
GPU-Accelerated Deep Convolutional Neural Network for Object Recognition

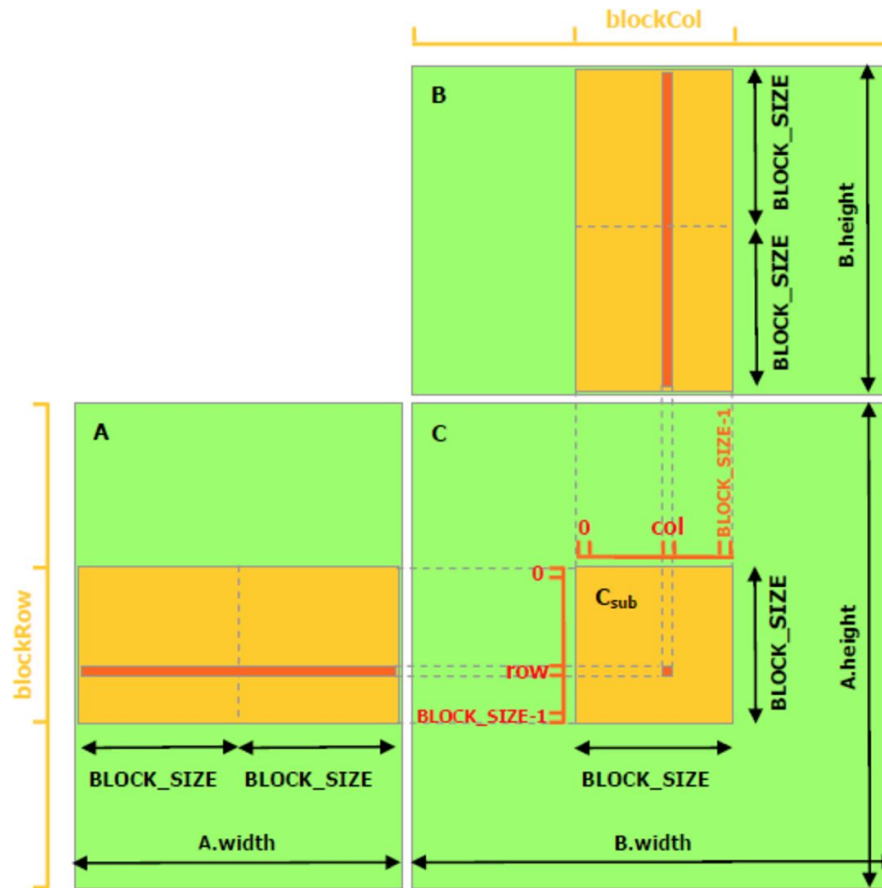
— Guan Sun —

Goal for this week

Continue working on optimization

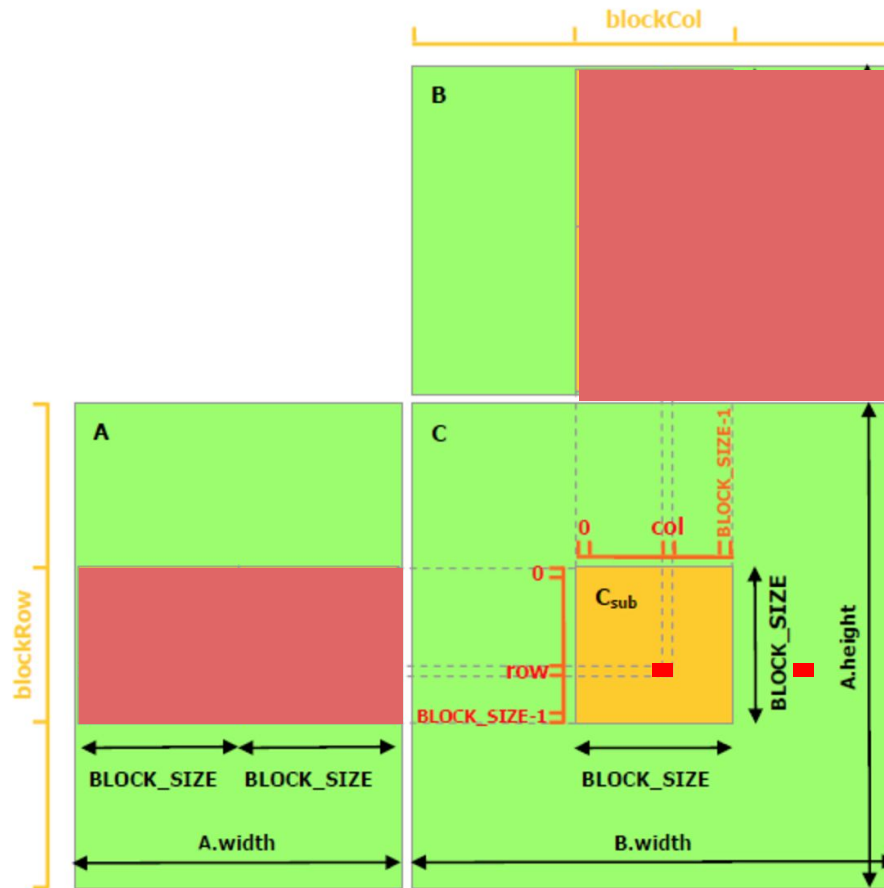
Run test on larger dataset

Matrix Multiplication Optimization



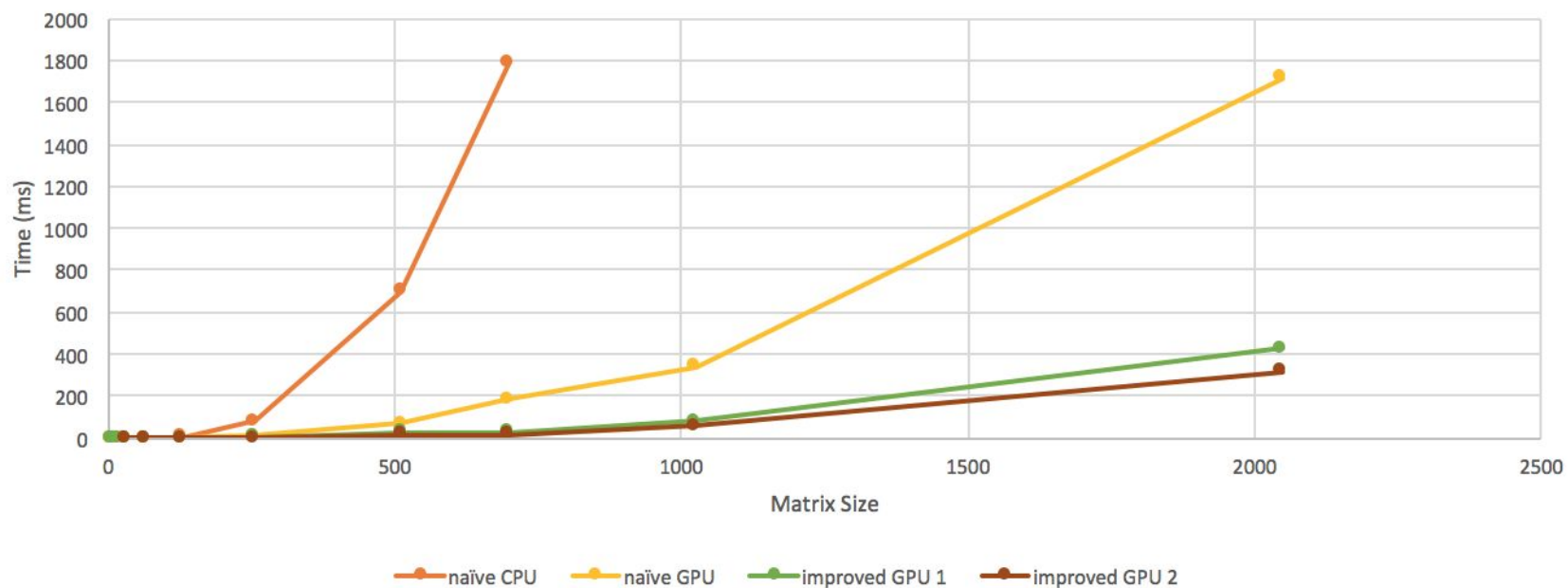
Each element is read $N/Block_size$ times from global memory.
Threads access elements from faster shared memory.

Matrix Multiplication Optimization



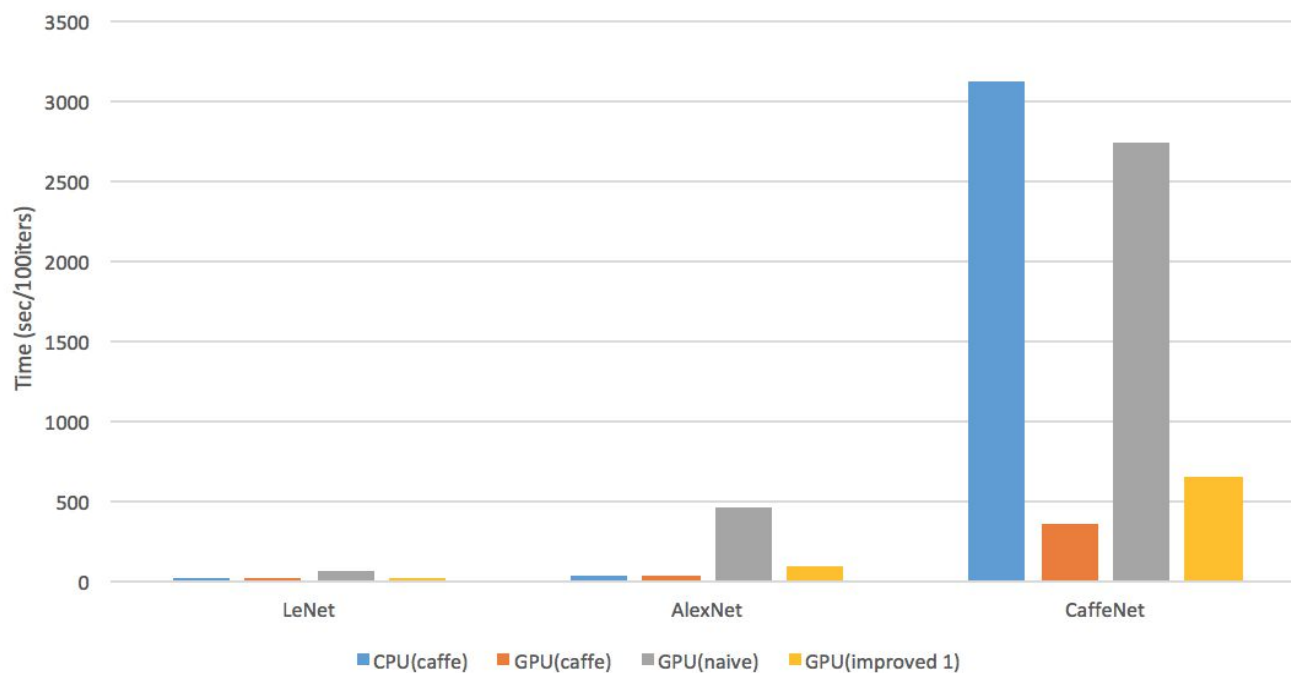
Each block stores more elements in the shared memory, so each thread can compute two elements. Number of threads is reduced to half!

Results after Optimization



(ms)

Performance

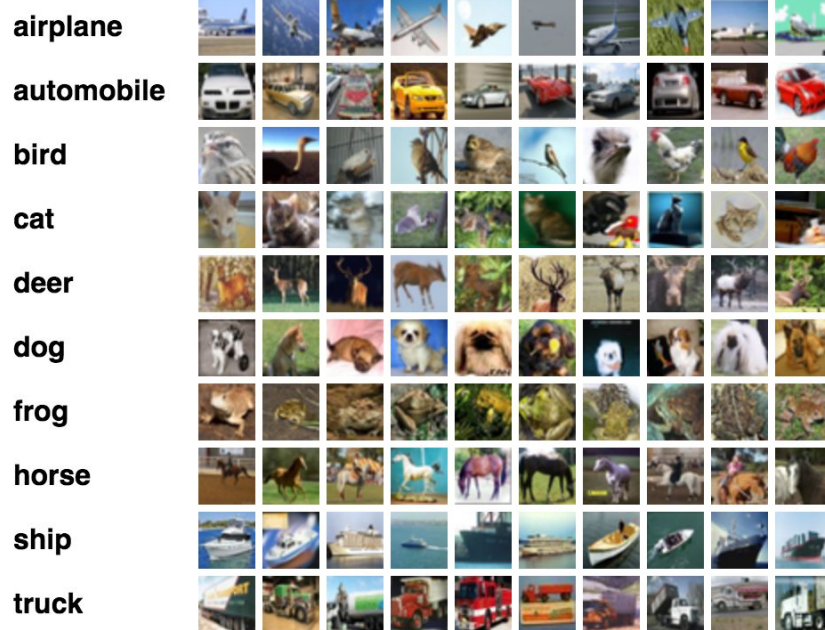


	LeNet	AlexNet	CaffeNet
CPU(caffe)	6.11	42.89	3119.48
GPU(caffe)	6.89	30.67	365.7
GPU(naive)	64.78	456.11	2742.8
GPU(improved 1)	15.82	89.06	653.04

(seconds/100 training iters)

Object Recognition

The CIFAR-10 dataset



input

conv1

pool1

relu1

conv2

relu2

pool2

conv3

relu3

pool3

inner product1

inner product2

softmax

output

```
I1207 17:07:34.491516 2071388928 solver.cpp:320] Iteration 5000, loss = 0.481036
I1207 17:07:34.491554 2071388928 solver.cpp:340] Iteration 5000, Testing net (#0)
I1207 17:07:44.542856 2071388928 solver.cpp:408] Test net output #0: accuracy = 0.7621
I1207 17:07:44.542901 2071388928 solver.cpp:408] Test net output #1: loss = 0.722371 (* 1 = 0.722371 loss)
I1207 17:07:44.542912 2071388928 solver.cpp:325] Optimization Done.
I1207 17:07:44.542920 2071388928 caffe.cpp:215] Optimization Done.
```

Accuracy = 76.21%

Final Results

- CUDA implementation of 5 neural network layers under Caffe framework
- A trained deep CNN model for object recognition
- Performance Analysis
- Optimized matrix multiplication

Thanks!