

EVALUATION SEGMENTATION ON GPUS

ABDUL DAKKAK

ABSTRACT

Image segmentation is an important and well studied problem in computer vision. The task is to label pixels, grouping them into the same category. In medical imaging, for example, one needs to classify a liver from the rest of the image, or discern the shape of cells to detect any diseases. Sadly, segmentation is slow — taking on the order of minutes. Speeding up segmentation would not only improve existing applications, but also open the possibility of new ones.

In this project we evaluate image segmentation on the GPU. We look at Cellular Automaton (CA) based approaches for binary segmentation (classifying pixels as either background or foreground). We try different optimizations on both CPUs and GPU and perform a comparison.

1 CELLULAR AUTOMATON (GROWCUT) BASED SEG- MENTATION

GrowCut [1] is an algorithm that uses ideas from cellular automaton to perform image segmentation. GrowCut flood neighbors based on an initial seed until it reaches a barrier. One can think of each pixel as a bacteria with an energy function, if a bacteria has higher energy than one of its neighbors, then it will devour them — otherwise it is the victim. The algorithm is iterated until either a fixed point is reached or the maximum number of iterations are reached. A pseudocode of the algorithm is shown in listing 2. For our experiments, we set MAX_ITERATIONS to 2048 and use a penalty function $g(x, y) = 1 - \|x - y\|_2$.

Listing 1: Pseudo Code for GrowCut Segmentation

```
1 Input: Image, Approximate Label
2 Output: Label Segmentation
3 -----
4 changed = true
5 do_while changed && iterations < MAX_ITERATIONS:
6   changed = false
7   for ii from 0 to height:
8     for jj from 0 to width:
9       cp = image[ii][jj]
10      nl = lp = label[ii][jj]
11      ns = sp = strength[ii][jj]
12      for ni from -1 to 1:
13        for nj from -1 to 1:
```

* dakkak@illinois.edu

```

14     cq = image[ii+ni][jj+nj]
15     lq = label[ii+ni][jj+nj]
16     sq = strength[ii+ni][jj+nj]
17     gc = g(cp, cq)
18     if gc * sq > sp:
19         nl = lq
20         ns = sq * gc
21         changed = true
22     end
23 end
24 end
25 nextLabel[ii][jj] = nl
26 nextStrength[ii][jj] = ns
27 end
28 end
29 iterations++
30 label = nextLabel
31 strength = nextStrength
32 end

```

1.1 Markov Random Field (GraphCut) Based Segmentation

MRF have been used in GPU computing with good success. The problem is that they do not map nicely to the hardware — use irregular data structures, require atomic updates, etc... Furthermore, the method is patented which limits its usefulness in real world applications. While we initially had an MRF implementation, we abandoned it in favor of a much more interesting algorithm that maps nicely to GPUs.

1.2 Optimizations

In this section we evaluate different optimizations on both CPUs and GPUs.

TASK BASED PARALLELIZATION We first parallelized the code to run on different CPU cores. This was done using a task based approach (using Thread Building Blocks [2]) which simplifies some of the programming.

BASE CUDA IMPLEMENTATION Once the threaded multi-core version was working, we ported the code to use CUDA. Each thread in CUDA processes a pixel, access its elements, and updates its state. A global synchronization is required after each iteration of the algorithm since boundary edges are not known after a thread group completes.

CUDA IMPLEMENTATION USING SHARED MEMORY Shared memory are memory residing on each core of the GPU. Placing reused memory there, instead of always accessing it of chip, is a common optimization. In place of explicitly computing the 2-norm, we use the `hypotf` function in this step as an added optimization. We also give compile hints so constant memory are to be placed in a constant memory space (which is a separate part of the hardware dedicated to read-only memory).

APPROXIMATING THE 2-NORM The 2-norm computation is the only complicated arithmetic operation in our code, we therefore substituting it by using the fact that $\|x - y\|_2$ can be approximated with $\alpha \text{Max}(x, y) + \beta \text{Min}(x, y)$, the choice of α and β depend on the error you can tolerate. We chose $\alpha = 1$ and $\beta = 1/2$ which gives a maximum error of 11.8% and a mean error of 8.86 on average.

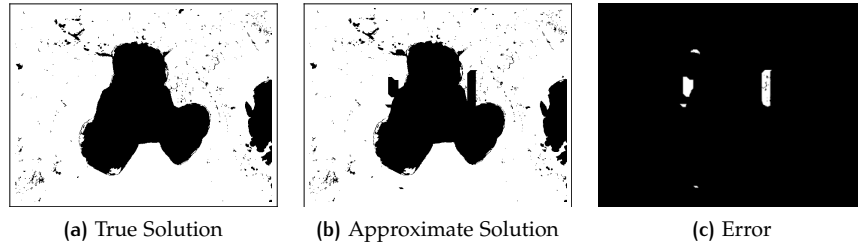


Figure 1: This shows the effects of approximating the algorithm by not having a global synchronization after each step. We rescale our differences so that even small errors are visible.

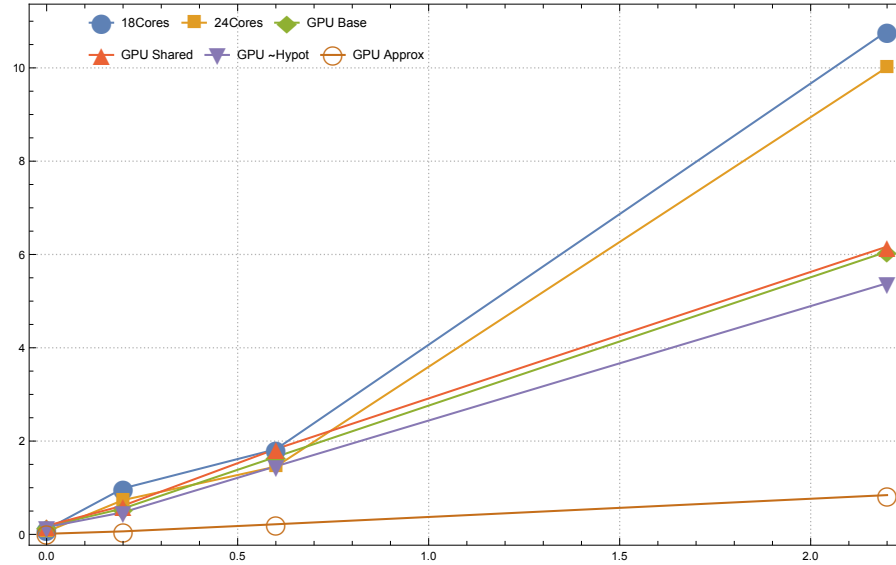


Figure 2: This shows the timing of GrowCut on a 24 and 32 core system. As expected, no difference is noticed between the two configurations due to hitting I/O bandwidth limits. The x axis is the number of megapixels in the image and the y axis is the execution time in seconds. It also shows the effects of different GPU optimizations.

APPROXIMATING THE ALGORITHM A synchronization is required between iterations different iterations of the algorithm. Since CUDA has no global synchronization capabilities, we essentially do a local synchronization (within the block). To not propagate error too much, we perform 32 iterations before doing a global synchronization.

Since the value of boundary pixels is not known within a local synchronization, we make the assumption that boundary pixels “look” like their neighbors or their previous state. If both the past history and the neighbors match, then we keep the value unchanged. Otherwise, we have a 50% probability to be correct and take the value of the neighborhood pixel. Figure 1 shows the effects of this approximation.

2 EVALUATION

In Figure 2 we see that having 24 core (Intel Xeon X5660) vs 32 core (Intel Xeon X5675) does not make a difference for streaming, since one hits the

memory and I/O bandwidth limits (and essentially encounters Amdahl's law). More advanced configurations, such as RAID, would allow us to scale better as the number of cores become large.

Figure 2 shows the effects of different GPU optimizations on a C2070 GPU. The GPU base version performs the computation using an 8×8 work group. The GPU shared version places the label, strength, and image data into shared memory before performing the computation. Since error can be tolerated, the GPU ~Hypot version approximates the 2-norm using the " α -max plus β -min algorithm" with $\alpha = 1$ and $\beta = 1/2$. The final version, GPU Approx, introduces more error by performing multiple iterations without doing a global synchronization. This shows that beyond program tuning for hardware, one can exploit the error tolerance for this class of algorithms to achieve great speedup.

3 FUTURE WORK

One key observation in segmentation is that the labeling is decided in the first few iterations of the algorithm. A simple optimization to reduce our error while approximating the algorithm is to start the algorithm with a global synchronization after each iterations (computing the true solution). At some time t we can more safely perform the approximation by having global synchronization after n iterations of the algorithm. At a later time $t + s$, we can perform the synchronization after $2n$ iterations.

REFERENCES

- [1] Vladimir Vezhnevets and Vadim Konouchine. Growcut: Interactive multi-label nd image segmentation by cellular automata. In *Proc. of Graphicon*, pages 150–156, 2005.
- [2] James Reinders. *Intel threaded building blocks*. O'Reilly, 2007.

4 APPENDIX

Listing 2: Final optimized and approximated code

```

1
2 #include <limits.h>
3 #include <cstdio>
4 #include <cstdlib>
5 #include <cstring>
6 #include <iostream>
7 #include <fstream>
8
9 #include "assert.h"
10 #include "dataset.h"
11 #include "mfi.h"
12 #include "timer.h"
13
14 #include "math.h"
15 #include "cuda.h"
16 #include "cuda_runtime.h"
17 #include "device_launch_parameters.h"
18
19 using namespace std;
20
21 #define RADIUS 1
22 #define MAX_ITERATIONS 2048 / 32
23 #define BLOCKDIM_X 8
24 #define BLOCKDIM_Y 8
25
26 template <typename T0, typename T1> static T0 zCeil(const T0 &n, const T1 &d) {
27     return (T0) ceil(static_cast<double>(n) / static_cast<double>(d));
28 }
29
30 #define cudaCheck(stmt) _cudaCheck(__LINE__, stmt)
31 #define _cudaCheck(line, stmt) \
32     do { \
33         cudaError_t p_err = stmt; \
34         assert(p_err == cudaSuccess); \
35         if (p_err != cudaSuccess) { \
36             printf("ERROR_on_line_%d(%s):_%s_--_%d\n", line, #stmt, \
37                 cudaGetErrorString(p_err), p_err); \
38         } \
39     } while (0)
40
41 #define _max(x, y) (((x) > (y)) ? x : y)
42 #define _min(x, y) (((x) < (y)) ? x : y)
43 #define g(x, y) (1 - _max(x, y) - _min(x, y) / 2)
44
45 __global__ void growcut(unsigned short const *__restrict__ img, char *label,
46                        float *strength, int height, int width) {
47
48     __shared__ unsigned short
49         imgShared[BLOCKDIM_Y + 2 * RADIUS][BLOCKDIM_X + 2 * RADIUS];
50     __shared__ char lShared[BLOCKDIM_Y + 2 * RADIUS][BLOCKDIM_X + 2 * RADIUS];
51     __shared__ float sShared[BLOCKDIM_Y + 2 * RADIUS][BLOCKDIM_X + 2 * RADIUS];
52     __shared__ char ltShared[BLOCKDIM_Y][BLOCKDIM_X];
53     __shared__ float stShared[BLOCKDIM_Y][BLOCKDIM_X];
54
55     int tidX = threadIdx.x;
56     int tidY = threadIdx.y;
57     int tx = tidX;
58     int ty = tidY;
59     int jj = tidX + BLOCKDIM_X * blockIdx.x;
60     int ii = tidY + BLOCKDIM_Y * blockIdx.y;
61
62     #define idx(array, y, x) \
63         ((y) >= 0 && (y) < height && (x) >= 0 && (x) < width) \
64         ? array[(y) * width + (x)] \
65         : 0
66
67     tidY += RADIUS;
68     tidX += RADIUS;
69
70     imgShared[tidY][tidX] = idx(img, ii, jj);
71     sShared[tidY][tidX] = idx(strength, ii, jj);
72     lShared[tidY][tidX] = idx(label, ii, jj);
73
74     if (tx < RADIUS) {

```

```

75     imgShared[tidY][tidX - RADIUS] = idx(img, ii, jj - RADIUS);
76     sShared[tidY][tidX - RADIUS] = idx(strength, ii, jj - RADIUS);
77     lShared[tidY][tidX - RADIUS] = idx(label, ii, jj - RADIUS);
78
79     imgShared[tidY][tidX + BLOCKDIM_X] = idx(img, ii, jj + BLOCKDIM_X);
80     sShared[tidY][tidX + BLOCKDIM_X] = idx(strength, ii, jj + BLOCKDIM_X);
81     lShared[tidY][tidX + BLOCKDIM_X] = idx(label, ii, jj + BLOCKDIM_X);
82 }
83
84 if (ty < RADIUS) {
85     imgShared[tidY - RADIUS][tidX] = idx(img, ii - RADIUS, jj);
86     sShared[tidY - RADIUS][tidX] = idx(strength, ii - RADIUS, jj);
87     lShared[tidY - RADIUS][tidX] = idx(label, ii - RADIUS, jj);
88
89     imgShared[tidY + BLOCKDIM_Y][tidX] = idx(img, ii + BLOCKDIM_Y, jj);
90     sShared[tidY + BLOCKDIM_Y][tidX] = idx(strength, ii + BLOCKDIM_Y, jj);
91     lShared[tidY + BLOCKDIM_Y][tidX] = idx(label, ii + BLOCKDIM_Y, jj);
92 }
93
94 for (int kk = 0; kk < 16; kk++) {
95     __syncthreads();
96     char nl = 0;
97     float ns = 0;
98     if (jj < width && ii < height) {
99         char lq;
100         float sq, gc;
101         unsigned short cq;
102
103         unsigned short cp = imgShared[tidY][tidX];
104         char lp = lShared[tidY][tidX];
105         float sp = sShared[tidY][tidX];
106         nl = lp;
107         ns = sp;
108
109         cq = imgShared[tidY - 1][tidX];
110         lq = lShared[tidY - 1][tidX];
111         sq = sShared[tidY - 1][tidX];
112         gc = g(cp, cq) * sq;
113         if (gc > sp) {
114             nl = lq;
115             ns = gc;
116         }
117
118         cq = imgShared[tidY + 1][tidX];
119         lq = lShared[tidY + 1][tidX];
120         sq = sShared[tidY + 1][tidX];
121         gc = g(cp, cq) * sq;
122         if (gc > sp) {
123             nl = lq;
124             ns = gc;
125         }
126
127         cq = imgShared[tidY][tidX - 1];
128         lq = lShared[tidY][tidX - 1];
129         sq = sShared[tidY][tidX - 1];
130         gc = g(cp, cq) * sq;
131         if (gc > sp) {
132             nl = lq;
133             ns = gc;
134         }
135
136         cq = imgShared[tidY][tidX + 1];
137         lq = lShared[tidY][tidX + 1];
138         sq = sShared[tidY][tidX + 1];
139         gc = g(cp, cq) * sq;
140         if (gc > sp) {
141             nl = lq;
142             ns = gc;
143         }
144     }
145
146     ltShared[ty][tx] = nl;
147     stShared[ty][tx] = ns;
148
149     __syncthreads();
150
151     lShared[tidY][tidX] = ltShared[ty][tx];
152     sShared[tidY][tidX] = stShared[ty][tx];
153
154 }

```

```

155     if (tx < RADIUS) {
156         sShared[tidY][tidX - RADIUS] =
157             (sShared[tidY][tidX - RADIUS] + stShared[ty][tx]) / 2;
158         lShared[tidY][tidX - RADIUS] =
159             (lShared[tidY][tidX - RADIUS] + ltShared[ty][tx]) / 2;
160
161         sShared[tidY][tidX + BLOCKDIM_X] =
162             (sShared[tidY][tidX + BLOCKDIM_X] +
163              stShared[ty][BLOCKDIM_X - tx - RADIUS]) / 2;
164         lShared[tidY][tidX + BLOCKDIM_X] =
165             (lShared[tidY][tidX + BLOCKDIM_X] +
166              ltShared[ty][BLOCKDIM_X - tx - RADIUS]) / 2;
167     }
168
169     if (ty < RADIUS) {
170         sShared[tidY - RADIUS][tidX] =
171             (sShared[tidY - RADIUS][tidX] + stShared[ty][tx]) / 2;
172         lShared[tidY - RADIUS][tidX] =
173             (lShared[tidY - RADIUS][tidX] + ltShared[ty][tx]) / 2;
174
175         sShared[tidY + BLOCKDIM_Y][tidX] =
176             (sShared[tidY + BLOCKDIM_Y][tidX] +
177              stShared[BLOCKDIM_Y - ty - RADIUS][tx]) / 2;
178         lShared[tidY + BLOCKDIM_Y][tidX] =
179             (lShared[tidY + BLOCKDIM_Y][tidX] +
180              ltShared[BLOCKDIM_Y - ty - RADIUS][tx]) / 2;
181     }
182 }
183
184 if (jj < width && ii < height) {
185     label[ii * width + jj] = ltShared[ty][tx];
186     strength[ii * width + jj] = stShared[ty][tx];
187 }
188 return;
189 }
190
191 int runGrowCut(MFI *mfi, char *label, int *iterations0) {
192     int iterations = 0;
193     int width = mfi->width;
194     int height = mfi->height;
195     char *nextLabel = (char *)malloc(sizeof(char) * width * height);
196     float *strength = (float *)malloc(sizeof(float) * width * height);
197     float *nextStrength = (float *)malloc(sizeof(float) * width * height);
198     unsigned short *cap_source = (unsigned short *)mfi->cap_source;
199     unsigned short *cap_sink = (unsigned short *)mfi->cap_sink;
200     int len = width * height;
201     for (int ii = 0; ii < len; ii++) {
202         float s = label[ii] != 0;
203         strength[ii] = s;
204     }
205
206     unsigned short *dcapsource;
207     float *dStrength;
208     char *dLabel;
209
210     cudaCheck(cudaMalloc(&dcapsource, sizeof(unsigned short) * len));
211     cudaCheck(cudaMalloc(&dStrength, sizeof(float) * len));
212     cudaCheck(cudaMalloc(&dLabel, sizeof(char) * len));
213
214     cudaCheck(cudaMemcpy(dcapsource, mfi->cap_source,
215                          sizeof(unsigned short) * len, cudaMemcpyHostToDevice));
216     cudaCheck(
217         cudaMemcpy(dLabel, label, sizeof(char) * len, cudaMemcpyHostToDevice));
218     cudaCheck(cudaMemcpy(dStrength, strength, sizeof(float) * len,
219                          cudaMemcpyHostToDevice));
220
221     dim3 blockDim(BLOCKDIM_X, BLOCKDIM_Y);
222     dim3 gridDim(zCeil(width, blockDim.x), zCeil(height, blockDim.y));
223
224     while (iterations++ < MAX_ITERATIONS) {
225         growcut << <gridDim, blockDim>>>
226             (dcapsource, dLabel, dStrength, height, width);
227         cudaCheck(cudaThreadSynchronize());
228     }
229
230     cudaCheck(
231         cudaMemcpy(label, dLabel, sizeof(char) * len, cudaMemcpyDeviceToHost));
232
233     cudaFree(dcapsource);
234     cudaFree(dStrength);

```

```

235     cudaFree(dLabel);
236
237     free(nextLabel);
238     free(strength);
239     free(nextStrength);
240
241     *iterations0 = iterations;
242     return -1;
243 }
244
245 int main(int argc, char **argv) {
246
247     const char *dataset_path =
248         argc == 2 ? argv[1] : "C:\\Users\\abduld\\Documents\\visual_studio_
249                             "2012\\Projects\\growcut\\x64\\Debug\\dataset";
250
251     int num_instances = (sizeof(instances) / sizeof(Instance));
252
253     ofstream timesFile;
254     string timeFileName = string(dataset_path);
255     timeFileName.append("\\\\.\\times_cuda_opt_c2070_alpha_max_itermore.data");
256     timesFile.open(timeFileName, ios::out);
257     timesFile << "instance,num,width,height,changes,iterations,time\n";
258
259     for (int i = 0; i < num_instances; i++) {
260
261         for (int j = 0; j < instances[i].count; j++) {
262             char fileName[1024];
263
264             ofstream myfile;
265             string sfileName;
266             int iterations;
267             unsigned short *cap_source;
268             unsigned short *cap_sink;
269             char *label;
270             int changes;
271             uint64_t tic, toc;
272             double compute_time;
273
274             sprintf(fileName, instances[i].filename, dataset_path, j);
275             MFI *mfi = mfi_read(fileName);
276
277             if (!mfi) {
278                 //printf("FAILED to read instance %s\n", fileName);
279                 goto skip;
280             }
281
282             if (mfi->connectivity != 4 || mfi->dimension != 2 ||
283                 mfi->type_terminal_cap != MFI::TYPE_UINT16 ||
284                 mfi->type_neighbor_cap != MFI::TYPE_UINT16) {
285                 goto skip;
286             }
287
288             cap_source = (unsigned short *)mfi->cap_source;
289             cap_sink = (unsigned short *)mfi->cap_sink;
290             label = (char *)calloc(mfi->width * mfi->height, sizeof(char));
291             for (int ii = 0; ii < mfi->width * mfi->height; ii++) {
292                 if (cap_source[ii] > 15000) {
293                     label[ii] = 1;
294                 }
295                 if (cap_sink[ii] > 15000) {
296                     label[ii] = -1;
297                 }
298             }
299
300             tic = _hrtime();
301             changes = runGrowCut(mfi, label, &iterations);
302             toc = _hrtime();
303
304             compute_time = (toc - tic) / 1000000000.0f;
305
306             sfileName = string(fileName);
307             sfileName.append(".dat");
308
309             myfile.open(sfileName, ios::out);
310
311             for (int ii = 0; ii < mfi->height; ii++) {
312                 for (int jj = 0; jj < mfi->width; jj++) {
313                     myfile << ((int) label[ii * mfi->width + jj] + 2) / 2 << " ";
314                     //myfile << ((int) cap_sink[ii*mfi->width + jj]) << " ";

```



```

315     }
316     myfile << "\n";
317 }
318 myfile.close();
319
320 printf("%s,%d,%d,%d,%d,%d,%0.9f_\n", instances[i].name, j, mfi->width,
321        mfi->height, changes, iterations, compute_time);
322 timesFile << instances[i].name << "," << j << "," << mfi->width << ","
323        << mfi->height << "," << changes << "," << iterations << ","
324        << compute_time << "\n";
325 timesFile.flush();
326 free(label);
327 skip:
328     if (mfi != NULL)
329         mfi_free(mfi);
330
331 }
332 }
333 timesFile.close();
334 return 0;
335 }

```