEProblem(basic_version!, r0, (0.0, 0.1), p)

ODEProblem with uType Array{Float64,3} and tType Float64. In-place: true
timespan: (0.0, 0.1)
u0: [11.067596791366832 11.049387382048813 ...@*( 11.096750972563049 11.024559555174504;
11.026898386485927 11.077605849059506 (*@...@*( 11.005020745576179 11.07782657582884;
(*@...@*( ; 11.094796724318485 11.090635445667045 (*@...@*(
11.08320744731535311.031644064912005; 11.089711890390022 11.026574702707167 (*@...@*(
11.01094076493651 11.046876982972178] [12.100000000000001 12.100000000000001 (*@...@*(
12.100000000000001 12.100000000000001; 12.100000000000001 12.100000000000001 (*@...@*(
12.100000000000001 12.100000000000001; (*@...@*( ; 12.100000000000001 12.100000000000001
(*@...@*( 12.100000000000001 12.100000000000001; 12.100000000000001 12.100000000000001
(*@...@*( 12.10000000000000112.100000000000001]

```

In this version we have encoded our initial condition to be a 3-dimensional array, with `u[:, :, 1]` being the A part and `u[:, :, 2]` being the B part.

```
@benchmark solve(prob, Tsit5())
```

```

BenchmarkTools.Trial:
 memory estimate: 186.88 MiB
  allocs estimate: 8551
  -----
 minimum time:      48.985 ms (4.07% GC)
 median time:       51.222 ms (7.58% GC)
 mean time:         50.884 ms (7.03% GC)
 maximum time:      51.489 ms (7.73% GC)
  -----
 samples:           99
 evals/sample:      1

```

While this version isn't very efficient,

We recommend writing the "high-level" code first, and iteratively optimizing it! The first thing that we can do is get rid of the slicing allocations. The operation `r[:, :, 1]` creates a temporary array instead of a "view", i.e. a pointer to the already existing memory. To make it a view, add `@view`. Note that we have to be careful with views because they point to the same memory, and thus changing a view changes the original values:

```

A = rand(4)
@show A
B = @view A[1:3]
B[2] = 2
@show A

```

```

A = [0.10681951152837588, 0.5992481459850683, 0.8083545274299431, 0.8565100
765399429]
A = [0.10681951152837588, 2.0, 0.8083545274299431, 0.8565100765399429]
4-element Array{Float64,1}:
 0.10681951152837588
 2.0
 0.8083545274299431
 0.8565100765399429

```

Notice that changing B changed A. This is something to be careful of, but at the same time we want to use this since we want to modify the output `dr`. Additionally, the last statement is a purely element-wise operation, and thus we can make use of broadcast fusion there. Let's rewrite `basic_version!` to `***avoid slicing allocations***` and to `***use broadcast fusion***`:

```

function gm2!(dr,r,p,t)
    a,α,ubar,β,D1,D2 = p
    u = @view r[:,1]
    v = @view r[:,2]
    du = @view dr[:,1]
    dv = @view dr[:,2]
    Du = D1*(Ay*u + u*Ax)
    Dv = D2*(Ay*v + v*Ax)
    @. du = Du + a.*u.*u./v + ubar - α*u
    @. dv = Dv + a.*u.*u - β*v
end
prob = ODEProblem(gm2!,r0,(0.0,0.1),p)
@benchmark solve(prob,Tsit5())

```

```

BenchmarkTools.Trial:
 memory estimate: 119.55 MiB
 allocs estimate: 7081
-----
 minimum time:      39.140 ms (5.13% GC)
 median time:      39.301 ms (5.09% GC)
 mean time:        39.801 ms (6.26% GC)
 maximum time:     41.433 ms (9.49% GC)
-----
 samples:          126
 evals/sample:     1

```

Now, most of the allocations are taking place in `Du = D1*(Ay*u + u*Ax)` since those operations are vectorized and not mutating. We should instead replace the matrix multiplications with `mul!`. When doing so, we will need to have cache variables to write into. This looks like:

```

Ayu = zeros(N,N)
uAx = zeros(N,N)
Du = zeros(N,N)
Ayv = zeros(N,N)
vAx = zeros(N,N)
Dv = zeros(N,N)
function gm3!(dr,r,p,t)
    a,α,ubar,β,D1,D2 = p
    u = @view r[:,1]
    v = @view r[:,2]
    du = @view dr[:,1]
    dv = @view dr[:,2]
    mul!(Ayu,Ay,u)

```

```

mul!(uAx,u,Ax)
mul!(Ayv,Ay,v)
mul!(vAx,v,Ax)
@. Du = D1*(Ayu + uAx)
@. Dv = D2*(Ayv + vAx)
@. du = Du + a*u*u./v + ubar -  $\alpha$ *u
@. dv = Dv + a*u*u -  $\beta$ *v
end
prob = ODEProblem(gm3!,r0,(0.0,0.1),p)
@benchmark solve(prob,Tsit5())

```

```

BenchmarkTools.Trial:
  memory estimate: 29.75 MiB
  allocs estimate: 5317
  -----
  minimum time:      33.814 ms (0.00% GC)
  median time:       34.143 ms (0.00% GC)
  mean time:         34.686 ms (1.72% GC)
  maximum time:      36.226 ms (5.36% GC)
  -----
  samples:           145
  evals/sample:      1

```

But our temporary variables are global variables. We need to either declare the caches as `const` or localize them. We can localize them by adding them to the parameters, `p`. It's easier for the compiler to reason about local variables than global variables. *****Localizing variables helps to ensure type stability*****.

```

p = (1.0,1.0,1.0,10.0,0.001,100.0,Ayu,uAx,Du,Ayv,vAx,Dv) # a, $\alpha$ ,ubar, $\beta$ ,D1,D2
function gm4!(dr,r,p,t)
  a, $\alpha$ ,ubar, $\beta$ ,D1,D2,Ayu,uAx,Du,Ayv,vAx,Dv = p
  u = @view r[:,1]
  v = @view r[:,2]
  du = @view dr[:,1]
  dv = @view dr[:,2]
  mul!(Ayu,Ay,u)
  mul!(uAx,u,Ax)
  mul!(Ayv,Ay,v)
  mul!(vAx,v,Ax)
  @. Du = D1*(Ayu + uAx)
  @. Dv = D2*(Ayv + vAx)
  @. du = Du + a*u*u./v + ubar -  $\alpha$ *u
  @. dv = Dv + a*u*u -  $\beta$ *v
end
prob = ODEProblem(gm4!,r0,(0.0,0.1),p)
@benchmark solve(prob,Tsit5())

```

```

BenchmarkTools.Trial:
  memory estimate: 29.66 MiB
  allocs estimate: 1053
  -----
  minimum time:      27.345 ms (0.00% GC)
  median time:       27.496 ms (0.00% GC)
  mean time:         28.078 ms (2.08% GC)
  maximum time:      29.597 ms (6.40% GC)
  -----
  samples:           179
  evals/sample:      1

```

We could then use the BLAS `gemmv` to optimize the matrix multiplications some more, but instead let's devectorize the stencil.

```
p = (1.0,1.0,1.0,10.0,0.001,100.0,N)
function fast_gm!(du,u,p,t)
    a,α,ubar,β,D1,D2,N = p

    @inbounds for j in 2:N-1, i in 2:N-1
        du[i,j,1] = D1*(u[i-1,j,1] + u[i+1,j,1] + u[i,j+1,1] + u[i,j-1,1] - 4u[i,j,1]) +
            a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
    end

    @inbounds for j in 2:N-1, i in 2:N-1
        du[i,j,2] = D2*(u[i-1,j,2] + u[i+1,j,2] + u[i,j+1,2] + u[i,j-1,2] - 4u[i,j,2]) +
            a*u[i,j,1]^2 - β*u[i,j,2]
    end

    @inbounds for j in 2:N-1
        i = 1
        du[1,j,1] = D1*(2u[i+1,j,1] + u[i,j+1,1] + u[i,j-1,1] - 4u[i,j,1]) +
            a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
    end
    @inbounds for j in 2:N-1
        i = 1
        du[1,j,2] = D2*(2u[i+1,j,2] + u[i,j+1,2] + u[i,j-1,2] - 4u[i,j,2]) +
            a*u[i,j,1]^2 - β*u[i,j,2]
    end

    @inbounds for j in 2:N-1
        i = N
        du[end,j,1] = D1*(2u[i-1,j,1] + u[i,j+1,1] + u[i,j-1,1] - 4u[i,j,1]) +
            a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
    end
    @inbounds for j in 2:N-1
        i = N
        du[end,j,2] = D2*(2u[i-1,j,2] + u[i,j+1,2] + u[i,j-1,2] - 4u[i,j,2]) +
            a*u[i,j,1]^2 - β*u[i,j,2]
    end

    @inbounds for i in 2:N-1
        j = 1
        du[i,1,1] = D1*(u[i-1,j,1] + u[i+1,j,1] + 2u[i,j+1,1] - 4u[i,j,1]) +
            a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
    end
    @inbounds for i in 2:N-1
        j = 1
        du[i,1,2] = D2*(u[i-1,j,2] + u[i+1,j,2] + 2u[i,j+1,2] - 4u[i,j,2]) +
            a*u[i,j,1]^2 - β*u[i,j,2]
    end

    @inbounds for i in 2:N-1
        j = N
        du[i,end,1] = D1*(u[i-1,j,1] + u[i+1,j,1] + 2u[i,j-1,1] - 4u[i,j,1]) +
            a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
    end
    @inbounds for i in 2:N-1
        j = N
        du[i,end,2] = D2*(u[i-1,j,2] + u[i+1,j,2] + 2u[i,j-1,2] - 4u[i,j,2]) +
            a*u[i,j,1]^2 - β*u[i,j,2]
    end

    @inbounds begin
```

```

i = 1; j = 1
du[1,1,1] = D1*(2u[i+1,j,1] + 2u[i,j+1,1] - 4u[i,j,1]) +
    a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
du[1,1,2] = D2*(2u[i+1,j,2] + 2u[i,j+1,2] - 4u[i,j,2]) +
    a*u[i,j,1]^2 - β*u[i,j,2]

i = 1; j = N
du[1,N,1] = D1*(2u[i+1,j,1] + 2u[i,j-1,1] - 4u[i,j,1]) +
    a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
du[1,N,2] = D2*(2u[i+1,j,2] + 2u[i,j-1,2] - 4u[i,j,2]) +
    a*u[i,j,1]^2 - β*u[i,j,2]

i = N; j = 1
du[N,1,1] = D1*(2u[i-1,j,1] + 2u[i,j+1,1] - 4u[i,j,1]) +
    a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
du[N,1,2] = D2*(2u[i-1,j,2] + 2u[i,j+1,2] - 4u[i,j,2]) +
    a*u[i,j,1]^2 - β*u[i,j,2]

i = N; j = N
du[end,end,1] = D1*(2u[i-1,j,1] + 2u[i,j-1,1] - 4u[i,j,1]) +
    a*u[i,j,1]^2/u[i,j,2] + ubar - α*u[i,j,1]
du[end,end,2] = D2*(2u[i-1,j,2] + 2u[i,j-1,2] - 4u[i,j,2]) +
    a*u[i,j,1]^2 - β*u[i,j,2]
end
end
prob = ODEProblem(fast_gm!,r0,(0.0,0.1),p)
@benchmark solve(prob,Tsit5())

```

```

BenchmarkTools.Trial:
 memory estimate: 29.62 MiB
  allocs estimate: 466
  -----
 minimum time:      8.409 ms (0.00% GC)
 median time:      8.524 ms (0.00% GC)
 mean time:        9.093 ms (6.24% GC)
 maximum time:     10.589 ms (17.45% GC)
  -----
 samples:          550
 evals/sample:     1

```

Lastly, we can do other things like multithread the main loops, but these optimizations get the last 2x-3x out. The main optimizations which apply everywhere are the ones we just performed (though the last one only works if your matrix is a stencil. This is known as a matrix-free implementation of the PDE discretization).

This gets us to about 8x faster than our original MATLAB/SciPy/R vectorized style code!

The last thing to do is then **optimize our algorithm choice**. We have been using `Tsit5()` as our test algorithm, but in reality this problem is a stiff PDE discretization and thus one recommendation is to use `CVODE_BDF()`. However, instead of using the default dense Jacobian, we should make use of the sparse Jacobian afforded by the problem. The Jacobian is the matrix $\frac{df_i}{dr_j}$, where r is read by the linear index (i.e. down columns). But since the u variables depend on the v , the band size here is large, and thus this will not do well with a Banded Jacobian solver. Instead, we utilize sparse Jacobian algorithms. `CVODE_BDF` allows us to use a sparse Newton-Krylov solver by setting `linear_solver = :GMRES` (see [the solver documentation](#), and thus we can solve this problem efficiently. Let's see how this scales as we increase the integration time.

```

prob = ODEProblem(fast_gm!,r0,(0.0,10.0),p)
@benchmark solve(prob,Tsit5())

BenchmarkTools.Trial:
  memory estimate: 2.76 GiB
  allocs estimate: 41650
  -----
  minimum time:      1.344 s (39.86% GC)
  median time:      1.616 s (36.71% GC)
  mean time:        1.567 s (40.08% GC)
  maximum time:      1.693 s (30.42% GC)
  -----
  samples:          4
  evals/sample:     1

using Sundials
@benchmark solve(prob,CVODE_BDF(linear_solver=:GMRES))

BenchmarkTools.Trial:
  memory estimate: 118.89 MiB
  allocs estimate: 20238
  -----
  minimum time:      805.338 ms (0.25% GC)
  median time:      808.489 ms (0.27% GC)
  mean time:        808.190 ms (0.26% GC)
  maximum time:      810.190 ms (0.27% GC)
  -----
  samples:          7
  evals/sample:     1

prob = ODEProblem(fast_gm!,r0,(0.0,100.0),p)
# Will go out of memory if we don't turn off `save_everystep`!
@benchmark solve(prob,Tsit5(),save_everystep=false)

BenchmarkTools.Trial:
  memory estimate: 2.90 MiB
  allocs estimate: 74
  -----
  minimum time:      5.987 s (0.00% GC)
  median time:      5.987 s (0.00% GC)
  mean time:        5.987 s (0.00% GC)
  maximum time:      5.987 s (0.00% GC)
  -----
  samples:          1
  evals/sample:     1

@benchmark solve(prob,CVODE_BDF(linear_solver=:GMRES))

BenchmarkTools.Trial:
  memory estimate: 311.31 MiB
  allocs estimate: 56248
  -----
  minimum time:      2.302 s (0.00% GC)
  median time:      2.338 s (1.27% GC)
  mean time:        2.360 s (2.21% GC)
  maximum time:      2.439 s (5.21% GC)
  -----
  samples:          3
  evals/sample:     1

```

Now let's check the allocation growth.

```
@benchmark solve(prob,CVODE_BDF(linear_solver=:GMRES),save_everystep=false)
```

```
BenchmarkTools.Trial:
  memory estimate:  3.60 MiB
  allocs estimate:  46825
  -----
  minimum time:     2.271 s (0.00% GC)
  median time:      2.271 s (0.00% GC)
  mean time:        2.271 s (0.00% GC)
  maximum time:     2.273 s (0.00% GC)
  -----
  samples:          3
  evals/sample:     1
```

```
prob = ODEProblem(fast_gm!,r0,(0.0,500.0),p)
```

```
@benchmark solve(prob,CVODE_BDF(linear_solver=:GMRES),save_everystep=false)
```

```
BenchmarkTools.Trial:
  memory estimate:  3.78 MiB
  allocs estimate:  50029
  -----
  minimum time:     2.432 s (0.00% GC)
  median time:      2.433 s (0.00% GC)
  mean time:        2.433 s (0.00% GC)
  maximum time:     2.434 s (0.00% GC)
  -----
  samples:          3
  evals/sample:     1
```

Notice that we've eliminated almost all allocations, allowing the code to grow without hitting garbage collection and slowing down.

Why is `CVODE_BDF` doing well? What's happening is that, because the problem is stiff, the number of steps required by the explicit Runge-Kutta method grows rapidly, whereas `CVODE_BDF` is taking large steps. Additionally, the `GMRES` linear solver form is quite an efficient way to solve the implicit system in this case. This is problem-dependent, and in many cases using a Krylov method effectively requires a preconditioner, so you need to play around with testing other algorithms and linear solvers to find out what works best with your problem.

0.3 Conclusion

Julia gives you the tools to optimize the solver "all the way", but you need to make use of it. The main thing to avoid is temporary allocations. For small systems, this is effectively done via static arrays. For large systems, this is done via in-place operations and cache arrays. Either way, the resulting solution can be immensely sped up over vectorized formulations by using these principles.

0.4 Appendix

This tutorial is part of the `SciMLTutorials.jl` repository, found at: <https://github.com/SciML/SciMLTutorials.jl>. For more information on doing scientific machine learning (SciML) with open source software, check out <https://sciml.ai/>.

To locally run this tutorial, do the following commands:

```
using SciMLTutorials
SciMLTutorials.weave_file("introduction","03-optimizing_diffeq_code.jmd")
```

Computer Information:

```
Julia Version 1.4.2
Commit 44fa15b150* (2020-05-23 18:35 UTC)
Platform Info:
  OS: Linux (x86_64-pc-linux-gnu)
  CPU: Intel(R) Core(TM) i7-9700K CPU @ 3.60GHz
  WORD_SIZE: 64
  LIBM: libopenlibm
  LLVM: libLLVM-8.0.1 (ORCJIT, skylake)
```

Environment:

```
JULIA_LOAD_PATH = /builds/JuliaGPU/DiffEqTutorials.jl:
JULIA_DEPOT_PATH = /builds/JuliaGPU/DiffEqTutorials.jl/.julia
JULIA_CUDA_MEMORY_LIMIT = 2147483648
JULIA_NUM_THREADS = 8
```

Package Information:

```
Status `~/builds/JuliaGPU/DiffEqTutorials.jl/tutorials/introduction/Project.toml`
[6e4b80f9-dd63-53aa-95a3-0cdb28fa8baf] BenchmarkTools 0.5.0
[0c46a032-eb83-5123-abaf-570d42b7fbaf] DifferentialEquations 6.15.0
[65888b18-ceab-5e60-b2b9-181511a3b968] ParameterizedFunctions 5.6.0
[91a5bcdd-55d7-5caf-9e0b-520d859cae80] Plots 1.6.7
[90137ffa-7385-5640-81b9-e52037218182] StaticArrays 0.12.4
[c3572dad-4567-51f8-b174-8c6c989267f4] Sundials 4.3.0
[37e2e46d-f89d-539d-b4ee-838fcccc9c8e] LinearAlgebra
```