

GPU-CPU Benchmarking and Analysis: Case Study Using Deep Learning Libraries

Clayton Turner

CSIS 612 - Advanced Computer Architecture

Fall 2015

Code Available at: https://github.com/ClaytonTurner/GPU-Benchmarking*

*References will be made to files here throughout the presentation

Outline

Theano

- Overview/Installation

- Example

Caffe

- Overview/Installation

- Example

Torch

- Overview/Installation

- Example

- Hardware/Specifications

- Deep Learning Overview

- Speedups by Comparison

- Demos

Hardware

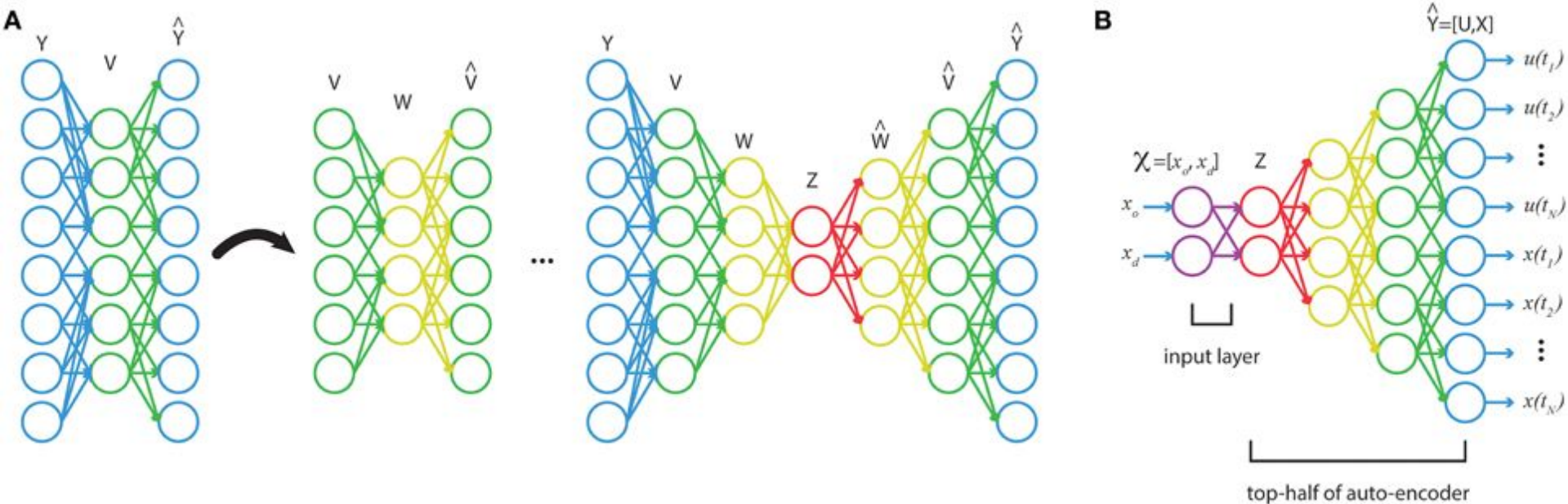
GPU: NVIDIA GeForce GTX 550 Ti
192 CUDA Cores
Graphics Clock: 900 MHz
Processor Clock: 1800 MHz

CPU: AMD FX™-8150 Eight-Core Processor (3.6 GHz)

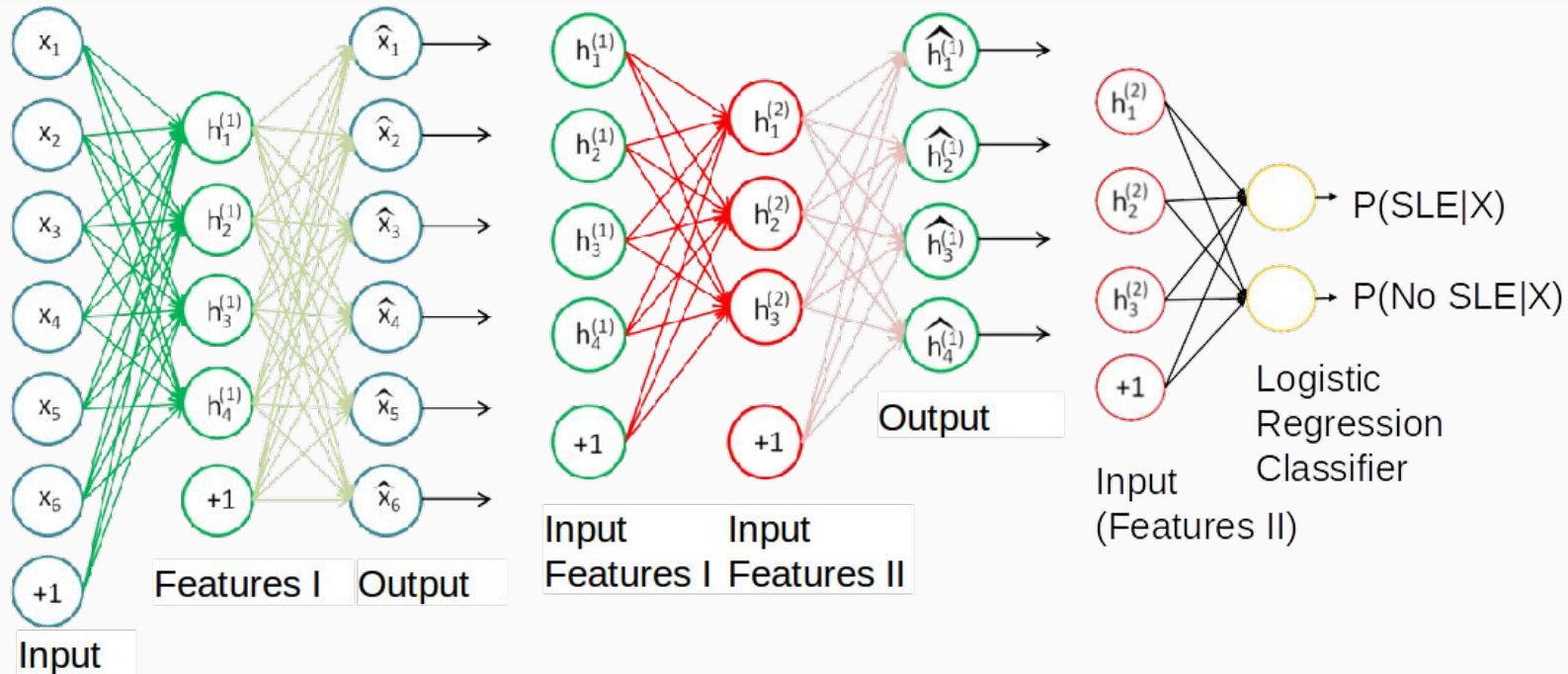


GeForce GTX 550 Ti

What Is Deep Learning? (SdA as use case)



What Is Deep Learning? (SdA as use case)



Theano: Overview

Python Deep Learning Library

Optimizes and evaluates user-defined multi-dimensional operations

Available through Git and PyPI (pip)

<https://github.com/Theano/Theano>

<http://www.deeplearning.net/software/theano>

Theano: Features

Transparent Use of GPU - up to 140x speedup over a CPU

Tightly integrated with NumPy

Efficient symbolic differentiation

Speed and stability optimizations - get the right answer for $\log(1+x)$ even when x is really tiny

Dynamic C code generation

Built-in unit-testing and self-verification

<http://www.deeplearning.net/software/theano>

Theano: Execution Modes

Runnable Modes

FAST_COMPILE ~ Python implementations where available with cheap graph transformations

FAST_RUN ~ C implementations where available with all transformations

DebugMode ~ Both implementations available with all transformations

MonitorMode ~ Debug mode with step-through (like pdb)

Example call using my file:

```
THEANO_FLAGS=mode=FAST_RUN,device=gpu,floatX=float32 python  
theano_exponentiate.py
```


Theano: Notes

Unused Optimizations

cuDNN <- NVIDIA library used by deep networks to optimize a subset of operations (such as a convolution)

CNMeM <- NVIDIA Deep Learning CUDA optimized memory management
Utilizes CUDA toolkit, pthreads/native Windows, and C++

Theano: Code

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x caturner3@caturner3-desktop: ~/GP... x

1 # Adapted from deeplearning.net
2
3 from theano import function, config, shared, sandbox
4 import theano.tensor as T
5 import numpy
6 from timeit import default_timer as timer # Since we need to use Python 2.7, this is a good alt
   ernative to Perfcounter
7
8 vlen = 10 * 30 * 768 # 10 x #cores x # threads per core
9 iters = 1000000
10
11 # Keep consistent results with random seed
12 rng = numpy.random.RandomState(42)
13
14 # Use Theano's shared variables in order to distribute computation to the GPU
15 # floatX == float32 must hold True for GPU to work - command line specification
16 # rng.rand(z) creates an array of size z
17 x = shared(numpy.asarray(rng.rand(vlen), config.floatX))
18
19 # Function syntax:
20 #   arg[0] = parameters to be passed in
21 #   arg[1] = the function to be performed
22 # We pass our shared array as part of the function definition
23 f = function([], T.exp(x))
24
25 begin = timer()
26 for i in xrange(iters):
27     r = f()
28 end = timer()
29 print("Looping %d times took %f seconds" % (iters, end - begin))
30 print("Result is %s" % (r,))
31
32 # Let's make sure we used the correct processing unit by printing and comparing to what we expect
33 if numpy.any([isinstance(x.op, T.Elemwise) for x in f.maker.fgraph.toposort()]):
34     print('Used the cpu')
35 else:
36     print('Used the gpu')
```

Theano: Digging Deeper ~ Optimization

<https://github.com/Theano/Theano/blob/master/theano/tensor/opt.py>

Optimization: Creating graphs for tensors before values are reached

```
def inplace_elemwise_optimizer_op(OP) # Parametrization for GPU | CPU
    def inplace_elemwise_optimizer(fgraph)
```

```
ex // x + y + z -> x += y += z
```

```
ex2// (x + y) * (x * y) -> (x += y) *= (x * y) or (x + y) *= (x *= y)
```

Caffe: Overview

Python Deep Learning framework
Focused on expression, speed, and modularity

Available on git and PyPI (requires several dependencies)

<https://github.com/BVLC/caffe>

<http://caffe.berkeleyvision.org/installation.html>

Caffe: Features

Expressive Architecture

switch between CPU and GPU through flags*

Extensible Code

Speed

can process over 60M images per day with a single NVIDIA K40 GPU

*We only show GPU results in later slides due to the nature of extracting functions from Caffe

<http://caffe.berkeleyvision.org>

Caffe:Notes

Works with Python and Matlab

Python version uses Boost.Python to integrate C++ and Python
http://www.boost.org/doc/libs/1_59_0/libs/python/doc/index.html

Hardware performance benchmarking: http://caffe.berkeleyvision.org/performance_hardware.html

Caffe: Code

Utils file located on Git:
https://github.com/ClaytonTurner/GPU-Benchmarking/blob/master/code/caffe_utils.py

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x caturner3@caturner3-desktop: ~/GP... x
1 from __future__ import print_function
2 from caffe import layers as L, params as P, to_proto
3 from caffe.proto import caffe_pb2
4 from timeit import default_timer as timer
5 import numpy
6
7 # helper function for common structures
8 vlen = 10 * 30 * 768
9 rng = numpy.random.RandomState(42)
10 blob = caffe_pb2.BlobProto()
11 blob.data.extend(list(rng.rand(vlen)))
12 data = blob.data
13
14 def conv_relu(bottom, ks, nout, stride=1, pad=0, group=1):
15     conv = L.Convolution(bottom, kernel_size=ks, stride=stride,
16                           num_output=nout, pad=pad, group=group)
17     return conv, L.ReLU(conv, in_place=True)
18
19 def fc_relu(bottom, nout):
20     fc = L.InnerProduct(bottom, num_output=nout)
21     return fc, L.ReLU(fc, in_place=True)
22
23 def max_pool(bottom, ks, stride=1):
24     return L.Pooling(bottom, pool=P.Pooling.MAX, kernel_size=ks, stride=stride)
25
26 def euc_loss(x):
27     return L.EuclideanLoss(x,x)
28
29 def caffe_net(lmdb, batch_size=256, include_acc=False):
30     euc1 = euc_loss(data)
31     #loss = L.SoftmaxWithLoss(euc1, label)
32     #return to_proto(loss)
33     return None
34
35 def make_net():
36     iters = 1000000
37     begin = timer()
38     for i in xrange(iters):
39         with open('../data/train.prototxt', 'w') as f:
40             print(caffe_net('/path/to/caffe-train-lmdb'), file=f)
41         with open('../data/test.prototxt', 'w') as f:
42             print(caffe_net('/path/to/caffe-val-lmdb', batch_size=50, include_acc=True), file=f)
43     end = timer()
44     print("Looping %d times took %f seconds" % (iters, end - begin))
45
46 if __name__ == '__main__':
47     make_net()
```

Caffe: Digging Deeper ~ Establishing GPU

<https://github.com/BVLC/caffe/blob/master/src/caffe/parallel.cpp>

```
67  
68  
69 template<typename Dtype>  
70 Params<Dtype>::Params(shared_ptr<Solver<Dtype> > root_solver)  
71 : size_(total_size<Dtype>(root_solver->net()->learnable_params())),  
72   data_(),  
73   diff_() {  
74 }  
75  
76 template<typename Dtype>  
77 GPUParams<Dtype>::GPUParams(shared_ptr<Solver<Dtype> > root_solver, int device)  
78 : Params<Dtype>(root_solver) {  
79 #ifndef CPU_ONLY  
80   int initial_device;  
81   CUDA_CHECK(cudaGetDevice(&initial_device));  
82  
83   // Allocate device buffers  
84   CUDA_CHECK(cudaSetDevice(device));  
85   CUDA_CHECK(cudaMalloc(&data_, size_ * sizeof(Dtype)));  
86  
87   // Copy blob values  
88   const vector<Blob<Dtype>*> &net =  
89     root_solver->net()->learnable_params();  
90   apply_buffers(net, data_, size_, copy);  
91  
92   CUDA_CHECK(cudaMalloc(&diff_, size_ * sizeof(Dtype)));  
93   caffe_gpu_set(size_, Dtype(0), diff_);  
94  
95   CUDA_CHECK(cudaSetDevice(initial_device));  
96 #else  
97   NO_GPU;  
98 #endif  
99 }  
100  
101 template<typename Dtype>  
102 GPUParams<Dtype>::~GPUParams() {  
103 #ifndef CPU_ONLY  
104   CUDA_CHECK(cudaFree(data_));  
105   CUDA_CHECK(cudaFree(diff_));  
106 #endif  
107 }  
108  
109 template<typename Dtype>  
110 void GPUParams<Dtype>::configure(Solver<Dtype>* solver) const {  
111   const vector<Blob<Dtype>*> &net =  
112     solver->net()->learnable_params();  
113   apply_buffers(net, data_, size_, replace_gpu);  
114   apply_buffers(net, diff_, size_, replace_gpu_diff);  
115 }  
116
```


Caffe: Digging Deeper ~ Updating Layers with Mutex Locking

[https://github.com/BVLC/caffe/blob/master/
src/caffe/layer.cpp](https://github.com/BVLC/caffe/blob/master/src/caffe/layer.cpp)

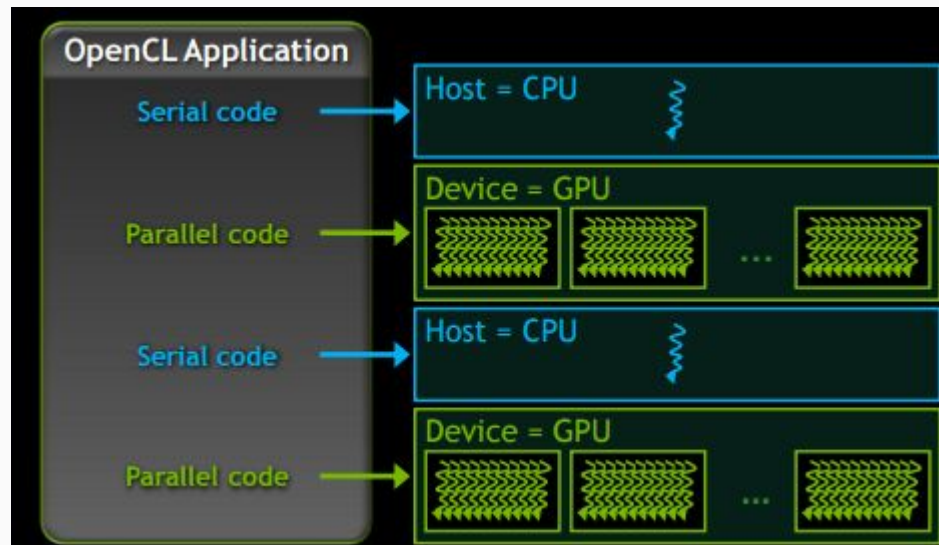
```
1  #include <boost/thread.hpp>
2  #include "caffe/layer.hpp"
3
4  namespace caffe {
5
6  template <typename Dtype>
7  void Layer<Dtype>::InitMutex() {
8      forward_mutex_.reset(new boost::mutex());
9  }
10
11 template <typename Dtype>
12 void Layer<Dtype>::Lock() {
13     if (IsShared()) {
14         forward_mutex_>lock();
15     }
16 }
17
18 template <typename Dtype>
19 void Layer<Dtype>::Unlock() {
20     if (IsShared()) {
21         forward_mutex_>unlock();
22     }
23 }
24
25 INSTANTIATE_CLASS(Layer);
26
27 } // namespace caffe
```

Torch: Overview

TorchCI Library for integrating with OpenCL

Flexible, fast framework operating
under LuaJIT and C/CUDA

Mostly neural network and optimization
libraries



Torch: Features

Powerful N-dimensional array

Efficient data munging routines

Interface to C, via LuaJIT

Linear algebra routines

Neural network and energy-based models at heart

Numeric optimization routines

Fast and efficient GPU support

Embeddable with iOS, Android, and FPGA backends

<http://torch.ch/>

Torch: Code

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x

1 x = torch.rand(1,16*30*768)
2
3 for i = 1, 1000 do
4     y = torch.exp(x)
5 end
6
7 time = os.clock()
8 print("Time elapsed: ", time)
```

CPU

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x

1 require 'cltorch'
2 require 'math'
3
4 x = torch.Tensor(1,16*30*768):uniform():cl()
5 y = torch.Tensor(1,16*30*768):zero():cl()
6
7 iters = 1000000
8
9 --begin = os.time()
10 for i=1,iters do
11     y:apply("x = exp(x)")
12     --x:map(y, "x = exp(x)")
13 end
14 --stop = os.time()
15
16 -- os.time() has second resolution - garbage for us
17 -- So we're going to use os.clock() even though it includes the
18 --print("Time elapsed: ",os.difftime(stop,begin))
19 print("Time elapsed: ", os.clock())
```

GPU

Torch: Digging Deeper ~

C <-> Lua integration

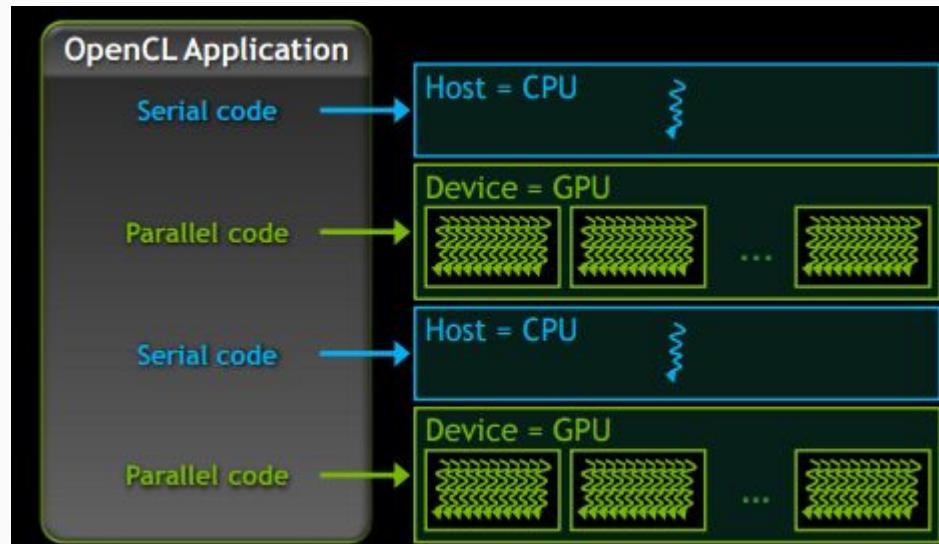
<https://github.com/torch/torch7/blob/master/lib/luaT/luaT.c>

```
1 #include <stdlib.h>
2 #include <string.h>
3
4 #include "luaT.h"
5
6 void* luaT_alloc(lua_State *L, long size)
7 {
8     void *ptr;
9
10    if(size == 0)
11        return NULL;
12
13    if(size < 0)
14        luaL_error(L, "$ Torch: invalid memory size -- maybe an overflow?");
15
16    ptr = malloc(size);
17    if(!ptr)
18        luaL_error(L, "$ Torch: not enough memory: you tried to allocate %dGB. Buy new RAM!", size/1073741824);
19
20    return ptr;
21 }
22
23 void* luaT_realloc(lua_State *L, void *ptr, long size)
24 {
25     if(!ptr)
26         return(luaT_alloc(L, size));
27
28     if(size == 0)
29     {
30         luaT_free(L, ptr);
31         return NULL;
32     }
33
34     if(size < 0)
35         luaL_error(L, "$ Torch: invalid memory size -- maybe an overflow?");
36
37     ptr = realloc(ptr, size);
38     if(!ptr)
39         luaL_error(L, "$ Torch: not enough memory: you tried to reallocate %dGB. Buy new RAM!", size/1073741824);
40     return ptr;
41 }
42
43 void luaT_free(lua_State *L, void *ptr)
44 {
45     free(ptr);
46 }
47
48 void luaT_setfuncs(lua_State *L, const luaL_Reg *l, int nup)
49 {
50     #if LUA_VERSION_NUM == 501
51         luaL_checkstack(L, nup+1, "too many upvalues");
```

PyOpenCL: OpenCL Refresher

Used for benchmarking standard OpenCL integration with Python for results

<http://mathema.tician.de/software/pyopencl/>
<http://documen.tician.de/pyopencl/>



PyOpenCL: Code

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x caturner3@caturner3-desktop: ~/GP... x

1 import pyopencl as cl
2 import numpy
3 import sys
4 from timeit import default_timer as timer
5
6 class CL(object):
7     def __init__(self, size=10*30*768):
8         self.size = size
9         self.ctx = cl.create_some_context()
10        self.queue = cl.CommandQueue(self.ctx)
11        self.iters = 1000000
12
13    def load_program(self):
14        # Using exponentiate instead of exp to avoid overloading/overriding
15        fstr="""
16        __kernel void exponentiate(__global float* a, __global float* c)
17        {
18            unsigned int i = get_global_id(0);
19
20            c[i] = a[i] * a[i];
21        }
22        """
23        self.program = cl.Program(self.ctx, fstr).build()
24
25    def popCorn(self):
26        mf = cl.mem_flags
27
28        # We use float32s here to have a fair comparison to the CPU
29        self.a = numpy.array(range(self.size), dtype=numpy.float32)
30
31        self.a_buf = cl.Buffer(self.ctx, mf.READ_ONLY | mf.COPY_HOST_PTR,
32                               hostbuf=self.a)
33        self.dest_buf = cl.Buffer(self.ctx, mf.WRITE_ONLY, self.a.nbytes)
34
35    def execute(self):
36        begin = timer()
37        for i in xrange(self.iters):
38            self.program.exponentiate(self.queue, self.a.shape, None, self.a_buf, self.dest_buf)
39        end = timer()
40        c = numpy.empty_like(self.a)
41        cl.enqueue_read_buffer(self.queue, self.dest_buf, c).wait()
42        print("Looping %d times took %f seconds" % (self.iters, end - begin))
43        print("Result is %s" % (c))
44
45 if __name__ == '__main__':
46     matrixmul = CL(10000000)
47     matrixmul.load_program()
48     matrixmul.popCorn()
49     matrixmul.execute()
```


Numpy Code

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x caturner3@caturner3-desktop: ~/GP... x

1 import numpy
2 from timeit import default_timer as timer # Since we need to use Python 2.7, this is a good alt
   ernative to Perfcounter
3
4 vlen = 10 * 30 * 768 # 10 x #cores x # threads per core
5 iters = 1000
6
7 # Keep consistent results with random seed
8 rng = numpy.random.RandomState(42)
9
10 # float32 so we have a fair comparison to the Theano times
11 # rng.rand(z) creates an array of size z
12 x = numpy.asarray(rng.rand(vlen), numpy.float32)
13
14 begin = timer()
15 for i in xrange(iters):
16     r = numpy.exp(x)
17 end = timer()
18 print("Looping %d times took %f seconds" % (iters, end - begin))
19 print("Result is %s" % (r,))
```

1,1 All

Baseline Code

```
caturner3@caturner3-desktop: ~/GPU-Benchmarking/code
caturner3@caturner3-desktop: ~ x caturner3@caturner3-desktop: ~/GP... x caturner3@cat

2  ## We convert back to Python native lists before timing to avoid timing bias
3
4  import numpy
5  from timeit import default_timer as timer # Since we need to use Python 2.
   alternative to Perfcounter
6
7  vlen = 10 * 30 * 768 # 10 x #cores x # threads per core
8  iters = 1000
9
10 # Keep consistent results with random seed
11 rng = numpy.random.RandomState(42)
12
13 # float32 so we have a fair comparison to the Theano times
14 # rng.rand(z) creates an array of size z
15 x = numpy.asarray(rng.rand(vlen), numpy.float32)
16 x = x.tolist()
17
18 def exp(exp_list):
19     for i in xrange(len(exp_list)):
20         exp_list[i] = exp_list[i] ** 2
21     return exp_list
22
23 begin = timer()
24 for i in xrange(iters):
25     r = exp(x)
26 end = timer()
27 print("Looping %d times took %f seconds" % (iters, end - begin))
28 # We're not printing the result because python doesn't, by default, have t
   en in Theano and Numpy by default
29 #print("Result is %s" % (r,))
```

Exponentiating Numbers

1000 iterations of exponentiating

Torch*

Timing included data migration
to GPU (see code)

Caffe**

No simple array interface
Utilized sum squared layers
to mimic exponentiating

Torch - TorchCL (GPU)*	0.156538 seconds
Caffe**	0.167299 seconds
Theano - FAST_RUN (GPU)	0.555048 seconds
PyOpenCL	0.976866seconds
Numpy	2.904454 seconds
Theano - FAST_COMPILE (CPU)	5.129712 seconds
Torch - (CPU)	6.009791 seconds
Normal Looping	22.315899 seconds
Theano - DebugMode	52.846434 seconds

Exponentiating Numbers

1000000 iterations of exponentiating

Torch*

Timing included data migration
to GPU (see code)

Caffe**

No simple array interface
Utilized sum squared layers
to mimic exponentiating

Torch - TorchCL (GPU)*	83.952904 seconds
Caffe**	305.870450 seconds
Theano - FAST_RUN (GPU)	525.722036 seconds
PyOpenCL	1359.186030 seconds
Numpy	2886.777740 seconds
Theano - FAST_COMPILE (CPU)	2921.630621 seconds
Torch - (CPU)	5375.984789 seconds
Normal Looping	seconds
Theano - DebugMode	seconds

Demos/Else

<https://github.com/lisa-lab/DeepLearningTutorials>

<http://demo.caffe.berkeleyvision.org/>

<https://github.com/torch/demos>

Code displayed using VIM with a custom .vimrc

Available at: <https://github.com/ClaytonTurner/UsefulScripts/tree/master/vim>

References

NVIDIA

<http://www.geforce.com/hardware/desktop-gpus/geforce-gtx-550ti>

Theano

<http://deeplearning.net/software/theano/index.html>

Caffe

<http://caffe.berkeleyvision.org/>

Torch

<http://torch.ch/docs/getting-started.html>

<https://github.com/hughperkins/cltorch>

PyOpenCL

<http://documen.tician.de/pyopencl/>