

HPC Python Programming

Ramses van Zon

SciNet HPC Consortium

Ontario HPCSS, July 2016

In this session...

- 1 Performance and Python
- 2 Profiling tools for Python
- 3 Numpy recap: Fast arrays for python
- 4 Multicore computations:
 - ▶ **Numexpr**
 - ▶ **Threading**
 - ▶ **Multiprocessing**
 - ▶ **Mpi4py**
- 5 Map/reduce approaches
 - ▶ **IPython Parallel**
 - ▶ **Apache Spark**

Getting started

Packages and code

Requirements for this session

If following along on your own laptop, you need the following packages:

- numpy
- psutil
- pyzmq
- scipy
- line_profiler
- mpi4py
- numexpr
- memory_profiler
- ipyparallel or
IPython.parallel
- matplotlib
- theano

Get the code and setup files on GPC

Code and installation can be copied from a GPC directory. It's in the directory `/scinet/course/ss2016/hpcpy`. The code accompanying this session is in the code subdirectory.

Setting up for today's class (GPC)

To get set up for today's class, perform the following steps.

1 Login to SciNet's GPC

```
$ ssh -Y USERNAME@login.scinet.utoronto.ca  
$ ssh -Y gpc
```

2 Install code and software to your own directory

```
$ cd $SCRATCH  
$ cp -r /scinet/course/ss2016/hpcpy .  
$ cd $SCRATCH/hpcpy/code  
$ source setup
```

The last command will install a few packages into your local account, so as to satisfy the requirements, and will load the correct modules.

3 Request an interactive session on a compute node

```
$ qsub -lnodes=1:ppn=8,walltime=3:00:00 -I -X -q teach
```



Introduction

Performance and Python

- Python is a high-level, interpreted language.
- Those defining features are often at odds with “high performance”.
- But the development in Python can be substantially easier (and thus faster) than compiled languages.
- In this session, we will explore when using Python still makes sense and how to get the most performance out of it, without losing the flexibility and ease of development.

What would make Python not “high performance”?

Interpreted language:

- Translation to machine language happens line-by-line as the script is read.
- Repeated lines are no faster.
- Cross-line optimizations are not possible.

What would make Python not “high performance” ?

Interpreted language:

- Translation to machine language happens line-by-line as the script is read.
- Repeated lines are no faster.
- Cross-line optimizations are not possible.

Dynamic language:

- Types are part of the data: extra overhead
- Memory management is automatic. Behind the scene that means reference counting and garbage collection.
- All this also interferes with optimal streaming of data to processor, which interferes with maximum performance.

Example: 2D diffusion equation

Suppose we are interested in the time evolution of the two-dimension diffusion equation:

$$\frac{\partial p(x, y, t)}{\partial t} = D \left(\frac{\partial^2 p(x, y, t)}{\partial x^2} + \frac{\partial^2 p(x, y, t)}{\partial y^2} \right),$$

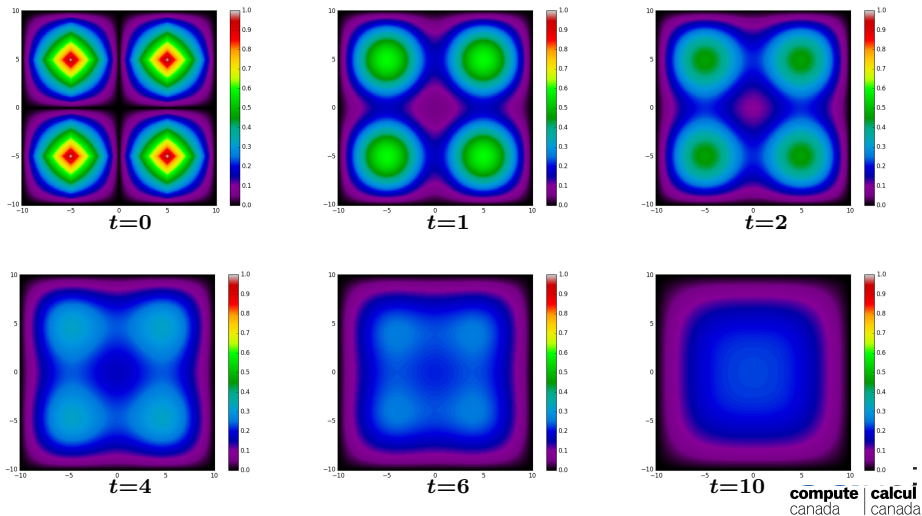
on domain $[x_1, x_2] \otimes [x_1, x_2]$,
with $P(x, y, t) = 0$ at all times for
all points on the domain boundary,
and for some given initial condition
 $p(x, y, t) = p_0(x, y)$.

Here:

- P : density
- x, y : spatial coordinates
- t : time
- D : diffusion constant

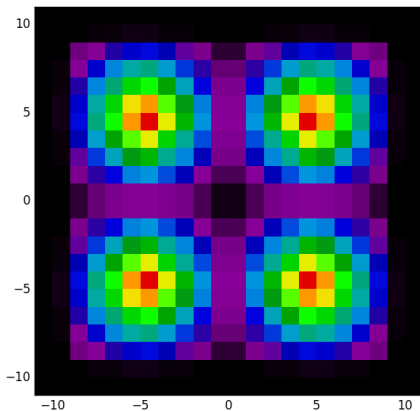
Example: 2D diffusion, result

$x_1 = -10, x_2 = 10, D = 1$, four-peak initial condition.



Example: 2D diffusion, algorithm

- Discretize space in both directions (points dx apart)
- Replace derivatives with finite differences.
- Explicit finite time stepping scheme (time step set by dx)
- For graphics: Matplotlib for python, pgplot for c++/fortran, every outtime time units



Parameters in file diff2dparams.py

Example: 2D diffusion, parameters

The fortran, C++ and python codes all read the same files (by some special tricks).

diff2dparams.py

```
D          = 1.0;
x1         = -10.0;
x2         = 10.0;
runtime    = 15.0;
dx         = 0.0667;
outtime    = 0.5;
graphics   = False;
```

Example: 2D diffusion, performance

The files `diff2d.cpp`, `diff2.f90` and `diff2d.py` contain the same algorithm, in C++, Fortran, and Python, respectively.

Example: 2D diffusion, performance

The files `diff2d.cpp`, `diff2.f90` and `diff2d.py` contain the same algorithm, in C++, Fortran, and Python, respectively.

```
$ etime() { /scinet/gpc/bin6/time -f "Elapsed: %e seconds" $@; }
$ etime make diff2d_cpp.ex diff2d_f90.ex
g++ -c -std=c++11 -O3 -o diff2d_cpp.o diff2d.cpp
gfortran -c -O3 -o pgplot90.o pgplot90.f90
...
Elapsed: 1.80 seconds
```

Example: 2D diffusion, performance

The files `diff2d.cpp`, `diff2.f90` and `diff2d.py` contain the same algorithm, in C++, Fortran, and Python, respectively.

```
$ etime() { /scinet/gpc/bin6/time -f "Elapsed: %e seconds" $@; }
$ etime make diff2d_cpp.ex diff2d_f90.ex
g++ -c -std=c++11 -O3 -o diff2d_cpp.o diff2d.cpp
gfortran -c -O3 -o pgplot90.o pgplot90.f90
...
Elapsed: 1.80 seconds
```

```
$ etime ./diff2d_cpp.ex > output_c.txt
Elapsed: 0.92 seconds
$ etime ./diff2d_f90.ex > output_f.txt
Elapsed: 0.60 seconds
$ etime python diff2d.py > output_n.txt
Elapsed: 337.20 seconds
```

This doesn't look too promising for Python for HPC...

Then why do we bother with Python?

```
import numpy as np
from diff2dplot import plotdens
D      = 1.0;
x1     = -10.0;
x2     = 10.0;
runtime = 15.0;
dx     = 0.0666;
outtime = 0.5;
nrows  = int((x2-x1)/dx)
npnts  = nrows + 2
xm     = (x1+x2)/2
dx     = (x2-x1)/nrows
dt     = 0.25*dx**2/D
nsteps = int(runtime/dt)
nper   = int(outtime/dt)
x=np.linspace(x1-dx,x2+dx,npnts)
dens1 = np.zeros((npnts,npnts))
dens2 = np.zeros((npnts,npnts))
lapl  = np.zeros((npnts,npnts))
simtime = 0
```

```
for i in xrange(1,npnts-1):
    a=1-abs(1-4*abs((x[i]-xm)/(x2-x1)))
    for j in xrange(1,npnts-1):
        b=1-abs(1-4*abs((x[j]-xm)/(x2-x1)))
        dens1[i][j]=a*b
print(simtime)
plotdens(dens1,x[0],x[-1],True)
for s in xrange(nsteps):
    lapl[1:nrows+1,1:nrows+1]=(
        dens1[2:nrows+2,1:nrows+1]
        +dens1[0:nrows+0,1:nrows+1]
        +dens1[1:nrows+1,2:nrows+2]
        +dens1[1:nrows+1,0:nrows+0]
        -4*dens1[1:nrows+1,1:nrows+1])
    dens2[:,:]=dens1+D/dx**2*dt*lapl
    dens1,dens2 = dens2,dens1
    simtime += dt
    if (s+1)%nper == 0:
        print(simtime)
        plotdens(dens1,x[0],x[-1])
```

Then why do we bother with Python?

- Python lends itself easily to writing clear, concise code. (2d diffusion fits on one slide!)
- Python is very flexible: large set of very useful packages.
- Easy of use → shorter development time
- Python's performance hit is most prominent on 'tightly coupled' calculation on fundamental data types that are known to the cpu (integers, doubles), which is exactly the case for the 2d diffusion.
- It does much less worse on file I/O, text comparisons, list manipularions etc.
- Hooks to compiled libraries to remove worst performance pitfalls.
- Once the performance isn't too bad, we can start thinking of parallelization, i.e., using more cpu cores working on the same problem.

Performance tuning tools for Python

CPU performance

- Performance is about maximizing the utility of a resource.
- This could be cpu processing power, memory, network, file I/O, etc.
- Let's focus on **cpu performance** first.

CPU Profiling by function

- To consider the cpu performance of functions, but not of individual lines in your code, there is the package called `cProfile`.

CPU Profiling by line

- To find cpu performance bottlenecks by line of code, there is package called `line_profiler`

cProfile

- Use cProfile or profile to know in which functions your script spends its time.
- You usually do this on a smaller but representative case.
- The code should be reasonably modular, i.e., with separate functions for different tasks, for cProfile to be useful.

Example

```
$ python -m cProfile -s cumulative diff2d_numpy.py
```

```
...
```

```
92333 function calls (92212 primitive calls) in 4.236 sec
```

```
Ordered by: cumulative time
```

ncalls	totttime	percall	cumtime	percall	filename:lineno(func
1	0.000	0.000	4.236	4.236	diff2d_numpy.py:9(<m
1	4.171	4.171	4.176	4.176	diff2d_numpy.py:15(m
3	0.010	0.003	0.078	0.026	__init__.py:1(<modul
1	0.003	0.003	0.059	0.059	__init__.py:106(<mod
1	0.000	0.000	0.045	0.045	add_newdocs.py:10(<m
1	0.000	0.000	0.034	0.034	type_check.py:3(<mod

line_profiler

- Use `line_profiler` to know, line-by-line, where your script spends its time.
- You usually do this on a smaller but representative case.
- First thing to do is to have your code be in a function.
- You also need to include modify your script slightly:
 - ▶ Decorate your function with `@profile`
 - ▶ Run your script on the command line with

```
$ kernprof -l -v SCRIPTNAME
```

line_profiler script instrumentation

Script before:

```
x=[1.0]*(2048*2048)
a=str(x[0])
a+="\nis a one\n"
del x
print(a)
```

line_profiler script instrumentation

Script before:

```
x=[1.0]*(2048*2048)
a=str(x[0])
a+="\nis a one\n"
del x
print(a)
```

Script after:

```
#file: profileme.py
@profile
def profilewrapper():
    x=[1.0]*(2048*2048)
    a=str(x[0])
    a+="\nis a one\n"
    del x
    print(a)
profilewrapper()
```

Run at the command line:

```
$ kernprof -l -v profileme.py
```


Output of line_profiler

```
1.0  
is a one
```

```
Wrote profile results to profileme.py.lprof  
Timer unit: 1e-06 s
```

```
Total time: 0.039049 s  
File: profileme.py  
Function: profilewrapper at line 2
```

Line #	Hits	Time	Per Hit	% Time	Line Contents
2					@profile
3					def profilewrapper
4	1	25832	25832.0	66.2	x=[1.0]*(2048*2048)
5	1	33	33.0	0.1	a=str(x[0])
6	1	3	3.0	0.0	a+="\nis a one\n"
7	1	13130	13130.0	33.6	del x
8	1	51	51.0	0.1	print(a)

Memory performance

Why worry about this?

Memory performance

Why worry about this?

Once your script runs out of memory, one of a number of things may happen:

- Computer may start using the harddrive as memory: **very slow**
- Your application crashes
- Your (compute) node crashes

Memory performance

Why worry about this?

Once your script runs out of memory, one of a number of things may happen:

- Computer may start using the harddrive as memory: **very slow**
- Your application crashes
- Your (compute) node crashes

How could you run out of memory?

- You're not quite sure how much memory your program takes.
- Python objects may take more memory than expected.
- Some functions may temporarily use extra memory.
- Python relies on a garbage collector to clean up unused variables

Garbage collector

- Python uses garbage collector to clean up un-needed variables
- You can force the garbage collection to run at any time by running:

```
>>> import gc
>>> collect = gc.collect()
```

- Running gc by hand should only be done in specific circumstances.
- You can also remove objects with del (if object larger than 32MB):

```
>>> x = [0,0,0,0]
>>> del x
>>> print (x)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'x' is not defined
```

But how would you know when the memory usage is problematic?

memory_profiler

- This module/utility monitors the python memory usage and its changes throughout the run.
- Good for catching memory leaks and unexpectedly large memory usage.
- Needs same instrumentation as line profiler.
- Requires the psutil module (at least on windows, but helps on linux/mac too).

memory_profiler, details

Your decorated script is usable by memory profiler.

You run your script through the profiler with the command

```
$ python -m memory_profiler profileme.py
```

memory_profiler, details

Your decorated script is usable by memory profiler.

You run your script through the profiler with the command

```
$ python -m memory_profiler profileme.py
```

```
1.0
is a one
```

Filename: profileme.py

Line #	Mem usage	Increment	Line Contents
2	19.223 MiB	0.000 MiB	@profile
3			def profilewrapper():
4	51.234 MiB	32.012 MiB	x=[1.0]*(2048*2048)
5	51.242 MiB	0.008 MiB	a=str(x[0])
6	51.242 MiB	0.000 MiB	a+="\nis a one\n"
7	19.238 MiB	-32.004 MiB	del x
8	19.242 MiB	0.004 MiB	print(a)

Hands-on

Profile the diff2d.py code

- Reduce the resolution in diff2dparams.py, i.e., increase dx to 0.1.
- In the same file, set `graphics=False`.
- Add `@profile` to the main function
- Run this through both the line and memory profilers.
 - ▶ What lines cause the most memory usage?
 - ▶ What lines cause the most cpu usage?

Numpy: faster numerical arrays for python

Numpy Recap

In the morning session we saw:

- Lists aren't the ideal data type
- Almost everything that you want to do starts with NumPy. Why?
 - 1 Lists are slow. They can contain any type, so the type has to be stored alongside every item, and when operations on lists are done, the type of the items must be checked each time.
 - 2 Secondly, lists do not allow elementwise algebra.

Nice, useful numpy features

- Numpy provides multi-dimensional arrays.
- Elements of numpy arrays are all of the same type.
- Numpy arrays do element-wise algebra
- `numpy.dot` does linear algebra multiplications (matrix-vector, vector-vector, inner product).
- Slicing etc.

Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #1: Array allocation/initialization

`diff2d.py` had:

```
x = [x1+((i-1)*(x2-x1))/nrows for i in xrange(npnts)]
dens = [[0.0]*npnts for i in xrange(npnts)]
densnext = [[0.0]*npnts for i in xrange(npnts)]
```

Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #1: Array allocation/initialization

`diff2d.py` had:

```
x = [x1+((i-1)*(x2-x1))/nrows for i in xrange(npnts)]
dens = [[0.0]*npnts for i in xrange(npnts)]
densnext = [[0.0]*npnts for i in xrange(npnts)]
```

In `diff2d_numpy`, this becomes

```
import numpy as np
x = np.linspace(x1-dx,x2+dx,num=npnts)
dens = np.zeros((npnts,npnts))
densnext = np.zeros((npnts,npnts))
```

Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #2: Updating the density field

`diff2d.py` had:

```
for i in range(1,nrows+1):
    for j in range(1,ncols+1):
        densnext[i][j] = dens[i][j] + (D/dx**2)*dt*laplacian[i][j]
```

Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #2: Updating the density field

`diff2d.py` had:

```
for i in range(1,nrows+1):
    for j in range(1,ncols+1):
        densnext[i][j] = dens[i][j] + (D/dx**2)*dt*laplacian[i][j]
```

In `diff2d_numpy`, this becomes

```
densnext[:,:] = dens + (D/dx**2)*dt*laplacian
```

- Element-wise algebra!
- Prevented a copy using `densnext[:,:]`

Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #3: Computing the Laplacian

`diff2d.py` had:

```
for i in range(1,nrows+1):
    for j in range(1,ncols+1):
        laplacian[i][j] = (dens[i+1][j]+dens[i-1][j]+dens[i][j+1]
                           +dens[i][j-1]-4*dens[i][j])
```


Switching to numpy

If you look at the `diff2d.py` and `diff2d_numpy.py` codes, you will see a number of differences.

Change #3: Computing the Laplacian

`diff2d.py` had:

```
for i in range(1,nrows+1):
    for j in range(1,ncols+1):
        laplacian[i][j] = (dens[i+1][j]+dens[i-1][j]+dens[i][j+1]
                           +dens[i][j-1]-4*dens[i][j])
```

In `diff2d_numpy`, this becomes

```
laplacian[1:nrows+1,1:ncols+1] = (dens[2:nrows+2,1:ncols+1] +
    dens[0:nrows+0,1:ncols+1] + dens[1:nrows+1,2:ncols+2] +
    dens[1:nrows+1,0:ncols+0] - 4*dens[1:nrows+1,1:ncols+1])
```

Data parallelism!

Does changing to numpy really help?

Let's return to our 2D diffusion example.

Pure python implementation:

```
$ etime python diff2d.py > output_p.txt  
Elapsed: 339.60 seconds
```

Does changing to numpy really help?

Let's return to our 2D diffusion example.

Pure python implementation:

```
$ etime python diff2d.py > output_p.txt  
Elapsed: 339.60 seconds
```

Numpy implementation:

```
$ etime python diff2d_numpy.py > output_n.txt  
Elapsed: 6.38 seconds
```

Yeah! About 30× faster.

Does changing to numpy really help?

Let's return to our 2D diffusion example.

Pure python implementation:

```
$ etime python diff2d.py > output_p.txt  
Elapsed: 339.60 seconds
```

Numpy implementation:

```
$ etime python diff2d_numpy.py > output_n.txt  
Elapsed: 6.38 seconds
```

Yeah! About 30× faster.

However, this is what the compiled versions do:

```
$ etime ./diff2d_cpp.ex > output_c.txt  
Elapsed: 1.09 seconds  
$ etime ./diff2d_f90.ex > output_f.txt  
Elapsed: 0.68 seconds
```

What about cython?

- Cython is a compiler for python code
- Almost all python is valid cython
- Typically used for packages, to be used in regular python scripts.
- It is important to realize that the compilation preserves the pythonic nature of the language, i.e, garbage collection, range checking, reference counting, etc, are still done: no performance enhancement.

```
$ etime python diff2d_numpy.py > output_n.txt  
Elapsed: 6.37 seconds  
$ etime python diff2d_numpy_cython.py > output_nc.txt  
Elapsed: 6.44 seconds
```

- If you want to get around that, you need to use Cython specific extensions that essentially use c types.
- From that point on, though, it isn't really python anymore, just a convenient way to write compiled python extensions.

Parallel Python

Parallel Python

We will look at a number of approaches to parallel programming with Python:

Package	Functionality
numexpr	threaded parallelization of certain numpy expressions
fork	create copies of existing process
threads	create threads sharing memory
multiprocessing	create processes that behave more like threads
mpi4py	message passing between processes
ipyparallel	framework of controlling parallel workers
spark	framework of controlling parallel workers

Numexpr

The numexpr package

The numexpr package is useful if you're doing matrix algebra:

- It is essentially a just-in-time compiler for NumPy.
- It takes matrix expressions, breaks things up into threads, and does the calculation in parallel.
- Somewhat awkwardly, it takes it's input in as a string.
- In some situations using numexpr can significantly speed up your calculations.

Numexpr in a nutshell

- Give it an array arithmetic expression, and it will compile and run it, and return or store the output.
- Supported operators:
+, -, *, /, **, %, <<, >>, <, <=, ==, !=, >=, >, &, |, ~
- Supported functions:
where, sin, cos, tan, arcsin, arccos, arctan, arctan2, sinh, cosh, tanh, arcsinh, arccosh, arctanh, log, log10, log1p, exp, expm1, sqrt, abs, conj, real, imag, complex, contains.
- Supported reductions:
sum, product

Using the numexpr package

Without numexpr:

```
>>> from etime import etime
>>> import numpy as np
>>> import numexpr as ne
>>> a = np.random.rand(1000000)
>>> b = np.random.rand(1000000)
>>> c = np.zeros(1000000)
>>> etime("c = a**2 + b**2 + 2*a*b", "a,b,c")
Elapsed: 0.0155924081802 seconds
```

Note: The python function etime measures the elapsed time. It is defined in the file etime.py that is part of the code of this session. The second argument should list the variables used (though some will be picked up automatically).

lpython has its own version of this, invoked (without quotes) as

```
In [10]: %time c = a**2 + b**2 +2*a*b
```



Using the numexpr package

With numexpr:

```
>>> from etime import etime
>>> import numpy as np
>>> import numexpr as ne
>>> a = np.random.rand(1000000)
>>> b = np.random.rand(1000000)
>>> c = np.zeros(1000000)
>>> etime("c = a**2 + b**2 + 2*a*b")
Elapsed: 0.021160197258 seconds
>>> old = ne.set_num_threads(1)
>>> etime("ne.evaluate('a**2 + b**2 + 2*a*b',out=c)", "a,b,c")
Elapsed: 0.0102033495903 seconds
>>> old = ne.set_num_threads(2)
>>> etime("ne.evaluate('a**2 + b**2 + 2*a*b',out=c)", "a,b,c")
Elapsed: 0.00552240610123 seconds
```

Numexpr for the diffusion example

- Annoyingly, numexpr has no facilities for slicing or offsets, etc.
- This is troubling for our diffusion code, in which we have to do something like

```
laplacian[1:nrows+1,1:ncols+1] = (dens[2:nrows+2,1:ncols+1]
                                   + dens[0:nrows+0,1:ncols+1]
                                   + dens[1:nrows+1,2:ncols+2]
                                   + dens[1:nrows+1,0:ncols+0]
                                   - 4*dens[1:nrows+1,1:ncols+1])
```

- We would need to make a copy of dens[2:nrows+2,1:ncols+1] etc. into a new numpy array before we can use numexpr, but copies are expensive.
- We want numexpr to use the same data as in dens, but *viewed* differently.

Numexpr for the diffusion example (cont.)

- We want numexpr to use the same data as in dens, but *viewed* differently.
- That is tricky, and requires knowledge of the data's memory structure.
- diff2d_numexpr shows one possible solution.

```
$ etime python diff2d_numpy.py > diff2d_numpy.out
```

```
Elapsed: 6.33 seconds
```

```
$ etime python diff2d_numexpr.py > diff2d_numexpr.out
```

```
Elapsed: 2.02 seconds
```

Theano (an alternative to numexpr)

- Theano is a numerical computation library.
- Much like numexpr, it takes an (array) expression and compiles it.
- Theano is frequently use in machine learning applications.
- Unlike numexpr, it can use multi-dimensional arrays and slices, like NumPy.
- Unlike numexpr, it does not natively use threads (though it may link to multithreaded blas libraries).
- But it can use GPUs (haven't tried).

Theano in the diffusion equation

For the diffusion code, only the computation of the laplacian changes

```
t_dens = theano.tensor.dmatrix('dens')
t_laplacian = (t_dens[2:nrows+2,1:ncols+1] +
               t_dens[0:nrows+0,1:ncols+1] +
               t_dens[1:nrows+1,2:ncols+2] +
               t_dens[1:nrows+1,0:ncols+0] -
               4*t_dens[1:nrows+1,1:ncols+1])
t_laplacian_fun = theano.function([t_dens], t_laplacian)
laplacian[1:nrows+1,1:ncols+1] = t_laplacian_fun(dens)
```


Theano in the diffusion equation

For the diffusion code, only the computation of the laplacian changes

```
t_dens = theano.tensor.dmatrix('dens')
t_laplacian = (t_dens[2:nrows+2,1:ncols+1] +
              t_dens[0:nrows+0,1:ncols+1] +
              t_dens[1:nrows+1,2:ncols+2] +
              t_dens[1:nrows+1,0:ncols+0] -
              4*t_dens[1:nrows+1,1:ncols+1])
t_laplacian_fun = theano.function([t_dens], t_laplacian)
laplacian[1:nrows+1,1:ncols+1] = t_laplacian_fun(dens)
```

Worth it, using 8 cores?

Theano in the diffusion equation

For the diffusion code, only the computation of the laplacian changes

```
t_dens = theano.tensor.dmatrix('dens')
t_laplacian = (t_dens[2:nrows+2,1:ncols+1] +
              t_dens[0:nrows+0,1:ncols+1] +
              t_dens[1:nrows+1,2:ncols+2] +
              t_dens[1:nrows+1,0:ncols+0] -
              4*t_dens[1:nrows+1,1:ncols+1])
t_laplacian_fun = theano.function([t_dens], t_laplacian)
laplacian[1:nrows+1,1:ncols+1] = t_laplacian_fun(dens)
```

Worth it, using 8 cores?

```
$ etime python diff2d_numpy.py
Elapsed: 6.25 seconds
$ etime python diff2d_numexpr.py
Elapsed: 1.99 seconds
$ etime python diff2d_theano.py
Elapsed: 13.26 seconds
```

Numexpr wins.

Theano in the diffusion equation

For the diffusion code, only the computation of the laplacian changes

```
t_dens = theano.tensor.dmatrix('dens')
t_laplacian = (t_dens[2:nrows+2,1:ncols+1] +
              t_dens[0:nrows+0,1:ncols+1] +
              t_dens[1:nrows+1,2:ncols+2] +
              t_dens[1:nrows+1,0:ncols+0] -
              4*t_dens[1:nrows+1,1:ncols+1])
t_laplacian_fun = theano.function([t_dens], t_laplacian)
laplacian[1:nrows+1,1:ncols+1] = t_laplacian_fun(dens)
```

Worth it, using 8 cores?

How about serially?

```
$ etime python diff2d_numpy.py
Elapsed: 6.25 seconds
$ etime python diff2d_numexpr.py
Elapsed: 1.99 seconds
$ etime python diff2d_theano.py
Elapsed: 13.26 seconds
```

Numexpr wins.

Theano in the diffusion equation

For the diffusion code, only the computation of the laplacian changes

```
t_dens = theano.tensor.dmatrix('dens')
t_laplacian = (t_dens[2:nrows+2,1:ncols+1] +
              t_dens[0:nrows+0,1:ncols+1] +
              t_dens[1:nrows+1,2:ncols+2] +
              t_dens[1:nrows+1,0:ncols+0] -
              4*t_dens[1:nrows+1,1:ncols+1])
t_laplacian_fun = theano.function([t_dens], t_laplacian)
laplacian[1:nrows+1,1:ncols+1] = t_laplacian_fun(dens)
```

Worth it, using 8 cores?

```
$ etime python diff2d_numpy.py
Elapsed: 6.25 seconds
$ etime python diff2d_numexpr.py
Elapsed: 1.99 seconds
$ etime python diff2d_theano.py
Elapsed: 13.26 seconds
```

Numexpr wins.

How about serially?

```
$ #with "ne.set_num_threads(1)"
$ etime python diff2d_numexpr.py
Elapsed: 6.97 seconds
$ etime python diff2d_theano.py
Elapsed: 13.26 seconds
```

numexpr still wins.

canada | canada

Forking

Forking (linux specific)

Another simple way to run code in parallel is to “fork” the process.

- The system call `fork()` creates a copy of the process that called it, and runs it as a child process.
- The child gets ALL the data of the parent process.
- The child gets its own process number (PID), and as such runs independently of the parent.
- We use the return value of `fork()` to determine which process we are; 0 means we're the child.

```
# firstfork.py
import os

# Our child process.
def child():
    print "Hello from", os.getpid()
    os._exit(0)

# The parent process.
while (True):
    newpid = os.fork()
    if newpid == 0:
        child()
    else:
        print "Hello from parent", os.getpid()

if raw_input() == "q": break
```

Process forking, continued

What does that look like?

```
$ python firstfork.py
Hello from parent 27089 27090
Hello from 27090
q
$
```

Forking/executing

What if we prefer to run a completely different code, rather than copying the existing code to the child?

- we can run one of the `os.exec` series of functions.
- The `os.execlp` call replaces the currently running program with the new one specified, in the child process only.
- If `os.execlp` is successful at launching the program, it never returns. Hence the `assert` statement is only invoked if something goes wrong.

```
# child.py
import os
print "Hello from", os.getpid()
os._exit(0)
```

```
# secondfork.py
import os

while (True):
    pid = os.fork()
    if pid == 0:
        os.execlp("python", "python",
                 "child.py")
        assert False, "Error starting"
    else:
        print "The child is", pid
        if raw_input() == "q": break
```

canada | canada

Notes about fork()

Fork was an early implementation used to spawn sub-processes, and is no longer commonly used. Some things to remember if you try to use this approach:

- use `os.waitpid(child_pid)` if you need to wait for the child process to finish. Otherwise the parent will exit and the child will live on.
- `fork()` is a Unix command. It doesn't work on Windows, except under Cygwin.
- This must be used very carefully, ALL the data is copied to the child process, including file handles, open sockets, database connections. . .
- Be sure to exit using `os._exit(0)` rather than `os.exit(0)`, or else the child process will try to clean up resources that the parent process is still using.
- Because of the above, `fork()` can lead to code that is difficult to maintain long-term.

Using fork in data analysis

Some notes about using forks in the context of data analysis:

- Something you may have noticed the about fork examples thus far is the lack of return from the functions.
- Forked processes, being processes and not threads, do not share anything with the parent process.
- As such, the only way they can return anything to the parent function is through inter-process communication.
- This is possible, though a bit tricky. We'll look at one way to do this later in the class.
- Your best bet, from a data processing point of view, is to just use fork for one-time functions that do not return anything to the parent.

Threads in Python

Processes versus threads

There is often confusion on the difference between threads and processes.

- A process provides the resources needed to execute a program. A thread is a path of execution within a process. As such, a process contains at least one thread, possibly many.
- A process contains a considerable amount of state information (handles to system objects, PID, address space, . . .). As such they are more resource-intensive to create. Threads are very light weight in comparison.
- Threads within the same process share the same address space. This means they can share the same memory and can easily communicate with each other.
- Different processes do not share the same address space. Different processes can only communicate with each other through OS-supplied mechanisms.

Notes about threads

Are there advantages to using threads, versus processes?

- As noted about, threads are light-weight compared to processes. As a result, they start up more quickly.
- Threads can be simpler to program, especially when the threads need to communicate with each other.
- Threads share memory, which can simplify (as well as obfuscate) programming.
- Threads are more portable than forked processes, as they are fully supported by Windows.

These points aside, there are downsides to using threads in a data-analysis application, as we'll see in a moment.

How much faster is it using threads?

```
# summer_threaded.py
import time, threading
from summer import my_summer

begin = time.time()
threads = []

for i in range(10):
    t = threading.Thread(
        target = my_summer,
        args = (0, 5000000))
    threads.append(t)
    t.start()

# Wait for all threads to finish.
for t in threads: t.join()
print ("Elapsed: %f"%
time.time() - begin,"seconds")
```

```
# summer.py - used in all summer*py
def my_summer(start, stop):
    tot = 0
    for i in xrange(start,stop):
        tot += i
```

```
# summer_serial.py
import time
from summer import my_summer
begin = time.time()
threads = []
for i in range(10):
    my_summer(0, 5000000)
print "Elapsed:", time.time() - begin
```

How much faster is it using threads?

```
# summer_threaded.py
import time, threading
from summer import my_summer

begin = time.time()
threads = []

for i in range(10):
    t = threading.Thread(
        target = my_summer,
        args = (0, 5000000))
    threads.append(t)
    t.start()

# Wait for all threads to finish.
for t in threads: t.join()
print ("Elapsed: %f"%
time.time() - begin,"seconds")
```

```
# summer.py - used in all summer*.py
def my_summer(start, stop):
    tot = 0
    for i in xrange(start,stop):
        tot += i
```

```
# summer_serial.py
import time
from summer import my_summer
begin = time.time()
threads = []
for i in range(10):
    my_summer(0, 5000000)
print "Elapsed:", time.time() - begin
```

Timings

```
$ python summer_serial.py
```

```
Elapsed: 11.58 seconds
```

```
$ python summer_threaded.py
```

```
Elapsed: 38.48 seconds
```

Not faster at all, slower!

The threading code is no faster than the serial code, even on my computer with two cores. Why?

- The Python Interpreter uses the Global Interpreter Lock (GIL).
- To prevent race conditions, the GIL prevents threads from the same Python program from running simultaneously. As such, only one core is used at any given time.
- Consequently the threaded code is no faster than the serial code, and is generally slower due to thread-creation overhead.
- As a general rule, threads are not used for most Python applications (GUIs being one important exception). This example is for demonstration purposes only.
- Instead, we will use one of several other modules, depending on the application in question. These modules will launch subprocesses rather than threads.

Multiprocessing

Multiprocessing

The multiprocessing module tries to strike a balance between forks and threads:

- Unlike fork, multiprocessing works on Windows (better portability).
- Slightly longer start-up time than threads.
- Multiprocessing spawns separate processes, like fork, and as such they each have their own memory.
- Multiprocessing requires pickleability for its processes on Windows, due to the way in which it is implemented. As such, passing non-pickleable objects, such as sockets, to spawned processes is not possible.

The multiprocessing module, continued

A few notes about the multiprocessing module:

- The Process function launches a separate process.
- The syntax is very similar to the threading module. This is intentional.
- The details under the hood depend strongly upon the system involved (Windows, Mac, Linux), thus the portability of code written with this module.

```
# summer_multiprocessing.py
import time, multiprocessing
from summer import my_summer
begin = time.time()
processes = []
for i in range(10):
    p = multiprocessing.Process(
        target = my_summer,
        args = (0, 5000000))
    processes.append(p)
    p.start()
# Wait for all processes to finish
for p in processes: p.join()
print ("Elapsed:%f"%
        time.time() - begin)
```

```
$ python summer_multiprocessing.py
Time:0.221297
```

Shared memory with multiprocessing

- multiprocessing allows one to seamlessly share memory between processes. This is done using 'Value' and 'Array'.
- Value is a wrapper around a strongly typed object called a ctype. When creating a Value, the first argument is the variable type, the second is that value.
- Code on the right has 10 processes add 50 increments of 1 to the Value v.

```
# multiprocessing_shared.py
from multiprocessing import Process
from multiprocessing import Value
def myfun(v):
    for i in range(50):
        time.sleep(0.001)
        v.value += 1

v = Value('i', 0);
procs = []
for i in range(10):
    p=Process(target=myfun,args=(v,))
    procs.append(p)
    p.start()
for proc in procs: proc.join()
print(v.value)
```

Shared memory with multiprocessing

- multiprocessing allows one to seamlessly share memory between processes. This is done using 'Value' and 'Array'.
- Value is a wrapper around a strongly typed object called a ctype. When creating a Value, the first argument is the variable type, the second is that value.
- Code on the right has 10 processes add 50 increments of 1 to the Value v

```
# multiprocessing_shared.py
from multiprocessing import Process
from multiprocessing import Value
def myfun(v):
    for i in range(50):
        time.sleep(0.001)
        v.value += 1

v = Value('i', 0);
procs = []
for i in range(10):
    p=Process(target=myfun,args=(v,))
    procs.append(p)
    p.start()
for proc in procs: proc.join()
print(v.value)
```

```
$ etime python multiprocessing_shared.py
499
```

Elapsed: 0.51 seconds

- Did the code behave as expect?

Race conditions

What went wrong?

- Race conditions occur when program instructions are executed in an order not intended by the programmer. The most common cause is when multiple processes are given access to a resource.
- In the example here, we've modified a location in memory that is being accessed by multiple processes.
- Note that it need not only be processes or threads that can modify a resource, anything can modify a resource, hardware or software.
- Bugs caused by race conditions are extremely hard to find.
- Disasters can occur.

Be very very careful when sharing resources between multiple processes or threads!

Using shared memory, continued

- The solution, of course, is to be more explicit in your locking.
- If you use shared memory, be sure to test everything thoroughly.

```
# multiprocessing_shared_fixed.py # multiprocessing_shared_fixed.py
from multiprocessing import Process # continued
from multiprocessing import Value v = Value('i', 0)
from multiprocessing import Lock lock = Lock()
proc = []
def myfun(v, lock):
    for i in range(50):
        time.sleep(0.001)
        with lock:
            v.value += 1
    for i in range(10):
        p=Process(target=myfun,
                  args=(v,lock))
        procs.append(p)
        p.start()
    for proc in procs: proc.join()
    print(v.value)
```

```
$ etime python multiprocessing_shared_fixed.py
500
Elapsed: 0.36 seconds
```

Using shared memory, arrays

- Multiprocessing also allows you to share a block of memory through the Array ctypes wrapper.
- Only 1-D arrays are permitted.
- Note that multiprocessing.Process must be used; shared memory does not work with multiprocessing.Pool.map.
- Note that, since arr is actually a ctypes object, you must print the contents of arr to see the result.

```
# multiprocessing_shared_array.py
from numpy import arange
from multiprocessing import Process
def myfun(a, i):
    a[i] = -a[i]
arr = Array('d', arange(10.))
procs = []
for i in range(10):
    p = Process(target=myfun,
                args=(arr, i))
    procs.append(p)
    p.start()
for proc in procs:
    proc.join()
print(arr[:])
```


Using shared memory, arrays

- Multiprocessing also allows you to share a block of memory through the Array ctypes wrapper.
- Only 1-D arrays are permitted.
- Note that multiprocessing.Process must be used; shared memory does not work with multiprocessing.Pool.map.
- Note that, since arr is actually a ctypes object, you must print the contents of arr to see the result.

```
# multiprocessing_shared_array.py
from numpy import arange
from multiprocessing import Process
def myfun(a, i):
    a[i] = -a[i]
arr = Array('d', arange(10.))
procs = []
for i in range(10):
    p = Process(target=myfun,
               args=(arr, i))
    procs.append(p)
    p.start()
for proc in procs:
    proc.join()
print(arr[:])
```

```
[-0.0, -1.0, -2.0, -3.0, -4.0, -5.0, -6.0, -7.0, -8.0, -9.0]
```

But there's more!

The multiprocessing module is loaded with functionality. Other features include:

- Inter-process communication, using Pipes and Queues.
- `multiprocessing.manager`, which allows jobs to be spread over multiple 'machines' (nodes).
- subclassing of the Process object, to allow further customization of the child process.
- `multiprocessing.Event`, which allows event-driven programming options.
- `multiprocess.condition`, which is used to synchronize processes.

We're not going to cover these features today.

MPI4PY

Message Passing Interface

The previous parallel techniques used processors on one node.

Using more than one node requires these nodes to communicate.

MPI is one way of doing that communication.

- MPI = Message Passing Interface.
- MPI is a C/Fortran Library API.
- Sending data = sending a message.
- Requires setup of processes through `mpirun/mpiexec`.
- Requires `MPI_Init(...)` in code to collect processes into a 'communicator'.
- Rather low level.

Mpi4py features

- mpi4py is a wrapper around the mpi library
- Point-to-point communication (sends, receives)
- Collective (broadcasts, scatters, gathers) communications of any picklable Python object,
- Optimized communications of Python object exposing the single-segment buffer interface (NumPy arrays, builtin bytes/string/array objects).
- Names of functions much the same as in C/Fortran, but are methods of the communicator (object-oriented).

MPI C/C++ recap

The following C++ code determines each process' rank and sends that rank to its left neighbor.

```
#include <mpi.h>
#include <iostream>

int main(int argc, char** argv) {
    int rank, size, rankr, right, left;
    MPI_Init(&argc, &argv);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    right = (rank+1)%size;
    left = (rank+size-1)%size;
    MPI_Sendrecv(&rank, 1, MPI_INT, left, 13,
                &rankr, 1, MPI_INT, right, 13,
                MPI_COMM_WORLD, MPI_STATUS_IGNORE);
    std::cout<<"I am rank "<<rank<<" ; my right neighbour is "<<rankr<<
    MPI_Finalize();
}
```

MPI Fortran recap

The following Fortran code determines each process' rank and sends that rank to its left neighbor.

```
program rightrank
  use mpi
  implicit none
  integer rank, size, rankr, right, left, e
  call MPI_Init(e)
  call MPI_Comm_rank(MPI_COMM_WORLD, rank, e)
  call MPI_Comm_size(MPI_COMM_WORLD, size, e)
  right = mod(rank+1, size)
  left  = mod(rank+size-1, size)
  call MPI_Sendrecv(rank, 1, MPI_INTEGER, left, 13, &
                    rankr, 1, MPI_INTEGER, right, 13, &
                    MPI_COMM_WORLD, MPI_STATUS_IGNORE, e)
  print *, "I am rank ", rank, "; my right neighbour is ", rankr
  call MPI_Finalize(e)
end program rightrank
```

Mpi4py

- One of the drudgeries of MPI is to have to express the binary layout of your data.
- The drudgery arises because C and Fortran do not have *introspection* and the MPI libraries cannot look inside your code.
- With Python, this is potentially different: we can investigate, within python, what the structure is.
- That means we should be able to express sending a piece of data without having to specify types and amounts.

```
from mpi4py import MPI
rank  = MPI.COMM_WORLD.Get_rank()
size  = MPI.COMM_WORLD.Get_size()
right = (rank+1)%size
left  = (rank+size-1)%size
rankr = MPI.COMM_WORLD.sendrecv(rank, left, source=right)
print "I am rank", rank, "; my right neighbour is", rankr
```


Mpi4py + numpy

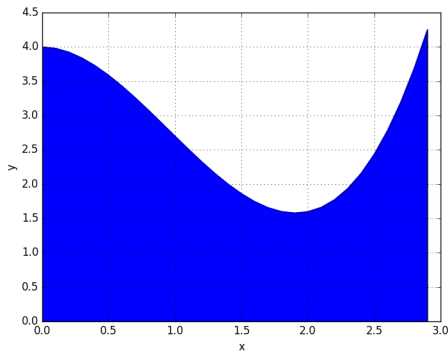
- It turns out that mpi4py's communication is pickle-based.
- Pickle is a *serialization* format which can convert any python object into a bytestream.
- Convenient as any python object can be sent, but conversion takes time.
- For numpy arrays, one can skip the pickling using Uppercase variants of the same communicator methods.
- However, this requires us to preallocate buffers to hold messages to be received.

Example: Area under the curve

- Let's consider a code that numerically computes the following integral:

$$b = \int_{x=0}^3 \left(\frac{7}{10}x^3 - 2x^2 + 4 \right) dx$$

- Exact answer $b = 8.175$
- It's the area under the curve on the right.



Method: sample $y = \frac{7}{10}x^3 - 2x^2 + 4$ at a uniform grid of x values (using `ntot` number of points), and add the y values.

Mpi4py+numpy: Upper/lowercase example

```
import sys
from mpi4py import MPI

rank = MPI.COMM_WORLD.Get_rank()
size = MPI.COMM_WORLD.Get_size()
ntot = int(sys.argv[1])
npnts = ntot/size
dx = 3.0/ntot
width = 3.0/size
x = rank*width
a = 0.0
for i in xrange(npnts):
    y = 0.7*x**3 - 2*x**2 + 4
    a += y*dx
    x += dx

b = MPI.COMM_WORLD.reduce(a)
if rank == 0:
    print "The area is", b
```

Mpi4py+numpy: Upper/lowercase example

```
import sys
from mpi4py import MPI

rank = MPI.COMM_WORLD.Get_rank()
size = MPI.COMM_WORLD.Get_size()
ntot = int(sys.argv[1])
npnts = ntot/size
dx = 3.0/ntot
width = 3.0/size
x = rank*width
a = 0.0
for i in xrange(npnts):
    y = 0.7*x**3 - 2*x**2 + 4
    a += y*dx
    x += dx

b = MPI.COMM_WORLD.reduce(a)
if rank == 0:
    print "The area is", b
```

```
import sys
from mpi4py import MPI
from numpy import zeros, asarray
rank = MPI.COMM_WORLD.Get_rank()
size = MPI.COMM_WORLD.Get_size()
ntot = int(sys.argv[1])
npnts = ntot/size
dx = 3.0/ntot
width = 3.0/size
x = rank*width
a = 0.0
for i in xrange(npnts):
    y = 0.7*x**3 - 2*x**2 + 4
    a += y*dx
    x += dx

b = np.zeros(1)
MPI.COMM_WORLD.Reduce(asarray(a), b)
if rank == 0:
    print "The area is", b[0]
```

Mpi4py Speedup?

```
$ etime mpirun -np 1 python auc.py 300000000  
The area is 8.175000  
Elapsed: 301.33 seconds
```

Mpi4py Speedup?

```
$ etime mpirun -np 1 python auc.py 300000000  
The area is 8.175000  
Elapsed: 301.33 seconds
```

```
$ etime mpirun -np 4 python auc.py 300000000  
The area is 8.175000  
Elapsed: 76.58 seconds
```

Mpi4py Speedup?

```
$ etime mpirun -np 1 python auc.py 300000000  
The area is 8.175000  
Elapsed: 301.33 seconds
```

```
$ etime mpirun -np 4 python auc.py 300000000  
The area is 8.175000  
Elapsed: 76.58 seconds
```

```
$ etime mpirun -np 4 python auc_numpy.py 300000000  
The area is 8.175000  
Elapsed: 77.19 seconds
```

Mpi4py Speedup?

```
$ etime mpirun -np 1 python auc.py 300000000  
The area is 8.175000  
Elapsed: 301.33 seconds
```

```
$ etime mpirun -np 4 python auc.py 300000000  
The area is 8.175000  
Elapsed: 76.58 seconds
```

```
$ etime mpirun -np 4 python auc_numpy.py 300000000  
The area is 8.175000  
Elapsed: 77.19 seconds
```

There simply isn't enough communication to see the difference between the pickled and non-pickled interface.

Hands-on

- 1 Use multiprocessing to parallelize the auc.py code.
- 2 Use numexpr to parallelize the auc.py code.
- 3 What else could we do to speed up the code?

Map/Reduce variations

Map/reduce

- The diffusion example is, as already admitted, a hard problem to get good performance out of with python.
- That was because it's a tightly coupled problem.
- Other problems aren't, e.g.:
 - ▶ Parameter sweeps
 - ▶ Reductions
 - ▶ Big data
- For such problems, there are some valuable frameworks of the **map/reduce** variety.
- We'll briefly consider two python-enabled map/reduce frameworks:
 - ▶ IPython Parallel
 - ▶ Apache Spark (in particular pyspark)

Common characteristics in map/reduce

- A master process + worker processes
- Master divides or requests work, and collects results
- Overall workflow is data based:
 - 1 Data is distributed over workers (or already resides there), and workers perform computation on their local data
 - 2 If reduction: data is moved between workers, and work is done by 'reducers'. This step is iterative.
 - 3 Result is reported to the master.
- Emphasis on distributing the work and bringing the work to the data. Works well if 'work chunks' take a good bit of time.

IPython's Parallel Architecture

Ipyparallel

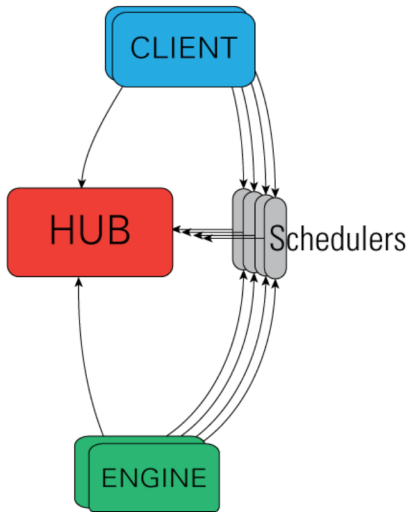
(formerly IPython Parallel)

Interestingly, IPython comes with a built-in parallel engine. It consists of four components:

- Engines: Do the work. One core, one engine.
- Schedulers: Deliver and divide the work.
- Hub: Coordinates and logs the engine and schedule activity.
- Clients: Request work to be done on engines.

Schedulers + Hub + Engine = Cluster

Schedulers + Hub = Controller



Starting an IPython cluster

Starting up an IPython cluster is not difficult. But there are some steps that you need to go through carefully:

```
$ hostname  
gpc-f123n063-ib0  
$ cd $SCRATCH/hpcpy/code  
$ source setup  
$ ipcluster start -n 4  
a bunch of messages...
```

Do this on the compute node, and remember the hostname. You can minimize this terminal but do not close it.

- Open a second terminal on your laptop
- Ssh into your compute node.

```
$ ssh -Y USERNAME@login.scinet.utoronto.ca  
$ ssh -Y gpc-f123n063-ib0  
$ cd $SCRATCH/hpcpy/code
```

Using the workers in python

Start python and grab the handles to the clients.

```
$ cd $SCRATCH/hpcpy/code  
$ source setup  
$ python
```

```
>>> from ipyparallel import Client  
>>> clients = Client()  
>>>
```

(older python version? try from IPython.parallel import Client)

Accessing the Clients

Let's see what we've got here.

```
>>>
>>> # use synchronous computations
>>> clients.block = True
>>> print(len(clients))
4
>>> print(clients.ids)
[0, 1, 2, 3]
>>>
```

Each client has been assigned an id, starting at 0.

Accessing the Clients, continued

Here's simple function to execute on the cores:

```
>>> def minus(a, b):  
...     return a - b  
>>> minus(5, 6)  
-1  
>>> #Execution on the first engine only:  
>>> clients[0].apply(minus, 5, 6)  
-1
```

Interfaces to the Engines

Some notes about the engines:

- The python environment that holds the 'Clients' part is completely separate from that of the 'Engines'.
- As such, you need to move data and code to the Engines.
- You also need to request to execute code on the Engines.
- The Controller (Schedulers + Hub) is the single point of contact for the clients.

Views

Views are a layer over sets of engines that allow access to engine variables through a dictionary, storing settings, and scheduling tasks.

There are two kinds of views:

- A Direct interface, where engines are addressed explicitly. You get this view by using square brackets. For example:

```
Client()[1:8:2]
```

- A LoadBalanced interface, where the Scheduler is trusted with assigning work to appropriate engines. You get this from

```
Client().load_balanced_view()
```

- View is selected by the client.

Parallel Execution

There are a number of ways to invoke the Engines:

- `clients[:].run` takes a script and runs in on the engine(s).
- `clients[:].execute` takes a command, as a string, to run on the engine(s).
- `clients[:].apply` takes a function and arguments, to run on the engine(s).
- `clients[:].map` takes a function and a list, to distribute over the engine(s).

In the last two, the function and arguments get shipped to the engine.

Blocking/Nonblocking

There are two modes in which execution of code can run:

- In blocking mode (“synchronous”), all execution must be finished before results are recorded.
- In non-blocking mode, an “AsyncResult” is returned, which we can ask if it is done (`.ready()`), and what the result is (`.get()`).

The latter is potentially faster, but requires a bit more ‘infrastructure’.

Examples

Execute minus in parallel on all the engines at once:

```
>>> clients[:].apply(minus, 5, 6)
[-1, -1, -1, -1]
```

What if we want different arguments to each engine? In normal Python we could use “map”:

```
>>> map(minus, [11, 10, 9, 8], [5, 6, 7, 8])
[6, 4, 2, 0]
```

Examples, continued

The client view's "map" function executes in parallel. Using a load-balanced view, this would look like this:

```
>>> view = clients.load_balanced_view()
>>> view.map(minus, [11, 10, 9, 8], [5, 6, 7, 8])
[6, 4, 2, 0]
```


Direct view

Recall that the “Direct View” allows you to directly command the engines. To execute a command on all the engines:

```
>>> clients.block = True
>>> dview = clients.direct_view()
>>> dview.block = True
>>> dview.apply(sum, [1, 2, 3])
[6, 6, 6, 6]
```

Direct view, continued

Slicing a Client's object gets you a Direct view as well:

```
>>> clients[::2]
<DirectView [0, 2]>
>>> clients[::2].apply(sum, [1, 2, 3])
[6, 6]
```

Which we saw previously, when we used `clients[:].apply(minus, 5, 6)`.

Load Balanced View

To execute a command on all the engines, using the Load Balanced View:

```
>>> dview = Client().load_balanced_view()
>>> dview.block = True
>>> dview.apply(sum, [1, 2, 3])
6
```

This view is useful if you're going to execute tasks one by one, or if the tasks take a varying amount of time.

We will focus on direct view in the remainder, which is a bit more flexible.

Data movement to and from engines

- Moving data around is straightforward.
- You can access variables through a dictionary-like interface.
- Indexing a client gives access to the dictionary for a particular engine.
- A view has a dictionary interface too, which gives you a list of the values in all the engines.

```
>>> v = clients[:]
>>> v.block = True
>>> v.execute('from os import getpid')
<AsyncResult: finished>
>>> v.execute('x = getpid()')
<AsyncResult: finished>
>>> v['x']
[24068, 24067, 24065, 24066]
>>> clients[3]['x']
24066
```

Scatter/Gather

- Some of the parallelization features we saw in the other parallel programming sessions are built-in as well.
- Sometimes you want to explicitly divide a list or array on the engines:

Scatter

- Or, reconstruct a larger list on the client from local lists on the engines: **Gather**

This is quite simple in IPython.parallel:

```
>>> v.scatter('a', np.arange(16))
>>> v['a']
[array([0, 1, 2, 3]),
 array([4, 5, 6, 7]),
 array([8, 9, 10, 11]),
 array([12, 13, 14, 15])]
>>> v.gather('a')
array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15])]
```

Spark

Why Spark?

Suppose you find yourself in a situation where your data is ridiculously huge:

- Facebook's daily logs: 60TB.
- 1,000 genomes project: 200TB.
- Google web index: > 10 PB.
- Cost of 1TB disk: \$65.
- Time to read 1TB from disk: 3 hours (100 MB/s).

The scale of these data sets means that we cannot analyse them on a single machine. The analysis must be done in a distributed-memory setting.

In the beginning

A few years ago, if you found yourself in this situation, you would have used Hadoop. Hadoop is essentially a distributed file system, that allows you to take your computation to the data.

But the original Hadoop had some serious issues:

- It was stuck with Map Reduce.
- This meant all data was mapped to a key-value pair and then 'reductions' were performed on this data.
- Each stage of any calculation involved:
 - 1 Read data from disk.
 - 2 Perform calculation.
 - 3 Write result to disk.
 - 4 Communicate the result.
- Very I/O heavy! And requires a disk available on all nodes.
- Slow! Inefficient!

Spark's approach

Apache Spark does not use I/O so heavily, and has ideal features:

- Spark keeps things in memory instead of disk (though on-disk storage is also supported).
- Map and Reduce are not your only available operations (though this is no longer true for Hadoop 2).
- Many language wrappers available (Scala, Java, R, Python).
- Uses lazy evaluation, which results in improved pipelining.
- Supports batch, interactive and streaming execution models.
- Spark has built-in redundancy; it keeps track of how the data are created, so that if a node fails the data can be rebuilt from scratch (like Hadoop).

Spark's anatomy

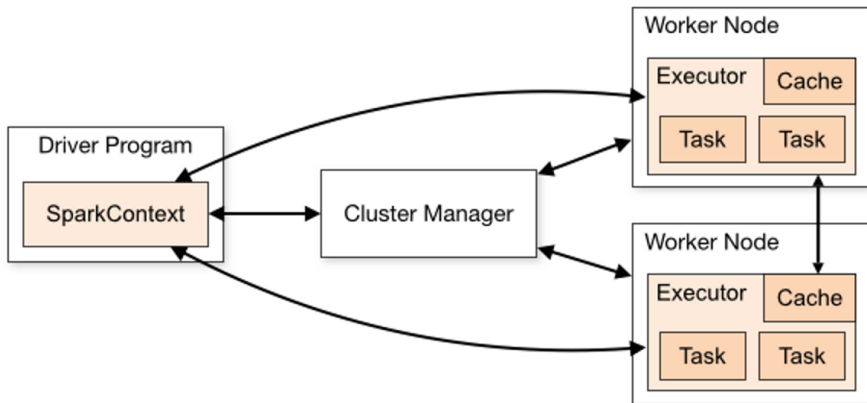
A Spark program consists of two parts:

- Driver program: runs on the driver machine. This is often called the 'master' program.
- Worker programs: run on cluster nodes or in local threads. These are sometimes called "Executors". When RDDs (Resilient Distributed Data sets) are created, they are distributed across the workers.

The first step in a Spark program is to create a SparkContext object.

- This tells Spark how and where to access a cluster.
- This is used to create our RDD.

Spark's anatomy, continued



“Driver program” == “Master”. Notice that only one worker program (executor) runs on each node, but that worker can launch multiple ‘tasks’.

About Spark and pySpark

We will use pySpark to access the running Spark machinery:

- Using Spark 1.6.1, which may still have some issues with Python 3.
- pySpark uses the standard CPython interpreter, so C libraries (like NumPy) can be used.
- Use the “pyspark” command to launch an interactive Python shell.
- The pySpark shell comes with a SparkContext built in, in the variable called “sc”.
- You can also invoke an IPython version of the pySpark shell, which comes with all the usual IPython functionality.

Setting up a Spark session

- 1 Make sure you're on an interactive node and have done source setup
- 2 Setup the Spark session.

```
$ source $SPARK_HOME/scripts/setup_spark.sh
starting org.apache.spark.deploy.master.Master, logging
to /scratch/s/scinet/rzon/spark-logs/spark-rzon-org.apache.spark.de
.
.
$
```

This script:

- Detects which nodes are part of your job.
- Starts a master process on the head node.
- Creates \$SCRATCH/temp, \$SCRATCH/spark-logs, \$SCRATCH/spark-workers directories, if necessary.
- Starts a worker process on all of your job's nodes. These are technically Java Virtual Machines.

Setting up a Spark session (cont.)

- 1 Move to the spark directory.

```
$ cd $SCRATCH/hpcpy/spark
```

- 2 Launch pySpark with IPython support.

```
$ IPYTHON=1 IPYTHON_OPTS="--pylab" pyspark
...
SparkContext available as sc.
In [1]:
```

We now have a Spark session running, which we access using pySpark.

Confirming your cluster

To confirm that your cluster is working:

```
>>>
>>> from socket import gethostname
>>>
>>> sc.parallelize(range(5)).map(lambda x:gethostname()).collect()
['gpc-f123n063-ib0', 'gpc-f123n063-ib0',
 'gpc-f123n063-ib0', 'gpc-f123n063-ib0',
 'gpc-f123n063-ib0']
>>>
```

Creating RDDs

Spark is all about RDDs (Resilient Distributed Data sets), which are what we use to hold our data:

- RDDs are objects; all operations on them consist of one of many existing functions.
- RDDs are immutable: they cannot be changed once created.
- To modify data you create new RDDs based on existing RDDs.
- Remember that RDDs are lazily evaluated, meaning not calculated until an 'action' is performed.
- RDDs can be 'cached' so that the results persist.
- RDDs contain data of two general types, regular (which are single values of any type), and (key, value).

RDD Pipelines

The general procedure for data analysis is a pipeline:

- create an RDD from a data source, file or list.
- apply various transformations to the RDD.
- apply an action to the RDD.

Creating RDDs, continued

So how do I create an RDD?

There are a few options in pySpark:

- `parallelize(x)`: create an RDD out of the data (list or Numpy vector) 'x'.
- `textFile(file)`: read a text file, and return each line as the data. Will also read all the files in a directory, if a path is given.
- There are also methods for reading in other types of files.

```
In [3]: myverbs = ['run', 'jump',  
                  'sit', 'laugh', 'run', 'smile']  
In [4]: verbRDD =  
        sc.parallelize(myverbs)  
In [5]:  
In [5]: myRDD =  
        sc.textFile('myfile.txt')  
In [6]:
```

Note that the file mentioned above doesn't (presumably) exist. Why isn't an error message thrown?

Manipulating RDDs: transformations

Ok, I've got my RDD. Now what? Well, the things you can do fall into two categories. The first is transformations:

- `map(func)`: map each value of the RDD to a function; return a new RDD that contains the return of the function on each RDD value.
- `filter(func)`: return a new dataset formed by applying 'func' to each element, and keeping those which return True.
- `flatMap(func)`: like map, but returns all elements in a single list.
- `distinct()`: reduce the RDD data points to distinct values only.
- `groupByKey()`: group key values as lists, return as (key, list) pairs.
- `reduceByKey(func)`: reduce the elements key-by-key, applying func to the elements.
- `sortBy(func)/sortByKey(func)`: sort the RDD, using func.

Transformations return an RDD, with the elements appropriately 'transformed.'

Manipulating RDDs: actions

The second category of RDD functions are actions:

- `reduce(func)`: reduce the elements back to the master process, by applying `func` to the elements.
- `count()`: counts the elements.
- `collect()`: bring all the elements back to the master.
- `take(n)`: bring in 'n' elements to the master.
- `takeOrdered(n, func)`: same as `take`, but re-order based on `func`.
`max()`, `min()`, `mean()`, `stdev()`, `sum()`

Actions return a value, or list. Read the API if you're not sure of the name of the function you need.

Manipulating RDDs, continued

```
In [6]: def pastTense(s):
...:     return s + 'ed'
...:
In [7]:
In [7]: myverbs = ['run', 'jump', 'sit', 'laugh', 'run', 'smile']
In [8]:
In [8]: verbRDD = sc.parallelize(myverbs)
In [9]:
In [9]: pastTenseRDD = verbRDD.map(pastTense)
In [10]:
In [10]: pastTenseRDD.collect()
Out[10]: ['runed', 'jumped', 'sited', 'laughed', 'runed', 'smileed']
In [11]:
In [11]: verbRDD.take(1)
Out[11]: ['run']
In [12]:
```

Note that nothing is actually calculated until the 'collect' is called.

Manipulating RDDs: partitioning

By default your RDD is sliced up into 'partitions' and spread across the workers.

- The number of partitions can greatly impact computational efficiency. When individual functions are applied to the data, one task is performed per partition.
- For load balancing you will want at least as many partitions as cores, possibly more.
- Spark documentation suggests 2-4 partitions per CPU. You will often see problems broken up into hundreds or thousands of partitions.
- Be aware that increasing the number of partitions before the data is read in can cause I/O issues. It's better to repartition after the data is read.

Manipulating RDDs: partitioning, cont.

There are built-in commands to adjust your partitions on the fly:

- You can get the system default minimum number of partitions using “sc.defaultMinPartitions”.
- You can set the number of partitions as an optional argument when you create your RDD:
 - ▶ “numSlices”: for parallelize.
 - ▶ “minPartitions”: for textFile.
- You can get the number of partitions for a given RDD using RDD.getNumPartitions().
- You can change the RDD’s number of partitions using “RDD.repartition(num)”. This can be an expensive operation, as the data is randomly shuffled.
- “RDD.coalesce(num)” also reduces the number of partitions without shuffling.

Manipulating RDDs: partitioning, cont.

```
In [12]:
In [12]: sc.defaultMinPartitions
Out[12]: 2
In [13]:
In [13]: myverbs = ['run', 'jump', 'sit', 'laugh', 'run', 'smile']
In [14]: verbRDD = sc.parallelize(myverbs, numSlices = 48)
In [15]: verbRDD.getNumPartitions()
Out[15]: 48
In [16]:
In [16]: myRDD=sc.textFile('data/stopwords.txt',minPartitions=100)
In [17]: myRDD.getNumPartitions()
Out[17]: 101
In [18]:
In [18]: newRDD = myRDD.repartition(200)
In [19]: newRDD.getNumPartitions()
Out[19]: 200
In [20]:
```


Python's lambda function

Python contains a function called 'lambda'. Lambda functions

- are small, anonymous functions (not bound to a name).
- have the format "lambda arguments: operation on arguments".
- are restricted to a single expression.
- example: `lambda a, b: a - b`.

Who cares? Well, if you combine lambda with the RDD 'map' command, you can do some pretty impressive things.

Using lambda

```
In [20]:  
In [20]: range(10)  
Out[20]: [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]  
In [21]:  
In [21]: myRDD = sc.parallelize(range(10))  
In [22]:  
In [22]: myRDD.map(lambda x: 2 * x)  
Out[22]: PythonRDD[1] at RDD at PythonRDD.scala:37  
In [23]:  
In [23]: myRDD.map(lambda x: 2 * x).collect()  
Out[23]: [0, 2, 4, 6, 8, 10, 12, 14, 16, 18]  
In [24]:  
In [24]: myRDD.filter(lambda x: x % 3 == 0).collect()  
Out[24]: [0, 3, 6, 9]  
In [25]:
```

Creating (key, value) pairs

The lambda function can be used to create compound variables, as well as to create (key, value) pairs.

```
In [25]: myRDD = sc.parallelize([3, 4, 5])
In [26]:
In [26]: myRDD.map(lambda x: [x, x - 2]).collect()
Out[26]: [[3, 1], [4, 2], [5, 3]]
In [27]:
In [27]: myRDD.flatMap(lambda x: [x, x - 2]).collect()
Out[27]: [3, 1, 4, 2, 5, 3]
In [28]:
In [28]: # create (key, value) pairs
In [28]: myRDD.map(lambda x: (x, 1)).collect()
Out[28]: [(3, 1), (4, 1), (5, 1)]
In [29]:
```

Note that once the data has been scattered to the workers, all operations occur locally, at the workers, until the 'collect' command.

(key, value) example

What is the frequency of words in “Moby Dick”? First prep the data:

```
In [29]: MD = sc.textFile("data/mobydick.txt")
In [30]:
In [30]: # the first 5 lines
In [30]: MD.take(5)
Out[30]: [u'MOBY DICK;', u'', u'or, THE WHALE', u'', u'']
In [31]:
In [31]: # We need to clean this data up.
In [31]: import string
In [32]:
In [32]: # convert from unicode to string.
In [32]: strMD = MD.map(lambda x: x.encode('ascii', 'ignore'))
In [33]:
In [33]: # Remove the punctuation.
In [33]: MDNoPu=strMD.map(
...:     lambda x:x.translate(None,string.punctuation))
In [34]:
In [34]: MDNoPu.take(5)
Out[34]: ['MOBY DICK', '', 'or, THE WHALE', '', '']
```

(key, value) example, data cleaning

```
In [35]: MDNoPu.take(5)
Out[35]: ['MOBY DICK', '', 'or THE WHALE', '', '']
In [36]:
In [36]: # We want the words, not the whole lines, so split.
In [36]: MDSplit = MDNoPu.flatMap(lambda x: x.split(' '))
In [37]:
In [37]: MDSplit.take(5)
Out[37]: ['MOBY', 'DICK', '', 'or', 'THE']
In [38]:
In [38]: # Make all the words lower case, and remove empty strings.
In [38]: MDlower = MDSplit.map(
...:     lambda x: x.lower()).filter(lambda x: x != '')
In [39]:
In [39]: MDlower.take(5)
Out[39]: ['moby', 'dick', 'or', 'the', 'whale']
In [40]:
```

(key, value) example, data cleaning, cont.

```
In [40]:
In [40]: MDlower.take(5)
Out[40]: ['moby', 'dick', 'or', 'the', 'whale']
In [41]:
In [41]: # stopwords are words which are considered 'uninteresting'
In [41]: stopwords = open('data/stopwords.txt',
...:                       'r').read().split('\n')
In [42]:
In [42]: stopwords[0:5]
Out[42]: ['a', 'about', 'above', 'across', 'after']
In [43]:
In [43]: MDReady = MDlower.filter(lambda x: x not in stopwords)
In [44]:
In [44]: MDReady.take(5)
Out[44]: ['moby', 'dick', 'whale', 'herman', 'melville']
In [45]:
```

Ready to start counting.

(key, value) example, counting occurrences

Using the (key, value) approach, 1) count the number of times each word appears and 2) give the top 5 words, in terms of frequency.

```
In [45]: kvMD = MDReady.map(lambda word: (word, 1))
In [46]:
In [46]: result = kvMD.reduceByKey(lambda x, y: x + y)
In [47]:
In [47]: result.take(5)
Out[47]:
[('knockers', 1),
 ('brevet', 1),
 ('yellow', 19),
 ('prefix', 1),
 ('looking', 56)]
In [48]:
In [48]: result.takeOrdered(5, lambda (k, v): -v)
Out[48]: [('whale', 892), ('like', 567), ('old', 436), ('man', 433),
 ('ahab', 417)]
```

Useful web sites for Spark

There are a number of useful web sites out there. Here are some of them.

<http://spark.apache.org/docs/latest/api/python/pyspark.html>

<http://spark.apache.org/docs/latest/programming-guide.html>

<http://spark.apache.org/docs/latest/api/python/pyspark.mllib.html>

<http://spark.apache.org/docs/latest/api/python/pyspark.sql.html>

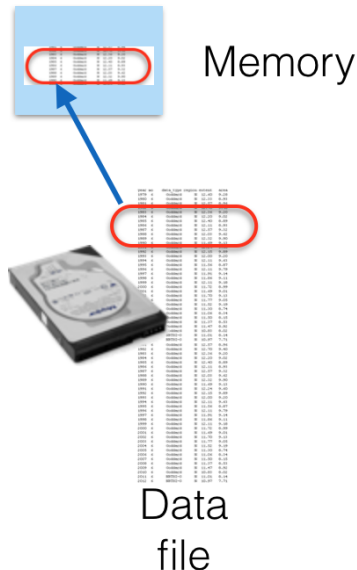
Out-of-core computations

Out-of-core computation

- Some problems require doing fairly simple analysis on data that is too large to fit into memory
 - ▶ Min/mean/max
 - ▶ Data cleaning
 - ▶ Even linear fitting is pretty simple
- In this case, one processor may be enough; you just want a way to not run out of memory.
- “Out of core” or “external memory” computation leaves the data on disk, bringing into memory only what is needed, or what fits, at any given time.
- For some computations, this works out well (but note: disk access is always much slower than memory access).

Out-of-core computation, continued

- The `numpy.memmap` class creates a memory-map to an array stored in a binary file on disk. This allows a file-backed out-of-memory computation, but only on numpy arrays.
- This approach works well when one's data access involves passing through an entire data sets a small number of times.
- There are other techniques for Python out-of-core computations, involving the combined use of `pytables`, `hdf5`, and `numpy`, but we won't cover them today.



Out-of-core computation, example

First, let us create a large array on file (don't actually perform these steps on GPC)

```
#file: create12GB.py
import numpy as np
n = 40000
f = np.memmap('bigfile', dtype='float64', mode='w+', shape=(n,n))
print ("initial fragment:"); print (f[0,0:5])
f[0,:] = np.random.rand(n)
print ("random fragment:"); print (f[0,0:5])
for i in xrange(n):
    f[i,:] = np.random.rand(n)
print("done")
```

```
initial fragment:
[ 0.  0.  0.  0.  0.]
random fragment:
[ 0.94379592  0.17967309  0.7169163  0.44854681  0.41266199]
done
```

Out-of-core computation, example

Exit the python prompt, and start over to calculate the mean of the array:

```
#file: average12GB.py
import numpy as np
n = 40000
f = np.memmap('bigfile', mode='r', shape=(n,n))
total = 0.0
for i in xrange(n): total += sum(f[i,:])
average = total / (n*n)
print(average)
```

```
130.186355131
```

Again, this is very hard on the file system!

Other good use of this technique:

- You have an array that is larger than would fit in memory
- You need only relatively few elements of the array
- But you do not know ahead of time which elements.