

Python for High Performance Computing

Pawel Pomorski
SHARCNET
University of Waterloo
ppomorsk@sharcnet.ca

June 14, 2019

Outline

- ▶ Speeding up Python code with NumPy
- ▶ Speeding up Python code with Cython
- ▶ Speeding up Python code with ctypes
- ▶ Using multiprocessing with Python
- ▶ Using MPI with Python, via mpi4py
- ▶ Using CUDA with Python, via PyCUDA

What is Python?

- ▶ Python is a programming language that appeared in 1991
- ▶ compare with Fortran (1957), C (1972), C++ (1983),
- ▶ While the older languages still dominate High Performance Computing (HPC), popularity of Python is growing

Python advantages

- ▶ Designed from the start for better code readability
- ▶ Allows expression of concepts in fewer lines of code
- ▶ Has dynamic type system, variables do not have to be declared
- ▶ Has automatic memory management
- ▶ Has large number of easily accessible, extensive libraries (eg. NumPy, SciPy)
- ▶ All this makes developing new codes easier

Python disadvantages

- ▶ Python is generally slower than compiled languages like C, C++ and Fortran
- ▶ Complex technical causes include dynamic typing and the fact that Python is interpreted, not compiled
- ▶ This does not matter much for a small desktop program that runs quickly.
- ▶ However, this will matter a lot in a High Performance Computing environment.
- ▶ Python use in HPC parallel environments is relatively recent, hence parallel techniques less well known
- ▶ Rest of this talk will describe approaches to ensure your Python code runs reasonably fast and in parallel

1D diffusion equation

To describe the dynamics of some quantity $u(x,t)$ (eg. heat) undergoing diffusion, use:

$$\frac{\partial u}{\partial t} = \kappa \frac{\partial^2 u}{\partial x^2}$$

Problem: given some initial condition $u(x,t=0)$, determine time evolution of u and obtain $u(x,t)$

Use finite difference with Euler method for time evolution

$$u(i\Delta x, (m+1)\Delta t) = u(i\Delta x, m\Delta t) + \frac{\kappa\Delta t}{\Delta x^2} \left[u((i+1)\Delta x, m\Delta t) + u((i-1)\Delta x, m\Delta t) - 2u(i\Delta x, m\Delta t) \right]$$

C code

```
1 #include <math.h>
2 #include <stdio.h>
3
4 int main() {
5
6     int const n=100000, niter=2000;
7
8     double x[n], u[n], udt[n];
9     int i, iter;
10    double dx=1.0;
11    double kappa=0.1;
12
13    for (i=0; i<n; i++){
14        u[i]=exp(-pow(dx*(i-n/2.0), 2.0)/100000.0);
15        udt[i]=0.0;
16    }
17
18    ...
```

C code continued :

```
1  ...
2  for( iter=0;iter<niter;iter++){
3      for ( i=1;i<n-1;i++){
4          udt[i]=u[i]+kappa*(u[i+1]+u[i-1]-2*u[i]);
5      }
6      for ( i=0;i<n;i++){
7          u[i]=udt[i];
8      }
9  }
10 return 0;
11 }
```


Program output

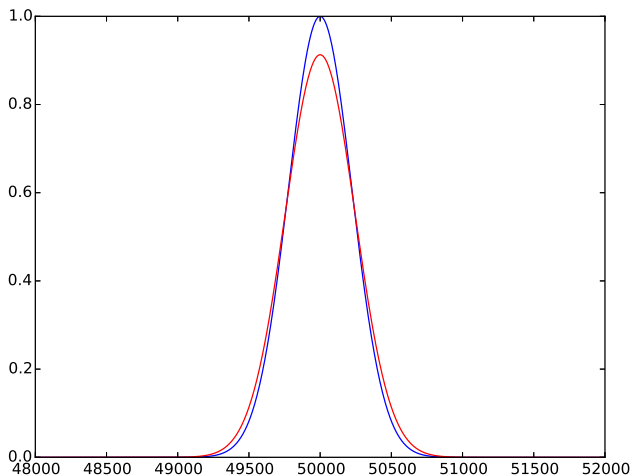


Figure: Evolution of $u(x)$ after 50,000 time steps (blue line initial, red line final)

"Vanilla" Python code

```
1 import math
2 n=100000 ; dx=1.0 ; niter=2000 ; kappa=0.1
3
4 x=n*[0.0 ,]
5 u=n*[0.0 ,]
6 udt=n*[0.0 ,]
7
8 for i in xrange(n):
9     u[i]=math.exp( -(dx*(i-n/2))**2/100000)
10
11 fac=1.0-2.0*kappa
12 for itern in xrange(niter):
13
14     for i in xrange(1,n-1):
15         udt[i]=fac*u[i]+kappa*(u[i+1]+u[i-1])
16
17     for i in xrange(n):
18         u[i]=udt[i]
```

Vanilla code performance

- ▶ 2000 iterations, tested on node in graham cluster
- ▶ C code compiled with Intel compiler (`icc -xCORE-AVX2 -fma -O2`) takes 0.30 seconds
- ▶ Python "vanilla" code takes 100.1 seconds
- ▶ Python is much slower (by factor 334)
- ▶ Code is slow because loops are explicit

NumPy

- ▶ To achieve reasonable efficiency in Python, will need support for efficient, large numerical arrays
- ▶ These provided by NumPy, an extension to Python
- ▶ NumPy (<http://www.numpy.org/>) along with SciPy (<http://www.scipy.org/>) provide a large set of easily accessible libraries which make Python so attractive to the scientific community
- ▶ The goal is to eliminate costly explicit loops and replace them with numpy operations instead
- ▶ Numpy functions invoke efficient libraries written in C
- ▶ The difficulty of eliminating costly explicit loops varies.

Slicing NumPy arrays :

```
1 sharcnet1:~ pawelpomorski$ python
2 Python 2.7.9 (default, Dec 12 2014, 12:40:21)
3 [GCC 4.2.1 Compatible Apple LLVM 6.0 (clang-600.0.56)]
   on darwin
4 Type "help", "copyright", "credits" or "license" for
   more information.
5 >>> import numpy as np
6 >>> a=np.arange(10)
7 >>> a
8 array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
9 >>> a[1:-1]
10 array([1, 2, 3, 4, 5, 6, 7, 8])
11 >>> a[0:-2]
12 array([0, 1, 2, 3, 4, 5, 6, 7])
13 >>> a[1:-1]+a[0:-2]
14 array([ 1,  3,  5,  7,  9, 11, 13, 15])
15 >>>
```

NumPy vector operations

Replace explicit loops

```
1  for i in xrange(1,n-1):  
2      udt[i]=fac*u[i]+kappa*(u[i+1]+u[i-1])
```

with NumPy vector operations using slicing

```
1  udt[1:-1]=fac*u[1:-1]+kappa*(u[0:-2]+u[2:])
```

Python code using Numpy operations instead of loops

```
1
2 import numpy as np
3
4 n=100000; dx=1.0; niter=50000; kappa=0.1
5
6 x=np.arange(n,dtype="float64")
7 u=np.empty(n,dtype="float64")
8 udt=np.empty(n,dtype="float64")
9
10 u_init = lambda x: np.exp( -(dx*(x-n/2))**2/100000)
11 u=u_init(x)
12 udt[:]=0.0
13 fac=1.0-2.0*kappa
14 for itern in xrange(niter):
15     udt[1:-1]=fac*u[1:-1]+kappa*(u[0:-2]+u[2:])
16     u[:]=udt[:]
```

Performance

- ▶ 50,000 iterations, tested on node in graham cluster
- ▶ C code compiled with `icc -xCORE-AVX2 -fma` - 7.32 s
- ▶ Python code with NumPy operations - 22.33 s
- ▶ Python 3.05 times slower than C code
- ▶ The compiler can achieve a higher degree of optimization than the numpy library

Numpy libraries

- ▶ For standard operations, eg. matrix multiply, will run at the speed of underlying numerical library
- ▶ Performance will strongly depend on which library is used, can see with **`numpy.show_config()`**
- ▶ If libraries are threaded, python will take advantage of multithreading, with no extra programming (free parallelism)

General approaches for code speedup

- ▶ NumPy does not help with all problems, some don't fit array operations
- ▶ Need a more general technique to speed up Python code
- ▶ As the problem is that Python is not a compiled language, one can try to compile it
- ▶ General compiler: **nuitka** (<http://nuitka.net/>) under active development
- ▶ **PyPy** (<http://pypy.org/>) - Just-in-Time (JIT) compiler
- ▶ **Cython** (<http://cython.org>) - turns Python program into C and compiles it

PyPy

- ▶ Python distribution with just in time compiling
- ▶ Under development, NumPy not fully supported yet
- ▶ PyPy is available on on graham as a module
- ▶ **module load pypy/5.8**
- ▶ Then, to run a python program with PyPy:
- ▶ **pypy yourprogram.py**
- ▶ Result on graham cluster for diffusion with 50,000 iterations of vanilla Python code: 30.5 s
- ▶ About 4.2 times slower than C code compiled with icc, 1.4 times slower than NumPy code

Euler problem

If p is the perimeter of a right angle triangle with integral length sides, a, b, c , there are exactly three solutions for $p = 120$.

$(20, 48, 52)$, $(24, 45, 51)$, $(30, 40, 50)$

For which value of $p < N$, is the number of solutions maximized?
Take $N=1000$ as starting point

(from <https://projecteuler.net>)

Get solutions at particular p

```
1 def find_num_solutions(p):  
2     n=0  
3     # a+b+c=p  
4     for a in range(1,p/2):  
5         for b in range(a,p):  
6  
7             c=p-a-b  
8             if (c>0):  
9                 if (a*a+b*b==c*c):  
10                     n=n+1  
11  
12     return n
```

Loop over possible value of p up to N

```
1 nmax=0 ; imax=0
2 N=1000
3
4 for i in range(1,N):
5     print i
6     nsols=find_num_solutions(i)
7     if (nsols>nmax):
8         nmax=nsols ; imax=i
9
10 print "maximum p , number of solutions",imax,nmax
```

Cython

- ▶ The goal is to identify functions in the code where it spends the most time. Python has profiler already built in
- ▶ **python -m cProfile euler39.py**
- ▶ Place those functions in a separate file so they are imported as module
- ▶ Cython will take a python module file, convert it into C code, and then compile it into a shared library
- ▶ Python will import that compiled library module at runtime just like it would import a standard Python module
- ▶ To make Cython work well, need to provide some hints to the compiler as to what the variables are, by defining some key variables

Invoking Cython

- ▶ Place module code (with Cython modifications) in `find_num_solutions.pyx`
- ▶ Create file `setup.py`

```
1 from distutils.core import setup
2 from Cython.Build import cythonize
3
4 setup(
5     ext_modules=cythonize("find_num_solutions.pyx"),
6 )
```

- ▶ Execute: `python setup.py build_ext -inplace`
- ▶ Creates `find_num_solutions.c`, C code implementation of the module
- ▶ From this creates `find_num_solutions.so` library which can be imported as Python module at runtime

Get solutions at particular p, cythonized

```
1 def find_num_solutions(int p): # note definition
2     cdef int a,b,c,n          # note definition
3     n=0
4     # a+b+c=p
5     for a in range(1,p/2):
6         for b in range(a,p):
7
8             c=p-a-b
9
10            if (c>0):
11                if (a*a+b*b==c*c):
12                    n=n+1
13
14    return n
```

This code in file find_num_solutions.pyx

Loop over possible value of p up to N , with Cython

Note changes at line 1 and line 7

```
1 import find_num_solutions
2
3 nmax=0 ; imax=0 ; N=1000
4
5 for i in range(1,N):
6     print i
7     nsols=find_num_solutions.find_num_solutions(i)
8     if(nsols>nmax):
9         nmax=nsols ; imax=i
10
11 print "maximum p and , number of solutions",imax,nmax
```

Speedup with Cython

For $N=1000$, tested on graham cluster

- ▶ vanilla python : 5.74 s
- ▶ Cython without variable definitions : 4.36 s, speedup factor 1.31
- ▶ Cython with integer variables defined : 0.056 s, speedup factor 102

ctypes - a foreign function library for Python

```
1 # compile C library with:  
2 # gcc -shared -o findnumsolutions.so findnumsolutions.c  
3 import ctypes  
4 findnumsolutions = ctypes.CDLL('./findnumsolutions.so')  
5 # ...  
6 nsols=findnumsolutions.findnumsolutions(i)
```

C code : findnumsolutions.c

```
1 int findnumsolutions(int p){
2     int a,b,c,n;
3     n=0;
4     for(a=1;a<p/2;a++){
5         for (b=a;b<p/2;b++){
6             c=p-a-b;
7             if (a*a+b*b==c*c) {
8                 n=n+1;
9             }
10        }
11    }
12    return n;
13 }
```

Speedup with ctypes

For N=1000, tested on graham node

- ▶ vanilla python : 5.74 s
- ▶ Cython without variable definitions : 4.36 s, speedup factor 1.31
- ▶ Cython with integer variables defined : 0.056 s, speedup factor 102
- ▶ ctypes (icc -O2) : 0.035 s, speedup factor 164, 1.6 times faster than cython
- ▶ ctypes (icc -xCORE-AVX2 -fma -O2) : 0.011 s , speedup factor 521, 5.1 times faster than cython
- ▶ pure C code (icc -xCORE-AVX2 -fma -O2): 0.0105 s (almost same as ctypes)
- ▶ It's important to choose the best compiler (Intel more efficient than GCC)
- ▶ pure C code (gcc -O2): 0.052 s

Parallelizing Python

- ▶ Once the serial version is optimized, need to parallelize Python to do true HPC
- ▶ Threading approach does not work due to Global Interpreter Lock
- ▶ In Python, you can have many threads, but only one executes at any one time, hence no speedup
- ▶ Have to use multiple processes instead
- ▶ Python has multiprocessing module but that only works within one node. Have to use MPI to achieve parallelism over many nodes

Multiprocessing - apply

```
1 import time,os
2 from multiprocessing import Pool
3
4 def f():
5     start=time.time()
6     time.sleep(2)
7     end=time.time()
8     print "inside f pid",os.getpid()
9     return end-start
10
11 p = Pool(processes=1)
12 result = p.apply(f)
13 print "apply is blocking, total time",result
14
15 result=p.apply_async(f)
16 print "apply_async is non-blocking"
17
18 while not result.ready():
19     time.sleep(0.5)
20     print "could work here while result computes"
21
22 print "total time",result.get()
```


Multiprocessing - Map

```
1 import time
2 from multiprocessing import Pool
3
4 def f(x):
5     return x**3
6
7 y = range(int(1e7))
8 p= Pool (processes=8)
9
10 start= time.time()
11 results = p.map(f,y)
12 end = time.time()
13
14 print "map blocks on launching process, time=",end-
    start
15
16 # map_async
17 start = time.time()
18 results = p.map_async(f,y)
19 print "map_async is non-blocking on launching process"
20 output = results.get()
21 end=time.time()
22 print "time",end-start
```

Euler problem with multiprocessing

```
1 from multiprocessing import Pool
2 import find_num_solutions
3
4 p= Pool (processes=4)
5
6 y=range(1,1000)
7 results=p.map( find_num_solutions.find_num_solutions , y)
8 print "answer",y[results.index(max(results))]
```

Multiprocessing performance

timing on graham node (32 cores)

n=10000 case

Number of processes	time(s)	speedup
1	55.0	1.0
2	33.0	1.7
4	18.3	3.0
8	9.56	5.8
16	5.13	10.7
32	2.99	18.4

Should scale better for larger values of N

MPI - Message Passing Interface

- ▶ Approach has multiple processors with independent memory running in parallel
- ▶ Since memory is not shared, data is exchanged via calls to MPI routines
- ▶ Each process runs same code, but can identify itself in the process set and execute code differently

Compare MPI in C and Python with mpi4py - MPI reduce

```
1 int main(int argc, char* argv[]) {  
2     int my_rank, imax, imax_in;  
3     MPI_Init(&argc, &argv);  
4     MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);  
5     imax_in=my_rank;  
6     MPI_Reduce(&imax_in, &imax, 1, MPI_INT, MPI_MAX, 0,  
7     MPI_COMM_WORLD);  
8     if (my_rank == 0) printf("%d \n", imax);  
9     MPI_Finalize();  
10    return 0;}
```

```
1 from mpi4py import MPI  
2 comm = MPI.COMM_WORLD  
3 myid = comm.Get_rank()  
4 imax_in = myid  
5 imax = comm.reduce(imax_in, op=MPI.MAX)  
6 if (myid==0):  
7     print imax  
8 MPI.Finalize
```

Loop over p values up to N distributed among MPI processes

```
1 from mpi4py import MPI
2 import find_num_solutions
3
4 comm = MPI.COMM_WORLD
5 myid = comm.Get_rank()
6 nprocs = comm.Get_size()
7
8 nmax=0 ; imax=0 ; N=5000
9
10 for i in range(1,N):
11
12     if (i%nprocs==myid):
13         nsols=find_num_solutions.find_num_solutions(i)
14         if (nsols>nmax):
15             nmax=nsols ; imax=i
16
17 nmax_global=comm.allreduce(nmax,op=MPI.MAX)
18 if (nmax_global==nmax):
19     print "process ",myid," found maximum at ",imax
21 MPI.Finalize
```

MPI performance

timing on graham

n=10000 case

MPI processes	time(s)	speedup
1	55.34	1.0
2	28.34	1.95
4	14.85	3.73
8	8.352	6.63
16	5.841	9.47
32	5.902	9.37
64	8.578	6.45

Will scale better for larger values of N. Main culprit for bad speedup is the long time it takes to import modules at the start.

Python on GPUs

- ▶ PyCUDA - Python wrapper for CUDA
(<https://mathemat.ician.de/software/pycuda/>)
- ▶ GPU Kernels must still be written in CUDA C
- ▶ Aside from that, more convenient to use than CUDA
- ▶ Popular software implemented in Python with GPU acceleration
- ▶ Theano (<http://deeplearning.net/software/theano/>)
- ▶ TensorFlow (<https://www.tensorflow.org/>)

PyCUDA example:

```
1 import pycuda.driver as drv
2 import pycuda.tools
3 import pycuda.autoinit
4 import numpy
5 import numpy.linalg as la
6 from pycuda.compiler import SourceModule
7
8 mod = SourceModule("""
9     __global__ void multiply_them(float *dest, float *a,
10         float *b)
11     {
12         const int i = threadIdx.x;
13         dest[i] = a[i] * b[i];
14     }
15 """)
16 multiply_them = mod.get_function("multiply_them")
```

...

PyCUDA example - continued:

```
1 a = numpy.random.randn(400).astype(numpy.float32)
2 b = numpy.random.randn(400).astype(numpy.float32)
3
4 dest = numpy.zeros_like(a)
5 multiply_them(
6     drv.Out(dest), drv.In(a), drv.In(b),
7     block=(400,1,1))
8
9 print dest-a*b
```

Conclusion

- ▶ Python is generally slower than compiled languages like C
- ▶ With a bit of effort, can take a Python code which is a great deal slower and make it only somewhat slower
- ▶ The tradeoff between slower code but faster development time is something the programmer has to decide
- ▶ Tools currently under development should make this problem less severe over time