

Introduction to OpenCL

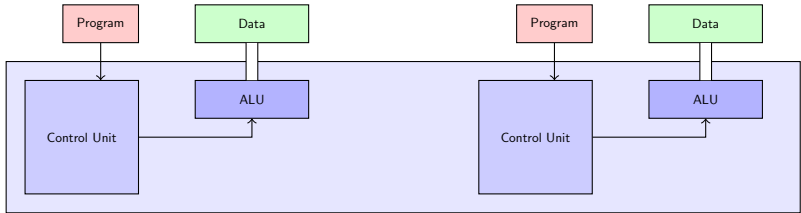
Computing on your Graphics Hardware

Maurice Leclaire

TumFUG
Linux / Unix
get-together

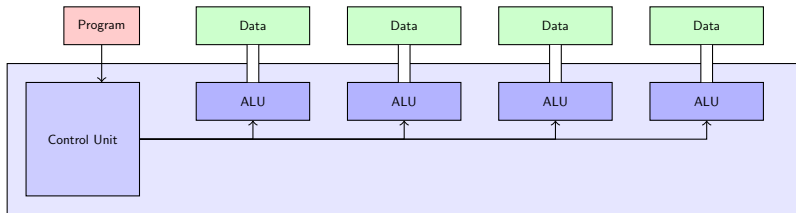
June 22, 2011

Normal CPU Architecture



- ▶ Separate control unit in each core
- ▶ Can execute different programs on each core

GPU / SIMD Architecture



- ▶ Single instruction multiple data
- ▶ Executes the same program on different data

When is this useful?

- ▶ The same computation on many data items
- ▶ Data items must be independent, e.g. matrix multiplication
- ▶ Not useful for totally different programs
- ▶ Not useful for dependent data items

OpenCL

- ▶ Device independent framework to use SIMD devices
- ▶ Runtime device scheduling
- ▶ Runtime compilation for different devices
- ▶ Host program on CPU executes a kernel on the target device

OpenCL C

- ▶ included language for programming OpenCL
- ▶ based on C99

Example: Matrix multiplication

$$A, B \in n \times n$$

$$C = A \times B \in n \times n$$

$$c_{i,j} = \sum_{k=0}^n a_{i,k} \cdot b_{k,j}$$

Sequential matrix multiplication

```
1 void matrixmul(int *A, int *B, int *C)
2 {
3     for (int i = 0; i < n; i++) {
4         for (int j = 0; j < n; j++) {
5             int m = 0;
6             for (int k = 0; k < n; k++) {
7                 m += A[i*n+k] * B[k*n+j];
8             }
9             C[i*n+j] = m;
10        }
11    }
12 }
```

Time $\mathcal{O}(n^3)$

Matrix multiplication with OpenCL

```
1  __kernel void matrixmul(__global int *A,  
2                          __global int *B,  
3                          __global int *C)  
4  {  
5      int n = get_global_size(0);  
6      int i = get_global_id(0);  
7      int j = get_global_id(1);  
8  
9      int m = 0;  
10     for (int k = 0; k < n; k++) {  
11         m += A[i*n+k] * B[k*n+j];  
12     }  
13     C[i*n+j] = m;  
14 }
```

Time $\mathcal{O}(n)$ (if enough execution units)

Executing the kernel

```
1 int matrixmul(int n, int *A, int*B, int *C)
2 {
3     cl_uint numPlatforms = 0;
4     clGetPlatformIDs(0, NULL, &numPlatforms);
5     if (numPlatforms == 0)
6         return -1;
7     cl_platform_id *platforms = (cl_platform_id *)
8         malloc(numPlatforms*sizeof(cl_platform_id));
9     clGetPlatformIDs(numPlatforms, platforms, NULL);
10
11     cl_uint numDevices = 0;
12     clGetDeviceIDs(platforms[0], CL_DEVICE_TYPE_GPU,
13         0, NULL, &numDevices);
14     if (numDevices == 0)
15         return -1;
16     cl_device_id *devices = (cl_device_id *)
17         malloc(numDevices*sizeof(cl_device_id));
18     clGetDeviceIds(platforms[0], CL_DEVICE_TYPE_GPU,
19         numDevices, devices, NULL);
```

Executing the kernel

```
20     cl_context context = clCreateContext(  
21         NULL, numDevices, devices, NULL, NULL, NULL);  
22  
23     cl_command_queue cmdQueue = clCreateCommandQueue(  
24         context, devices[0], 0, NULL);  
25  
26     cl_mem a = clCreateBuffer(context,  
27         CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,  
28         n*n*sizeof(int), A, NULL);  
29     cl_mem b = clCreateBuffer(context,  
30         CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,  
31         n*n*sizeof(int), B, NULL);  
32     cl_mem c = clCreateBuffer(context,  
33         CL_MEM_READ_WRITE,  
34         n*n*sizeof(int), NULL, NULL);
```

Executing the kernel

```
35     const char *source; //contains OpenCl source
36
37     program = clCreateProgramWithSource(context, 1,
38         &source, NULL, NULL);
39
40     cl_int buildErr = clBuilProgram(program,
41         numDevices, devices, NULL, NULL, NULL);
42     if (buildErr != CL_SUCCESS)
43         return -1;
44
45     cl_kernel kernel = clCreateKernel(program,
46         "matrixmul", NULL);
```

Executing the kernel

```
47     size_t globalWorkSize [2];
48     globalWorkSize [0] = n;
49     globalWorkSize [1] = n;
50
51     clSetKernelArg(kernel, 0, sizeof(cl_mem), &a);
52     clSetKernelArg(kernel, 1, sizeof(cl_mem), &b);
53     clSetKernelArg(kernel, 2, sizeof(cl_mem), &b);
54
55     clEnqueueNDRangeKernel(cmdQueue, kernel, 2,
56         NULL, globalWorkSize, NULL, 0, NULL, NULL);
57
58     clEnqueueReadBuffer(cmdQueue, c, CL_TRUE, 0,
59         n*n*sizeof(int), C, 0, NULL, NULL);
```

Executing the kernel

```
60     clReleaseKernel(kernel);
61     clReleaseProgram(program);
62     clReleaseCommandQueue(cmdQueue);
63     clReleaseMemObj(a);
64     clReleaseMemObj(b);
65     clReleaseMemObj(c);
66     clReleaseContext(context);
67
68     free(devices);
69     free(platforms);
70
71     return 0;
72 }
```

File Edit Help

Source Code

```

1
2 __kernel void matmult( __global int *A,
3                       __global int *B,
4                       __global int *C)
5 {
6     int n = get_global_size(0);
7     int i = get_global_id(1);
8     int j = get_global_id(0);
9
10    int m = 0;
11
12    for (int k = 0; k < n; k++) {
13        m += A[i*n+k] + B[k*n+j];
14    }
15    C[i*n+j] = m;
16 }
17

```

Compile

Source type: OpenCL

OpenCL Compiler

Options

Macro Definitions

Symbol	Value
Right-click to add macros.	

Object Code

Format: Radeon HD 5770 (Uniper) Assembly

```

1; ----- Disassembly -----
2 00 ALU: PUSH_BEFORE: ADDR(12) CNT(13) KCACHE(CB0:0-15)
3     0 w: MOV     R8.x, RCO[0].x
4     1 w: MOV     R2.y, RCO[0].y
5     2 w: MOV     R2.x, RCO[0].x
6     3 w: MOV     R2.w, RCO[0].w
7     4 w: MOV     R4.x, 0.0f
8     1 v: NULLO_INT    __, R1.x, RCO[1].x
9     2 w: ADD_INT    __, R0.w, P81
10    1 v: NULLO_INT    __, R1.y, RCO[1].y
11    3 w: ADD_INT    R0.x, P12.w, RCO[4].x
12    w: ADD_INT    __, R0.y, P82
13    4 w: ADD_INT    R0.y, P10.w, RCO[6].y
14    8 v: NULLO_INT    R1.x, P14.y, R8.x
15    6 w: PREDOT_INT    __, R8.x, 0.0f UPDATE_EXECP
16 01 JMP: POP_CNT(1) ADDR(9)
17 02 ALU: ADDR(14) CNT(9) KCACHE(CB0:0-15) KCACHE(CB1:0-15)
18     7 w: LSHL    __, R0.x, (0x00000002, 2.8028962)
19     w: LSHL    R2.w, RCO[0].x, (0x00000002, 2.802)
20     0 w: MOV     R4.x, 0.0f
21     9 w: ADD_INT    R8.w, RCL[1].x, P17.w
22     1 v: NULLO_INT    __, RCO[0].x, R0.y
23     9 y: LSHL    __, P8, (0x00000002, 2.8028962)
24    10 w: ADD_INT    R2.x, RCL[0].x, P19.y
25 03 LOOP: FILL 10 FILL_JUMP_ADDR(8)
26 04 ALU: ADDR(14) CNT(8)
27     11 w: ADD_INT    R8.x, -1, R8.x
28     y: LSHR    R0.y, R8.w, (0x00000002, 2.802)
29     w: LSHR    R0.x, R2.x, (0x00000002, 2.802)
30     w: ADD_INT    R3.w, R2.w, R3.w VEC_110
31     v: ADD_INT    R2.x, R2.x, (0x00000004, 0.605)
32 05 TEX: ADDR(10) CNT(2)
33     12 VFETCH R0 __, R0.x, f0178 MESA(4)
34     VFETCH_TYPE(NO_INDEX_OFFSET)
35     13 VFETCH R1 __, R0.y, f0178 MESA(4)
36     VFETCH_TYPE(NO_INDEX_OFFSET)
37 06 ALU: BREAK: ADDR(0) CNT(4)
38     14 w: ADD_INT    __, R4.x, R0.w

```

Compiler Statistics (Using CAL 11.3)

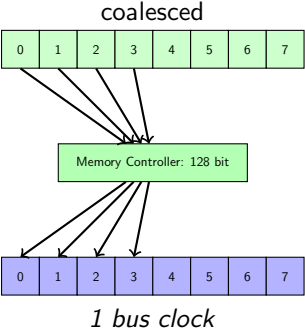
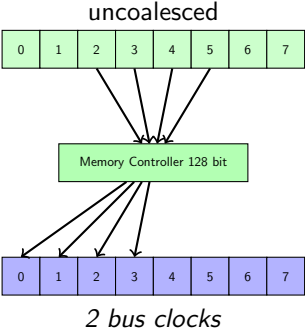
Name	GPR	Scratch Reg	Min	Max	Avg	ALU	Fetch	Write	Est Cycles	ALU/Fetch	Bottle/ck	Thread/Clock	Throughput
Radeon HD 2600	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Radeon HD 3870	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Radeon HD 4890	S	0	3.20	285.60	34.95	40	2	1	30.84	0.61	Global Fetch	0.52	441M Threads/Sec
Radeon HD 4770	S	0	4.00	357.00	43.69	40	2	1	38.55	0.61	Global Fetch	0.42	311M Threads/Sec
Radeon HD 4870	S	0	3.20	285.60	25.22	40	2	1	22.26	0.65	Global Fetch	0.72	539M Threads/Sec
Radeon HD 4670	S	0	4.00	357.00	23.92	40	2	1	21.11	0.70	Global Fetch	0.38	284M Threads/Sec
Radeon HD 4550	S	0	16.00	714.00	59.43	40	2	1	59.43	1.41	ALU Ops	0.13	81M Threads/Sec
Radeon HD 5870	S	0	1.20	285.60	34.95	19	2	1	30.84	0.54	Global Fetch	1.04	882M Threads/Sec
Radeon HD 5770	S	0	1.20	285.60	34.95	19	2	1	19.57	1.16	Global Fetch	0.72	317M Threads/Sec
Radeon HD 5670	S	0	2.20	285.60	19.57	19	2	1	19.57	1.16	ALU Ops	0.41	317M Threads/Sec
Radeon HD 5450	S	0	5.50	517.50	31.32	19	2	1	31.32	2.54	ALU Ops	0.13	83M Threads/Sec
FireStream 9170	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A

Compiler Output

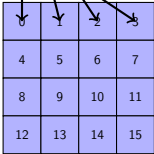
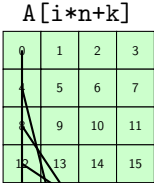
Compiler Output

Coalescing memory accesses

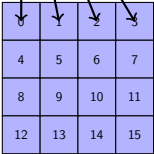
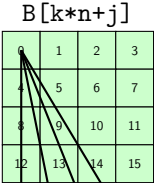
ATI Radeon HD 5770 / Memory bus with: 128 bit



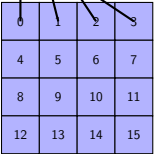
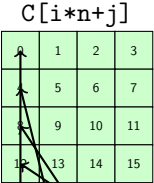
Memory accesses



uncoalesced



broadcast

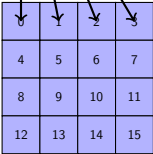
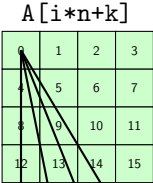


uncoalesced

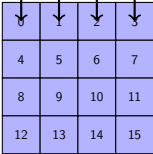
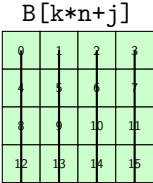
Matrix multiplication with different thread mapping

```
1  __kernel void matrixmul(__global int *A,  
2                               __global int *B,  
3                               __global int *C)  
4  {  
5      int n = get_global_size(0);  
6      int i = get_global_id(1);  
7      int j = get_global_id(0);  
8  
9      int m = 0;  
10     for (int k = 0; k < n; k++) {  
11         m += A[i*n+k] * B[k*n+j];  
12     }  
13     C[i*n+j] = m;  
14 }
```

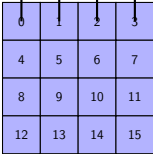
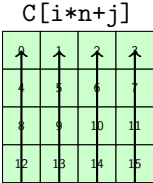
Memory accesses with different thread mapping



broadcast

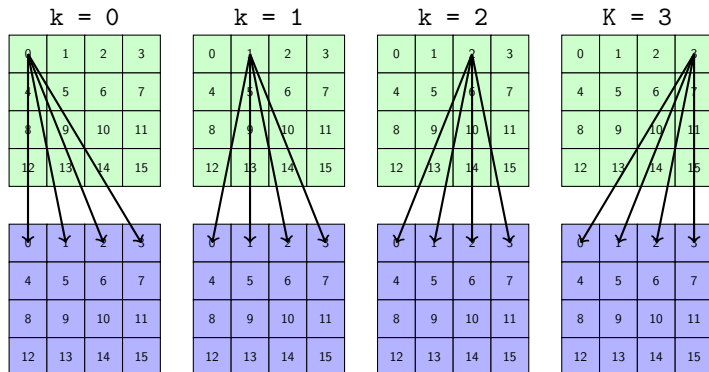


coalesced

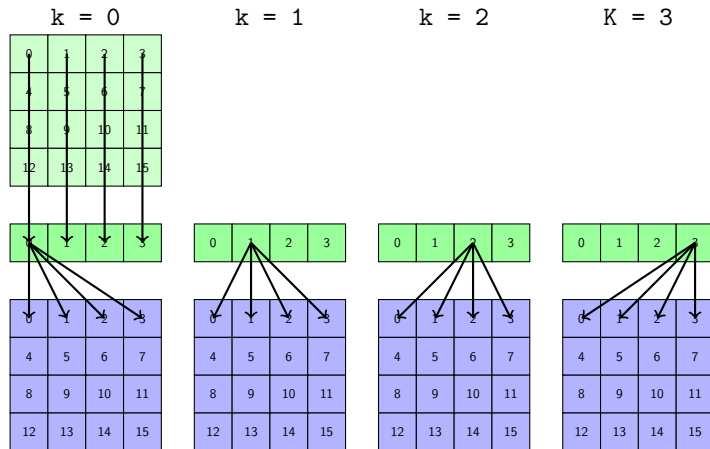


coalesced

Broadcast memory accesses during loop



Memory accesses using local memory



Matrix multiplication using local memory

How can this be done?

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

- ▶ Dividing the global space into local parts
- ▶ Using a local row number of 1
- ▶ Local column size should be as big as possible

Matrix multiplication using local memory

```
1  __kernel void matrixmul(__global int *A,
2  __global int *B, __global int *C,  __local int *tmp)
3  {
4      int n = get_global_size(0);
5      int i = get_global_id(1);
6      int j = get_global_id(0);
7      int nl = get_local_size(0);
8      int jl = get_local_id(0);
9      int num = get_num_groups(0);
10
11     int m = 0;
12     for (int k = 0; k < num; k++) {
13         tmp[jl] = A[i*n+k*nl+jl];
14         barrier(CLK_LOCAL_MEM_FENCE);
15         for (int l = 0; l < nl; l++) {
16             m += tmp[l] * B[(k*nl + l)*n+j];
17         }
18         barrier(CLK_LOCAL_MEM_FENCE);
19     }
20     C[i*n+j] = m;
21 }
```

Executing the kernel

```
22     size_t globalWorkSize [2];
23     globalWorkSize [0] = n;
24     globalWorkSize [1] = n;
25
26     size_t localWorkSize [2];
27     clGetDeviceInfo (devices [0],
28                     CL_DEVICE_MAX_WORK_GROUP_SIZE,
29                     sizeof (size_t), localWorkSize, NULL);
30     localWorkSize [1] = 1;
31
32     clSetKernelArg (kernel, 0, sizeof (cl_mem), &a);
33     clSetKernelArg (kernel, 1, sizeof (cl_mem), &b);
34     clSetKernelArg (kernel, 2, sizeof (cl_mem), &b);
35     clSetKernelArg (kernel, 3,
36                     localWorkSize [0] * sizeof (int), NULL);
37
38     clEnqueueNDRangeKernel (cmdQueue, kernel, 2,
39                             NULL, globalWorkSize, NULL, 0, NULL, NULL);
```

Speed comparison

AMD Phenom II x4 965 vs. ATI Radeon HD 5770

Matrix multiplication with 4096×4096 elements

Environment	Time	Speed factor
CPU (sequential)	4020 sec	1
CPU (OpenMP)	990 sec	4
GPU	100 sec	40
GPU (coalesced)	5,5 sec	700
GPU (local memory)	5 sec	800