

 Features Explore Pricing

This repository

serban / kmeans

<> Code

! Issues 3

🔗 Pull requests 1


📁 Projects 0

📈 Pulse

📊 Graphs

Branch: master ▾

kmeans / README

 **serban** Update README with -DBLOCK_SHARED_MEM_OPTIMIZATION=0

1 contributor

147 lines (115 sloc) 6.31 KB

```
1 This software is dervied from Professor Wei-keng Liao's parallel k-means
2 clustering code obtained on November 21, 2010 from
3 http://users.eecs.northwestern.edu/~wkliao/Kmeans/index.html
4 (http://users.eecs.northwestern.edu/~wkliao/Kmeans/simple_kmeans.tar.gz).
5
6 With his permission, I am publishing my CUDA implementation based on his code
7 under the open-source MIT license. See the LICENSE file for more details.
8
9 For starters, run the benchmark.sh script to see how fast this code runs.
10
11 It's pretty fast! Depending on your hardware, data set, and k, you should see
12 dramatic improvements in performance over CPU implementations.
13
14 On an 8-core 2.4 GHz Intel Xeon E5620 machine with an NVIDIA Tesla C1060 card,
15 the CUDA implementation runs almost 50 times faster than the sequential version
16 on the color17695.bin data set (for k = 128)!
17
18 Here's some sample output for the same data set on an 8-core 2.67 GHz Intel Core
19 i7 920 machine with an NVIDIA GeForce GTX 480 card. For k = 128, the speedup is
20 over 75x!
21
22 -----
23 k = 2 seqTime = 0.2870s ompTime = 0.1497s cudaTime = 0.1483s speedup = 1.9x
24 k = 4 seqTime = 0.0945s ompTime = 0.0351s cudaTime = 0.0652s speedup = 1.4x
25 k = 8 seqTime = 0.4379s ompTime = 0.1439s cudaTime = 0.0919s speedup = 4.7x
```

```
26 k = 16 seqTime = 2.0539s ompTime = 0.6628s cudaTime = 0.1611s speedup = 12.7x
27 k = 32 seqTime = 3.0699s ompTime = 1.0319s cudaTime = 0.1461s speedup = 21.0x
28 k = 64 seqTime = 8.4675s ompTime = 3.1946s cudaTime = 0.2487s speedup = 34.0x
29 k = 128 seqTime = 22.9694s ompTime = 10.0000s cudaTime = 0.3031s speedup = 75.7x
```

```
30 -----
31
```

```
32 The CUDA implementation may need some tweaking to work with larger data sets,
33 but the basic functionality is there. I've optimized the code with the general
34 assumption that k and the number of clusters will be relatively small compared
35 to the number of data points.
```

```
36
```

```
37 In fact, you may run into problems if k is too big or if your data has a large
38 number of dimensions because there won't be enough space in block shared memory
39 to hold all the clusters (the C1060 and GTX 480 have 16 and 48 KiB of block
40 shared memory, respectively). If you hit this limitation, you should be able to
41 get around it easily. Do the following:
```

```
42
```

```
43 1) Run 'make clean'
```

```
44 2) Edit the Makefile. Find the line at the top of the file that looks like this:
```

```
45
```

```
46         CFLAGS = $(OPTFLAGS) $(DFLAGS) $(INCFLAGS) -DBLOCK_SHARED_MEM_OPTIMIZATION=1
```

```
47
```

```
48 3) Set -DBLOCK_SHARED_MEM_OPTIMIZATION=0
```

```
49 4) Run 'make cuda' and try again
```

```
50
```

```
51 Please don't hesitate to contact me with any questions you may have. I'd love to
52 help you out if you run into a problem. Of course, the more information you give
53 me about your CUDA hardware and your data set (number of data points,
54 dimensionality, number of clusters), the more helpful I can be.
```

```
55
```

```
56 The original README, with some additions, is reproduced below.
```

```
57
```

```
58 Cheers!
```

```
59
```

```
60 Serban Giuroiu
```

```
61 http://serban.org
```

```
62
```

```
63 # -----
```

```
64
```

```
65 Parallel K-Means Data Clustering
```

```
66
```

```
67 The software package of parallel K-means data clustering contains the
68 followings:
```

```
69
```

```
70  * A parallel implementation using OpenMP and C
71  * A parallel implementation using MPI and C
72  * A parallel implementation using CUDA and C
73  * A sequential version in C
74
75  To compile:
76  Although I used Intel C compiler, icc, version 7.1 during the code
77  development, there is no particular features required except for OpenMP.
78  Thus, the implementation should be fairly portable. Please modify
79  Makefile to change the compiler if needed.
80
81  You will need the NVIDIA CUDA toolkit, which contains nvcc, to build the CUDA
82  version. It works fine in concert with gcc.
83
84  To run:
85  * The Makefile will produce executables
86    o "omp_main" for OpenMP version
87    o "mpi_main" for MPI version
88    o "cuda_main" for CUDA version
89    o "seq_main" for sequential version
90
91  * The list of available command-line arguments can be obtained by
92  running -h option
93    o For example, running command "omp_main -h" will produce:
94      Usage: main [switches] -i filename -n num_clusters
95          -i filename      : file containing data to be clustered
96          -b                : input file is in binary format (default no)
97          -n num_clusters: number of clusters (K must > 1)
98          -t threshold     : threshold value (default 0.0010)
99          -p nproc         : number of threads (default system allocated)
100          -a               : perform atomic OpenMP pragma (default no)
101          -o               : output timing results (default no)
102          -d               : enable debug mode
103
104  Input file format:
105  The executables read an input file that stores the data points to be
106  clustered. A few example files are provided in the sub-directory
107  ./Image_data. The input files can be in two formats: ASCII text and raw
108  binary.
109
110  * ASCII text format:
111    o Each line contains the coordinates of a single data point
112    o The number of coordinates must be equal for all data points
113  * Raw binary format:
```

```
114     o There is a header of 2 integers.
115     o The first 4-byte integer must be the number of data points.
116     o The second integer must be the number of coordinates.
117     o The rest of the file contains the coordinates of all data
118       points and each coordinate is of type 4-byte float.
119
120 Output files: There are two output files:
121     * Coordinates of cluster centers
122     o The file name is the input file name appended with ".cluster_centres".
123     o It is in ASCII text format.
124     o Each line contains an integer indicating the cluster id and the
125       coordinates of the cluster center.
126     * Membership of all data points to the clusters
127     o The file name is the input file name appended with ".membership".
128     o It is in ASCII text format.
129     o Each line contains two integers: data point index (from 0 to
130       the number of points) and the cluster id indicating the membership of
131       the point.
132
133 Limitations:
134     * Data type -- This implementation uses C float data type for all
135       coordinates and other real numbers.
136     * Large number of data points -- The number of data points cannot
137       exceed 2G due to the 4-byte integers used in the programs. (But do
138       let me know if it is desired.)
139
140
141 Wei-keng Liao (wkliao@ece.northwestern.edu)
142 EECS Department
143 Northwestern University
144
145 Sep. 17, 2005
146
```

