

CS 406 Project Final Report

Efe Şencan, Gökberk Yar, Ahmet Enes Deveci
Alper Giray, Zeynep Özge Ergin

28 April 2021

Abstract

Depth first search (DFS) is a common algorithm that is used to traverse graph data structures. The algorithm starts by selecting an arbitrary root node as a starting point and it keeps traversing nodes starting from the root node as deep as possible until some stopping conditions are met. The goal of this project is to count the number of cycles ($k = 3, 4, 5$) in a cyclic undirected graph using the DFS algorithm as efficiently as possible by utilizing parallelization. For this purpose, we implemented several algorithms using OpenMP and CUDA and experimented with various parallelization methods on high performance computing (HPC) clusters. Based on our observations, we improved the sequential version of the implementation significantly both on CPU and GPU. In Amazon dataset, we improved sequential CPU runtime for $k=3$ case by a factor of 19 and in Dblp dataset we got more than 6 times speed up in $k=4$.

Keywords— Openmp, Cuda, Multicore, Cycle Count

Contents

1	Introduction	3
2	Data Description and Preprocess	3
3	General DFS	5
4	Methods and Results	5
4.1	Conversion of inputs to CSR format	5
4.2	Multi-core CPU implementation with OpenMP	6
4.3	GPU Parallelization	7
4.3.1	Single GPU Recursive Implementation	7
4.3.2	Single GPU Non-Recursive Implementation	8
4.3.3	Multi-GPU Non-Recursive Implementation	9
4.4	Combining GPU and CPU	10
4.4.1	Multi-GPU + CPU Non-Recursive Implementation	10
4.4.2	Multi-GPU + CPU Non-Recursive Dynamic Workqueue Implementation	11
5	How to reproduce results	12
6	Future Work	12
7	Appendix	13
7.1	File Reading (Spare Matrix Creation)	13
7.2	OpenMP Implementation	13
7.3	Recursive Cuda (Single GPU) Kernel	14
7.4	Non-recursive Cuda	15
7.5	Multi GPU (Non-recursive) Kernel	17
7.6	Multi GPU (Non-recursive)+ CPU Kernel	19
7.7	Multi GPU + CPU + Dynamic Workload	22

1 Introduction

Un-directed graph implies $u, v \in V$ and $(u, v) \in E \implies (v, u) \in E$. A k -cycle in G defined as $p = (v_0, v_1, v_2 \dots v_{k-1}, v_k)$ where $v_0 = v_k$ and $(v_i, v_{i+1}) \in E$ for $i = 0, 1, \dots, k$ and $v_i \neq v_j$ if $i \neq j$ for $i = 0, 1, \dots, k$. Given an un-directed **graph** G and a **number** k where $2 < k < 6$. $G = (V, E)$, where V is **vertices set**, and E is the **edges set**, our goal is to find the number of k -cycles for $\forall v \in V$.

Depth first search (DFS) is a common algorithm that is used to traverse graph data structures. The DFS algorithm contains many sub-tasks, in particular graph traversals which can be initialized with multiple processors. Since, one traversal starting from a particular root node j , does not depend on another traversal starting from j , DFS can be implemented in parallel. Moreover, as the size of the graph G grows, the runtime of the single thread approach can cost a significant amount of computation time (exponential time). For these reasons, HPC is applicable to DFS to reduce the runtime.

2 Data Description and Preprocess

A matrix is a two-dimensional data object made of m rows and n columns, therefore having total $m \times n$ values. If most of the elements of the matrix have 0 value, then it is called a sparse matrix. In this project, we have 2 sparse matrices that we can conduct tests on, named Dblp and Amazon. These matrices are taken as inputs where input format is like:

$$u_1 \ v_1$$

$$u_2 \ v_2$$

$$\dots$$

$$u_m \ v_m$$

$$\text{where } u_i < v_i \text{ and } (u_i, v_i) \in E, \forall i.$$

To be able to do the DFS, we need to have an adjacency matrix(list). So, we take the inputs and transform it into a sparse adjacency matrix. However, in both cases, we have two sparse matrices with Dblp having 425956 columns and rows, and Amazon having 548551 rows and columns. Since storing these

kinds of matrices consumes too much memory and is nearly impossible, we used a method called Compressed Sparse Row (CSR) representations. The CSR representation uses two arrays, `adj` and `xadj` for storing the column indices of the nonzeros within each row of the matrix adjacently. In CSR, the length of the array `adj` is equal to the number of nonzeros, and it keeps the column indices for the rows. The array `xadj` keeps the starting index for each row in `adj`. Hence, for a row `i`, the sub-array of column indices of nonzeros at row `i` starts with the entry `adj[xadj[i]]` and ends with the entry `adj[xadj[i+1]] - 1`. To simplify the implementations and make the previous statement correct for the last row, the length of `xadj` is set to `n + 1`, and the last element, `xadj[n]`, is set to the number of nonzeros. An example CSR representation for a toy sparse matrix is given in Figure 1.

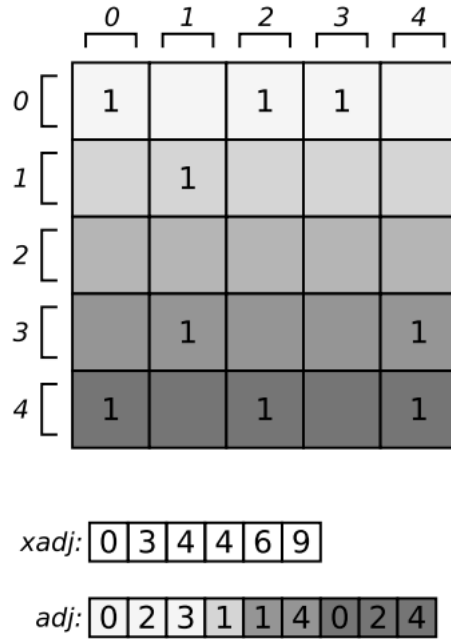


Figure 1: CSR Matrix Visualization

3 General DFS

There are approaches that utilize Breadth First Search (BFS) algorithm for detection of cycles in an undirected graph. However, this approach does not count number of cycles for each vertex. For this purpose(counting cycles of length k from a starting vertex i), Depth First Search(DFS) algorithm is more suitable.

The DFS algorithm starts by selecting an arbitrary root node as a starting point and it keeps traversing nodes starting from the root node as deep as possible until some stopping conditions are met. The essential part is, graphs may contain cycles and a particular node in a graph may be visited more than once. To avoid processing a node more than once, a boolean array is used to keep track of the visited nodes. Different from the regular DFS approach, since we are trying to find the cycles of length k , k -depth search is enough to determine whether the subgraph starting from node i is a cycle of length k . Therefore, we implemented DFS algorithm with k level nested for loops, since k is small.

Algorithm 1 Recursive_DFS(G, u):

Require: G , undirected graph**Require:** u , starting vertex

```
1: if  $visited[u] = \text{true}$  then  
2:   return  
3: end if  
4:  $visited[u] \leftarrow \text{true}$   
5: for  $v \in G[u].neighbors$  do  
   Recursive_DFS( $G, v$ )  
6: end for
```

4 Methods and Results

4.1 Conversion of inputs to CSR format

As it mentioned in the “Data Description & Preprocess” section, the graph that we provided as input is mostly sparse. In other words, most of the vertices do not have any neighbours. Therefore, it is more efficient to store matrices with CSR format to be memory efficient. To begin with our input data preprocessing method, we first read the input file line by line. For each vertex, we store it’s neighbours as key-value pairs by using `unordered_map` in C++. The major

reason for using unordered map is to perform find and insert operations in constant time. After filling the unordered map, the keys would be the vertices in our graph and the values are the corresponding neighbours of these vertices.

After reading the inputs, adj and xadj vectors must be filled correctly. For this purpose, there is a for loop that iterates number of vertices many times. First, the next element of the unordered map is taken. Then, we checked whether this element has connection to any vertex or not. If the element is not connected, its index of the xadj is filled with the previous element's value. If it is connected, we take the number of neighbors it has and fill the xadj with it. After filling the xadj, we started to visit every neighbor of the element and fill the adj array accordingly.

4.2 Multi-core CPU implementation with OpenMP

As it mentioned earlier in the “General DFS Algorithm”, our implementation was based on the DFS. For each vertex in the graph, we applied the DFS algorithm but only with depth k where k is the length of the cycle. When the algorithm reaches k th depth, it checks whether the starting vertex is the neighbor of it. If that is the case, the count for this vertex is incremented by one.

To parallelize this algorithm, we distributed the vertices to threads. We were calling the DFS function in a loop so we parallelized this loop. In each iteration, a thread arrives and takes an available vertex. After traversing and detecting the k -cycles, the thread takes another vertex until no available vertex is left. We used different scheduling options and guided options worked best as expected because the number of neighbours of each vertex is different and it leads to unbalanced workload. This unbalanced workload is more problematic in Dblp case, most of the CPU threads becomes idle and wait for the responsible thread which executes vertex with the most neighbours. In Table 1, efficiencies for Amazon is almost 1, our algorithm scales with number of processors. In Dblp, efficiency is less than Amazon since it is denser but speed up still keeps increasing with number of threads used.

OPENMP				
K=3	Thread Count	Duration	Speedup	Efficiency
amazon	1	0.603501	1	1
	2	0.310298	1.944907798	0.9724538992
	4	0.156428	3.858011353	0.9645028384
	8	0.0786051	7.677631604	0.9597039505
	16	0.0402396	14.99768884	0.9373555527
dblp	1	1.8533	1	1
	2	1.12785	1.643214967	0.8216074833
	4	0.659231	2.8113059	0.7028264751
	8	0.353933	5.236301786	0.6545377232
	16	0.178739	10.36874997	0.6480468728

K=4	Thread Count	Duration	Speedup	Efficiency
amazon	1	4.55978	1	1
	2	2.29525	1.986615837	0.9933079185
	4	1.19381	3.819519019	0.9548797547
	8	0.578946	7.876002252	0.9845002815
	16	0.287246	15.8741288	0.9921330497
dblp	1	66.5606	1	1
	2	34.1796	1.947377968	0.9736889841
	4	20.5717	3.235542031	0.8088855078
	8	12.5448	5.305831898	0.6632289873
	16	9.59927	6.9339231	0.4333701938

K=5	Thread Count	Duration	Speedup	Efficiency
amazon	1	43.1156	1	1
	2	21.9676	1.962690508	0.9813452539
	4	10.9534	3.936275494	0.9840688736
	8	5.42065	7.953953862	0.9942442327
	16	2.86459	15.05122897	0.9407018107
dblp	1	-	-	-
	2	-	-	-
	4	-	-	-
	8	-	-	-
	32	977	-	-

Table 1: Multicore OpenMP Implementation

4.3 GPU Parallelization

4.3.1 Single GPU Recursive Implementation

After implementing the sequential and OpenMP algorithms, to run these algorithms in the GPU, we tried to give each vertex to a thread. To do that, we

used a single dimensional grid and block. These threads are sent to a single kernel that gives each thread a startVertex. When a thread takes this startVertex, it starts to dive k-depth and checks if the kth vertex is a neighbor of the startVertex. For each cycle, the count is incremented and at the end, results are written a global array according to the global id of the thread.

One the other hand, there is a drawback of using recursion operations in CUDA. First of all, the stack size that is allocated for each recursive call is limited, and the programmer should be careful about not exceeding the memory limit. In addition, we realized that recursive calls become costly in CUDA which brings us inefficiency in our implementation. Hence, we abandoned the recursive implementation and decided to implement a non-recursive algorithm.

4.3.2 Single GPU Non-Recursive Implementation

The most common method for converting a recursive algorithm to a non recursive version is to use stack data structure and simulate the recursion by a first come last out manner. However, since the cycle lengths in our problem definition is limited within the range 3 to 5, instead of utilizing stack, we wrote nested for loops of k many times. We implemented three different kernels for each cycle length scenario. After determining the root node, we computed its neighbour by looking at the xadj array. Then, for every neighbour of the starting point, this time we computed the neighbours of that neighbour. This process continues until we reach the k^{th} neighbour of the root node. Finally, we checked whether the k^{th} neighbour is equal to the starting node. If they are equal, then this means that we find a cycle of length k. In Table 2, one can see that single GPU with non-recursive implementation produces 13 threads equivalent work for Amazon and nearly 4 threads worth work for Dblp case. When k=4, performance nearly dropped by 50 percent. When k = 5, performance dropped by 25 more percent on k=4 case for Amazon. Dblp requires too much time that's why it is not reported. As an observation GPU performs best when sparsity of matrix is large and performance drops dramatically as density of graph increases. Increasing k or using dblp are ways to increase density of graph. The reason why performance drops is most of the GPU threads become idle and wait for the bottleneck in the calculation (vertex with most neighbours).

Non Recursive		
K=3	Duration	Speedup
amazon	0.045381	13.29853904
dblp	0.487654	3.800440476

K=4	Duration	Speedup
amazon	0.514845	8.856607328
dblp	44.543335	1.494288652

K=5	Duration	Speedup
amazon	6.447344	6.687342881
dblp	-	-

Table 2: GPU Non-Recursive

4.3.3 Multi-GPU Non-Recursive Implementation

With Multi-GPU Non-Recursive implementation, we designed the algorithm such that it could run on up to 4 GPUs concurrently. In fact, it can work as much as GPUs that the computer allows and our testing environment has 4 GPUs available. The algorithm is essentially very similar to the single GPU Non-Recursive implementation. The appropriate kernel function is found by the main processor and every thread starts working on a vertex. Main difference of this approach is that since there are multiple GPUs, it first divides the workload among the GPUs. In other words, the total number of vertices was divided into the GPUs equally. GPU 1 takes the first forth of the data, GPU 2 takes the second forth data and continues. GPUs in Nebula are heterogenous, there are 2 faster and 2 slower GPUs. This method’s performance is directly affected by this fact. For run time, we took the slowest GPU among 4 GPUs. Our observation is generally other GPUs complete their work significantly early and become idle for a duration. This idle time shows there is still room for improvement with custom workloading the GPUs. For example giving more work to better GPUs or applying a preprocessing step for computation order of the vertices can improve performance. In Table 3, we achieve 19 threads worth speed up for the Amazon dataset. This is about 50 percent faster than single GPU implementation. This speed up pattern continues for k=4 and k=5.

Multigpu - Non Recursive		
K=3	Duration	Speedup
amazon	0.030755	19.62285807
dblp	0.49196	3.767176193

K=4	Duration	Speedup
amazon	0.372918	12.2272993
dblp	57.406826	1.159454452

K=5	Duration	Speedup
amazon	5.285642	8.157116959
dblp	-	-

Table 3: Multi-GPU Non-Recursive

4.4 Combining GPU and CPU

4.4.1 Multi-GPU + CPU Non-Recursive Implementation

With the inspiration we got from Multi-GPU Non-Recursive Implementation, we thought that we could run the code with multiple GPUs and CPUs together. In this implementation, we used 4 GPU and 28 CPU together. Similar to the multi-GPU case, it can work as much as GPUs and CPUs that the computer allows and our testing environment has 4 GPUs and 28 CPUs available. For this purpose, we created 32 threads initially and assigned the first 4 of them to the GPU and rest to the CPU. First, we divided the total work to the threads chunk by chunk and distributed these chunks to the processors equally. Based on our experiments, threads in the GPU finished their job, whereas threads in the CPU still have work to do. Since, the GPU has more available threads than the CPU, overall, GPU can handle more vertices for a particular time period. Hence, it has a better throughput compared to CPU. Therefore, we distributed more chunks to the GPU threads. After conducting some runs, we figured out that GPU threads can handle 10 times more vertices than CPU threads. Then, GPU threads found the appropriate kernel according to cycle length and started to work on this kernel like in the previous non-recursive implementations. On the other hand, CPU threads started to perform DFS like Multi-core CPU implementation with OpenMP. In Table 4, except the most sparse case Amazon dataset k=3, adding CPU threads improved the performance. For instance in Dblp k=3, speedup doubled compared to multigpu case. Also for Amazon k=4 and k=5 speedups are better when CPU is added on top.

Multigpu + CPU Non Recursive		
K=3	Duration	Speedup
amazon	0.037041	16.29278367
dblp	0.239933	7.724239684

K=4	Duration	Speedup
amazon	0.274205	16.62909137
dblp	22.658237	2.93758954

K=5	Duration	Speedup
amazon	4.201431	10.2621226
dblp	2200	-

Table 4: Multi-GPU + CPU Non-Recursive

4.4.2 Multi-GPU + CPU Non-Recursive Dynamic Workqueue Implementation

After running some tests on the Multi-GPU + CPU Non-Recursive Implementation, despite the GPU threads having 10 times more workload, there were still CPU and GPU threads waiting whereas a single thread works so much longer. To optimize the usage of the threads, we did not divide all the work into the threads initially. Instead, we had created chunks and calculated the total number of chunks that needed to be completed. After that, until we reached the total number of chunks, the available threads both in GPU and CPU came forward and took the next available chunk. We were aware that picking up the next available chunk caused a race condition. That's why we took those parts into the critical region. But this critical region created an overhead for threads and that caused the algorithm to run slower. A solution to this problem would be the incrementing the chunk size, since for larger chunks threads will spend more time on that, hence the threads would need less to get the next chunk. However, for larger chunk sizes, it would be similar to Multi-GPU + CPU Non-Recursive Implementation, a non-dynamic method, so there would be idle threads and that would also increase the run time. In Table 5, only Dplb k=4 case performed better and it is actually the best case among all GPU versions. However, speedup dropped under 1 for Amazon cases. For this method, parameters like chunk size and GPU multiplier are very important, one should find the best pair that is consistent for in all datasets and cycle lengths.

Multigpu + CPU Non Recursive Dynamic Workload		
K=3	Duration	Speedup
amazon	2.076866	0.2905825412
dblp	3.700613	0.50080892

K=4	Duration	Speedup
amazon	5.658017	0.8058971898
dblp	10.617888	6.268723121

K=5	Duration	Speedup
amazon	8.319901	5.182225125
dblp	-	-

Table 5: Multi-GPU + CPU Non-Recursive Dynamic Workload

5 How to reproduce results

A makefile is shared with the project. One can reproduce the results by running different commands in the makefile. For all testcase, there is run that only measures the run time and another run that produces the desired output.

6 Future Work

After observing GPU performs best when sparsity is large and idling increase when workload is imbalanced, one should preprocess the input and sort vertices accordingly their required work amount. One heuristic that can be used to estimate required work for a vertex is its neighbour count (one vertex could have less neighbours but its neighbours could potentially have many neighbours, so it is only a heuristic). Also, there is no need to re-compute same cycles for each vertex in the path, by using a better algorithm, shared memory and atomic memory accesses redundant calculation could be minimized. For GPU implementation, one could implement different kernels for different neighbour count vertices. For example, neighbour count space could be splitted into 3 : less than 32, 32 to 256 and more than 256. For less than 32 case, one can pack different vertices into single GPU warp and run together. For 32 to 256, one can use block level implementation and vertices more than 256 one can use grid level implementations to reduce idling.

7 Appendix

7.1 File Reading (Spare Matrix Creation)

```
1 void read_mtxbin(string fname, int k){
2     ifstream infile(fname);
3     int a, b;
4     int nnv = 0;
5     unordered_map<int, vector<int> > hashmap;
6
7     int maxElement = -1;
8     while (infile >> a >> b)
9     {
10         nnv+=2;
11         hashmap[a].push_back(b);
12         hashmap[b].push_back(a);
13
14         if(b > maxElement){
15             maxElement = b;
16         }
17     }
18
19     int nov = maxElement + 1;
20
21     int * adj = new int[nnv];
22     int * xadj = new int[nov+1];
23     xadj[0]=0;
24
25     int j = 0;
26     int maxSize = -1;
27
28     for(int i=0; i < nov ; i++){
29         auto current = hashmap.find(i);
30         if (current == hashmap.end()){
31             xadj[i+1] = xadj[i];
32         }
33         else{
34             int size = current->second.size();
35             maxSize = max(size,maxSize);
36
37             xadj[i+1] = xadj[i] + size;
38             for(auto val : current->second) {
39                 adj[j] = val;
40                 j++;
41             }
42         }
43     }
44
45     wrapper(xadj,adj,k,nov,nnv);
46 }
47
48 }
```

7.2 OpenMP Implementation

```
1 void DFS_sparse(int xadj[], int adj[], bool marked[], int n,
2     int vert, int start, int &count)
3     //vert: bulunduđu konum //start: baslangıc noktası
4 {
5     marked[vert] = true;
6     int start_index = xadj[vert];
7     int path_length = xadj[vert+1];
8     if (n == 0){
9         marked[vert] = false;
10         for(int i = start_index; i < path_length; i++){
11             if(adj[i] == start){
12                 count++;
13                 break;
14             }
15         }
16     }
17 }
```

```

13     }
14     }
15     return;
16 }
17 for(int i=start_index; i < path_length; i++){
18     if(!marked[adj[i]]){
19         DFS_sparse(xadj, adj, marked, n-1, adj[i], start, count);
20     }
21 }
22 marked[vert] = false;
23 }
24
25 void countCycles_sparse(int *xadj, int *adj, int n, int nov)
26 {
27     double start, end;
28     start = omp_get_wtime();
29     int *arr = new int[nov];
30
31     #pragma omp parallel num_threads(16)
32     {
33         bool *marked = new bool[nov];
34         memset(marked, false, nov * sizeof(bool)); // bu belki silinebilir
35
36         #pragma omp for schedule(guided)
37         for(int i = 0; i < nov; i++){
38             int localcount = 0;
39             DFS_sparse(xadj, adj, marked, n - 1, i, i, localcount);
40             arr[i] = localcount;
41         }
42     }
43     end = omp_get_wtime();
44 }

```

7.3 Recursive Cuda (Single GPU) Kernel

```

1  __device__ bool check(int marked[], int round, int val){
2  for(int i = 0; i < round; i++){
3      if(marked[i] == val){return false;}
4  }
5  return true;
6  }
7
8  __device__ void DFS_sparse(int xadj[], int adj[], int marked[], int n,
9      int vert, int start, int &count, int round)
10     //vert: bulunduğu konum //start: başlangıç noktası
11 {
12     marked[round] = vert;
13     int start_index = xadj[vert];
14     int path_length = xadj[vert+1];
15     if (n == 0){
16         marked[round] = -1;
17         for(int i = start_index; i < path_length; i++){
18             if(adj[i] == start){
19                 count++;
20                 break;
21             }
22         }
23         return;
24     }
25     for(int i=start_index; i < path_length; i++){
26         if(check(marked, round, adj[i])){
27             DFS_sparse(xadj, adj, marked, n-1, adj[i], start, count, round +
28                 1);
29         }
30     }
31     marked[round] = -1;
32 }
33
34 __global__ void kernel(int* adj, int* xadj, int* output, int n, int nov){

```

```

35     int index = threadIdx.x + (blockIdx.x * blockDim.x);
36     __shared__ int marked[THREADS_PER_BLOCK][10];
37     if(index < nov){
38         //int *marked = new int[n];
39         //memset(marked, -1, n * sizeof(int)); // bu belki silinebilir
40         int localcount = 0;
41         int round = 0;
42         DFS_sparse(xadj, adj, marked[threadIdx.x], n - 1, index, index,
43                 localcount, round);
44         output[index] = localcount;
45     }
46 }
47 void wrapper(int *xadj, int *adj, int n, int nov, int nnz){
48     cudaSetDevice(0);
49     int *adj_d;
50     int *xadj_d;
51     int *output_d;
52     int *output_h = new int[nov];
53     int numBlock = (nov + THREADS_PER_BLOCK - 1) / THREADS_PER_BLOCK;
54     cudaEvent_t start, stop;
55     float elapsedTime;
56
57     gpuErrchk(cudaMalloc((void**)&adj_d, (nnz) * sizeof(int)));
58     gpuErrchk(cudaMalloc((void**)&xadj_d, (nov + 1) * sizeof(int)));
59
60     gpuErrchk(cudaMalloc((void**)&output_d, (nov) * sizeof(int)));
61     //gpuErrchk(cudaMallocHost((void **)&output_h, (nov) * sizeof(int)));
62     gpuErrchk(cudaMemcpy(adj_d, adj, (nnz) * sizeof(int),
63                         cudaMemcpyHostToDevice));
64     gpuErrchk(cudaMemcpy(xadj_d, xadj, (nov + 1) * sizeof(int),
65                         cudaMemcpyHostToDevice));
66
67     cudaEventCreate(&start);
68     cudaEventRecord(start, 0);
69
70     kernel<<<numBlock, THREADS_PER_BLOCK>>>(adj_d, xadj_d, output_d, n, nov);
71
72     gpuErrchk(cudaDeviceSynchronize());
73
74     gpuErrchk(cudaMemcpy(output_h, output_d, (nov) * sizeof(int),
75                         cudaMemcpyDeviceToHost));
76     //printArray(output_h, nov);
77     cudaEventCreate(&stop);
78     cudaEventRecord(stop, 0);
79     cudaEventSynchronize(stop);
80
81     cudaEventElapsedTime(&elapsedTime, start, stop);
82     //printf("GPU scale took: %f s\n", elapsedTime/1000);
83
84     cudaFree(adj_d);
85     cudaFree(xadj_d);
86 }

```

7.4 Non-recursive Cuda

```

1  __global__ void kernel3(int* adj, int* xadj, int* output, int nov){
2      int index = threadIdx.x + (blockIdx.x * blockDim.x);
3      if(index < nov){
4          //int *marked = new int[n];
5          //memset(marked, -1, n * sizeof(int)); // bu belki silinebilir
6          int localcount = 0;
7          // int round = 0;
8          // 0==>
9          int s0 = xadj[index];
10         int e0 = xadj[index+1];
11
12         for(int i=s0; i < e0; i++){
13

```

```

14 // 0 --> 1
15
16 int neighbour_1 = adj[i];
17 int s1 = xadj[neighbour_1];
18 int e1 = xadj[neighbour_1+1];
19
20 for(int j=s1; j < e1; j++){
21 // 0 --> 1 --> 2
22
23 int neighbour_2 = adj[j];
24 if (neighbour_2 == index) continue;
25 int s2 = xadj[neighbour_2];
26 int e2 = xadj[neighbour_2+1];
27
28 for(int k=s2; k < e2; k++){
29 // 0 --> 1 --> 2 --> 3
30 int neighbour_3 = adj[k];
31 if (neighbour_3 == index){
32 localcount+=1;
33 break;
34 }
35 }
36 }
37 }
38 output[index] = localcount;
39 }
40 }
41 }
42
43 void wrapper(int *xadj, int *adj, int n, int nov, int nnz){
44
45 // int X = nov;
46 // int Y = maxSize;
47 // int Z = maxSize;
48
49 dim3 threadsPerBlock(8, 8, 8);
50 dim3 numBlocks(X/threadsPerBlock.x, /* for instance 512/8 = 64 */
51 Y/threadsPerBlock.y,
52 Z/threadsPerBlock.z);
53
54
55
56
57
58
59
60
61 cudaSetDevice(0);
62 int *adj_d;
63 int *xadj_d;
64 int *output_d;
65 int *output_h = new int[nov];
66 int numBlock = (nov + THREADS_PER_BLOCK - 1) / THREADS_PER_BLOCK;
67 cudaEvent_t start, stop;
68 float elapsedTime;
69
70 gpuErrchk(cudaMalloc((void**)&adj_d, (nnz) * sizeof(int)));
71 gpuErrchk(cudaMalloc((void**)&xadj_d, (nov + 1) * sizeof(int)));
72 gpuErrchk(cudaMalloc((void**)&output_d, (nov) * sizeof(int)));
73 //gpuErrchk(cudaMallocHost((void **)&output_h, (nov) * sizeof(int)));
74
75 gpuErrchk(cudaMemcpy(adj_d, adj, (nnz) * sizeof(int),
76 cudaMemcpyHostToDevice));
77
78 gpuErrchk(cudaMemcpy(xadj_d, xadj, (nov + 1) * sizeof(int),
79 cudaMemcpyHostToDevice));
80
81 cudaEventCreate(&start);
82 cudaEventRecord(start, 0);
83
84 kernel5<<<numBlock, THREADS_PER_BLOCK>>>(adj_d, xadj_d, output_d, nov);
85
86 //combination<<<numBlocks, threadsPerBlock>>>(adj_d, xadj_d, output_d, n, nov);
87
88 gpuErrchk(cudaDeviceSynchronize());
89
90 gpuErrchk(cudaMemcpy(output_h, output_d, (nov) * sizeof(int),
91 cudaMemcpyDeviceToHost));
92
93 cudaEventCreate(&stop);
94 cudaEventRecord(stop, 0);
95 cudaEventSynchronize(stop);
96
97 cudaEventElapsedTime(&elapsedTime, start, stop);
98 // printf("GPU scale took: %f s\n", elapsedTime/1000);

```



```

96     printArray(output_h,nov);
97     cudaFree(adj_d);
98     cudaFree(xadj_d);
99 }

```

7.5 Multi GPU (Non-recursive) Kernel

```

1  void DFS_sparse(int xadj[], int adj[], bool marked[], int n,
2      int vert, int start, int &count)
3      //vert: bulundugu konum //start: baslangic noktası
4  {
5      marked[vert] = true;
6      int start_index = xadj[vert];
7      int path_length = xadj[vert+1];
8      if (n == 0){
9          marked[vert] = false;
10         for(int i = start_index; i < path_length; i++){
11             if(adj[i] == start){
12                 count++;
13                 break;
14             }
15         }
16         return;
17     }
18 }
19 // if(path_length-start_index <=1) return;
20
21 for(int i=start_index; i < path_length; i++){
22     if(!marked[adj[i]]){
23         DFS_sparse(xadj, adj, marked, n-1, adj[i], start, count);
24     }
25 }
26 marked[vert] = false;
27 }
28
29 __global__ void kernel3(int* adj, int* xadj, int* output, int nov, int
30     novStart){
31     int index = novStart + threadIdx.x + (blockIdx.x * blockDim.x);
32     if(index < nov){
33         // if(index ==0)printf("called gpu \n");
34         //int *marked = new int[n];
35         //memset(marked, -1, n * sizeof(int)); // bu belki silinebilir
36         int localcount = 0;
37         // int round = 0;
38
39         // 0-->
40         int s0 = xadj[index];
41         int e0 = xadj[index+1];
42
43         for(int i=s0; i < e0; i++){
44             // 0 --> 1
45
46             int neighbour_1 = adj[i];
47             int s1 = xadj[neighbour_1];
48             int e1 = xadj[neighbour_1+1];
49
50             for(int j=s1; j < e1; j++){
51                 // 0 --> 1 --> 2
52
53                 int neighbour_2 = adj[j];
54                 if (neighbour_2 == index) continue;
55                 int s2 = xadj[neighbour_2];
56                 int e2 = xadj[neighbour_2+1];
57
58                 for(int k=s2; k < e2; k++){
59                     // 0 --> 1 --> 2 --> 3
60                     int neighbour_3 = adj[k];
61                     if (neighbour_3 == index){
62                         localcount+=1;
63                     }

```

```

64         break;
65     }
66 }
67 }
68 }
69     output[index-novStart] = localcount;
70 }
71 }
72 }
73 }
74 void wrapper(int *xadj, int *adj, int n, int nov, int nnz){
75     int *output_h = new int[nov];
76     double start_cpu, end_cpu;
77     start_cpu = omp_get_wtime();
78
79     #pragma omp parallel num_threads(PARALEL_THREAD_COUNT)
80     {
81         int threadId=omp_get_thread_num ();
82         // cout<< threadId<<endl;
83
84         int virtual_thread_count = GPU_MULTIPLIER *4 + PARALEL_CPU;
85         int novForThread = (nov+virtual_thread_count-1)/virtual_thread_count;
86
87         if(threadId <=3)
88         {
89             int novStart = GPU_MULTIPLIER * novForThread * threadId;
90             int novEnd = GPU_MULTIPLIER * novForThread * (threadId+1);
91             if (novEnd > nov) novEnd = nov;
92             int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) /
93                 THREADS_PER_BLOCK;
94             // printf("nov s %d e %d \n", novStart,novEnd);
95
96             cudaSetDevice(threadId);
97
98             int *adj_d;
99             int *xadj_d;
100             int *output_d;
101             cudaEvent_t start, stop;
102             float elapsedTime;
103
104             gpuErrchk(cudaMalloc((void**)&adj_d, (nnz) * sizeof(int)));
105             gpuErrchk(cudaMalloc((void**)&xadj_d, (nov + 1) * sizeof(int)));
106             gpuErrchk(cudaMalloc((void**)&output_d, (novEnd-novStart) *
107                 sizeof(int)));
108             //gpuErrchk(cudaMallocHost((void **)&output_h, (nov) * sizeof(int)));
109
110             gpuErrchk(cudaMemcpy(adj_d, adj, (nnz) * sizeof(int),
111                 cudaMemcpyHostToDevice));
112             gpuErrchk(cudaMemcpy(xadj_d, xadj, (nov + 1) * sizeof(int),
113                 cudaMemcpyHostToDevice));
114
115             cudaEventCreate(&start);
116             cudaEventRecord(start, 0);
117             double start_gpu = omp_get_wtime();
118             cudaStream_t stream1;
119             cudaStreamCreate (&stream1);
120
121             // printf("threadId entry to kernel %d GPU \n", threadId);
122             if (n==3)kernel3<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
123                 xadj_d, output_d, novEnd,novStart);
124             else if (n==4)kernel4<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
125                 xadj_d, output_d, novEnd,novStart);
126             else if (n==5)kernel5<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
127                 xadj_d, output_d, novEnd,novStart);
128
129             //combination<<<numBlocks, threadsPerBlock>>>(adj_d, xadj_d, output_d, n, nov);
130             // printf("threadId exit to kernel %d GPU \n", threadId);
131             double end_gpu = omp_get_wtime();
132
133             gpuErrchk(cudaDeviceSynchronize());
134             cudaEventCreate(&stop);
135             cudaEventRecord(stop, 0);
136             cudaEventSynchronize(stop);
137             cudaEventElapsedTime(&elapsedTime, start, stop);
138         }
139     }
140     end_cpu = omp_get_wtime();

```

```

141     printf("GPU scale took: %f s on gpu  %d \n", elapsedTime/1000, threadId);
142
143     gpuErrchk(cudaMemcpy(output_h+novStart, output_d, (novEnd-novStart) *
144                       sizeof(int), cudaMemcpyDeviceToHost));
145     cudaFree(adj_d);
146     cudaFree(xadj_d);
147
148 }
149 else{
150     // printf("Entered \n");
151
152     // int novForThread = (nov+PARALEL_THREAD_COUNT-1)/PARALEL_THREAD_COUNT;
153     // int novStart = novForThread * threadId;
154     // int novEnd = novForThread * (threadId+1);
155     // if (novEnd > nov) novEnd = nov;
156
157     // int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) / THREADS_PER_BLOCK;
158     int novStart = 4 * GPU_MULTIPLIER*novForThread + 1 * novForThread *
159               (threadId-4);
160     int novEnd = novStart + 1* novForThread ;
161     if (novEnd > nov) novEnd = nov;
162
163     bool *marked = new bool[nov];
164     memset(marked, false, nov * sizeof(bool)); // bu belki silinebilir
165
166     double start_thread = omp_get_wtime();
167     for(int i = novStart; i < novEnd; i++){
168         int localcount = 0;
169         DFS_sparse(xadj, adj, marked, n - 1, i, i, localcount);
170         output_h[i] = localcount;
171     }
172
173     double end_thread = omp_get_wtime();
174     printf("Took %f secs \n", end_thread -start_thread );
175
176 }
177
178 }
179
180 }
181
182 end_cpu = omp_get_wtime();
183
184 }

```

7.6 Multi GPU (Non-recursive)+ CPU Kernel

```

1 void DFS_sparse(int xadj[], int adj[], bool marked[], int n,
2               int vert, int start, int &count)
3 {
4     //vert: bulundugu konum //start: baslangic noktasi
5     marked[vert] = true;
6     int start_index = xadj[vert];
7     int path_length = xadj[vert+1];
8     if (n == 0){
9         marked[vert] = false;
10        for(int i = start_index; i < path_length; i++){
11            if(adj[i] == start){
12                count++;
13                break;
14            }
15        }
16        return;
17    }
18
19    // if(path_length-start_index <=1) return;
20
21    for(int i=start_index; i < path_length; i++){
22        if(!marked[adj[i]]){
23            DFS_sparse(xadj, adj, marked, n-1, adj[i], start, count);
24        }
25    }
26 }

```

```

25     }
26     marked[vert] = false;
27 }
28
29
30
31 __global__ void kernel3(int* adj, int* xadj, int* output, int nov, int
    novStart){
32
33     int index = novStart + threadIdx.x + (blockIdx.x * blockDim.x);
34     if(index < nov){
35         // if(index ==0)printf("called gpu \n");
36         //int *marked = new int[n];
37         //memset(marked, -1, n * sizeof(int)); // bu belki silinebilir
38         int localcount = 0;
39         // int round = 0;
40
41         // 0-->
42         int s0 = xadj[index];
43         int e0 = xadj[index+1];
44
45         for(int i=s0; i < e0; i++){
46             // 0 --> 1
47
48             int neighbour_1 = adj[i];
49             int s1 = xadj[neighbour_1];
50             int e1 = xadj[neighbour_1+1];
51
52             for(int j=s1; j < e1; j++){
53                 // 0 --> 1 --> 2
54
55                 int neighbour_2 = adj[j];
56                 if (neighbour_2 == index) continue;
57                 int s2 = xadj[neighbour_2];
58                 int e2 = xadj[neighbour_2+1];
59
60                 for(int k=s2; k < e2; k++){
61                     // 0 --> 1 --> 2 --> 3
62                     int neighbour_3 = adj[k];
63                     if (neighbour_3 == index){
64                         localcount+=1;
65                         break;
66                     }
67                 }
68             }
69         }
70         output[index-novStart] = localcount;
71     }
72 }
73
74
75
76 void wrapper(int *xadj, int *adj, int n, int nov, int nnz){
77
78     int *output_h = new int[nov];
79
80     double start_cpu, end_cpu;
81     start_cpu = omp_get_wtime();
82
83
84     #pragma omp parallel num_threads(PARALEL_THREAD_COUNT)
85     {
86
87         int threadId=omp_get_thread_num ();
88         // cout<< threadId<<endl;
89
90         int virtual_thread_count = GPU_MULTIPLIER *4 + PARALEL_CPU;
91         int novForThread = (nov+virtual_thread_count-1)/virtual_thread_count;
92
93
94
95         if(threadId <=3)
96         {
97             int novStart = GPU_MULTIPLIER * novForThread * threadId;
98             int novEnd = GPU_MULTIPLIER * novForThread * (threadId+1);
99             if (novEnd > nov) novEnd = nov;
100             int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) /
                THREADS_PER_BLOCK;
101             // printf("nov s %d e %d \n", novStart,novEnd);
102
103
104             cudaSetDevice(threadId);
105
106             int *adj_d;
107             int *xadj_d;
108             int *output_d;
109             cudaEvent_t start, stop;

```

```

110 float elapsedTime;
111
112 gpuErrchk(cudaMalloc((void**)&adj_d, (nnz) * sizeof(int)));
113 gpuErrchk(cudaMalloc((void**)&xadj_d, (nov + 1) * sizeof(int)));
114
115 gpuErrchk(cudaMalloc((void**)&output_d, (novEnd-novStart) *
116     sizeof(int)));
117 //gpuErrchk(cudaMallocHost((void **)&output_h, (nov) * sizeof(int)));
118
119 gpuErrchk(cudaMemcpy(adj_d, adj, (nnz) * sizeof(int),
120     cudaMemcpyHostToDevice));
121 gpuErrchk(cudaMemcpy(xadj_d, xadj, (nov + 1) * sizeof(int),
122     cudaMemcpyHostToDevice));
123
124 cudaEventCreate(&start);
125 cudaEventRecord(start, 0);
126 double start_gpu = omp_get_wtime();
127 cudaStream_t stream1;
128 cudaStreamCreate (&stream1);
129 // printf("threadId entry to kernel %d GPU \n", threadId);
130 if (n==3)kernel3<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
131     xadj_d, output_d, novEnd,novStart);
132 else if (n==4)kernel4<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
133     xadj_d, output_d, novEnd,novStart);
134 else if (n==5)kernel5<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
135     xadj_d, output_d, novEnd,novStart);
136 //combination<<<numBlocks, threadsPerBlock>>>(adj_d, xadj_d, output_d, n, nov);
137 // printf("threadId exit to kernel %d GPU \n", threadId);
138 double end_gpu = omp_get_wtime();
139
140 gpuErrchk(cudaDeviceSynchronize());
141 cudaEventCreate(&stop);
142 cudaEventRecord(stop, 0);
143 cudaEventSynchronize(stop);
144 cudaEventElapsedTime(&elapsedTime, start, stop);
145 printf("GPU scale took: %f s on gpu %d \n", elapsedTime/1000, threadId);
146
147 gpuErrchk(cudaMemcpy(output_h+novStart, output_d, (novEnd-novStart) *
148     sizeof(int), cudaMemcpyDeviceToHost));
149 cudaFree(adj_d);
150 cudaFree(xadj_d);
151
152 }
153 else{
154     // printf("Entered \n");
155     // int novForThread = (nov+PARALEL_THREAD_COUNT-1)/PARALEL_THREAD_COUNT;
156     // int novStart = novForThread * threadId;
157     // int novEnd = novForThread * (threadId+1);
158     // if (novEnd> nov) novEnd = nov;
159     // int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) / THREADS_PER_BLOCK;
160     int novStart = 4 * GPU_MULTIPLIER*novForThread + 1 * novForThread *
161         (threadId-4);
162     int novEnd = novStart + 1* novForThread ;
163     if (novEnd> nov) novEnd = nov;
164     // printf("nov s %d e %d -cpu \n", novStart,novEnd);
165
166     bool *marked = new bool[nov];
167     memset(marked, false, nov * sizeof(bool)); // bu belki silinebilir
168
169     double start_thread = omp_get_wtime();
170     for(int i = novStart; i < novEnd; i++){
171         int localcount = 0;
172         DFS_sparse(xadj, adj, marked, n - 1, i, i, localcount);
173         output_h[i] = localcount;
174     }
175
176     double end_thread = omp_get_wtime();
177     printf("Took %f secs \n", end_thread -start_thread );
178
179 }
180
181 }
182
183

```

```

184     }
185     end_cpu = omp_get_wtime();
186
187     // printf("Took %f secs \n", end_cpu - start_cpu);
188     // printArray(output_h, nov);
189 }

```

7.7 Multi GPU + CPU + Dynamic Workload

```

1  __host__ void DFS_sparse(int xadj[], int adj[], bool marked[], int n,
2  int vert, int start, int &count)
3      //vert: bulunduđu konum //start: baslangıc noktası
4  {
5      marked[vert] = true;
6      int start_index = xadj[vert];
7      int path_length = xadj[vert+1];
8      if (n == 0){
9          marked[vert] = false;
10         for(int i = start_index; i < path_length; i++){
11             if(adj[i] == start){
12                 count++;
13                 break;
14             }
15         }
16         return;
17     }
18
19     // if(path_length-start_index <=1) return;
20
21     for(int i=start_index; i < path_length; i++){
22         if(!marked[adj[i]]){
23             DFS_sparse(xadj, adj, marked, n-1, adj[i], start, count);
24         }
25     }
26     marked[vert] = false;
27 }
28
29
30 __global__ void kernel3(int* adj, int* xadj, int* output, int nov, int
31 novStart){
32     int index = novStart + threadIdx.x + (blockIdx.x * blockDim.x);
33     if(index < nov){
34         // if(index ==0)printf("called gpu \n");
35         // int *marked = new int[n];
36         // memset(marked, -1, n * sizeof(int)); // bu belki silinebilir
37         int localcount = 0;
38         // int round = 0;
39
40         // 0-->
41         int s0 = xadj[index];
42         int e0 = xadj[index+1];
43
44         for(int i=s0; i < e0; i++){
45             // 0 --> 1
46
47             int neighbour_1 = adj[i];
48             int s1 = xadj[neighbour_1];
49             int e1 = xadj[neighbour_1+1];
50
51             for(int j=s1; j < e1; j++){
52                 // 0 --> 1 --> 2
53
54                 int neighbour_2 = adj[j];
55                 if (neighbour_2 == index) continue;
56                 int s2 = xadj[neighbour_2];
57                 int e2 = xadj[neighbour_2+1];
58
59                 for(int k=s2; k < e2; k++){
60                     // 0 --> 1 --> 2 --> 3
61                     int neighbour_3 = adj[k];
62

```

```

63         if (neighbour_3 == index){
64             localcount+=1;
65             break;
66         }
67     }
68 }
69 }
70 output[index-novStart] = localcount;
71 }
72 }
73
74 void wrapper(int *xadj, int *adj, int n, int nov, int nnz){
75     int *output_h = new int[nov];
76     double start_cpu, end_cpu;
77     start_cpu = omp_get_wtime();
78
79     int totalChunk = (nov + CHUNK_SIZE_OUR -1) /CHUNK_SIZE_OUR;
80     int currentChunk = 0; //mask accesed atomicly
81
82     #pragma omp parallel num_threads(PARALEL_THREAD_COUNT)
83     {
84
85         int threadId=omp_get_thread_num ();
86         // cout<< threadId<<endl;
87         // int virtual_thread_count = GPU_MULTIPLIER *4 + PARALEL_CPU;
88         // int novForThread = (nov+virtual_thread_count-1)/virtual_thread_count;
89
90         if(threadId <=(OUR_GPU_COUNT-1))
91         {
92             // int novStart = GPU_MULTIPLIER * novForThread * threadId;
93             // int novEnd = GPU_MULTIPLIER * novForThread * (threadId+1);
94             // if (novEnd > nov) novEnd = nov;
95             // int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) / THREADS_PER_BLOCK;
96             // printf("nov s %d e %d \n", novStart,novEnd);
97
98             cudaSetDevice(threadId);
99             int *adj_d;
100             int *xadj_d;
101             int *output_d;
102             cudaEvent_t start, stop;
103             float elapsedTime;
104
105             gpuErrchk(cudaMalloc((void**)&adj_d, (nnz) * sizeof(int)));
106             gpuErrchk(cudaMalloc((void**)&xadj_d, (nov + 1) * sizeof(int)));
107             gpuErrchk(cudaMalloc((void**)&output_d, (GPU_MULTIPLIER *CHUNK_SIZE_OUR)
108                 * sizeof(int)));
109             //gpuErrchk(cudaMallocHost((void **)&output_h, (nov) * sizeof(int)));
110             gpuErrchk(cudaMemcpy(adj_d, adj, (nnz) * sizeof(int),
111                 cudaMemcpyHostToDevice));
112             gpuErrchk(cudaMemcpy(xadj_d, xadj, (nov + 1) * sizeof(int),
113                 cudaMemcpyHostToDevice));
114
115             cudaEventCreate(&start);
116             cudaEventRecord(start, 0);
117             double start_gpu = omp_get_wtime();
118             cudaStream_t stream1;
119             cudaStreamCreate (&stream1);
120
121             while(true){
122                 int thisChunk;
123                 #pragma omp critical
124                 {
125                     thisChunk=currentChunk; //thisChunk is different on everyone.
126                     currentChunk+= GPU_MULTIPLIER;
127                 }
128                 if (thisChunk >= totalChunk) break;
129                 int novStart = thisChunk * CHUNK_SIZE_OUR;
130                 int novEnd = novStart + GPU_MULTIPLIER * CHUNK_SIZE_OUR;
131
132             }
133         }
134     }
135     end_cpu = omp_get_wtime();
136 }

```

```

147 // printf("nov s %d e %d \n", novStart,novEnd);
148 if(novEnd>nov) novEnd=nov;
149
150 int numBlock = (novEnd-novStart + THREADS_PER_BLOCK-1) /
    THREADS_PER_BLOCK;
151
152
153
154 if (n==3)kernel3<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
    xadj_d, output_d, novEnd,novStart);
155 else if (n==4)kernel4<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
    xadj_d, output_d, novEnd,novStart);
156 else if (n==5)kernel5<<<numBlock, THREADS_PER_BLOCK,0,stream1>>>(adj_d,
    xadj_d, output_d, novEnd,novStart);
157
158 gpuErrchk(cudaDeviceSynchronize());
159 gpuErrchk(cudaMemcpy(output_h+novStart, output_d, (novEnd-novStart) *
160     sizeof(int), cudaMemcpyDeviceToHost));
161
162 }
163
164
165 double end_gpu = omp_get_wtime();
166
167
168
169 cudaEventCreate(&stop);
170 cudaEventRecord(stop, 0);
171 cudaEventSynchronize(stop);
172 cudaEventElapsedTime(&elapsedTime, start, stop);
173
174 printf("GPU scale took: %f s on gpu %d \n", elapsedTime/1000, threadId);
175
176 cudaFree(adj_d);
177 cudaFree(xadj_d);
178
179
180 }
181
182 else{
183
184 double start_thread = omp_get_wtime();
185
186 while(true){
187     int thisChunk;
188     #pragma omp critical
189     {
190         thisChunk=currentChunk; //thisChunk is different on everyone.
191         currentChunk++;
192     }
193
194     if(thisChunk >= totalChunck) break;
195
196     int novStart = thisChunk * CHUNK_SIZE_OUR;
197     int novEnd = (thisChunk+1)* CHUNK_SIZE_OUR;
198     if(novEnd>nov) novEnd=nov;
199
200
201     bool *marked = new bool[nov];
202     memset(marked, false, nov * sizeof(bool)); // bu belki silinebilir
203
204
205     for(int i = novStart; i < novEnd; i++){
206         int localcount = 0;
207         DFS_sparse(xadj, adj, marked, n - 1, i, i, localcount);
208         output_h[i] = localcount;
209
210     }
211
212 }
213 double end_thread = omp_get_wtime();
214 printf("Took %f secs \n", end_thread -start_thread );
215
216 }
217
218 }
219 end_cpu = omp_get_wtime();
220
221 printf("Took %f secs \n", end_cpu - start_cpu);
222

```