# Advanced OpenCL Programming

Pangfeng Liu
National Taiwan University

June 7, 2015

## Multiple Devices

- We have been using one device for computation.
- Now we will use multiple device to solve a problem.
- We will use matrix multiplication as an example.

## Partition

- Since we are very poor and cannot afford four GPUs, we will partition the data by rows (not by block) among devices.
- We will partition matrix by rows, so that the memory assigned to a device is contiguous.
- Each kernel will run on a device, and compute part of the answers.

## Partition

- We will partition $A$ and $C$ by rows, and do *not* partition $B$.
- For example, if we use two devices, then the first kernel will compute the top half of $C$, and another kernel will compute the bottom half of $C$.
- The reason for not partitioning $B$ is that all kernels need the entire $B$.

## Constants

- MAXLOG is the maximum number of bytes for storing compilation log. More on this later.
- The number of device is DEVICENUM. We will use 2 in our humble installation.
- ITEMPERDEVICE is the number of work item in NDRange of a kernel.

## Constants

**Example 1:   (matrixMul-time-copy-local-multidevice.c)**

```
1   #define CL_USE_DEPRECATED_OPENCL_2_0_APIS
2   #include <stdio.h>
3   #include <assert.h>
4   #include <CL/cl.h>
5
6   #define N 1024
7   #define Blk 64
8   #define BSIDE (N / Blk)
9   #define MAXGPU 10
10  #define MAXK 1024
11  #define MAXLOG 4096
12  #define DEVICENUM 2
13  #define ITEMPERDEVICE (N * N / DEVICENUM)
14  #define NANO2SECOND 1000000000.0
15
16  cl_uint A[N][N], B[N][N], C[N][N];
```

## Discussion

- If $N$ is 1024, and the number of device is 2, then how many works items are there in a kernel?

## Select Devices

- We will select the first DEVICENUM GPU device from our humble installation.
- The code is similar to those earlier version, except that now we need to make sure that the number of GPUs found is at least DEVICENUM.

## Devices

**Example 2: (matrixMul-time-copy-local-multidevice.c)**

```
29    cl_device_id GPU[MAXGPU];
30    cl_uint GPU_id_got;
31    status =
32      clGetDeviceIDs(platform_id, CL_DEVICE_TYPE_GPU,
33                     MAXGPU, GPU, &GPU_id_got);
34    assert(status == CL_SUCCESS &&
35           GPU_id_got >= DEVICENUM);
36    printf("There are %d GPU devices\n", GPU_id_got);
```

## Context

- Since a context can consist of multiple devices, we need only one context.
- Here we include the first DEVICENUM GPUs in the context.

## Context

**Example 3:   (matrixMul-time-copy-local-multidevice.c)**

```
38    cl_context context =
39      clCreateContext(NULL, DEVICENUM, GPU, NULL, NULL,
40                      &status);
41    assert(status == CL_SUCCESS);
```

## Multiple Queue

- Recall that a command queue connects to a device.
- Since we have multiple devices, we need multiple command queues.
- We put the command queues in the commandQueue array.
- Note that we set the CL_QUEUE_PROFILING_ENABLE to enable profiling.

## Command Queue

**Example 4:   (matrixMul-time-copy-local-multidevice.c)**

```
43    cl_command_queue commandQueue[DEVICENUM];
44    for (int device = 0; device < DEVICENUM; device++) {
45      commandQueue[device] =
46        clCreateCommandQueue(context, GPU[device],
47                             CL_QUEUE_PROFILING_ENABLE,
48                             &status);
49      assert(status == CL_SUCCESS);
50    }
```

## One Context, Multiple Devices

- Up to this point, we can image there is only one context, which has multiple devices.
- Within this context we will have one command queue to connect to each device.
- Later we will send commands to these devices. The command sent into a command queue will run on the device this command queue connects to.

## Discussion

- What will be sent into these command queues as commands?

## Program

- The kernel will be sent into the command queue as commands.
- Before this can happen we need to compile the kernel.
- Recall that the "kernel" at this point is only a set of strings, which need to be compiled.

## Build Program

- We call `clBuildProgram` to build executable.
- Note that we pass *all* devices we want to use as parameters, so OpenCL will build executable for all our `DEVICENUM` devices.

## Compilation Log

- The kernel source is compiled when an OpenCL program runs.
- It is very inconvenient because when we run an OpenCL program the execution will not display compilation errors of the kernel.
- We will use `clGetProgramBuildInfo` to show the error message if the compilation of the kernel fails.

**Prototype 5: clGetProgramBuildInfo.h**

```
1   cl_int clGetProgramBuildInfo(cl_program program,
2                                cl_device_id device,
3                                cl_program_build_info param_name,
4                                size_t param_value_size,
5                                void *param_value,
6                                size_t *param_value_size_ret);
```

## Parameters

program The compiled program.

device The device you want to query.

param_name The information you want to query.

param_value_size The length of param_value in bytes.

param_value The location for the query answer.

param_value_size_ret The number of bytes returned from the query.

## Error

- When there is an error in clBuildProgram, we will call clGetProgramBuildInfo to find out.
- Since we are not sure about which device causes the compilation error, we list the information from all of them.
- We query with CL PROGRAM BUILD LOG for the compilation log, and place it into the log buffer we prepared.

## Build Program

**Example 6:    (matrixMul-time-copy-local-multidevice.c)**

```
64    status =
65      clBuildProgram ( program , DEVICENUM , GPU , NULL ,
66                       NULL , NULL );
67    if ( status != CL_SUCCESS ) {
68      char log [ MAXLOG ];
69      size_t logLength ;
70      for ( int device = 0; device < DEVICENUM ; device ++) {
71        clGetProgramBuildInfo ( program , GPU [ device ],
72                               CL_PROGRAM_BUILD_LOG ,
73                               MAXLOG , log , & logLength );
74        puts ( log );
75      }
76      exit ( -1);
77    }
78    printf ( " Build program completes \n " );
```

## Discussion

- What is the type of the compilation log?
- What does puts do?

## Partition

- Now we are ready to partition host buffers $A$, $B$ and $C$.
- We will create DEVICENUM buffers for $A$, and each of which will be given to a kernel as a parameter for computations.
- We partition $A$ into sub-matrix of N / DEVICENUM rows of $A$ each, and assign each partition to an OpenCL buffer.

## Buffer A

- ITEMPERDEVICE means the number of items per device.
- The index device determines the starting position of a buffer, i.e., where the matrix $A$ will be partitioned.

## Buffer for A

**Example 7:   (matrixMul-time-copy-local-multidevice.c)**

```
90    cl_mem bufferA[DEVICENUM];
91    for (int device = 0; device < DEVICENUM; device++) {
92      bufferA[device] =
93        clCreateBuffer(context,
94          CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
95          ITEMPERDEVICE * sizeof(cl_uint),
96          ((cl_uint *)A) + device * ITEMPERDEVICE,
97          &status);
98      assert(status == CL_SUCCESS);
99    }
```

## Buffer B

- We do not need to partition $B$ because it will be used by all kernels.

## Buffer for B

**Example 8:   (matrixMul-time-copy-local-multidevice.c)**

```
101    cl_mem bufferB =
102      clCreateBuffer ( context ,
103        CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR ,
104        N * N * sizeof ( cl_uint ), B, & status );
105    assert ( status == CL_SUCCESS );
```

# Buffer C

- Buffers for $C$ are similarly built as for $A$.
- Unlike $A$ and $B$ that will copy host buffers to device buffers, $C$ will use host memory directly.

## Buffer for C

**Example 9:  (matrixMul-time-copy-local-multidevice.c)**

```
107    cl_mem bufferC[DEVICENUM];
108    for (int device = 0; device < DEVICENUM; device++) {
109      bufferC[device] =
110        clCreateBuffer(context,
111          CL_MEM_WRITE_ONLY | CL_MEM_USE_HOST_PTR,
112          ITEMPERDEVICE * sizeof(cl_uint),
113          ((cl_uint *) C) + device * ITEMPERDEVICE,
114          &status);
115      assert(status == CL_SUCCESS);
116    }
117    printf("Build buffers completes\n");
```

## Discussion

- Make sure that you understand the offset calculation in $A$ and $C$.

## NDRange

- According to our partition we declare NDRange as a two dimension array with N / DEVICENUM rows and N columns.
- This is consistent with $C$, where each work item will compute an element of $C$.
- The work group is still BSIDE by BSIDE since we will use the previous local memory algorithm.
- We will have DEVICENUM kernels, so we will need DEVICENUM events for synchronization and profiling.

## NDRange

**Example 10: (matrixMul-time-copy-local-multidevice.c)**

```
119    size_t globalThreads[] =
120      {(size_t)(N / DEVICENUM), (size_t)N};
121    size_t localThreads[] = {BSIDE, BSIDE};
122    cl_event events[DEVICENUM];
```

## Set Argument

- We will launch DEVICENUM kernels. Each of them will compute N / DEVICENUM rows and N columns of $C$.

- We loop through all devices, and assign the argument order as $A$, $B$, and $C$.

- Note that we need to supply the correct $A$ and $C$ since there are DEVICENUM of them.

## Set Arguments

**Example 11: (matrixMul-time-copy-local-multidevice.c)**

```
124    for (int device = 0; device < DEVICENUM; device++) {
125       status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
126                               (void*)(&bufferA[device]));
127       assert(status == CL_SUCCESS);
128       status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
129                               (void*)&bufferB);
130       assert(status == CL_SUCCESS);
131       status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
132                               (void*)(&bufferC[device]));
133       assert(status == CL_SUCCESS);
134       printf("Set kernel arguments completes\n");
```

## Kernel

- Now we are ready to launch kernel.
- Note that we need to supply the command queue to the device as a parameters.
- The same kernel is launched again and again for DEVICENUM times. The only difference is the buffers $A$ and $C$ used in launching the kernel.
- We associate an event with each kernel launching.

## Start Kernel

**Example 12:  (matrixMul-time-copy-local-multidevice.c)**

```
136        status =
137          clEnqueueNDRangeKernel ( commandQueue [ device ] ,
138                                   kernel , 2 , NULL ,
139                                   globalThreads , localThreads ,
140                                   0 , NULL , &( events [ device ] ) );
141        assert ( status == CL_SUCCESS );
142      }
```

## Wait for Completion

- Since kernel launching is non-blocking. we need to wait for them to complete.
- We simply use `clWaitForEvents`, and wait for all events to complete.

## Wait for Kernels

**Example 13:   (matrixMul-time-copy-local-multidevice.c)**

```
144    clWaitForEvents ( DEVICENUM , events );
145    printf ( "Kernel execution completes .\n" );
```

## Discussion

- How many events do we need to wait for the whole thing to complete?

## Kernel

- The kernel function is very similar to the previous local memory version.
- We also use the constant DEVICENUM to denote the number of devices.

## Kernel

**Example 14: (mul-local-multidevice-kernel.cl)**

```
1  #define N 1024
2  #define Blk 64
3  #define DEVICENUM 2
4  #define BSIDE (N / Blk)
```

## Kernel

Example 15: (mul-local-multidevice-kernel.cl)

```
6  __kernel void mul(__global int A[N/DEVICENUM][N],
7                    __global int B[N][N],
8                    __global int C[N/DEVICENUM][N])
9  {
10   int globalRow = get_global_id(0);
11   int globalCol = get_global_id(1);
12   int localRow = get_local_id(0);
13   int localCol = get_local_id(1);
14
15   __local int ALocal[BSIDE][BSIDE];
16   __local int BLocal[BSIDE][BSIDE];
```

## NDRange

- We declare $A$ and $C$ as N / DEVICENUM by N arrays, because the domain is partitioned among devices.
- The rest of the code is *not* changed at all!

## Reason

- The reason for this "no need to change" is that a kernel function is acting from a local view, i.e., it is a work item within a work group, so all the operations are still the same.
- The kernel still thinks from the point of view of blocks, and the related operations are done in local indices.
- We only need the global index while accessing $A$, $B$, and $C$.
- Since the kernel is given the corresponding parts of $A$ and $B$, so the operation is correct.

## Discussion

- Convince yourself that the code is correct.

## Timing

- We would like to know the detailed timing of the kernels.
- In particular we would like to prove the two kernel runs *in parallel*.
- To do so we need the absolute time of the events from all devices.

## Base

- The time returned by the event is difficult to interpret.
- We would like to establish a relative time for inspection.
- We choose the time the first kernel entering the queue as the base time, and report the relative time to it for ease understanding.

## Get Time

**Example 16:   (matrixMul-time-copy-local-multidevice.c)**

```
147    cl_ulong base;
148    status =
149      clGetEventProfilingInfo(events[0],
150        CL_PROFILING_COMMAND_QUEUED,
151        sizeof(cl_ulong), &base, NULL);
152    assert(status == CL_SUCCESS);
```

## Get Time

- We loop through all devices and get timing information from all of them.
- This part is the same as before.

## Get Time

**Example 17:  (matrixMul-time-copy-local-multidevice.c)**

```
154    for (int device = 0; device < DEVICENUM; device++) {
155      cl_ulong timeEnterQueue, timeSubmit, timeStart,
156        timeEnd;
157      status =
158        clGetEventProfilingInfo(events[device],
159          CL_PROFILING_COMMAND_QUEUED,
160          sizeof(cl_ulong), &timeEnterQueue, NULL);
161      assert(status == CL_SUCCESS);
162      status =
163        clGetEventProfilingInfo(events[device],
164          CL_PROFILING_COMMAND_SUBMIT,
165          sizeof(cl_ulong), &timeSubmit, NULL);
166      assert(status == CL_SUCCESS);
```

## Get Time

**Example 18:   (matrixMul-time-copy-local-multidevice.c)**

```
168        status =
169          clGetEventProfilingInfo(events[device],
170            CL_PROFILING_COMMAND_START,
171            sizeof(cl_ulong), &timeStart, NULL);
172        assert(status == CL_SUCCESS);
173        status =
174          clGetEventProfilingInfo(events[device],
175            CL_PROFILING_COMMAND_END,
176            sizeof(cl_ulong), &timeEnd, NULL);
177        assert(status == CL_SUCCESS);
```

## Relative Time

- In addition to queuing time, submission time, and execution, we also report the relative time when the four events happened.
- The relative time is in nanosecond.

## Print Time

**Example 19: (matrixMul-time-copy-local-multidevice.c)**

```
179        printf("\nkernel entered queue at %f\n",
180                (timeEnterQueue - base) / NANO2SECOND);
181        printf("kernel submitted to device at %f\n",
182                (timeSubmit - base) / NANO2SECOND);
183        printf("kernel started at %f\n",
184                (timeStart - base) / NANO2SECOND);
185        printf("kernel ended  at %f\n",
186                (timeEnd - base) / NANO2SECOND);
187        printf("kernel queued time %f seconds\n",
188                (timeSubmit - timeEnterQueue) / NANO2SECOND);
189        printf("kernel submission time %f seconds\n",
190                (timeStart - timeSubmit) / NANO2SECOND);
191        printf("kernel execution time %f seconds\n",
192            (timeEnd - timeStart) / NANO2SECOND);
193    }
```

## Demonstration

- Run the `matrixMul-time-copy-local-multidevice-cl` program.

## Discussion

- Does the kernels run in parallel? Observe the results and give your answer.

# Dependency

- We may launch multiple kernels into multiple devices.
- Kernels may have dependency, and we need to ensure dependency among kernels.

## Example

- In the following example we compute
  $G = (A + B) + (D + E)$, where all vector has length $N$.
- We compute $C = A + B$, then $F = D + E$, then we compute
  $G = C + F$. The last computation must wait for the first two
  to complete.
- We will use one device for *each* of the three additions.

## Constants

- We first declare the constants and variables.

**Example 20:   (vectorAdd-dependency.c)**

```
 1  #define CL_USE_DEPRECATED_OPENCL_2_0_APIS
 2  #include <stdio.h>
 3  #include <assert.h>
 4  #include <CL/cl.h>
 5  #define N (65536 * 4)
 6  #define MAXGPU 10
 7  #define MAXK 1024
 8  #define MAXLOG 4096
 9  #define DEVICENUM 3
10  #define NANO2SECOND 1000000000.0
11
12  cl_uint A[N], B[N], C[N], D[N], E[N], F[N], G[N];
```

## Initialization

- We initialize the $A$, $B$, $D$ and $E$ vectors.

**Example 21: (vectorAdd-dependency.c)**

```
81      for (int i = 0; i < N; i++) {
82        A[i] = i;
83        B[i] = N - i;
84        D[i] = i;
85        E[i] = N - i;
86      }
```

## Buffers

- We then create buffers for kernel parameters.
- $A$ and $B$ are read only, and need to be copied into device.
- $C$ is both read and write enable.

**Example 22: (vectorAdd-dependency.c)**

```
88    cl_mem bufferA =
89      clCreateBuffer(context,
90                     CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
91                     N * sizeof(cl_uint), A, &status);
92    assert(status == CL_SUCCESS);
93    cl_mem bufferB =
94      clCreateBuffer(context,
95                     CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
96                     N * sizeof(cl_uint), B, &status);
97    assert(status == CL_SUCCESS);
98    cl_mem bufferC =
99      clCreateBuffer(context,
100                    CL_MEM_READ_WRITE,
101                    N * sizeof(cl_uint), NULL, &status);
102   assert(status == CL_SUCCESS);
```

## Buffers

- Similarly we create buffers for $D$, $E$ and $F$.

**Example 23: (vectorAdd-dependency.c)**

```
104    cl_mem bufferD =
105      clCreateBuffer(context,
106                     CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
107                     N * sizeof(cl_uint), D, &status);
108    assert(status == CL_SUCCESS);
109    cl_mem bufferE =
110      clCreateBuffer(context,
111                     CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
112                     N * sizeof(cl_uint), E, &status);
113    assert(status == CL_SUCCESS);
114    cl_mem bufferF =
115      clCreateBuffer(context,
116                     CL_MEM_READ_WRITE,
117                     N * sizeof(cl_uint), NULL, &status);
118    assert(status == CL_SUCCESS);
```

## Buffers

- We create buffers for $G$.
- This buffer is write only and use the host memory directly, so the host can print it directly.

**Example 24: (vectorAdd-dependency.c)**

```
120    cl_mem bufferG =
121      clCreateBuffer(context,
122                     CL_MEM_WRITE_ONLY | CL_MEM_USE_HOST_PTR,
123                     N * sizeof(cl_uint), G, &status);
124    assert(status == CL_SUCCESS);
125    printf("Build buffers completes\n");
```

## Discussion

- What are the characteristic for read only, read and write, and write only buffers respectively?

## Buffers

- Both NDRange and the work group are one dimensional.
- The size of NDRange is $N$.
- The size of a work group is 256.

**Example 25: (vectorAdd-dependency.c)**

```
127   size_t globalThreads[] = {(size_t)N};
128   size_t localThreads[] = {256};
```

# $C = A + B$

- Now we launch the first kernel for $C = A + B$.
- We have already built three command queues for three devices.
- We use commandQueue[0] for the first addition.
- We also declared three events to denote the completion of the three kernels that we will launch.

**Example 26: (vectorAdd-dependency.c)**

```
130    cl_event events[3];
131    status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
132                            (void*)&bufferA);
133    assert(status == CL_SUCCESS);
134    status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
135                            (void*)&bufferB);
136    assert(status == CL_SUCCESS);
137    status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
138                            (void*)&bufferC);
139    assert(status == CL_SUCCESS);
140
141    status =
142      clEnqueueNDRangeKernel(commandQueue[0], kernel, 1, NULL,
143                             globalThreads, localThreads,
144                             0, NULL, &(events[0]));
145    assert(status == CL_SUCCESS);
```

# $F = D + E$

- Similarly we launch the second kernel for $F = D + E$.
- Note that we used commandQueue[1] and events[1].

**Example 27: (vectorAdd-dependency.c)**

```
147    status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
148                            (void*)&bufferD);
149    assert(status == CL_SUCCESS);
150    status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
151                            (void*)&bufferE);
152    assert(status == CL_SUCCESS);
153    status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
154                            (void*)&bufferF);
155    assert(status == CL_SUCCESS);
156    status =
157      clEnqueueNDRangeKernel(commandQueue[1], kernel, 1, NULL,
158                             globalThreads, localThreads,
159                             0, NULL, &(events[1]));
160    assert(status == CL_SUCCESS);
```

# $G = C + F$

- Finally we launch the third kernel for $G = C + F$.
- Note that we used commandQueue[2] and events[2].

**Example 28:   (vectorAdd-dependency.c)**

```
162    status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
163                            (void*)&bufferC);
164    assert(status == CL_SUCCESS);
165    status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
166                            (void*)&bufferF);
167    assert(status == CL_SUCCESS);
168    status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
169                            (void*)&bufferG);
170    assert(status == CL_SUCCESS);
171    status =
172      clEnqueueNDRangeKernel(commandQueue[2], kernel, 1, NULL,
173                             globalThreads, localThreads,
174                             0, NULL, &(events[2]));
175    assert(status == CL_SUCCESS);
```

# Wait for Events

- We wait for all three kernel to finish.

**Example 29: (vectorAdd-dependency.c)**

```
177    status = clWaitForEvents(DEVICENUM, events);
178    assert(status == CL_SUCCESS);
179    printf("All three kernels complete.\n");
```

# Wait for Events

- The rest of the code just prints the timing information, checks for correctness, and releases resources.
- Please refer to previous discussion.

# Demonstration

- Run the `vectorAdd-dependency-cl`.

## Discussion

- Does the program produce the correct answer?
- If not what is the possible reason?

## Problem

- The previous program waits for all three kernels to complete, but the third kernel did *not* wait for the first two.
- We will use the *wait for events* mechanism to launch the third kernel.

**Prototype 30: clEnqueueNDRangeKernel.h**

```
1   cl_int
2   clEnqueueNDRangeKernel (cl_command_queue command_queue ,
3                           cl_kernel kernel ,
4                           cl_uint work_dim ,
5                           const size_t *global_work_offset ,
6                           const size_t *global_work_size ,
7                           const size_t *local_work_size ,
8                           cl_uint num_events_in_wait_list ,
9                           const cl_event *event_wait_list ,
10                          cl_event *event );
```

# $G = C + F$

- We launch the third kernel and wait for the first two events in the events array.
- Set num_events_in_wait_list to 2 and event_wait_list to events.

**Example 31: (vectorAdd-dependency-correct.c)**

```
162   status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
163                           (void*)&bufferC);
164   assert(status == CL_SUCCESS);
165   status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
166                           (void*)&bufferF);
167   assert(status == CL_SUCCESS);
168   status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
169                           (void*)&bufferG);
170   assert(status == CL_SUCCESS);
171   status =
172     clEnqueueNDRangeKernel(commandQueue[2], kernel, 1, NULL,
173                            globalThreads, localThreads,
174                            2, events, &(events[2]));
175   assert(status == CL_SUCCESS);
```

## Demonstration

- Run the `vectorAdd-dependency-correct-cl`.

## Discussion

- Does the program produce the correct answer?
- Observe the timing and make sure the first two kernels run in parallel, and the third kernel will wait for the first two.

## flFinish

- The host can also explicitly wait for the first two kernels to finish before launching the third kernel by clFinish.
- We just need to wait for the command queues to the first two devices to become empty.

**Example 32:  (vectorAdd-dependency-clfinish.c)**

```
162    clFinish(commandQueue[0]);
163    clFinish(commandQueue[1]);
164    status = clSetKernelArg(kernel, 0, sizeof(cl_mem),
165                            (void*)&bufferC);
166    assert(status == CL_SUCCESS);
167    status = clSetKernelArg(kernel, 1, sizeof(cl_mem),
168                            (void*)&bufferF);
169    assert(status == CL_SUCCESS);
170    status = clSetKernelArg(kernel, 2, sizeof(cl_mem),
171                            (void*)&bufferG);
172    assert(status == CL_SUCCESS);
173    status =
174      clEnqueueNDRangeKernel(commandQueue[2], kernel, 1, NULL,
175                              globalThreads, localThreads,
176                              0, NULL, &(events[2]));
177    assert(status == CL_SUCCESS);
```

## Demonstration

- Run the `vectorAdd-dependency-clfinish-cl` program.

## Discussion

- Does the program produce the correct answer?
- Observe the timing and make sure the first two kernels run in parallel, and the third kernel will wait for the first two.
- Compare the timing results between vectorAdd-dependency-correct-cl and vectorAdd-dependency-clfinish-cl, especially in when the last kernel joined the command queue.

## Group Size

- The previous matrix multiplication program has extremely good performance due to two reasons.
    - It uses local memory to speed up data access.
    - It has a larger group size (256).
- Now we would like to study the effects of group size on performance.

## Compute Units

- A device has only a limited number of compute units.
- The work items of a work group will occupy a compute unit.
- For a given number of work items, if the work group size is small, then the number of work groups becomes large, and some work groups will wait for compute unit.
- In other words, a small work group size limits the the number of work items that we can process in parallel.

## Classroom

- A school has only a limited number of classrooms.
- The students are divided into group, and one group will occupy a classroom.
- For a given number of students, if the group size is small, then the number of groups becomes large, and some groups will wait for classroom.
- In other words, a small group size limits the the number of students that can study in parallel.

## Discussion

- Give your own example of this phenomenon.

# Compute Unit Number

- We need only one GPU device.
- We call clGetDeviceInfo with CL_DEVICE_MAX_COMPUTE_UNITS to get the number of compute units on GPU.
- We need the number of compute units to determine the work group size.

**Example 33: (vectorAdd-groupsize.c)**

```
23    cl_device_id GPU[MAXGPU];
24    cl_uint GPU_id_got;
25    status =
26      clGetDeviceIDs(platform_id, CL_DEVICE_TYPE_GPU,
27                     MAXGPU, GPU, &GPU_id_got);
28    assert(status == CL_SUCCESS && GPU_id_got >= 1);
29    printf("There are %d GPU devices\n", GPU_id_got);
30    cl_uint unit;
31    status =
32      clGetDeviceInfo(GPU[0], CL_DEVICE_MAX_COMPUTE_UNITS,
33                      sizeof(cl_uint), &unit, NULL);
34    assert(status == CL_SUCCESS);
35    printf("# of compute units is %d\n", unit);
```

## Group Size

- We fix the number of work items to $N$, and vary the size of the work group from 1 to 256.

- We wait for the event before launching the next kernel.

- We will record the timing information into a file vectorAdd-grouopsize.dat.

**Example 34: (vectorAdd-groupsize.c)**

```
108    size_t workItem[] = {(size_t)N};
109    FILE *timefp = fopen("vectorAdd-groupsize.dat", "w");
110    assert(timefp != NULL);
111    for (int groupSize = 1; groupSize <= 256; groupSize *= 2) {
112      cl_event event;
113      size_t localSize[1];
114      localSize[0] = groupSize;
115      status =
116        clEnqueueNDRangeKernel(commandQueue, kernel, 1, NULL,
117                               workItem, localSize,
118                               0, NULL, &event);
119      assert(status == CL_SUCCESS);
120      /* waitforevent */
121      status = clWaitForEvents(1, &event);
122      assert(status == CL_SUCCESS);
123      printf("The kernel with group size %d completes.\n",
124              groupSize);
```

# Get Timing Information

- The rest of the code will retrieve timing information from the event.
- We output the group size and the execution time into a file vectorAdd-grouopsize.dat.

**Example 35:   (vectorAdd-groupsize.c)**

```
157        printf("\nkernel entered queue at %f\n",
158                (timeEnterQueue - base) / NANO2SECOND);
159        printf("kernel submitted to device at %f\n",
160                (timeSubmit - base) / NANO2SECOND);
161        printf("kernel started at %f\n",
162                (timeStart - base) / NANO2SECOND);
163        printf("kernel ended  at %f\n",
164                (timeEnd - base) / NANO2SECOND);
165        printf("kernel queued time %f seconds\n",
166                (timeSubmit - timeEnterQueue) / NANO2SECOND);
167        printf("kernel submission time %f seconds\n",
168                (timeStart - timeSubmit) / NANO2SECOND);
169        printf("kernel execution time %f seconds\n",
170                (timeEnd - timeStart) / NANO2SECOND);
171        fprintf(timefp, "%d %f\n", groupSize ,
172                (timeEnd - timeStart) / NANO2SECOND);
173     }
```
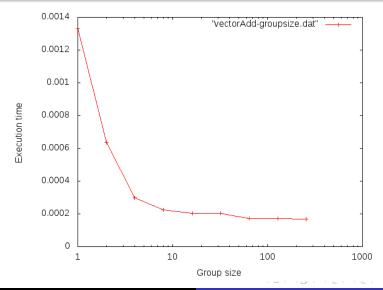
## Demonstration

- We plot the kernel execution time for different group sizes.
- The x-xis is in log scale.

# Execution Time

## Discussion

- Observe the execution time and draw you conclusion about the effects of group size on performance.