## a5/cuda_kmeans_transpose.cu

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3
4   #include "kmeans.h"
5   #include "alloc.h"
6   #include "error.h"
7
8   #ifdef __CUDACC__
9   inline void checkCuda(cudaError_t e) {
10      if (e != cudaSuccess) {
11          // cudaGetErrorString() isn't always very helpful. Look up the error
12          // number in the cudaError enum in driver_types.h in the CUDA includes
13          // directory for a better explanation.
14          error("CUDA Error %d: %s\n", e, cudaGetErrorString(e));
15      }
16  }
17
18  inline void checkLastCudaError() {
19      checkCuda(cudaGetLastError());
20  }
21  #endif
22
23  __device__ int get_tid() {
24    return blockIdx.x * blockDim.x + threadIdx.x;
25  }
26
27  /* square of Euclid distance between two multi-dimensional points using column-base format
    */
28  __host__ __device__ inline static
29  double euclid_dist_2_transpose(int numCoords,
30                                 int numObjs,
31                                 int numClusters,
32                                 double *objects,     // [numCoords][numObjs]
33                                 double *clusters,    // [numCoords][numClusters]
34                                 int objectId,
35                                 int clusterId) {
36    int i;
37    double ans = 0.0;
38
39    /* TODO: Calculate the euclid_dist of elem=objectId of objects from elem=clusterId from
    clusters, but for column-base format!!! */
40    for (i = 0; i < numCoords; i++) {
41      double objectVal = objects[i * numObjs + objectId];
42      double clusterVal = clusters[i * numClusters + clusterId];
43
44      double diff = objectVal - clusterVal;
45      ans += diff * diff;
46    }
47
48    return (ans);
49  }
50
```

```
 51  __global__ static
 52  void find_nearest_cluster(int numCoords,
 53                            int numObjs,
 54                            int numClusters,
 55                            double *objects,          //  [numCoords][numObjs]
 56                            double *deviceClusters,   //  [numCoords][numClusters]
 57                            int *membership,          //  [numObjs]
 58                            double *devdelta) {
 59    /* Get the global ID of the thread. */
 60    int tid = get_tid();
 61
 62    if (tid < numObjs) {
 63      int index, i;
 64      double dist, min_dist;
 65
 66      /* find the cluster id that has min distance to object */
 67      index = 0;
 68
 69      min_dist = euclid_dist_2_transpose(numCoords, numObjs, numClusters,
 70                            objects, deviceClusters,
 71                            tid, index);
 72
 73      for (i = 1; i < numClusters; i++) {
 74
 75        dist = euclid_dist_2_transpose(numCoords, numObjs, numClusters,
 76                            objects, deviceClusters,
 77                            tid, i);
 78        /* no need square root */
 79        if (dist < min_dist) { /* find the min and its array index */
 80          min_dist = dist;
 81          index = i;
 82        }
 83      }
 84
 85      if (membership[tid] != index) {
 86
 87        atomicAdd(devdelta, 1.0);
 88      }
 89
 90      /* assign the deviceMembership to object objectId */
 91      membership[tid] = index;
 92    }
 93  }
 94
 95  //
 96  //  ---------------------------------------
 97  //  DATA LAYOUT
 98  //
 99  //  objects          [numObjs][numCoords]
100  //  clusters         [numClusters][numCoords]
101  //  dimObjects       [numCoords][numObjs]
102  //  dimClusters      [numCoords][numClusters]
103  //  newClusters      [numCoords][numClusters]
104  //  deviceObjects    [numCoords][numObjs]
```

```
105    //  deviceClusters  [numCoords][numClusters]
106    //  --------------------------------------
107    //
108    /* return an array of cluster centers of size [numClusters][numCoords]      */
109    void kmeans_gpu(double *objects,      /* in: [numObjs][numCoords] */
110                     int numCoords,    /* no. features */
111                     int numObjs,       /* no. objects */
112                     int numClusters,  /* no. clusters */
113                     double threshold,    /* % objects change membership */
114                     long loop_threshold,   /* maximum number of iterations */
115                     int *membership,   /* out: [numObjs] */
116                     double *clusters,   /* out: [numClusters][numCoords] */
117                     int blockSize) {
118       double timing = wtime(), timing_internal, timer_min = 1e42, timer_max = 0;
119       double timing_gpu, timing_cpu, timing_transfers, transfers_time = 0.0, cpu_time = 0.0,
       gpu_time = 0.0;
120       int loop_iterations = 0;
121       int i, j, index, loop = 0;
122       int *newClusterSize; /* [numClusters]: no. objects assigned in each
123                              new cluster */
124       double delta = 0, *dev_delta_ptr;          /* % of objects change their clusters */
125
126       /* TODO: Transpose dims */
127       double **dimObjects = (double **) calloc_2d(numCoords, numObjs, sizeof(double));
       //calloc_2d(...) -> [numCoords][numObjs]
128       double **dimClusters = (double **) calloc_2d(numCoords, numClusters, sizeof(double));
       //calloc_2d(...) -> [numCoords][numClusters]
129       double **newClusters = (double **) calloc_2d(numCoords, numClusters, sizeof(double));
       //calloc_2d(...) -> [numCoords][numClusters]
130
131       double *deviceObjects;
132       double *deviceClusters;
133       int *deviceMembership;
134
135       printf("\n|-----------Transpose GPU Kmeans------------|\n\n");
136
137       //  TODO: Copy objects given in [numObjs][numCoords] layout to new
138       //  [numCoords][numObjs] layout
139       for (i=0 ; i <numObjs; i++){
140         for (j=0; j<numCoords; j++){
141           dimObjects[j][i]=objects[i*numCoords + j];
142         }
143       }
144
145       /* pick first numClusters elements of objects[] as initial cluster centers*/
146       for (i = 0; i < numCoords; i++) {
147         for (j = 0; j < numClusters; j++) {
148           dimClusters[i][j] = dimObjects[i][j];
149         }
150       }
151
152       /* initialize membership[] */
153       for (i = 0; i < numObjs; i++) membership[i] = -1;
154
```

```
155      /* need to initialize newClusterSize and newClusters[0] to all 0 */
156      newClusterSize = (int *) calloc(numClusters, sizeof(int));
157      assert(newClusterSize != NULL);
158
159      timing = wtime() - timing;
160      printf("t_alloc: %lf ms\n\n", 1000 * timing);
161      timing = wtime();
162
163      const unsigned int numThreadsPerClusterBlock = (numObjs > blockSize) ? blockSize :
     numObjs;
164      const unsigned int numClusterBlocks = (numObjs + numThreadsPerClusterBlock - 1) /
     numThreadsPerClusterBlock;
165      const unsigned int clusterBlockSharedDataSize = 0;
166
167      checkCuda(cudaMalloc(&deviceObjects, numObjs * numCoords * sizeof(double)));
168      checkCuda(cudaMalloc(&deviceClusters, numClusters * numCoords * sizeof(double)));
169      checkCuda(cudaMalloc(&deviceMembership, numObjs * sizeof(int)));
170      checkCuda(cudaMalloc(&dev_delta_ptr, sizeof(double)));
171      timing = wtime() - timing;
172      printf("t_alloc_gpu: %lf ms\n\n", 1000 * timing);
173      timing = wtime();
174
175      checkCuda(cudaMemcpy(deviceObjects, dimObjects[0],
176                           numObjs * numCoords * sizeof(double), cudaMemcpyHostToDevice));
177      checkCuda(cudaMemcpy(deviceMembership, membership,
178                           numObjs * sizeof(int), cudaMemcpyHostToDevice));
179      timing = wtime() - timing;
180      printf("t_get_gpu: %lf ms\n\n", 1000 * timing);
181      timing = wtime();
182
183      do {
184        timing_internal = wtime();
185
186        /* GPU part: calculate new memberships */
187
188        timing_transfers = wtime();
189        // TODO: Copy clusters to deviceClusters
190        checkCuda(cudaMemcpy(deviceClusters, dimClusters[0],
191                             numClusters * numCoords * sizeof(double),
192                             cudaMemcpyHostToDevice));
193
194        transfers_time += wtime() - timing_transfers;
195
196        checkCuda(cudaMemset(dev_delta_ptr, 0, sizeof(double)));
197
198        //printf("Launching find_nearest_cluster Kernel with grid_size = %d, block_size = %d,
     shared_mem = %d KB\n", numClusterBlocks, numThreadsPerClusterBlock, clusterBlockSharedDa-
     taSize/1000);
199        timing_gpu = wtime();
200        find_nearest_cluster
201        <<< numClusterBlocks, numThreadsPerClusterBlock, clusterBlockSharedDataSize >>>
202                (numCoords, numObjs, numClusters,
203                 deviceObjects, deviceClusters, deviceMembership, dev_delta_ptr);
204
```

```
205        cudaDeviceSynchronize();
206        checkLastCudaError();
207        gpu_time += wtime() - timing_gpu;
208        //printf("Kernels complete for itter %d, updating data in CPU\n", loop);
209
210        timing_transfers = wtime();
211
212        checkCuda(cudaMemcpy(membership, deviceMembership,
213                        numObjs * sizeof(int),
214                        cudaMemcpyDeviceToHost));
215
216        checkCuda(cudaMemcpy(&delta, dev_delta_ptr,
217                        sizeof(double),
218                        cudaMemcpyDeviceToHost));
219        transfers_time += wtime() - timing_transfers;
220
221        /* CPU part: Update cluster centers*/
222
223        timing_cpu = wtime();
224        for (i = 0; i < numObjs; i++) {
225          /* find the array index of nestest cluster center */
226          index = membership[i];
227
228          /* update new cluster centers : sum of objects located within */
229          newClusterSize[index]++;
230          for (j = 0; j < numCoords; j++)
231            newClusters[j][index] += objects[i * numCoords + j];
232        }
233
234        /* average the sum and replace old cluster centers with newClusters */
235        for (i = 0; i < numClusters; i++) {
236          for (j = 0; j < numCoords; j++) {
237            if (newClusterSize[i] > 0)
238              dimClusters[j][i] = newClusters[j][i] / newClusterSize[i];
239            newClusters[j][i] = 0.0;   /* set back to 0 */
240          }
241          newClusterSize[i] = 0;   /* set back to 0 */
242        }
243
244        delta /= numObjs;
245        //printf("delta is %f - ", delta);
246        loop++;
247        //printf("completed loop %d\n", loop);
248        cpu_time += wtime() - timing_cpu;
249
250        timing_internal = wtime() - timing_internal;
251        if (timing_internal < timer_min) timer_min = timing_internal;
252        if (timing_internal > timer_max) timer_max = timing_internal;
253    } while (delta > threshold && loop < loop_threshold);
254
255    /*TODO: Update clusters using dimClusters. Be carefull of layout!!!
    clusters[numClusters][numCoords] vs dimClusters[numCoords][numClusters] */
256    for (i = 0; i < numClusters; i++) {
257        for (j = 0; j < numCoords; j++) {
```

```
258                 clusters[i * numCoords + j] = dimClusters[j][i];
259             }
260         }
261
262     timing = wtime() - timing;
263     printf("nloops = %d  : total = %lf ms\n\t-> t_loop_avg = %lf ms\n\t-> t_loop_min = %lf
        ms\n\t-> t_loop_max = %lf ms\n\t"
264             "-> t_cpu_avg = %lf ms\n\t-> t_gpu_avg = %lf ms\n\t-> t_transfers_avg = %lf
        ms\n\n|----------------------------------------|\n",
265             loop, 1000 * timing, 1000 * timing / loop, 1000 * timer_min, 1000 * timer_max,
266             1000 * cpu_time / loop, 1000 * gpu_time / loop, 1000 * transfers_time / loop);
267
268     char outfile_name[1024] = {0};
269     sprintf(outfile_name, "Execution_logs/silver1-V100_Sz-%lu_Coo-%d_Cl-%d.csv",
270             numObjs * numCoords * sizeof(double) / (1024 * 1024), numCoords, numClusters);
271     FILE *fp = fopen(outfile_name, "a+");
272     if (!fp) error("Filename %s did not open succesfully, no logging performed\n",
        outfile_name);
273     fprintf(fp, "%s,%d,%lf,%lf,%lf\n", "Transpose", blockSize, timing / loop, timer_min,
        timer_max);
274     fclose(fp);
275
276     checkCuda(cudaFree(deviceObjects));
277     checkCuda(cudaFree(deviceClusters));
278     checkCuda(cudaFree(deviceMembership));
279
280     free(dimObjects[0]);
281     free(dimObjects);
282     free(dimClusters[0]);
283     free(dimClusters);
284     free(newClusters[0]);
285     free(newClusters);
286     free(newClusterSize);
287
288     return;
289 }
290
291
```