

The parallel algorithm (source code provided on the following pages) differs from the obvious sequential solution in a few key ways:

- “Closest centroid” cluster assignments are computed one point per thread.
- A sum reduction is used to combine the per-thread computations into per-block values.
- Two more block sum reductions are performed to reduce the values into a single block for the purposes of updating centroid positions.
- Branch prediction optimizations were removed since the CUDA architecture does not benefit from branches being taken more or less frequently.

Let $C = \frac{k \cdot n \cdot h}{t}$ be the computation rate where $1 \leq k \leq 1000$ is the number of clusters (means), $1 \leq n \leq 10^8$ is the number of points (in 2D), h is the number of iterations until convergence, and t is elapsed time (in seconds). Then, on a system with an E5-2650v3 CPU and TITAN V GPU,

- Our sequential code’s performance peaked at ~ 950 million computations per second (across all parameter choices in these ranges).
- Our parallel code’s performance peaked at ~ 65 billion computations per second (across all parameter choices in these ranges).

This constitutes approximately a 70x speedup (keeping in mind that a single core of the E5-2650v3 CPU and the entirety of a TITAN V are drastically different architectures).

You can gain a significant speedup when k is large by just using atomic CUDA operations (e.g. we’ve observed up to ~ 575 billion computations per second for $k = 1000$), but in practice k is usually small, so we’ve intentionally avoided this solution.

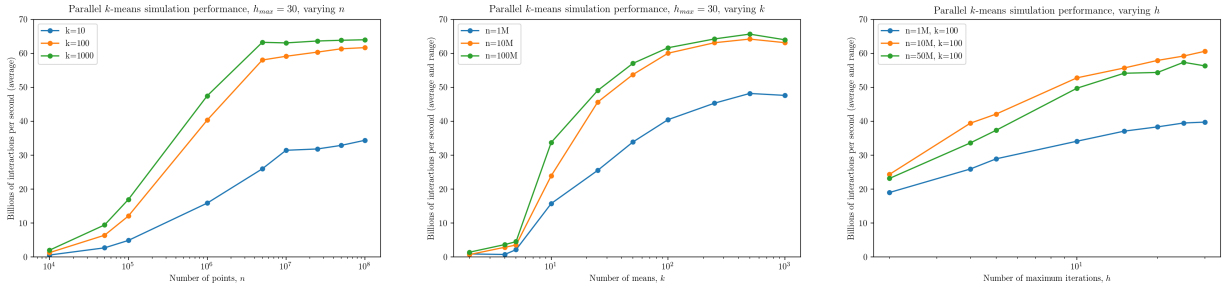


Figure 1: Performance of the parallel algorithm (in log-scale) when compiled with `nvcc`. Left-Center-Right: Performance when varying n , k , and h , respectively.

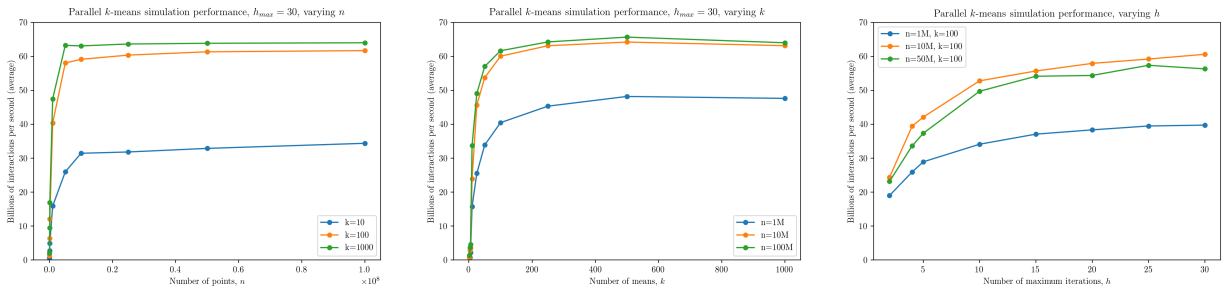


Figure 2: Performance of the parallel algorithm (in linear-scale) when compiled with `nvcc`. Left-Center-Right: Performance when varying n , k , and h , respectively.

Sample Output

We wrote a short data generator in Python that draws points from 2D Gaussians on a radial tree.

```

1  import numpy as np
2  import sys
3
4  # 10 20 5 2 2000000 0.1 generates 100M points
5  num_branches1 = int(sys.argv[1])
6  dist_branches1 = float(sys.argv[2])
7  num_branches2 = int(sys.argv[3])
8  dist_branches2 = float(sys.argv[4])
9  cluster_size = int(sys.argv[5])
10 cluster_scale = float(sys.argv[6])
11
12 for i in range(num_branches1):
13     for j in range(num_branches2):
14         val = dist_branches1 * 1j ** (i * 4 / num_branches1) + \
15             dist_branches2 * 1j ** (j * 4 / num_branches2)
16         for _ in range(cluster_size):
17             x = val.real + np.random.normal(scale=cluster_scale)
18             y = val.imag + np.random.normal(scale=cluster_scale)
19             print(x, y)

```

For example, the 1 million point data sets were generated with

```
python3 gen_radial_points.py 10 20 5 2 20000 0.1
```

yielding

- 10 “large clusters” 20 units away from the origin
- 5 “small clusters” in each large one, 5 units away from the centers of the large clusters.
- 20 thousand points per cluster, drawn from a 2D Gaussian distribution with standard deviation 0.1.

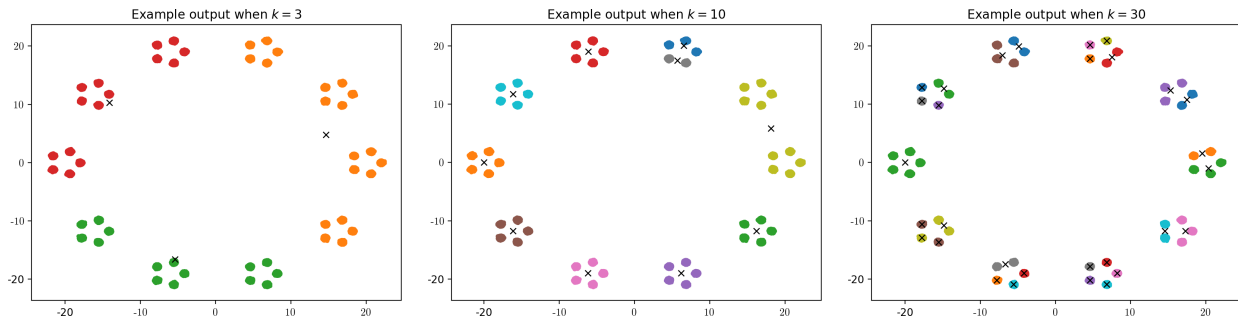


Figure 3: Example output from our parallel k -means algorithm. Left-Center-Right: Output with $k = 3, 10, 30$, respectively.

Strictly speaking, the CUDA/C++ code simply prints a performance reading, followed by lines of the form

```

=====cluster 0 centered at 14.635420 4.755381 has size 400000=====
22.014330 -0.275197
21.905111 0.078775
21.974599 0.088757
<OMIT>
=====cluster 1 centered at -5.393562 -16.599389 has size 300000=====
-14.175596 -11.695731
-14.225406 -11.717935
-14.166853 -11.820264
<OMIT>

```

but this could be easily altered to have the results exported in other formats.

CUDA Code Listing

```

1  #include <assert.h>
2  #include <math.h>
3  #include <stdio.h>
4  #include <stdlib.h>
5  #include <string.h>
6  #include <time.h>
7
8  #define BLOCK_SIZE 512
9  #define MAX_POINTS 100000000 // 100M points
10 #define MAX_MEANS 1000
11 #define MAX_ITER 30
12
13 // CUDA prefers struct-of-arrays style here (for cache purposes)
14 typedef struct {
15     double *x, *y;
16     int *membership;
17 } points;
18
19 typedef struct {
20     double *x, *y;
21 } centroids;
22
23 typedef struct {
24     double *x_sum, *y_sum;
25     int *size;
26 } temp_centroids;
27
28 // algorithm termination flag
29 __managed__ int assignment_changed = 1;
30
31 // reads n data points from input file
32 __host__ void read_data(int n, char *file_name, points P) {
33     unsigned int i = 0;
34     double x, y;
35     FILE *file = fopen(file_name, "r");
36     assert(file != NULL);
37
38     while (!feof(file) && i < n) {
39         if (fscanf(file, "%lf%lf", &x, &y) != 2)
40             break;
41         P.x[i] = x;
42         P.y[i] = y;
43         P.membership[i++] = -1;
44     }
45 }
46
47 // selects k centers at random from n points
48 __host__ void init_centers(int n, int k, points P, centroids C) {
49     srand(time(NULL));
50     for (int i = 0; i < k; ++i) {
51         // not actually uniform random sampling, but very close
52         int rand_idx = rand() % n;
53         C.x[i] = P.x[rand_idx];
54         C.y[i] = P.y[rand_idx];
55     }
56 }
57
58 // computes ||p-c||^2 for a point p and center c
59 __device__ inline double norm_2D_sqr(double x1, double y1,
60                                     double x2, double y2) {
61     // sqrt is monotonic, so we may omit it in the distance calculation
62     // i.e. application of sqrt does not change the order of distances
63     return (x1 - x2) * (x1 - x2) + (y1 - y2) * (y1 - y2);
64 }
65
66 // assign each point to the cluster given by the closest centroid
67 // NVIDIA suggests const and restrict here to improve compiler optimization
68 __global__ void
69 assign_clusters(int n, int k,
70               const double *__restrict__ Px,
71               const double *__restrict__ Py,
72               int *__restrict__ Pmembership,
73               double *__restrict__ Cx,
74               double *__restrict__ Cy,
75               double *__restrict__ Ox_sum,
76               double *__restrict__ Oy_sum,
77               int *__restrict__ Osize) {
78     int index = blockIdx.x * blockDim.x + threadIdx.x;
79     int tid = threadIdx.x;
80
81     // thread-local values that will be reduced
82     __shared__ double x_sum[BLOCK_SIZE];

```

```

83  __shared__ double y_sum[BLOCK_SIZE];
84  __shared__ int size[BLOCK_SIZE];
85
86  int membership = -1;
87
88  if (index < n) {
89      double min_dist = INFINITY;
90      for (int i = 0; i < k; ++i) {
91          double current_dist = norm_2D_sqr(Px[index], Py[index], Cx[i], Cy[i]);
92          if (current_dist < min_dist) {
93              min_dist = current_dist;
94              membership = i;
95          }
96      }
97
98      // arbitrary concurrent write is valid since all
99      // threads write the same value
100     if (membership != Pmembership[index])
101         assignment_changed = 1;
102     Pmembership[index] = membership;
103 }
104 __syncthreads();
105
106 // k reductions (one per centroid)
107 for (int c = 0; c < k; ++c) {
108     x_sum[tid] = (membership == c) ? Px[index] : 0;
109     y_sum[tid] = (membership == c) ? Py[index] : 0;
110     size[tid] = (membership == c) ? 1 : 0;
111     __syncthreads();
112
113     // reduce block's sums into one value (in thread 0)
114     for (int offset = BLOCK_SIZE >> 1; offset > 0; offset >>= 1) {
115         if (tid < offset) {
116             x_sum[tid] += x_sum[tid + offset];
117             y_sum[tid] += y_sum[tid + offset];
118             size[tid] += size[tid + offset];
119         }
120         __syncthreads();
121     }
122
123     // save block's sums to output arrays
124     if (tid == 0) {
125         Ox_sum[blockIdx.x * k + c] = x_sum[tid];
126         Oy_sum[blockIdx.x * k + c] = y_sum[tid];
127         Osize[blockIdx.x * k + c] = size[tid];
128     }
129     __syncthreads();
130 }
131 }
132
133 // reduce temporary cluster sizes and centroid x/y sums to smaller arrays
134 __global__ void
135 reduce_temp_clusters(int n, int k,
136                     const double *__restrict__ Ix_sum,
137                     const double *__restrict__ Iy_sum,
138                     const int *__restrict__ Isize,
139                     double *__restrict__ Ox_sum,
140                     double *__restrict__ Oy_sum,
141                     int *__restrict__ Osize) {
142     int index = blockIdx.x * blockDim.x + threadIdx.x;
143     int stride = blockDim.x * gridDim.x;
144     int tid = threadIdx.x;
145
146     // thread-local values that will be reduced
147     __shared__ double x_sum[BLOCK_SIZE];
148     __shared__ double y_sum[BLOCK_SIZE];
149     __shared__ int size[BLOCK_SIZE];
150
151     for (int c = 0; c < k; ++c) {
152         x_sum[tid] = 0;
153         y_sum[tid] = 0;
154         size[tid] = 0;
155
156         // if necessary, sum multiple items per thread
157         for (int b = index; b < n; b += stride) {
158             x_sum[tid] += Ix_sum[b * k + c];
159             y_sum[tid] += Iy_sum[b * k + c];
160             size[tid] += Isize[b * k + c];
161         }
162         __syncthreads();
163
164         // reduce block's sums into one value (in thread 0)
165         for (int offset = BLOCK_SIZE >> 1; offset > 0; offset >>= 1) {
166             if (tid < offset) {

```

```

167         x_sum[tid] += x_sum[tid + offset];
168         y_sum[tid] += y_sum[tid + offset];
169         size[tid] += size[tid + offset];
170     }
171     __syncthreads();
172 }
173
174 // save block's sums to output arrays
175 if (tid == 0) {
176     Ox_sum[blockIdx.x * k + c] = x_sum[tid];
177     Oy_sum[blockIdx.x * k + c] = y_sum[tid];
178     Osize[blockIdx.x * k + c] = size[tid];
179 }
180 __syncthreads();
181 }
182 }
183
184 // update cluster centroid positions
185 __global__ void update_clusters(int n, int k,
186                                double *__restrict__ Cx,
187                                double *__restrict__ Cy,
188                                const double *__restrict__ Ix_sum,
189                                const double *__restrict__ Iy_sum,
190                                const int *__restrict__ Isize) {
191     int index = blockIdx.x * blockDim.x + threadIdx.x;
192
193     if (index < k && Isize[index]) {
194         Cx[index] = Ix_sum[index] / Isize[index];
195         Cy[index] = Iy_sum[index] / Isize[index];
196     }
197 }
198
199 /*
200 * prints results and performance where
201 * k = number of clusters (means)
202 * n = number of points (in 2D)
203 * h = number of iterations until convergence
204 * t = elapsed time (in seconds)
205 *
206 * P contains the input points
207 * C contains the final cluster centroids
208 * T contains (in part) the final cluster sizes
209 */
210 __host__ void print_results(int k, int n, int h, double t,
211                             points P, centroids C, temp_centroids T) {
212     printf("performed %d iterations in %.2f s, perf: %.2f billion\n", h, t,
213           (double)k * n * h / t * 1e-9);
214
215     double *xs = (double *)malloc(sizeof(double) * n);
216     double *ys = (double *)malloc(sizeof(double) * n);
217     int offsets[k + 1];
218
219     offsets[0] = 0;
220     for (int i = 0; i < k; ++i) {
221         offsets[i + 1] = offsets[i] + T.size[i];
222     }
223
224     // pack permutation of input points into clusters in a single pass by using
225     // prefix-sum on the cluster sizes as offsets into our output arrays
226     for (int i = 0; i < n; ++i) {
227         int m = P.membership[i];
228         xs[offsets[m]] = P.x[i];
229         ys[offsets[m]++] = P.y[i];
230     }
231
232     for (int c = 0; c < k; ++c) {
233         printf("====cluster %d, centered at %lf %lf, has %d size====\n", c, C.x[c],
234               C.y[c], T.size[c]);
235         for (int i = offsets[c] - T.size[c]; i < offsets[c]; ++i) {
236             printf("%lf %lf\n", xs[i], ys[i]);
237         }
238     }
239
240     free(xs);
241     free(ys);
242 }
243
244 int main(int argc, char **argv) {
245     int k, n, h;
246     char *file_name;
247     points P;
248     centroids C;
249     temp_centroids T1;
250     temp_centroids T2;

```

```

251  cudaEvent_t start, stop;
252  float time;
253
254  // read in number of points and means
255  assert(argc >= 4);
256  n = atoi(argv[1]);
257  k = atoi(argv[2]);
258  file_name = argv[3];
259  assert(n <= MAX_POINTS && k <= MAX_MEANS);
260
261  int blockSize = BLOCK_SIZE;
262  int numBlocks = (n + blockSize - 1) / blockSize;
263  int reductionBlockSize = BLOCK_SIZE;
264  int reductionNumBlocks =
265      (numBlocks + reductionBlockSize - 1) / reductionBlockSize;
266
267  // make sure that we can support the number of points with our two block
268  // reductions. with BLOCK_SIZE = 512, this limit is ~250M points
269  assert(reductionNumBlocks <= 1024);
270
271  // malloc memory and set up GPU timers
272  cudaMallocManaged(&P.x, sizeof(double) * n);
273  cudaMallocManaged(&P.y, sizeof(double) * n);
274  cudaMallocManaged(&P.membership, sizeof(int) * n);
275  cudaMallocManaged(&C.x, sizeof(double) * k);
276  cudaMallocManaged(&C.y, sizeof(double) * k);
277  cudaMallocManaged(&T1.x_sum, sizeof(double) * numBlocks * k);
278  cudaMallocManaged(&T1.y_sum, sizeof(double) * numBlocks * k);
279  cudaMallocManaged(&T1.size, sizeof(int) * numBlocks * k);
280  cudaMallocManaged(&T2.x_sum, sizeof(double) * reductionNumBlocks * k);
281  cudaMallocManaged(&T2.y_sum, sizeof(double) * reductionNumBlocks * k);
282  cudaMallocManaged(&T2.size, sizeof(int) * reductionNumBlocks * k);
283  cudaEventCreate(&start);
284  cudaEventCreate(&stop);
285
286  read_data(n, file_name, P);
287  init_centers(n, k, P, C);
288
289  cudaEventRecord(start, 0);
290  for (h = 0; assignment_changed && h < MAX_ITER; ++h) {
291      // assign points to nearest clusters
292      assignment_changed = 0;
293      assign_clusters<<<numBlocks, blockSize>>>
294          (n, k,
295           P.x, P.y, P.membership,
296           C.x, C.y,
297           T1.x_sum, T1.y_sum, T1.size);
298      cudaDeviceSynchronize();
299
300      // two block reductions of cluster sizes and centroid x/y sums
301      reduce_temp_clusters<<<reductionNumBlocks, reductionBlockSize>>>
302          (numBlocks, k,
303           T1.x_sum, T1.y_sum, T1.size, // input values to reduce
304           T2.x_sum, T2.y_sum, T2.size); // reduced output values
305      cudaDeviceSynchronize();
306      reduce_temp_clusters<<<1, reductionBlockSize>>>
307          (reductionNumBlocks, k,
308           T2.x_sum, T2.y_sum, T2.size, // reduce values from T2
309           T1.x_sum, T1.y_sum, T1.size); // back into T1
310      cudaDeviceSynchronize();
311
312      // update centroid positions
313      update_clusters<<<1, k>>>
314          (n, k,
315           C.x, C.y,
316           T1.x_sum, T1.y_sum, T1.size);
317      cudaDeviceSynchronize();
318  }
319
320  cudaEventRecord(stop, 0);
321  cudaEventSynchronize(stop);
322  cudaEventElapsedTime(&time, start, stop);
323  cudaEventDestroy(start);
324  cudaEventDestroy(stop);
325
326  print_results(k, n, h, time * 1e-3, P, C, T1);
327
328  // CUDA automatically frees and resets device on program exit
329  }

```