

Assignment 4 Report

By Khyateeswar Naidu Nalla
2019CS10376

Approach 1:

I first stored the input matrices in the form of blocks i.e had an array which contains the nonzero blocks continuously and every block here is indexed by another array. Now I sent these 4 matrices to gpu, here the multiplication is done such that each block calculates the final block of matrix so $n/m * n/m$ blocks are used for computation and each thread in that block distributes the multiplication into nonzero matrix block multiplication which adds to that block. And then we send the output matrix to cpu and we output to the output file. I used OPENMP for speeding up the output pushing to output file.

```
__global__ void matmul(int *a, int *b, int *c, int m, int u, int n, int *p1, int* p2, int* p3) {  
    //p1 has the index to access i,j block from a  
    //p2 has the index to access i,j block from b  
    //p3 contains information about non zero blocks  
    //c is the final matrix  
    int i = blockIdx.x/u;  
    int j = blockIdx.x%u;  
    int l = threadIdx.x;  
    int s = u/(blockDim.x)+1;  
    for(int k=s*l; k<s*(l+1); k++){  
        if(i<u && j<u && k<u){  
            if(p1[i*u+k]>=0 && p2[k*u+j]>=0){  
                atomicOr(p3+i*u+j,1);  
                int ik = p1[i*u+k];  
                int kj = p2[k*u+j];  
                for(int p=0; p<m; p++){  
                    for(int q=0; q<m; q++){  
                        for(int z=0; z<m; z++){  
                            atomicAdd(c+(i*m+p)*n+j*m+q, b[ kj*m+m+z*m+q]*a[ ik*m+m+p*m+z ] );  
                        }  
                    }  
                }  
            }  
        }  
    }  
}
```

Approach 2:

I made some changes to the above approach. Instead of the sending the whole output matrix from gpu ram to cpu ram. I used shared memory to calculate the final block corresponding that block and then converted the non zero blocks to char array and then send the char array to cpu ram where we can directly output it. This helps in reducing output time.

```
32 __global__ void matmul(int *a, int *b, char* cf, int m, int u, int n, int *p1, int *p2, int *nz, int* fnz) {
33 //p1 has the index to access i,j block from a
34 //p2 has the index to access i,j block from b
35 // nz contains the no of nonzero blocks
36 //cf is the char array to send to output
37
38 __shared__ int c[64]; // array for storing the block computed by this block
39 for(int h=0; h<64; h++){
40     c[h]=0;
41 }
42 __shared__ int p3;
43 p3=0;
44
45 __syncthreads();
46
47 int i = blockIdx.x/u;
48 int j = blockIdx.x%u;
49 int l = threadIdx.x;
50 int s = u/(blockDim.x)+1;
51
52 for(int k=s*l; k<s*l+s; k++){
53
54     if(i<u && j<u && k<u){
55
56         if(p1[i*u+k]>=0 && p2[k*u+j]>=0){
57             atomicOr(&p3, 1);
58             int ik = p1[i*u+k];
59             int kj = p2[k*u+j];
60             for(int p=0; p<m; p++){
61                 for(int q=0; q<m; q++){
62                     for(int z=0; z<m; z++){
63
64                         atomicAdd( &c[p+m*q], b[ kj*m+m+z+m*q]*a[ ik*m+m+p+m*z ] );
65
66                     }
67                 }
68             }
69         }
70     }
71 }
72
73
74 }
75
76
77 __syncthreads();
78
79
80 if(l==0 && i<u && j<u){
81     if( p3>0 ){
82         p3 = atomicAdd(nz, 1);
83         int off = (p3-1)*(4*m*m+8)+12;
84         conv(i, cf, off);
85         conv(j, cf, off+4);
86         int sd=0;
87         atomicAdd(fnz, 1);
88         for(int a=0; a<m; a++){
89             for(int b=0; b<m; b++){
90                 conv(c[a*m+b], cf, off+8+sd);
91                 sd=sd+4;
92             }
93         }
94     }
95     if(i==0 && j==0){
96         conv(n, cf, 0);
97         conv(m, cf, 4);
98     }
99 }
100
101
102
103 }
```

N=30000

M=8

Varying k

K(No of nonzero blocks of each matrix)	Time taken(s)
20000	6
50000	46
100000	61
500000	90
1000000	