# Python on GPUs (work in progress!)







## Laurie Stephey GPUs for Science Day, July 3, 2019 Rollin Thomas, NERSC Lawrence Berkeley National Laboratory



# Python is friendly and popular



## **TIOBE Index for June 2019**

#### June Headline: Python continues to soar in the TIOBE index

This month Python has reached again an all time high in TIOBE index of 8.5%. If Python can keep this pace, it will probably replace C and Java in 3 to 4 years time, thus becoming the most popular programming language of the world. The main reason for this is that software engineering is booming. It attracts lots of newcomers to the field. Java's way of programming is too verbose for beginners. In order to fully understand and run a simple program such as "hello world" in Java you need to have knowledge of classes, static methods and packages. In C this is a bit easier, but then you will be hit in the face with explicit memory management. In Python this is just a one-liner. Enough said.

| Jun 2019 | Jun 2018 | Change | Programming Language | Ratings | Change |
|----------|----------|--------|----------------------|---------|--------|
| 1        | 1        |        | Java                 | 15.004% | -0.36% |
| 2        | 2        |        | С                    | 13.300% | -1.64% |
| 3        | 4        | *      | Python               | 8.530%  | +2.77% |
| 4        | 3        | *      | C++                  | 7.384%  | -0.95% |
| 5        | 6        | *      | Visual Basic .NET    | 4.624%  | +0.86% |



#### Screenshots from: <a href="https://www.tiobe.com/tiobe-index/">https://www.tiobe.com/tiobe-index/</a>









You have some Python code you like. Can you just run it on a GPU?

import numpy as np
from scipy import special
import qpu ?

## Unfortunately no.







## Right now, there is no "right" answer

- CuPy
- Numba
- pyCUDA (<u>https://mathema.tician.de/software/pycuda/</u>)
- pyOpenCL (<u>https://mathema.tician.de/software/pyopencl/</u>)
- Rewrite kernels in C, Fortran, CUDA...





## **DESI: Our case study**

Wavelength





# CuPy (https://cupy.chainer.org/)



- Developed by Chainer, supported in RAPIDS
- Meant to be a drop-in replacement for NumPy
- Some, but not all, NumPy coverage

| NumPy                             | СиРу          |
|-----------------------------------|---------------|
| numpy.abs                         | cupy.abs      |
| numpy.absolute                    | cupy.absolute |
| numpy.add                         | cupy.add      |
| <pre>numpy.add_docstring</pre>    | -             |
| numpy.add_newdoc                  | -             |
| <pre>numpy.add_newdoc_ufunc</pre> | -             |
| numpy.alen                        | -             |
| numpy.all                         | cupy.all      |
| numpy.allclose                    | cupy.allclose |

Screenshot from:

import numpy as np
import cupy as cp

```
cpu ans = np.abs(data)
```

#same thing on gpu
gpu\_data = cp.asarray(data)
gpu\_temp = cp.abs(gpu\_data)
gpu\_ans = cp.asnumpy(gpu\_temp)



Office of Science https://docs-cupy.chainer.org/en/stable/reference/comparison.htm



# eigh in CuPy



- Important function for DESI
- Compared CuPy eigh on Cori Volta GPU to Cori Haswell and Cori KNL
- Tried "divide-and-conquer" 10°
   approach on both CPU and GPU (1, 2, 5, 10 divisions)
- Volta wins only at very large<sup>10<sup>-2</sup></sup> matrix sizes
- Major pro: eigh really easy to use in CuPy!







# legval in CuPy



```
import cupy as cp
#expects cupy arrays as input
def legval_cupy(x, c):
```

```
#x and c are cupy arrays
```

```
ndd = c.shape[0] #will be an int
#and now change nd into a cupy array
nd = cp.array(ndd)
```

```
xlen = x.shape[0] #will be an int
```

```
c0=c[-2]*cp.ones(xlen) #cupy
c1=c[-1]*cp.ones(xlen) #cupy
```

```
for i in range(3, ndd, + 1):
    tmp = c0
    nd = nd - 1
    nd_inv = 1/nd
    c0 = c[-i] - (c1*(nd - 1))*nd_inv
    c1 = tmp + (c1*x*(2*nd - 1))*nd_inv
    cupy_result = c0 + c1*x
return cupy_result
```

- Easy to convert from NumPy arrays to CuPy arrays
- This function is ~150x slower than the cpu version!

10 iterations of legval\_numba ran in 0.0012524127960205078 s 10 iterations of legval\_cupy ran in 0.1504673957824707 s

- This implies there is probably some undesirable data movement between the cpu and gpu
- Maybe I'm just doing it wrong









- We like Numba because it has worked well for JIT-compiling code for a CPU
- Unlike CPU Numba, CUDA Numba doesn't allow most NumPy (which is a problem for scientific users like DESI)
- Numba for GPUs doesn't look very much like "normal" Python, looks a lot more like CUDA







#### Step 1: Invoke kernel with thread information

threadsperblock = 32
blockspergrid = (xx.size + (threadsperblock - 1)) // threadsperblock
for i in range(niter):
 legval\_numba\_gpu[blockspergrid, threadsperblock](xx, cc, c0, c1, results)

#### Step 3: Troubleshoot Numba type errors! (Common)

This error is usually caused by passing an argument of a type that is unsupported by the named f unction.

[1] During: typing of intrinsic-call at /global/cscratch1/sd/stephey/git\_repo/specter/py/specter
/util/util.py (299)

#### Step 2: Numba gpu function

from numba import cuda @cuda.jit def legval\_numba\_gpu(x, c, c0, c1, results): nd = len(c) ndd = nd xlen = x.size #c0=c[-2]\*np.ones(xlen) #c1=c[-1]\*np.ones(xlen) #need a way to do this without numpy for i in range(3, ndd + 1): #instead use i from cuda.grid #i = cuda.grid(1) tmp = c0

> c0 = c[-i] - (c1\*(nd - 1))\*nd\_inv c1 = tmp + (c1\*x\*(2\*nd - 1))\*nd inv

#### Step 4: Profit!

results = c0 + c1 \* x

nd = nd - 1

 $nd_inv = 1/nd$ 





## • Screenshot from:

https://numba.pydata.org/numba-doc/dev/cuda/exa

### mples.html

```
@cuda.jit
def matmul(A, B, C):
    """Perform square matrix multiplication of C = A * B
    """
    i, j = cuda.grid(2)
    if i < C.shape[0] and j < C.shape[1]:
        tmp = 0.
        for k in range(A.shape[1]):
            tmp += A[i, k] * B[k, j]
        C[i, j] = tmp</pre>
```







- Code that runs well on a CPU might not be good for a GPU
- More than just porting some kernels/functions, it could require a substantial rewrite
- How to avoid doing this every few years? How to be able to run on many architectures?
- Unfortunately there are no easy answers





## What have we learned?



- Python on GPUs is still evolving
- We have tried:
  - $\circ$  CuPy  $\rightarrow$  difficulty easy, but not every NumPy/SciPy function
  - $\circ$  Numba  $\rightarrow$  difficulty hard, looks less like Python, but more flexible
- Our job is to help DESI and our users figure out the best strategy (performance + maintainability + portability)
- Stay tuned!





# Thank you!







