

# What's wrong with this



Kenta Sato (佐藤建太)

Twitter/GitHub: @bicycle1885



# A Sad Story

Is Julia as fast as C/C++? Where is the download link!?

Oh, the syntax looks really neat!

Ok, it works. Now let's compare it with my Python script.

Huh? not so fast as he said, maybe comparable with Python?

Julia? Yes, I tried it last month but blahblah...

# A Sad Story

Is Julia as fast as C/C++? Where is the download link!?

Oh, the syntax looks really neat!

Ok, it works. Now let's compare it with my Python script.

Huh? not so fast as he said, maybe comparable with Python?

Julia? Yes, I tried it last month but blahblah...

Prob(their code does something wrong | newbies complain about performance)  $> 0.95$

We are going to learn very basic skills to write (nearly) optimal Julia code from three cases.

# Case 1

```
s = 0
n = 500000000
for i in 1:n
    s += i
end
println(s)
```

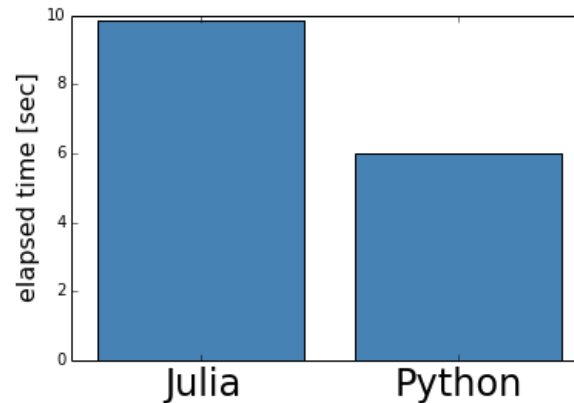
Julia

```
$ time julia case1.jl
real    0m9.859s
...
```

```
s = 0
n = 500000000
for i in xrange(1, n + 1):
    s += i
print(s)
```

Python

```
$ time python case1.py
real    0m5.998s
...
```



Is Julia slower than Python? Really??

# Case 1

## Do Not Write a Heavy Loop at the Top Level!

- The top-level code is interpreted, not compiled.
- Top level is not intended to run numeric computation.
- Should be written in a `function` / `let` block.

```
s = 0
n = 500000000
for i in 1:n
    s += i
end
println(s)
```

Julia

```
$ time julia case1.jl
real    0m9.859s
...
```

**let**

```
s = 0
n = 500000000
for i in 1:n
    s += i
end
println(s)
```

Julia

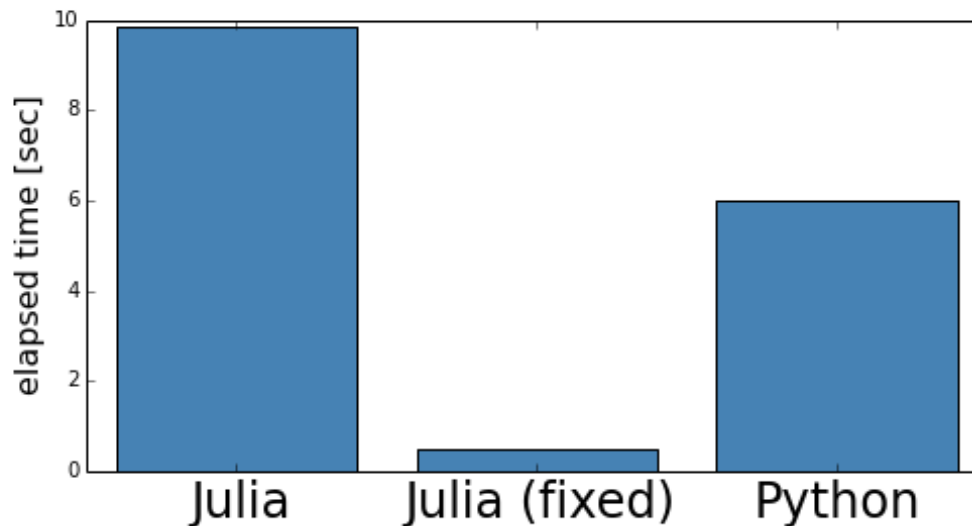
**end**

```
$ time julia case1.fix.jl
real    0m0.494s
...
```

# Case 1

## Do Not Write a Heavy Loop at the Top Level!

- The top-level code is interpreted, not compiled.
- Top level is not intended to run numeric computation.
- Should be written in a `function` / `let` block.



9.859s / 0.494s = 20x faster

# Case 2

```
function dot(x, y)          Julia
    s = 0
    for i in 1:length(x)
        s += x[i] * y[i]
    end
    s
end
```

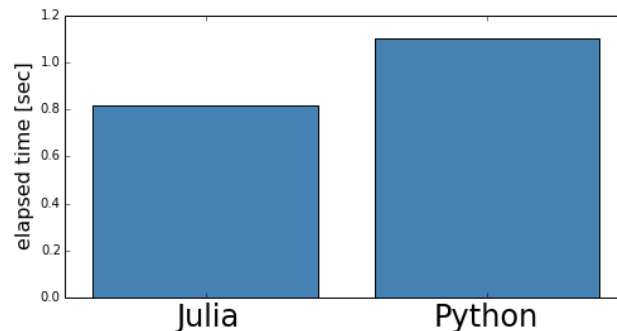
```
n = 10000000
x = zeros(n)
y = zeros(n)
dot(x[1:3], y[1:3]) # compile
```

```
def dot(x, y):              Python
    s = 0
    for i in xrange(len(x)):
        s += x[i] * y[i]
    return s
```

```
n = 100000000
x = [0] * n
y = [0] * n
```

```
% timeit dot(x, y)
1 loops, best of 3: 1.1 s per lo
```

```
$ julia -L case2.jl -e '@time do
elapsed time: 0.814667538 second
```



# Case 2

## Make Variable Types Stable!

- Literal `0` is typed as `Int`.
- Hence Julia assumes the type of the variable `s` is `Int`.
- The element type of `x` and `y` is `Float64`.
- In Julia, `Int + Float64 * Float64` is `Float64`.
- `s += x[i] * y[i]` disturbs the type stability of `s`.

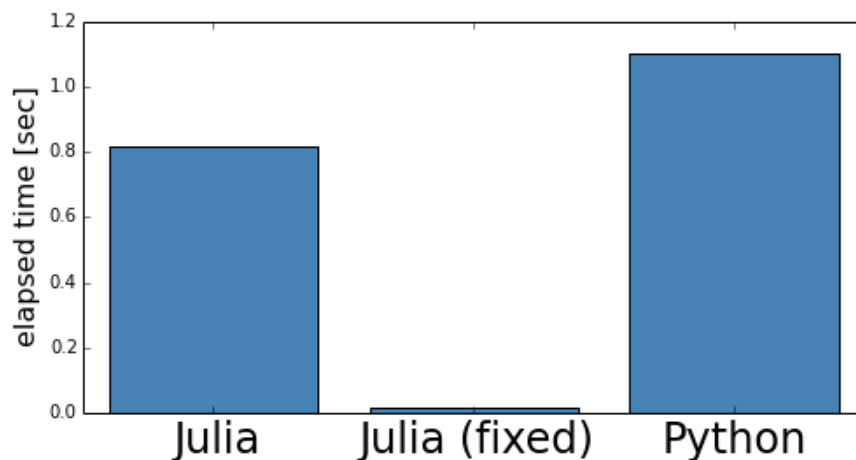
```
function dot(x, y)      Julia
    s = 0
    for i in 1:length(x)
        s += x[i] * y[i]
    end
    s
end
...
```

```
function dot(x, y)      Julia
    s = zero(eltype(x))
    for i in 1:length(x)
        s += x[i] * y[i]
    end
    s
end
...
```



## Case 2

- Type stability of variables dramatically improves the performance.



- Much less memory allocation and GC time

```
$ julia -L case2.jl -e '@time dot(x, y)'  
elapsed time: 0.814667538 seconds (320013880 bytes allocated, 17.5
```

```
$ julia -L case2.fix.jl -e '@time dot(x, y)'  
elapsed time: 0.017035491 seconds (13896 bytes allocated)
```

0.814667538s / 0.017035491s = 48x faster

# Case 2

- The emitted LLVM code is cleaner when types are stable.

[illegible]

```

In [143]: | julia -L cas3.fix.jl -e '@code_livm_dot(x, y)'

define double @julia_dot_19813(%jl_value_t*, %jl_value_t*) {
top:
  %2 = getelementptr inbounds %jl_value_t*, @0, 164 2, idbg 1287
  %3 = bitcast %jl_value_t* %2 to i64*, idbg 1287
  %4 = load i64*, %3, align 8, idbg 1287, !tbaa %jtbaa_arraylen
  %5 = icmp eqg 164 %4, 0, idbg 1287
  br i1 %5, label %L.preheader, label %L5, idbg 1287

L.preheader:                                ; preds = %top
  %6 = load i64*, %3, align 8, idbg 1292, !tbaa %jtbaa_arraylen
  %7 = getelementptr inbounds %jl_value_t*, %1, 164 2, idbg 1292
  %8 = bitcast %jl_value_t* %7 to i64*, idbg 1292
  %9 = getelementptr inbounds %jl_value_t*, @0, 164 1, idbg 1292
  %10 = bitcast %jl_value_t* %9 to i64*, idbg 1292
  %11 = getelementptr inbounds %jl_value_t*, %1, 164 1, idbg 1292
  %12 = bitcast %jl_value_t* %11 to i8**, idbg 1292
  br label %L, idbg 1287

L:                                            ; preds = %L.preheader, %idxend2
  %"#2.0" = phi i64 [ %19, %idxend2 ], [ 1, %L.preheader ]
  %8.0 = phi double [ %29, %idxend2 ], [ 0.000000e+00, %L.preheader ]
  %13 = add i64 %4 %13, -1, idbg 1292
  %14 = icmp ult 164 %13, %6, idbg 1292
  br i1 %14, label %idxend, label %toob, idbg 1292

oob:
  %15 = load %jl_value_t** %jl_bounds_exception, align 8, idbg 1292, !tbaa %jtbaa_const
  call void @j_lthrow_with_superfluous_argument(%jl_value_t* %15, 132 4), idbg 1292
  unreachable, idbg 1292

idxend:                                    ; preds = %L
  %16 = load i64*, %8, align 8, idbg 1292, !tbaa %jtbaa_arraylen
  %17 = icmp ult 164 %13, %16, idbg 1292
  br i1 %17, label %idxend, label %toob1, idbg 1292

oob1:
  %18 = load %jl_value_t** %jl_bounds_exception, align 8, idbg 1292, !tbaa %jtbaa_const
  call void @j_lthrow_with_superfluous_argument(%jl_value_t* %18, 132 4), idbg 1292
  unreachable, idbg 1292

idxend2:                                   ; preds = %idxend
  %19 = add i64 %"#2.0", 1, idbg 1287
  %20 = load i8**, %10, align 8, idbg 1292, !tbaa %jtbaa_arrayptr
  %21 = bitcast 18** %20 to double*, idbg 1292
  %22 = getelementptr double*, %21, 164 %13, idbg 1292
  %23 = load double*, %22, align 8, idbg 1292, !tbaa %jtbaa_user
  %24 = load i8**, %12, align 8, idbg 1292, !tbaa %jtbaa_arrayptr
  %25 = bitcast 18** %24 to double*, idbg 1292
  %26 = getelementptr double*, %25, 164 %13, idbg 1292
  %27 = load double*, %26, align 8, idbg 1292, !tbaa %jtbaa_user
  %28 = fmul double %23, %27, idbg 1292
  %29 = fadd double %8.0, %28, idbg 1292
  %30 = icmp eq 164 %"#2.0", %4, idbg 1292
  br i1 %30, label %L5, label %L, idbg 1292

L5:
  %31 = phi double [ 0.000000e+00, %top ], [ %29, %idxend2 ]
  ret double %31, idbg 1296
}

```

## Type-unstable code

# Type-stable code

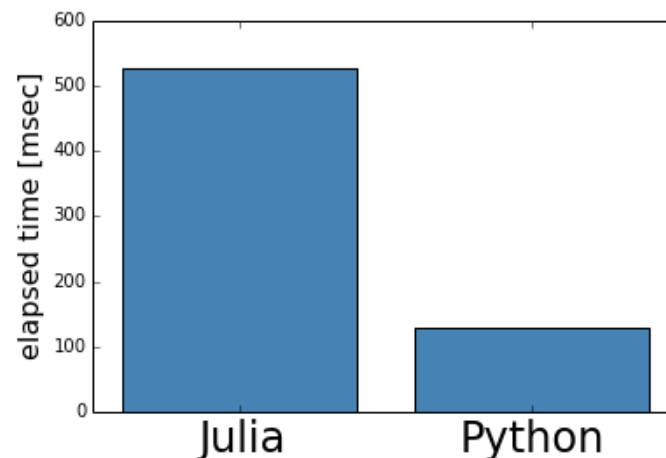
# Case 3

Julia

```
...  
function optimize(x0, n_iter)  
     $\eta$  = 0.1  
    x = copy(x0)  
    g = zeros(length(x0))  
    for i in 1:n_iter  
        grad!(g, x)  
        x -=  $\eta$  * g  
    end  
    x  
end
```

Python

```
...  
def optimize(x0, n_iter):  
    eta = 0.1  
    x = np.copy(x0)  
    g = np.zeros(len(x0))  
    for i in xrange(n_iter):  
        set_grad(g, x)  
        x -= eta * g  
    return x
```



# Case 3

## `-=` for array types is not an in-place operator!

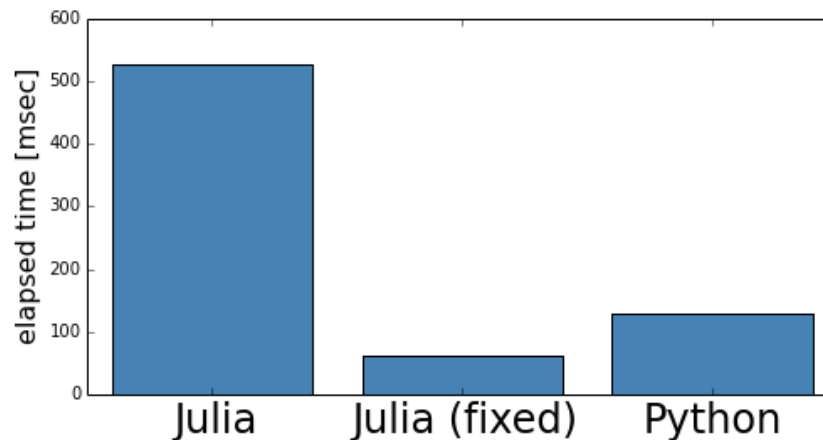
- `x -= η * g` works like:
  - `y = η * g`: product
  - `z = x - y`: subtract
  - `x = z`: bind
- For each iteration:
  - Create two arrays (`y` and `z`) and
  - Discard two arrays (`x` and `y`)
- Lots of memories are wasted!
  - 53% of time is consumed in GC!

```
$ julia -L case3.jl -e '@time optimize(x0, 50000)'  
elapsed time: 0.527045846 seconds  
(812429992 bytes allocated, 53.60% gc time)
```

# Case 3

```
Julia
...
function optimize(x0, n_iter)
    ...
    for i in 1:n_iter
        grad!(g, x)
        x -=  $\eta$  * g
    end
    x
end
```

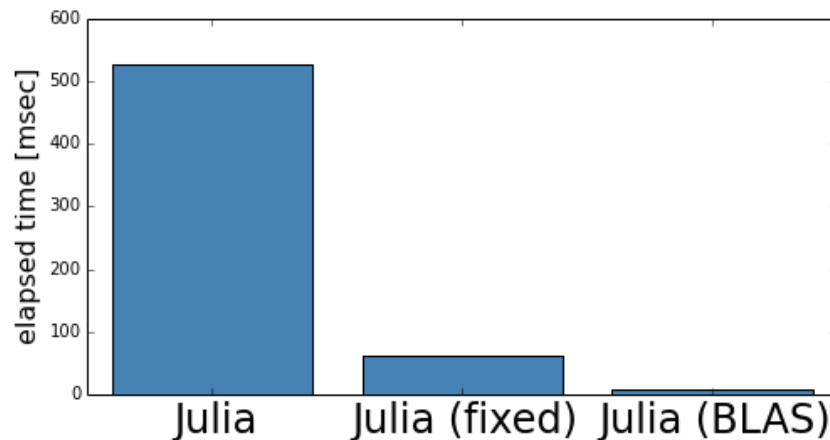
```
Julia
...
function optimize(x0, n_iter)
    ...
    for i in 1:n_iter
        grad!(g, x)
        for j in 1:n
            x[j] -=  $\eta$  * g[j]
        end
    end
    x
end
```



# Case 3

```
...                                     Julia
function optimize(x0, n_iter)
    ...
    for i in 1:n_iter
        grad!(g, x)
        for j in 1:n
            x[j] -=  $\eta$  * g[j]
        end
    end
    x
end
```

```
...                                     Julia
function optimize(x0, n_iter)
    ...
    for i in 1:n_iter
        grad!(g, x)
        BLAS.scal!(n,  $\eta$ , g, 1)
    end
    x
end
```



# Take-home Points

- Julia is definitely fast unless you do it wrong.
- You should know what code is *smooth* in Julia.
- [Manual/Performance Tips](#) is a definitive guide to know the smoothness of Julia.
- Utilize standard tools to analyze bottlenecks:
  - `@time` macro
  - `@profile` macro
  - `--track-allocation` option
  - `@code_typed` / `@code_lowered` / `@code_llvm` / `@code_native` macros (from higher level to lower)
  - `@code_warntype` (v0.4)