

Lyon Kee

[keel@oregonstate.edu](mailto:keel@oregonstate.edu)

Project #6: [OpenCL](#)

## Key snippets of code

1. Program header
2. Allocate the host memory buffers
3. Create an OpenCL context
4. Create an OpenCL command queue
5. Allocate the device memory buffers
6. Write the data from the host buffers to the device buffers
7. Read the kernel code from a file
8. Compile and link the kernel code
9. Create the kernel object
10. Setup the arguments to the kernel object
11. Enqueue the kernel object for execution
12. Read the results buffer back from the device to the host
13. Clean everything up

### 1. Program header

```
1 // 1. Program header
2
3 #include <stdio.h>
4 #include <math.h>
5 #include <string.h>
6 #include <stdlib.h>
7 #include <omp.h>
8
9 #include "cl.h"
10 #include "cl_platform.h"
```

We observe that `cl.h` and `cl_platform.h` is added which is needed for OpenCL, OpenMP for time measurement, and the rest are the standard libraries for operations.

### 2. Allocate the host memory buffers

```
106 // 2. create the host memory buffers:
107
108 // read the data file:
109
110 FILE* fdata;
111 #ifdef _WIN32
112 err = fopen_s(&fdata, DATAFILE, "r");
113 if (err != 0)
114 #else
115 fdata = fopen(DATAFILE, "r");
116 if (fdata == NULL)
117 #endif
118 {
119     fprintf(stderr, "Cannot open data file '%s'\n", DATAFILE);
120     return -1;
121 }
122
123 float x, y;
124 for( int i = 0; i < DATASIZE; i++ )
125 {
126     #ifdef _WIN32
127         fscanf_s(fdata, "%f %f", &x, &y);
128     #else
129         fscanf( fdata, "%f %f", &x, &y );
130     #endif
131     hX[i] = x;
132     hY[i] = y;
133 }
134 fclose( fdata );
```

Before this, we check if kernel file is readable and also select the best OpenCL device. We read in data and store it into host device arrays of `hX` and `hY`.

### 3. Create an OpenCL context

```
137 // 3. create an opengl context:
138
139 Context = clCreateContext( NULL, 1, &Device, NULL, NULL, &status );
140 if( status != CL_SUCCESS )
141     fprintf( stderr, "clCreateContext failed\n" );
142
```

Created an OpenCL context and set it to `cl\_context context` and check for CL\_SUCCESS. Check for CL\_SUCCESS before we proceed.

### 4. Create an OpenCL command queue

```
144 // 4. create an opengl command queue:
145
146 CmdQueue = clCreateCommandQueue( Context, Device, 0, &status );
147 if( status != CL_SUCCESS )
148     fprintf( stderr, "clCreateCommandQueue failed\n" );
```

Create a cl\_command\_queue to store commands for OpenCL to quickly grab commands to execute. Check for CL\_SUCCESS before we proceed.

### 5. Allocate the device memory buffers

```
151 // 5. allocate the device memory buffers:
152
153 size_t xySize = DATASIZE * sizeof(float);
154
155 cl_mem dx = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
156 cl_mem dy = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
157 cl_mem dSumx2 = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
158 cl_mem dSumx = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
159 cl_mem dSumxy = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
160 cl_mem dSumy = clCreateBuffer( Context, CL_MEM_READ_ONLY, xySize, NULL, &status );
161
162 if( status != CL_SUCCESS )
163     fprintf( stderr, "clCreateBuffer failed\n" );
```

Used clCreateBuffer to allocate memory on the OpenCL device to allow the device to read and write to when it is doing computation. Check for CL\_SUCCESS before we proceed.

### 6. Write the data from the host buffers to the device buffers

```
166 // 6. enqueue the 2 commands to write the data from the host buffers to the device buffers:
167
168 status = clEnqueueWriteBuffer( CmdQueue, dx, CL_FALSE, 0, DATASIZE, hX, 0, NULL, NULL );
169 if( status != CL_SUCCESS )
170     fprintf( stderr, "clEnqueueWriteBuffer failed (1)\n" );
171
172 status = clEnqueueWriteBuffer( CmdQueue, dy, CL_FALSE, 0, DATASIZE, hY, 0, NULL, NULL );
173 if( status != CL_SUCCESS )
174     fprintf( stderr, "clEnqueueWriteBuffer failed (2)\n" );
175
176 Wait( CmdQueue );
```

Writing data loaded from file -> host memory to device memory and waiting for it to finish writing.

### 7. Read the kernel code from a file

```

179 // 7. read the kernel code from a file ...
180
181 fseek( fp, 0, SEEK_END );
182 size_t fileSize = ftell( fp );
183 fseek( fp, 0, SEEK_SET );
184 char *clProgramText = new char[ fileSize+1 ]; // leave room for '\0'
185 size_t n = fread( clProgramText, 1, fileSize, fp );
186 clProgramText[fileSize] = '\0';
187 fclose( fp );
188 if( n != fileSize )
189     fprintf( stderr, "Expected to read %d bytes read from '%s' -- actually read %d.\n", (int)fileSize, CL_FILE_NAME, (int)n );
190

```

Reading our OpenCL code to create a kernel program

## 8. Compile and link the kernel code

```

201 // 8. compile and link the kernel code:
202
203 char *options = { (char *)"" };
204 status = clBuildProgram( Program, 1, &Device, options, NULL, NULL );
205 if( status != CL_SUCCESS )
206 {
207     size_t size;
208     clGetProgramBuildInfo( Program, Device, CL_PROGRAM_BUILD_LOG, 0, NULL, &size );
209     cl_char *log = new cl_char[ size ];
210     clGetProgramBuildInfo( Program, Device, CL_PROGRAM_BUILD_LOG, size, log, NULL );
211     fprintf( stderr, "clBuildProgram failed:\n%s\n", log );
212     delete [ ] log;
213 }
214

```

Compile and link the kernel code that is loaded to a string in step 7

## 9. Create the kernel object

```

216 // 9. create the kernel object:
217
218 Kernel = clCreateKernel( Program, "Regression", &status );
219 if( status != CL_SUCCESS )
220     fprintf( stderr, "clCreateKernel failed\n" );
221

```

Create a kernel from the compiled and linked program in step 8.

## 10. Setup the arguments to the kernel object

```

223 // 10. setup the arguments to the kernel object:
224
225 status = clSetKernelArg( Kernel, 0, sizeof(cl_mem), &dX );
226 status = clSetKernelArg( Kernel, 1, sizeof(cl_mem), &dY );
227
228 status = clSetKernelArg( Kernel, 2, sizeof(cl_mem), &dSumx2 );
229 status = clSetKernelArg( Kernel, 3, sizeof(cl_mem), &dSumx );
230 status = clSetKernelArg( Kernel, 4, sizeof(cl_mem), &dSumxy );
231 status = clSetKernelArg( Kernel, 5, sizeof(cl_mem), &dSumy );
232 if( status != CL_SUCCESS )
233     fprintf( stderr, "clSetKernelArg failed: %d\n", status );
234

```

Set kernel arguments which should align with the kernel program arguments below:

```
4  kernel
5  void
6  Regression( IN global const float *dX,
7              IN global const float *dY,
8              OUT global float *dSumx2,
9              OUT global float *dSumx,
10             OUT global float *dSumxy,
11             OUT global float *dSumy )
12  {
13      int gid = get_global_id( 0 );
14
15      float x = dX[gid];
16      float y = dY[gid];
17      dSumx2[ gid ] = x * x;
18      dSumx[ gid ] = x;
19      dSumxy[ gid ] = x * y;
20      dSumy[ gid ] = y;
21  }
22
```

### 11. Enqueue the kernel object for execution

```
235      // 11. enqueue the kernel object for execution:
236
237      size_t globalWorkSize[3] = { DATASIZE, 1, 1 };
238      size_t localWorkSize[3] = { LOCALSIZE, 1, 1 };
239
240      Wait( CmdQueue );
241
242      double time0 = omp_get_wtime( );
243
244      status = clEnqueueNDRangeKernel( CmdQueue, Kernel, 1, NULL, globalWorkSize, localWorkSize, 0, NULL, NULL );
245      if( status != CL_SUCCESS )
246          fprintf( stderr, "clEnqueueNDRangeKernel failed: %d\n", status );
247
248      Wait( CmdQueue );
249      double time1 = omp_get_wtime( );
```

Enqueue the Kernel and time the execution after waiting for it to finish running.

### 12. Read the results buffer back from the device to the host

```
252      // 12. read the results buffer back from the device to the host:
253
254      status = clEnqueueReadBuffer( CmdQueue, dSumx2, CL_FALSE, 0, xySize, hSumx2, 0, NULL, NULL );
255      status = clEnqueueReadBuffer( CmdQueue, dSumx, CL_FALSE, 0, xySize, hSumx, 0, NULL, NULL );
256      status = clEnqueueReadBuffer( CmdQueue, dSumxy, CL_FALSE, 0, xySize, hSumxy, 0, NULL, NULL );
257      status = clEnqueueReadBuffer( CmdQueue, dSumy, CL_FALSE, 0, xySize, hSumy, 0, NULL, NULL );
258
259      Wait( CmdQueue );
```

Read the output of the Kernel back to host device to compute sum.

### 13. Clean everything up

```
288      // 13. clean everything up:
289
290      clReleaseKernel( Kernel );
291      clReleaseProgram( Program );
292      clReleaseCommandQueue( CmdQueue );
293      clReleaseMemObject( dSumx2 );
294      clReleaseMemObject( dSumx );
295      clReleaseMemObject( dSumxy );
296      clReleaseMemObject( dSumy );
297
```

Clean up created objects relating to the OpenCL device.

## Tables of data

Columns: Array Size

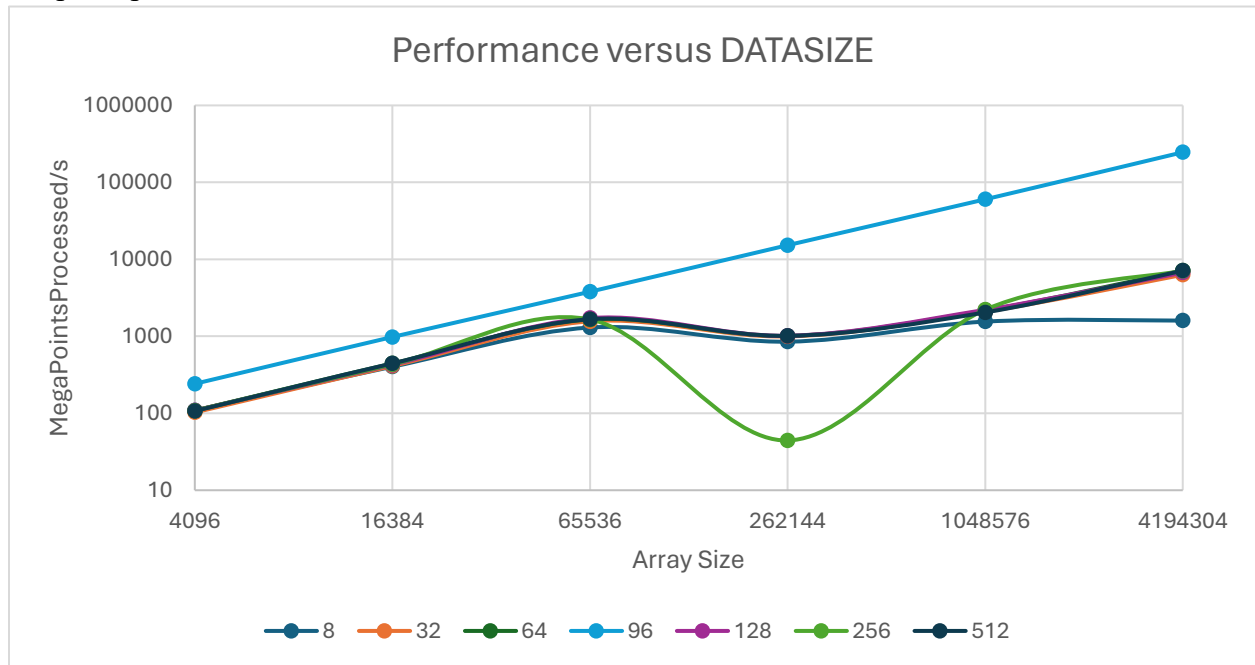
Rows: Work Elements

Values: megaPointsProcessedPerSecond

Sum of MegaPointsProcessedPerSecond	Array Size					
Work Elements	4096	16384	65536	262144	1048576	4194304
8	108.68	402.71	1301.2	851.84	1553.27	1599.62
32	102.55	411.84	1567.16	997.16	2084.37	6297.17
64	108.28	443.89	1657.76	1011.41	2162	7150.43
96	241.49	979.85	3793.46	15259.4	60221.43	246152.14
128	108.87	426.04	1725.4	1016.4	2207.97	6726.27
256	108.27	433.16	1656.98	44.35	2220.25	7087.84
512	107.01	444.34	1669.25	1014.2	2042.92	7129.99

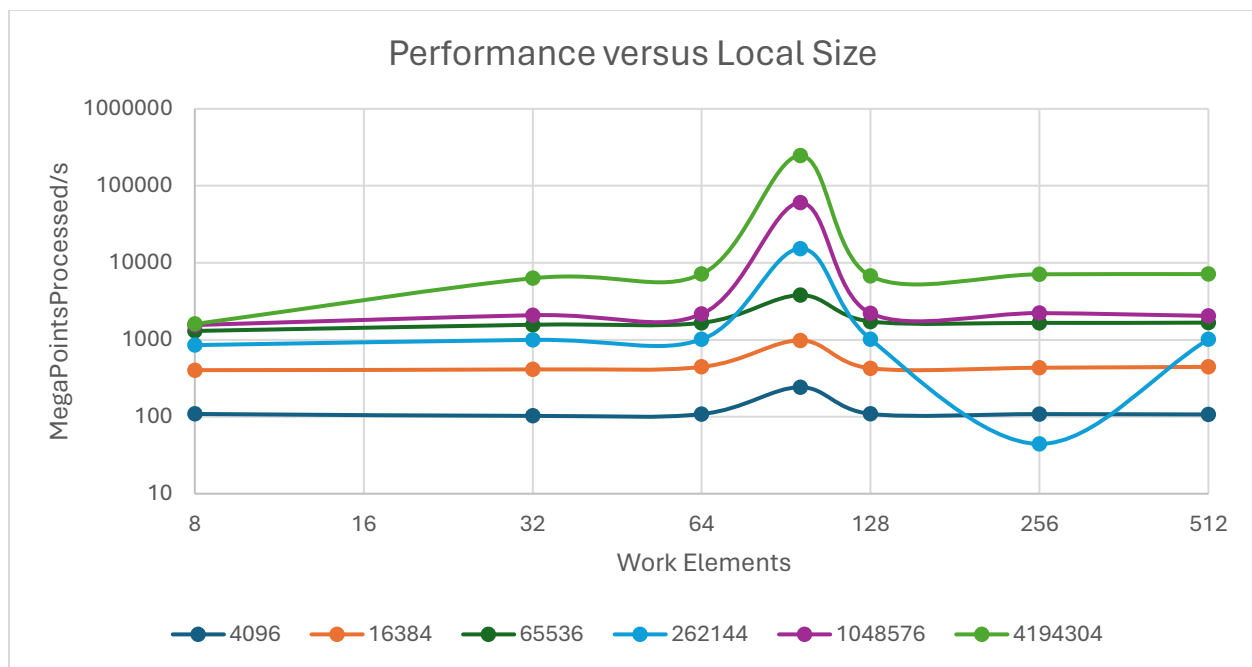
## Graphs of data

Graph of performance versus DATASIZE with colored curves of constant LOCALSIZE



In the graph above, we observe that as the Array Size increases, the performance increases as well, although we do observe that given specific “Work Elements” or constant local sizes, the performance peaks and stays constant even when we increase DATASIZE, this is likely the peak performance of the device with data being abundant so there is little inefficiency.

Graph of performance versus LOCALSIZE with colored curves of constant DATASIZE



The graph above shows that there is a sweet spot in the number of work elements, and a number above or below it will yield a similar performance. We observe that the best performance is at 96 work elements, from lecture it seems that each “gripper” has 32 work elements, and at 64, we have 2 sets which allows the scheduler to swap out memory when it is performing tasks that require waiting and it just seems like in this system, 96 was the perfect number which yielded the best performance.

## PDF Commentary

1. What machine you ran this on

I ran this on the one Tesla V100 of the A100s that the class has access to, therefore I get dedicated CPU&GPU time and it is very reliable.

2. Show the table and the graphs

Shown in the above sections.

3. What patterns are you seeing in the performance curves? What difference does the size of data make? What difference does the size of each work-group make?

Graph analysis and evaluations are shown in the above sections.