

Implementing Lenet5 by using CUDA extension

2014190150

Jung Woo Lee (mgp10067)

1. Introduction

The critical problem of initial multilayered neural network was less flexibility toward topology of data. Even the date set with little change, the model consider data as all different data to train. Since the raw fully connected multi layered neural network did not consider the topology but only consider raw data, it requires extreme numbers of training data and high time cost to train and test.

LeNet is a model that first introduced Convolutional Neural Network(CNN), which solved above problem of fully connected neural network. The model includes convolution process that use kernel filter to extract feature maps. Also, by using max-pooling to subsample feature maps for greatest correlation. As a result, by adding convolution and pooling process before training fully connected network, model can be trained with smaller data set with more flexibility. Since the model uses features of data to train model, it can consider topology without new data set to test another data set.

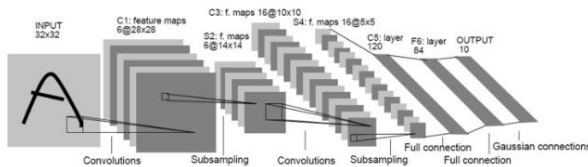


Figure 1. LeNet5 Architecture

Although LeNet5 uses smaller size data set, the training time and testing time still costs high when the data is more complex. To reduce training time of model, general purposed GPU technique is applied. Not only for LeNet, but also for every machine learning models, gpGPU skill greatly reduced time to both train and testing.

The GPU has pretty simple architecture, containing a lot of ALU (or cores) that can run operations in parallel. Therefore, by calculating training process in optimized parallel operations, the time cost to perform machine learning will be greatly decreased.

In this report, I will introduce the characteristics of GPU programming (CUDA), implementation of LeNet5 by using CUDA, and optimization to speed up testing performance.

2. Implementation of CUDA

To run GPU as a general purpose, user must allocate certain memory space of GPU and copy data

from main memory to the GPU memory. Since CPU and GPU do not share their memory, data must be copied from host to device and then, after all calculation, result must be copied back to host from device. The data transfer is one of the greatest bottle neck of gpGPU processing and thus, user should not frequently call data from CPU in the Cuda kernel function but copy before calculation.

```
cudaMalloc((void**)&d_output, sizeof(double) * batch * output_size);  
  
// Copy Parameters  
  
cudaMemcpy(d_conv1_weight, conv1_weight,  
           sizeof(double) * conv1_in_channel * conv1_out_channel *  
           conv1_kernel_size * conv1_kernel_size,  
           cudaMemcpyHostToDevice);
```

Figure 2. Prepare device memory

A. Normalization

Normalization of input data increases final accuracy after the training. I implemented normalize function with global and applied same as cpu version with mean 0.5 and var 0.5. The difference is an use of dim3 variables to allocate block size and grid size to GPU. I used 2d grid with input channel and batch, and 2d block with input size and input size to fit input data. By running the function, each thread will normalize input values in parallel instead of using naïve iteration.

B. Convolution

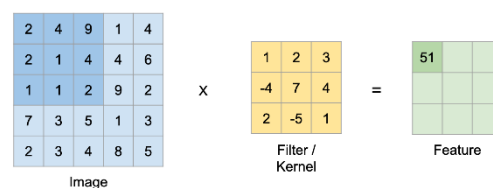


Figure 3. Convolution technique example

In the convolution process, I applied grid dimension as conv1 output channel x batch number and block dimension as input size-(kernel size -1). Operations in the function is same as CPU version, for instance in convolution 1, 5x5 size kernel iterate the 32x32 input to extract 28x28 feature maps. First 4 loops were replaced by each parallel thread while, other 3 to calculate feature map still remains.

The key difference is that in CPU version code, the calculated value with kernel, continually added to output channel data directly while in CUDA version, I used temporal variable to sum up all calculation and applied to output channel at last to minimize direct


```
int output_index = blockIdx.y * (gridDim.x * blockDim.y * blockDim.x)
                  + blockIdx.x * (blockDim.y * blockDim.x)
                  + (int)fmaf(threadIdx.y, blockDim.x, threadIdx.x);
double tmp = bias[output_index];
```

I applied the `fmaf` function to `conv` function and fully connected functions to reduce calculation times. As a result, the predict speed reduced from 1.55 to about 1.49ms.

One of the bottle neck of CUDA version is internal loops in the cuda functions. In the conv, pool, and fc function, there are internal loop to calculate output from each layer. Although each thread in the block runs in parallel, the internal loop runs as sequential way. To minimize the bottle neck due to internal sequential loop, I applied `#pragma unroll` above the for loop. The command will unroll the loop with branch prediction.

```
#pragma unroll
for(int ic = 0; ic < IC; ic++){
    tmp += weight[(int)fmaf(threadIdx.x, IC, ic)] * input[(int)fmaf(blockIdx.x, IC, ic)];
}
```

As a result, I was able to reduce the speed of prediction from 1.49ms to about 1.47ms.

By using `gpu trace` shell command, I was able to check the runtime of each kernel functions and the size of grids and blocks.

===== Preparing application: c:\program files\hp\hpsetup.exe 1.71616121 1.71616121 1.71616121 1.71616121													
===== Preparing release:													
File Name	Duration	File Size	Block Size	Age*	SPW*	DPW*	Size Transferred	Percentage	Software	Device	Index	Content	Stream
321.396	1.440000	1.440000	1.440000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	1.440000	1.440000	1.440000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	800.000000	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0	1.515040	100.000000	Pagefile	Device	defragment H3A_206	1	7 [EISA memory H3D]
321.396	0.720000	0.720000	0.720000	0	0	0							

According to the figure 10, conv, fc layer functions and memory copy cause most delays. First and second convolution took 23.6us and 19.2 us respectively. Also, the first fully connected layer ran for 40.2us which is the largest time consuming. The reason of fc function latency is an internal 1d loops that calculate 120 feature maps in 1 grid.

To prevent errors from prediction, the last memcpy is occurred to synchronize others. In my first guess, concurrent running of memcpy functions will reduce the total time of running prepare function. However, time duration of each memcpy increased as figure 11 shows and thus could not solve the memcpy bottle neck.

[illegible]

According to reference [3], calculation of converted number from 32-bit floating point to 16-bit fixed point has only small change. Thus, if there was no limitation in conversion of precision, data conversion of input can increase the speed performance.

[2] Yann LeCun, Leon Bottou, Gradient-Based Learning Applied to Document Recognition

[3] Tianshi Chen, Zidong Du, DianNao: a small-footprint high-throughput accelerator for ubiquitous machine-learning.