

A gentle tutorial on programming scientific applications on NVIDIA GPU's with Python and CUDA

Daniel Wysocki

Rochester Institute of Technology

RIT Scientific Computing Group
October 18th, 2018

R·I·T

Scientific computing on a *graphics* processor?



© Rob Boudon – <https://www.flickr.com/people/88562024@N00>

Graphics requires lots of linear algebra and trig – really fast

RAY TRACING

(for one pixel up to first bounce)

The diagram illustrates the ray tracing process. A light source at the top left emits rays towards a green sphere. From the eye point on a viewing plane, a primary ray is cast towards the sphere. Upon hitting the sphere, it reflects a primary ray back towards the eye and emits shadow rays to determine if other objects are in the way. The sphere's surface is shown with a normal vector N and reflection/transmission vectors L, R, and T.

① Sphere equation: $(\vec{p} - \vec{c}) \cdot (\vec{p} - \vec{c}) = r^2$ **Intersection:** $(\vec{o} + t\vec{d} - \vec{c}) \cdot (\vec{o} + t\vec{d} - \vec{c}) = r^2$
Ray equation: $\vec{r}(t) = \vec{o} + t\vec{d}$
 $t^2(\vec{d} \cdot \vec{d}) + 2(\vec{o} - \vec{c}) \cdot t\vec{d} + (\vec{o} - \vec{c}) \cdot (\vec{o} - \vec{c}) - r^2 = 0$

② Illumination Equation (Blinn-Phong) with recursive Transmitted and Reflected Intensity:

$I = k_a I_a + I_i \left(k_d (\vec{L} \cdot \vec{N}) + k_s (\vec{V} \cdot \vec{R})^n \right) + \underbrace{k_t I_t + k_r I_r}_{recursion}$

③ Snell's law: $\frac{\sin \theta_1}{\sin \theta_2} = \frac{v_1}{v_2} = \frac{n_2}{n_1}$ $n_{air} \sin \theta_i = n_{glass} \sin \theta_t$ **refraction coefficients:** $n_{air} = 1, n_{glass} = 1.5$

④ Area Light Simulation: $I_{light} \frac{\# \text{ (visible shadow rays)}}{\# \text{ (all shadow rays)}}$

© Niklaus Leopold

Single Instruction, Multiple Data (SIMD)



Ford assembly line – public domain

NVIDIA CUDA



- Currently the best performing GPUs in the world are made by NVIDIA, and are most efficiently programmed using their proprietary (freeware) language CUDA
- CUDA is basically an extension to C/C++
- Very low-level, requiring understanding of how GPU's operate to get the most efficient code possible

R·I·T

CuPy – CUDA programming in NumPy style



- CuPy is a free and open source Python library, meant as a replacement for NumPy, but using CUDA under the hood
- It includes a large subset of NumPy's features, along with additional tools for low-level GPU stuff, and for converting data between CuPy and NumPy
- Available at <https://cupy.chainer.org/>

R·I·T

CuPy – Basic example

```
1 | >>> import cupy as cp
2 | >>> x = cp.arange(6).reshape(2, 3).astype('f')
3 | >>> x
4 | array([[ 0.,  1.,  2.],
5 |        [ 3.,  4.,  5.]], dtype=float32)
6 | >>> x.sum(axis=1)
7 | array([ 3., 12.], dtype=float32)
```

CuPy installation

Assuming you already have CUDA installed, you can install CuPy with the standard Python package manager pip.

```
# (For CUDA 8.0)
$ pip install cupy-cuda80

# (For CUDA 9.0)
$ pip install cupy-cuda90

# (For CUDA 9.1)
$ pip install cupy-cuda91

# (Install CuPy from source)
$ pip install cupy
```

CuPy Arrays

- The core data structure in CuPy is the N -dimensional array, `cupy.ndarray`
- Looks and behaves almost exactly like the N -dimensional array in NumPy, except it is allocated in the GPU's memory instead

```
1 |   1>>> import numpy, cupy
2 |
3 |   2>>> cupy.arange(5)
4 |   array([0, 1, 2, 3, 4])
5 |   >>> numpy.arange(5)
6 |   array([0, 1, 2, 3, 4])
7 |
8 |   7>>> cupy.ones((2,2), dtype=float)
9 |   array([[ 1.,  1.],
10 |          [ 1.,  1.]])
11 |   >>> numpy.ones((2,2), dtype=float)
12 |   array([[ 1.,  1.],
13 |          [ 1.,  1.]])
```

Basic arithmetic in CuPy

- CuPy arrays support all basic arithmetic NumPy arrays do

```
1 | >>> x = cupy.arange(4)
2 | >>> x
3 | array([0, 1, 2, 3])
```

```
1 | >>> x + x
2 | array([0, 2, 4, 6])
3 | >>> x * x
4 | array([0, 1, 4, 9])
5 | >>> x / x
6 | array([0, 1, 1, 1])
7 | >>> x - x
8 | array([0, 0, 0, 0])
9 | >>> 2 * x
10 | array([0, 2, 4, 6])
11 | >>> x**2
12 | array([0, 1, 4, 9])
13 | >>> 3*(x**2 + 4*x + 1)
14 | array([ 3, 18, 39, 66])
```

More math functions in CuPy

```
1  >>> cupy.sin(x)
2  array([ 0.           ,  0.84147098,  0.90929743,  0.14112001])
3  >>> cupy.cos(x)
4  array([ 1.           ,  0.54030231, -0.41614684, -0.9899925 ])
5  >>> cupy.exp(x)
6  array([ 1.           ,   2.71828183,   7.3890561 ,  20.08553692])
7  >>> cupy.sqrt(x)
8  array([ 0.           ,  1.           ,  1.41421356,  1.73205081])
9  >>> cupy.outer(x, x)
10 array([[0, 0, 0, 0],
11        [0, 1, 2, 3],
12        [0, 2, 4, 6],
13        [0, 3, 6, 9]])
```

Mixing CuPy and NumPy arrays (wrong way)

```
1  >>> cupy.arange(0, 5) + cupy.arange(3, 8)
2  array([ 3,  5,  7,  9, 11])
3
4  >>> cupy.arange(0, 5) + numpy.arange(3, 8)
5  -----
6  TypeError                                Traceback (most recent call last)
7  <ipython-input-8-9951f4a4a370> in <module>()
8  ----> 1 cupy.arange(0, 5) + numpy.arange(3, 8)
9
10 cupy/core/core.pyx in cupy.core.core.ndarray.__add__()
11
12 cupy/core/elementwise.pxi in cupy.core.core.ufunc.__call__()
13
14 cupy/core/elementwise.pxi in cupy.core.core._preprocess_args()
15
16 TypeError: Unsupported type <type 'numpy.ndarray'>
```

Mixing CuPy and NumPy arrays (right way)

- must convert between CuPy and NumPy arrays to mix
 - `cupy.asarray()`: CuPy→NumPy
 - `cupy.asnumpy()`: NuMPy→CuPy

```
1  >>> cupy.asnumpy(cupy.arange(0, 5)) + numpy.arange(3, 8)
2  array([ 3,  5,  7,  9, 11])
3
4  >>> type(cupy.asnumpy(cupy.arange(0, 5)) + numpy.arange(3, 8))
5  numpy.ndarray
6
7  >>> type(cupy.arange(0, 5) + cupy.asarray(numpy.arange(3, 8)))
8  cupy.core.core.ndarray
```

- Beware: slow process, should avoid everywhere possible

Writing CuPy/NumPy agnostic code

- Can write functions that work on both CuPy and NumPy

```
1  >>> def euler_formula(x):
2      ...     "exp(i*x) = cos(x) + i*sin(x)"
3      ...     # Can be either `numpy` or `cupy`.
4      ...     xpy = cupy.get_array_module(x)
5      ...     # Compute the result with the right library.
6      ...     return xpy.cos(x) + 1j*xpy.sin(x)
7
8  >>> type(euler_formula(cupy.arange(10)))
9  cupy.core.core.ndarray
10
11 >>> type(euler_formula(numpy.arange(10)))
12 numpy.ndarray
```

Writing CuPy/NumPy agnostic code (optimized)

- For extra speed, you can save the call to `cupy.get_array_module`

```
1  >>> def euler_formula(x, xpy=cupy):
2      ...     "exp(i*x) = cos(x) + i*sin(x)"
3      ...     # Compute the result with the right library.
4      ...     return xpy.cos(x) + 1j*xpy.sin(x)
5
6  >>> type(euler_formula(cupy.arange(10), xpy=cupy))
7  cupy.core.core.ndarray
8
9  >>> type(euler_formula(numpy.arange(10), xpy=numpy))
10 numpy.ndarray
11
12 >>> type(euler_formula(numpy.arange(10), xpy=cupy))
13 TypeError: Unsupported type <type 'numpy.ndarray'>
```

More power – kernels

```
1  >>> squared_diff = cupy.ElementwiseKernel(
2      ...     'float32 x, float32 y', # Input arrays
3      ...     'float32 z', # Output array
4      ...     'z = (x - y) * (x - y)', # Compute the result and store in output array.
5      ...     'squared_diff', # Name the kernel
6      ... )
7
8  >>> x = cupy.arange(10, dtype=np.float32).reshape(2, 5)
9  >>> y = cupy.arange(5, dtype=np.float32)
10 >>> squared_diff(x, y)
11 array([[ 0.,  0.,  0.,  0.,  0.],
12        [25., 25., 25., 25., 25.]], dtype=float32)
13 >>> squared_diff(x, 5)
14 array([[25., 16.,  9.,  4.,  1.],
15        [ 0.,  1.,  4.,  9., 16.]], dtype=float32)
```

Type-generic kernels

```
1  >>> squared_diff_generic = cupy.ElementwiseKernel(
2      ...     'T x, T y',
3      ...     'T z',
4      ...     'z = (x - y) * (x - y)',
5      ...     'squared_diff_generic',
6      ...     )
```

Map/Reduce kernels

```
1  >>> l2norm_kernel = cupy.ReductionKernel(
2      ...     'T x',    # input params
3      ...     'T y',    # output params
4      ...     'x * x',  # map
5      ...     'a + b',  # reduce,
6      ...     'y = sqrt(a)', # post-reduction map
7      ...     '0',       # identity value
8      ...     'l2norm'   # kernel name
9      ... )
10 >>> x = cp.arange(10, dtype=np.float32).reshape(2, 5)
11 >>> l2norm_kernel(x, axis=1)
12 array([ 5.477226 , 15.9687195], dtype=float32)
```

Profiling – cProfile

- `cProfile` records function-level timing statistics
- Simply run `python` command as usual, but with `-m cProfile`
 - e.g., `python my_script.py` becomes `python -m cProfile my_script.py`

Test program – naive version

```
1 import numpy
2
3 seed = 1; random = numpy.random.RandomState(seed)
4
5 x = random.uniform(size=(5000,5000)); y = numpy.zeros(5000)
6
7 def sum_2d():
8     for i in range(5000):
9         for j in range(5000):
10             y[i] += x[i,j]
11
12 def sum_1d():
13     global z
14     z = 0
15     for i in range(5000):
16         z += y[i]
17
sum_2d()
sum_1d()
```

Profiling – naive version

```
$ python -m cProfile big_calculation_naive.py  
12502566.4643  
    17098 function calls (16977 primitive calls) in 9.976 seconds
```

Ordered by: standard name

ncalls	tottime	percall	cumtime	percall	filename:lineno(function)
1	0.000	0.000	0.000	0.000	<string>:1(<module>)
.....					
1	0.001	0.001	9.976	9.976	big_calculation_naive.py:1(<module>)
1	0.001	0.001	0.001	0.001	big_calculation_naive.py:11(sum_1d)
1	9.378	9.378	9.482	9.482	big_calculation_naive.py:7(sum_2d)
.....					

Test program – efficiently with NumPy

```
import numpy

seed = 1; random = numpy.random.RandomState(seed)

x = random.uniform(size=(5000,5000)); y = numpy.zeros(5000)

def sum_2d():
    x.sum(axis=1, out=y)
def sum_1d():
    global z
    z = y.sum()
sum_2d()
sum_1d()

print(z)
```

Profiling – efficiently with NumPy

```
$ python -m cProfile big_calculation_numpy.py
12502566.4643
    12102 function calls (11981 primitive calls) in 0.508 seconds

Ordered by: standard name

      ncalls  tottime  percall  cumtime  percall filename:lineno(function)
          1    0.000    0.000    0.000    0.000 <string>:1(<module>)
. . .
          1    0.001    0.001    0.508    0.508 big_calculation_numpy.py:1(<module>)
          1    0.000    0.000    0.015    0.015 big_calculation_numpy.py:7(sum_2d)
          1    0.000    0.000    0.000    0.000 big_calculation_numpy.py:9(sum_1d)
. . .
```

Test program – efficiently with CuPy

```
import cupy

seed = 1; random = cupy.random.RandomState(seed)

x = random.uniform(size=(5000,5000)); y = cupy.zeros(5000)

def sum_2d():
    x.sum(axis=1, out=y)
def sum_1d():
    global z
    z = y.sum()
sum_2d()
sum_1d()

print(z)
```

Profiling – efficiently with CuPy

```
$ python -m cProfile big_calculation_cupy.py  
12499843.3008  
    247065 function calls (242929 primitive calls) in 1.514 seconds
```

Ordered by: standard name

ncalls	tottime	percall	cumtime	percall	filename:lineno(function)
1	0.000	0.000	0.000	0.000	<string>:1(<module>)
.
1	0.001	0.001	1.516	1.516	big_calculation_cupy.py:1(<module>)
1	0.000	0.000	0.020	0.020	big_calculation_cupy.py:7(sum_2d)
1	0.000	0.000	0.033	0.033	big_calculation_cupy.py:9(sum_1d)
.

Profiling – kernprof

- `kernprof` records line-level timing statistics
- Need to pip install `line_profiler`
- Add `@profile` before functions you want profiled
- Run script with `kernprof -l` instead of `python`
 - e.g., `python my_script.py` becomes `kernprof -l my_script.py`
- Then read profiling summary with `python -m line_profiler my_script.py.lprof`

Test program – naive version

```
1 import numpy
2
3 seed = 1; random = numpy.random.RandomState(seed)
4
5 @profile
6 def main():
7     x = random.uniform(size=(1000,1000))
8     z = numpy.empty((1000,1000))
9     for i in range(1000):
10         for j in range(1000):
11             z[i,j] = x[i,j] + x[j,i]
12     y = z.sum()
13     print(y)
14
15 main()
```

Profiling – naive version (I)

```
$ kernprof -l kernprof_demo_slow.py
999896.508257
Wrote profile results to kernprof_demo_slow.py.lprof
$ python -m line_profiler kernprof_demo_slow.py.lprof
...
```

Profiling – naive version (II)

Timer unit: 1e-06 s

Total time: 3.19258 s

File: kernprof_demo_slow.py

Function: main at line 5

Line #	Hits	Time	Per Hit	% Time	Line Contents
=====					
5					@profile
6					def main():
7	1	16749.0	16749.0	0.5	x = random.uniform(size=(1000,1000))
8	1	13.0	13.0	0.0	z = numpy.empty((1000,1000))
9	1001	1238.0	1.2	0.0	for i in range(1000):
10	1001000	1235929.0	1.2	38.7	for j in range(1000):
11	1000000	1937855.0	1.9	60.7	z[i,j] = x[i,j] + x[j,i]
12	1	728.0	728.0	0.0	y = z.sum()
13	1	64.0	64.0	0.0	print(y)

Test program – fast version

```
1 import numpy  
2  
3 seed = 1; random = numpy.random.RandomState(seed)  
4  
5 @profile  
6 def main():  
7     x = random.uniform(size=(1000,1000))  
8     z = x + x.T  
9     y = z.sum()  
10    print(y)  
11  
12 main()
```

Profiling – fast version (I)

```
$ kernprof -l kernprof_demo_fast.py
999896.508257
Wrote profile results to kernprof_demo_fast.py.lprof
$ python -m line_profiler kernprof_demo_fast.py.lprof
...
```

Profiling – fast version (II)

```
Timer unit: 1e-06 s
```

```
Total time: 0.022384 s
```

```
File: kernprof_demo_fast.py
```

```
Function: main at line 5
```

Line #	Hits	Time	Per Hit	% Time	Line Contents
=====					
5					@profile
6					def main():
7	1	16635.0	16635.0	74.3	x = random.uniform(size=(1000,1000))
8	1	5027.0	5027.0	22.5	z = x + x.T
9	1	645.0	645.0	2.9	y = z.sum()
10	1	77.0	77.0	0.3	print(y)

Profiling GPU code with NVIDIA Profiler (nvprof)

```
$ nvprof python big_calculation_cupy.py
==3127650== NVPROF is profiling process 3127650, command: python big_calculation_cupy.py
12499843.3008
==3127650== Profiling application: python big_calculation_cupy.py
==3127650== Profiling result:
      Type  Time(%)     Time    Calls      Avg      Min      Max  Name
GPU activities:  33.59%  52.765ms      1  52.765ms  52.765ms  52.765ms  generate_seed_pse
                  21.08%  33.114ms      2  16.557ms  12.287us  33.102ms  cupy_sum
                  13.30%  20.890ms      1  20.890ms  20.890ms  20.890ms  cupy_multiply
                  13.29%  20.882ms      1  20.882ms  20.882ms  20.882ms  cupy_add
                  11.90%  18.695ms      1  18.695ms  18.695ms  18.695ms  cupy_random_1_min
                  6.84%  10.752ms      1  10.752ms  10.752ms  10.752ms  void gen_sequence
                  0.00%  6.3350us      1  6.3350us  6.3350us  6.3350us  [CUDA memset]
                  0.00%  1.4400us      1  1.4400us  1.4400us  1.4400us  [CUDA memcpy DtoH
API calls:    52.39%  169.80ms      6  28.299ms  1.1880us  116.57ms  cudaFree
.....
```

Profiling GPU code with NVIDIA Visual Profiler (nvvp)

```
$ nvcc python big_calculation_cupy.py
```

Will add this if I can get it to work in time.