

EE-147 Lab 2

Tiled Matrix Multiplication

Cody Simons
861177050

April 28, 2018

1 Introduction

This lab our goal was to implement a tiled matrix multiplication algorithm. The goal of a tiled algorithm is to increase the the memory locality of an operation. We do this by batching together operations that access the same memory. This gives almost a 4 times speed up on large matrices.

2 Questions

1. In your kernel implementation, how many threads can be simultaneously executing? Assume a GPU which has 30 streaming multiprocessors.
Each multiprocessor executes 32 threads at a time, so there can be at a maximum 30×32 excuting at a time or 960 threads.
2. Use `nvcc --ptxas-options="-v"` to report the resource usage of your implementation. Note that the compilation will fail but you will still get a report of the relevant information. Experiment with the Nvidia visual profiler, which is part of the CUDA toolkit, and use it to further understand the resource usage. In particular, report your branch divergence behavior and whether your memory accesses are coalesced.
My kernel used 25 registers, 2048 bytes of smem, and 360 bytes of cmem. Using nvprof I found that there are 7.8 Gloabel Load Transactions Per Request. I think this means that it is loading 7 bytes per request. There are 5.9 Global store Transactions Per Request. Nvprof also classified the Control-Flow Function Unit Utilization as low which I think indicates that there isnt a lot of branch divergence.
3. How many times is each element of the input matrices loaded during the execution of the kernel?
While it will very depending on the tile size and the size of the input matrices. If A is an $m \times k$ matrix and B is and $k \times n$ matrix then every element in A will be loaded $\left\lceil \frac{n}{tile_size} \right\rceil$ times and every element in B will be loaded $\left\lceil \frac{m}{tile_size} \right\rceil$ times.

3 Code

3.1 Main.cu

The main file for this lab is very similar to the last lab. The main difference is that instead of taking the size of the vector as a parameter, we must calculate the length which is the total number of elements in the matrix.

```
1 /*****
2 *
3 *           (C) Copyright 2010 The Board of Trustees of the
4 *           University of Illinois
5 *           All Rights Reserved
6 *
7 *****/
```

```

8
9 #include <stdio.h>
10 #include <stdlib.h>
11 #include "kernel.cu"
12 #include "support.h"
13
14 int main (int argc, char *argv[])
15 {
16
17     Timer timer;
18     cudaError_t cuda_ret;
19
20     // Initialize host variables -----
21
22     printf("\nSetting up the problem..."); fflush(stdout);
23     startTime(&timer);
24
25     float *A_h, *B_h, *C_h;
26     float *A_d, *B_d, *C_d;
27     size_t A_sz, B_sz, C_sz;
28     unsigned matArow, matAcol;
29     unsigned matBrow, matBcol;
30     dim3 dim_grid, dim_block;
31
32     if (argc == 1) {
33         matArow = 1000;
34         matAcol = matBrow = 1000;
35         matBcol = 1000;
36     } else if (argc == 2) {
37         matArow = atoi(argv[1]);
38         matAcol = matBrow = atoi(argv[1]);
39         matBcol = atoi(argv[1]);
40     } else if (argc == 4) {
41         matArow = atoi(argv[1]);
42         matAcol = matBrow = atoi(argv[2]);
43         matBcol = atoi(argv[3]);
44     } else {
45         printf("\n    Invalid input parameters!"
46             "\n    Usage: ./sgemm-tiled          # All matrices are 1000 x 1000"
47             "\n    Usage: ./sgemm-tiled <n>        # All matrices are m x m"
48             "\n    Usage: ./sgemm-tiled <n> <k> <n> # A: m x k, B: k x n, C: m x n"
49             "\n");
50         exit(0);
51     }
52
53     A_sz = matArow*matAcol;
54     B_sz = matBrow*matBcol;
55     C_sz = matArow*matBcol;
56
57     A_h = (float*) malloc( sizeof(float)*A_sz );
58     for (unsigned int i=0; i < A_sz; i++) { A_h[i] = (rand()%100)/100.00; }
59
60     B_h = (float*) malloc( sizeof(float)*B_sz );
61     for (unsigned int i=0; i < B_sz; i++) { B_h[i] = (rand()%100)/100.00; }
62
63     C_h = (float*) malloc( sizeof(float)*C_sz );
64
65     stopTime(&timer); printf("%f s\n", elapsedTime(timer));
66     printf("    A: %u x %u\n    B: %u x %u\n    C: %u x %u\n", matArow, matAcol,
67         matBrow, matBcol, matArow, matBcol);
68
69     // Allocate device variables -----
70
71     printf("Allocating device variables..."); fflush(stdout);
72     startTime(&timer);
73
74     //INSERT CODE HERE
75     cudaMalloc((void**) &A_d, sizeof(float)*A_sz);
76     cudaMalloc((void**) &B_d, sizeof(float)*B_sz);
77     cudaMalloc((void**) &C_d, sizeof(float)*C_sz);
78

```

```

79  cudaDeviceSynchronize();
80  stopTime(&timer); printf("%f s\n", elapsedTime(timer));
81
82  // Copy host variables to device -----
83
84  printf("Copying data from host to device..."); fflush(stdout);
85  startTime(&timer);
86
87  //INSERT CODE HERE
88  cudaMemcpy(A_d, A_h, sizeof(float)*A_sz, cudaMemcpyHostToDevice);
89  cudaMemcpy(B_d, B_h, sizeof(float)*B_sz, cudaMemcpyHostToDevice);
90
91  cudaDeviceSynchronize();
92  stopTime(&timer); printf("%f s\n", elapsedTime(timer));
93
94  // Launch kernel using standard sgemv interface -----
95  printf("Launching kernel..."); fflush(stdout);
96  startTime(&timer);
97  basicSgemv('N', 'N', matArow, matBcol, matBrow, 1.0f, \
98            A_d, matArow, B_d, matBrow, 0.0f, C_d, matBrow);
99
100  cuda_ret = cudaDeviceSynchronize();
101  if(cuda_ret != cudaSuccess) FATAL("Unable to launch kernel");
102  stopTime(&timer); printf("%f s\n", elapsedTime(timer));
103
104  // Copy device variables from host -----
105
106  printf("Copying data from device to host..."); fflush(stdout);
107  startTime(&timer);
108
109  //INSERT CODE HERE
110  cudaMemcpy(C_h, C_d, sizeof(float)*C_sz, cudaMemcpyDeviceToHost);
111
112  cudaDeviceSynchronize();
113  stopTime(&timer); printf("%f s\n", elapsedTime(timer));
114
115  // Verify correctness -----
116
117  printf("Verifying results..."); fflush(stdout);
118
119  verify(A_h, B_h, C_h, matArow, matAcol, matBcol);
120
121
122  // Free memory -----
123
124  free(A_h);
125  free(B_h);
126  free(C_h);
127
128  //INSERT CODE HERE
129  cudaFree(A_d);
130  cudaFree(B_d);
131  cudaFree(C_d);
132
133  return 0;
134
135 }

```

3.2 Kernel.cu

In the kernel file I implemented the naive implementation and the tiled algorithm. The naive one is fairly straight forward. The tiled one takes a bit more work because we must load each element of the inputs in chunks. We must make sure to take care of the edge cases as well.

```

1  /*****
2  *cr
3  *cr          (C) Copyright 2010 The Board of Trustees of the
4  *cr          University of Illinois
5  *cr          All Rights Reserved

```

```

6  *CR
7  *****/
8
9  #include <stdio.h>
10 #define TILE_SIZE 16
11
12 __global__ void mysgemm(int m, int n, int k, const float *A, const float *B, float* C) {
13
14     /******
15     *
16     * Compute C = A x B
17     *   where A is a (m x k) matrix
18     *   where B is a (k x n) matrix
19     *   where C is a (m x n) matrix
20     *
21     * Use shared memory for tiling
22     *
23     *****/
24
25     // INSERT KERNEL CODE HERE
26
27     // Calculate position of row and col
28     int row = blockIdx.y*blockDim.y+threadIdx.y;
29     int col = blockIdx.x*blockDim.x+threadIdx.x;
30
31     // Make sure it is inside range
32     if (row<m && col<n) {
33         float val=0;
34         for(int i=0;i<k;i++) {
35             val += A[row*k+i]*B[i*n+col];
36         }
37         C[row*n+col] = val;
38     }
39 }
40
41 __global__ void mysgemm_tiled(int m, int n, int k,
42                             const float *A, const float *B, float *C) {
43     __shared__ float ds_M[TILE_SIZE][TILE_SIZE];
44     __shared__ float ds_N[TILE_SIZE][TILE_SIZE];
45
46     int bx = blockIdx.x; int by = blockIdx.y;
47     int tx = threadIdx.x; int ty = threadIdx.y;
48
49     int row = by*blockDim.y+ty;
50     int col = bx*blockDim.x+tx;
51
52     float val = 0;
53
54     for(int i=0;i<(k-1)/TILE_SIZE+1;++i) {
55         if(row<m && TILE_SIZE*i+tx<k) {
56             ds_M[ty][tx] = A[row*k+tx+i*TILE_SIZE];
57         } else {
58             ds_M[ty][tx] = 0;
59         }
60         if(col<n && TILE_SIZE*i+ty<k) {
61             ds_N[ty][tx] = B[(ty+i*TILE_SIZE)*n+col];
62         } else {
63             ds_N[ty][tx] = 0;
64         }
65
66         __syncthreads();
67
68         if(row<m && col<n)
69             for(int j=0;j<TILE_SIZE;++j)
70                 val+=ds_M[ty][j]*ds_N[j][tx];
71
72         __syncthreads();
73     }
74
75     if(row<m && col<n)
76         C[row*n+col] = val;

```

```

77 }
78
79 void basicSgemm(char transa, char transb, int m, int n, int k,
80               float alpha, const float *A, int lda,
81               const float *B, int ldb, float beta,
82               float *C, int ldc)
83 {
84     if ((transa != 'N') && (transa != 'n')) {
85         printf("unsupported value of 'transa'\n");
86         return;
87     }
88
89     if ((transb != 'N') && (transb != 'n')) {
90         printf("unsupported value of 'transb'\n");
91         return;
92     }
93
94     if ((alpha - 1.0f > 1e-10) || (alpha - 1.0f < -1e-10)) {
95         printf("unsupported value of alpha\n");
96         return;
97     }
98
99     if ((beta - 0.0f > 1e-10) || (beta - 0.0f < -1e-10)) {
100         printf("unsupported value of beta\n");
101         return;
102     }
103
104     // Initialize thread block and kernel grid dimensions
105
106     const unsigned int BLOCK_SIZE = TILE_SIZE;
107
108     //INSERT CODE HERE
109     dim3 dimGrid((n-1)/BLOCK_SIZE+1, (m-1)/BLOCK_SIZE+1, 1);
110     dim3 dimBlock(BLOCK_SIZE, BLOCK_SIZE, 1);
111
112     // Invoke CUDA kernel
113     //INSERT CODE HERE
114     mysgemm_tiled<<<dimGrid, dimBlock>>>(m, n, k, A, B, C);
115 }

```

4 Program Output

Below is the program output. This one is with zero parameters so it defaults to m=n=k=1000.

```

./sgemm-tiled

Setting up the problem...0.026839 s
  A: 1000 x 1000
  B: 1000 x 1000
  C: 1000 x 1000
Allocating device variables...0.086983 s
Copying data from host to device...0.003272 s
Launching kernel...0.007087 s
Copying data from device to host...0.002299 s
Verifying results...TEST PASSED

```