Programming in OpenCL

Timo Stich, NVIDIA



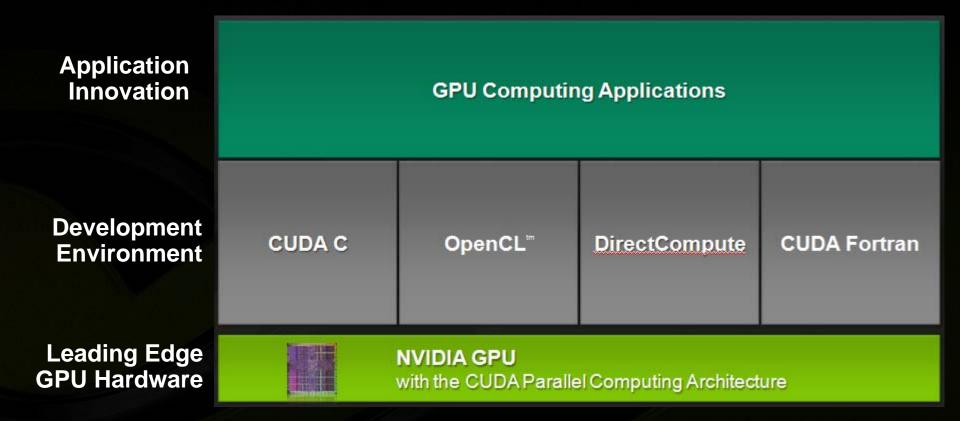
Outline



- Introduction to OpenCL
- OpenCL API Overview
- Performance Tuning on NVIDIA GPUs
- OpenCL Programming Tools & Resources

OpenCL and the CUDA Architecture





OpenCL Portability



Portable code across multiple devices
 GPU, CPU, Cell, mobiles, embedded systems, ...

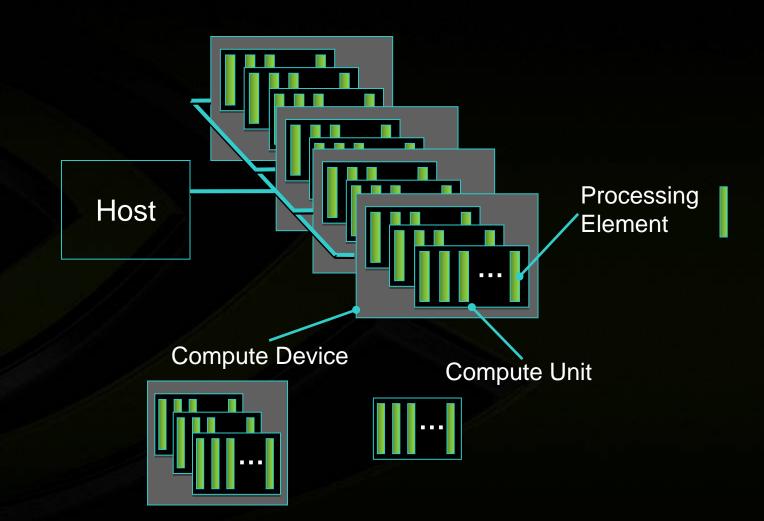
• NOTE:

functional portability != performance portability

- Different code for each device is necessary to get good performance
- Even for GPUs from different vendors!

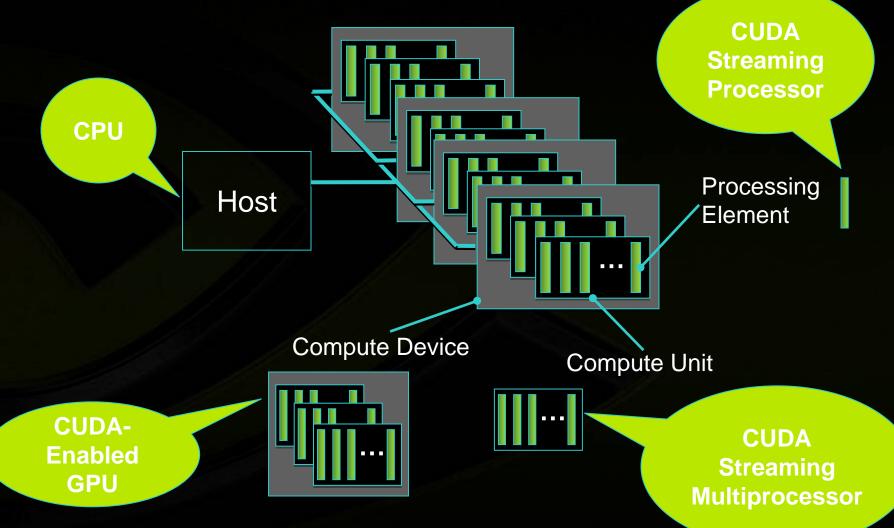
OpenCL Platform Model





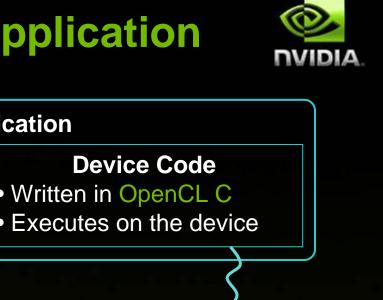
OpenCL Platform Model on the CUDA Architecture

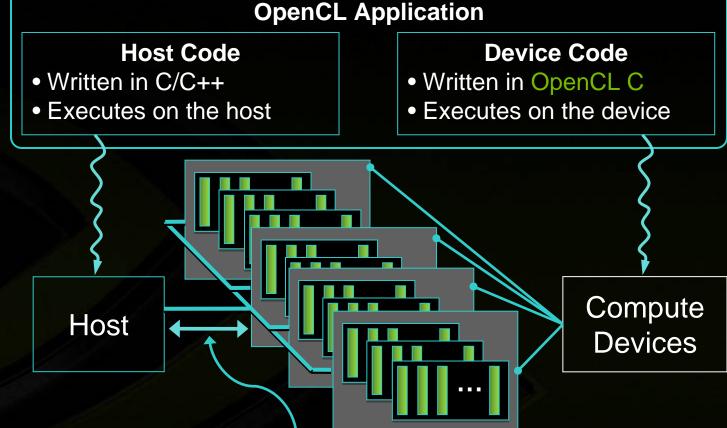




NVIDIA GPU Computing Master Class

Anatomy of an OpenCL Application





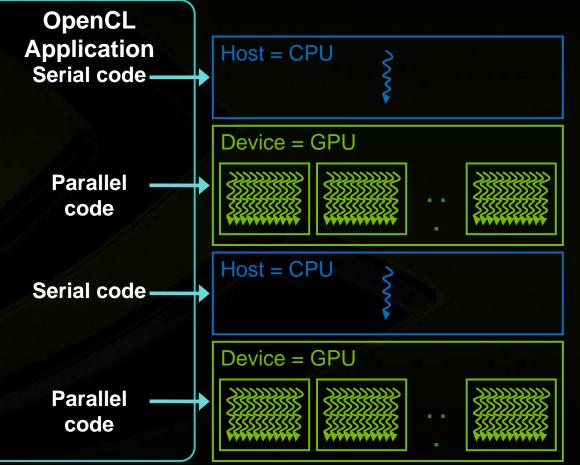
Host code sends commands to the devices:

- to transfer data between host memory and device memories
- to execute device code

Heterogeneous Computing



- Serial code executes in a CPU thread
- Parallel code executes in many GPU threads across multiple processing elements



OpenCL Framework



Platform layer

- Discover OpenCL devices and their capabilities and create contexts
- Runtime layer
 - Memory management and command execution within a context
- OpenCL C Compiler
 - Creates program executables that contain OpenCL kernels

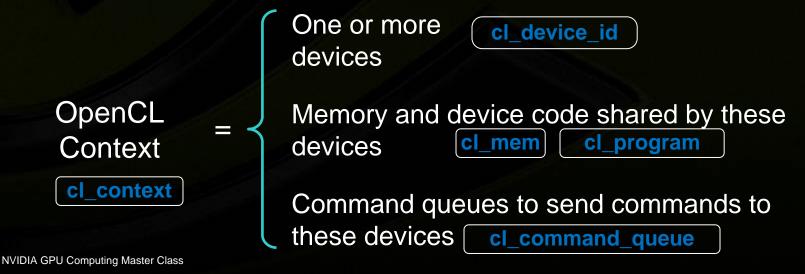
Platform Layer



Query platform information

- clGetPlatformIDs(): list of platforms
- clGetPlatformInfo(): profile, version, vendor, extensions
- clGetDeviceIDs(): list of devices
- clGetDeviceInfo(): type, capabilities

Create OpenCL context on one or more devices of one platform



Error Handling, Resource Deallocation



Error handling:

- All host functions return an error code
- Context error callback can be specified
- Resource deallocation
 - Reference counting API: clRetain*(), clRelease*()
- Both are removed from code samples for clarity
 - Please see SDK samples for complete code

Context Creation



// Create an OpenCL context for all GPU devices on the first Platform cl_context* CreateContext() { cl_platform_id platform_id; clGetPlatformIDs(1, &platform_id, NULL);

return clCreateContextFromType({CL_CONTEXT_PLATFORM, platform_id, 0}, CL_DEVICE_TYPE_GPU, NULL, NULL, NULL);

}
// Get the list of GPU devices associated with a context
cl_device_id* GetDevices(cl_context context) {
 size_t size;
 clGetContextInfo(context, CL_CONTEXT_DEVICES, 0, NULL, &size);
 cl_device_id* device_id = malloc(size);
 clGetContextInfo(context, CL_CONTEXT_DEVICES, size, device_id, NULL);
 return device_id;
}

Runtime



- Command queues creation and management
- Memory allocation and management
- Device code compilation and execution
- Event creation and management (synchronization, profiling)

Command Queue



 Sequence of commands scheduled for execution on a specific device

- Enqueuing functions: clEnqueue*()
- Multiple queues can execute on the same device
- Two modes of execution:
 - In-order: Each command in the queue executes only when the preceding command has completed
 - Including all memory writes, so memory is consistent with all prior command executions
 - Out-of-order: No guaranteed order of completion for commands

Commands



- Memory copy or mapping
- Device code execution
- Synchronization point

Command Queue Creation



code

// Create a command-queue for a specific device cl_command_queue CreateCommandQueue(cl_context context, cl_device_id device_id) { return clCreateCommandQueue(context, device_id, 0, NULL); } Properties Error

Command Synchronization



- Some clEnqueue*() calls can be optionally blocking
- Queue barrier command
 - Any commands after the barrier start executing only after all commands before the barrier have completed
- An event object can be associated to each enqueued command
 - Any commands (or clWaitForEvents()) can wait on events before executing
 - Can be queried to track execution status and get profiling information

Memory Objects



• Two types of memory objects (cl_mem):

- Buffer objects
- Image objects
- Associated with context, only implicitly with device
- Memory objects can be copied to host memory, from host memory, or to other memory objects
- Regions of a memory object can be accessed from host by mapping them into the host address space

Buffer Object



- One-dimensional array
- Elements are scalars, vectors, or any user-defined structures
- Accessed within device code via pointers

```
_kernel void myKernel(__global int* buffer) {
  <...>
  // Access element in buffer object
  int v = buffer[get_global_id(0)];
  <...>
```

Image Object



- Two- or three-dimensional array
- Elements are 4-component vectors from a list of predefined formats
- Accessed within device code via built-in functions (storage format not exposed to application)
 - Sampler objects are used to configure how built-in functions sample images (addressing modes, filtering modes)
- Can be created from OpenGL texture or renderbuffer

Data Transfer between Host and Device



int main() {

cl context context = CreateContext(); cl_device_id* device_id = GetDevices(context); cl command queue command queue = CreateCommandQueue(context, device_id[0]); size_t size = 100000 * sizeof(int); int* h_buffer = (int*)malloc(size); cl_mem* d_buffer = clCreateBuffer(context, CL_MEM_READ_WRITE, size, NULL, NULL); // Initialize host buffer h buffer clEnqueueWriteBuffer(command_queue, d_buffer, CL_FALSE, 0, size, h_buffer, 0, NULL, NULL); // Process device buffer d buffer clEngueueReadBuffer(command_gueue, d buffer, CL TRUE, 0, size, h buffer, 0, NULL, NULL);

Device Code in OpenCL C



Derived from ISO C99

- A few restrictions: recursion, function pointers, functions in C99 standard headers
- Some extensions: built-in variables and functions, function qualifiers, address space qualifiers, e.g:

_global float* a; // Pointer to device memory

 Functions qualified by __kernel keyword (a.k.a kernels) can be invoked by host code

kernel void MyKernel() { ... }

Kernel Execution: NDRange and Work-Items



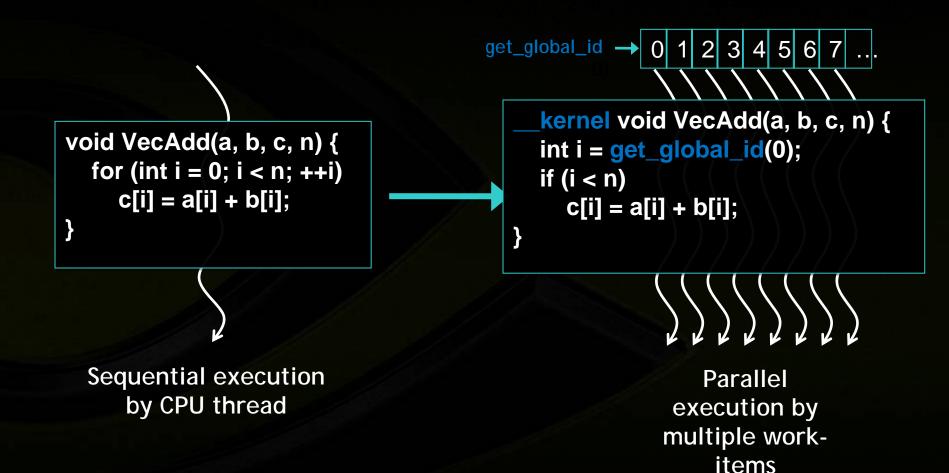
 Host code invokes a kernel over an index space called an NDRange

- NDRange = "N-Dimensional Range"
- NDRange can be a 1-, 2-, or 3-dimensional space
- A single kernel instance at a point in the index space is called a *work-item*
 - Each work-item has a unique global ID within the index space (accessible from device code via get_global_id())

Each work-item is free to execute a unique Code path NVIDIA GPU Computing Master Class

Example: Vector Addition



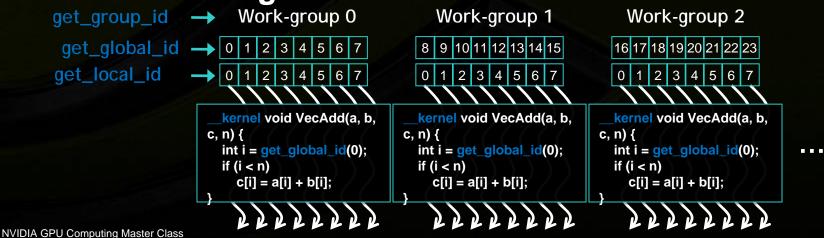


Kernel Execution: Work-Groups



Work-items are grouped into work-groups

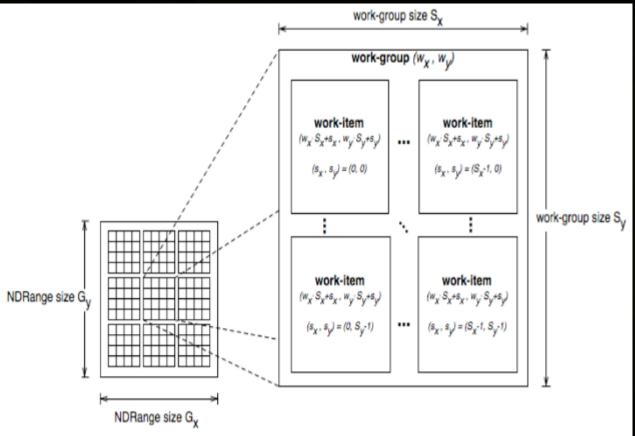
- Each work-group has a unique work-group ID (accessible from device code via get_group_id())
- Each work-item has a unique local ID within a work-group (accessible from device code via get_local_id())
- Work-group has same dimensionality as NDRange



Example of 2D NDRange

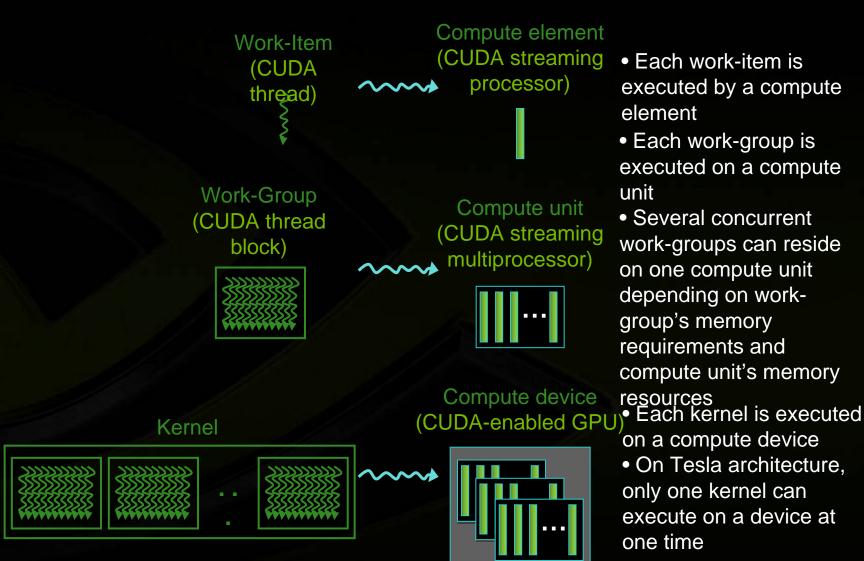


- Total number of work-items = $G_x \times G_y$
- Size of each work-group = $S_x \times S_y$
- Number of work-groups = $(G_x / S_x) \times (G_y / S_y)$ (must be dividable)



Kernel Execution on Platform Model





Benefits of Work-Groups

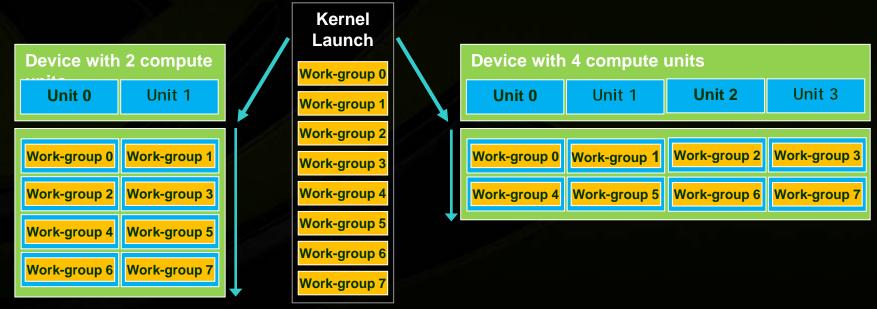


- Automatic scalability across devices with different numbers of compute units
- Efficient cooperation between work-items of same work-group
 - Fast shared memory and synchronization

Scalability

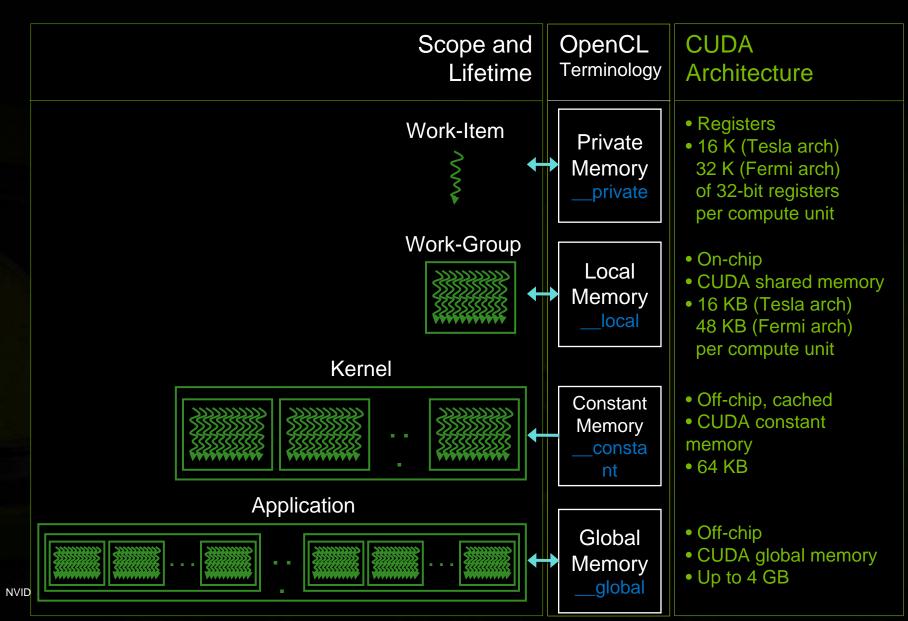


- Work-groups can execute in any order, concurrently or sequentially
- This independence between work-groups gives scalability:
 - A kernel scales across any number of compute units



Memory Spaces





Cooperation between Work-Items of same Work-Group



- Built-in functions to order memory operations and synchronize execution:
 - mem_fence(CLK_LOCAL_MEM_FENCE and/or CLK_GLOBAL_MEM_FENCE): waits until all reads/writes to local and/or global memory made by the calling work-item prior to mem_fence() are visible to all threads in the work-group
 - barrier(CLK_LOCAL_MEM_FENCE and/or CLK_GLOBAL_MEM_FENCE): waits until all workitems in the work-group have reached this point and calls mem_fence(CLK_LOCAL_MEM_FENCE and/or CLK_GLOBAL_MEM_FENCE)
- Used to coordinate accesses to local or global memory shared among work-items

Program and Kernel Objects



- A program object encapsulates some source code (with potentially several kernel functions) and its last successful build
 - clCreateProgramWithSource() // Create program from source
 - **clBuildProgram()** // Compile program
- A kernel object encapsulates the values of the kernel's arguments used when the kernel is executed
 - clCreateKernel() // Create kernel from successfully compiled
 // program
 - clSetKernelArg() // Set values of kernel's arguments

Kernel Invocation



int main() { ... // Create context and command gueue, allocate host and device buffers of N elements char* source = "__kernel void MyKernel(__global int* buffer, int N) {\n" if (get_global_id(0) < N) buffer[get_global_id(0)] = 7;\n" "}\n ": cl_program program = clCreateProgramWithSource(context, 1, &source, NULL, NULL); clBuildProgram(program, 0, NULL, NULL, NULL, NULL); ci kernel kernel = ciCreateKernel(program, "MyKernel", NULL); clSetKernelArg(kernel, 0, sizeof(cl_mem), (void*)&d_buffer); clSetKernelArg(kernel, 1, sizeof(int), (void*)&N); size_t localWorkSize = 256; // Number of work-items in a work-group int numWorkGroups = (N + localWorkSize - 1) / localWorkSize;size t globalWorkSize = numWorkGroups * localWorkSize; clEngueueNDRangeKernel(command_gueue, kernel, 1, NULL, &globalWorkSize, &localWorkSize, 0, NULL, NULL); ... // Read back buffer **NDRange**

dimension

OpenCL Local Memory on the CUDA Architecture



On-chip memory (CUDA shared memory)

- 2 orders of magnitude lower latency than global memory
- Order of magnitude higher bandwidth than global memory
- 16 KB per compute unit on Tesla architecture (up to 30 compute units)
- 48 KB per compute unit on Fermi architecture (up to 16 compute units)
- Acts as a user-managed cache to reduce global memory accesses
- Typical usage pattern for work-items within a workgroup:
 - Read data from global memory to local memory; synchronize with barrier()
 - Process data within local memory; synchronize with barrier()

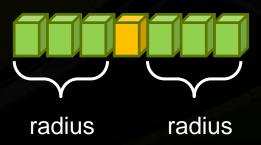
NVIDIA GPU Computing aste Write result to global memory

Example of Using Local Memory



Applying a 1D stencil to a 1D array of elements:

- Each output element is the sum of all elements within a radius
- For example, for radius = 3, each output element is the sum of 7 input elements:



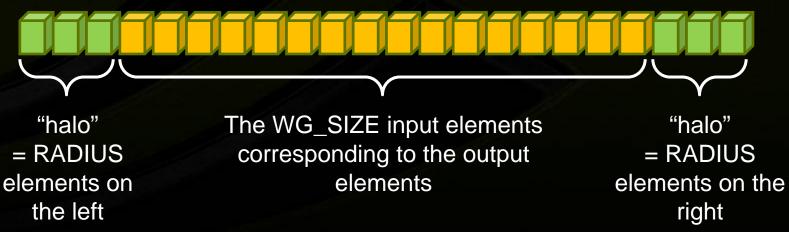
Implementation with Local Memory



 Each work-group outputs one element per workitem, so a total of WG_SIZE output elements

(WG_SIZE = number of work-items per work-group):

- Read (WG_SIZE + 2 * RADIUS) elements from global memory to local memory
- Compute WG_SIZE output elements in local memory
- Write WG_SIZE output elements to global memory





Kernel Code

RADIUS = 3kernel void stencil(__global int* input, WG SIZE = 16global int* output) { Local ID = 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 _local int local[WG_SIZE + 2 * RADIUS]; **i =** 0 1 2 3 4 5 6 7 8 9 101112131415161718192021 int i = get_local_id(0) + RADIUS; local[i] = input[get_global_id(0)]; if (get_local_id(0) < RADIUS) {</pre> local[i - RADIUS] = input[get_global_id(0) - RADIUS]; local[i + WG_SIZE] = input[get_global_id(0) + WG_SIZE]; } barrier(CLK_LOCAL_MEM_FENCE); // Blocks until work-items are done writing to local memory int value = 0;

for (offset = - RADIUS; offset <= RADIUS; ++offset) value += local[i + offset]; // Sum
output[get_global_id(0)] = value; }</pre>

OpenCL C Language Restrictions



- Pointers to functions are not allowed
- Pointers to pointers allowed within a kernel, but not as an argument
- Bit-fields are not supported
- Variable length arrays and structures are not supported
- Recursion is not supported
- Writes to a pointer of types less than 32-bit are not supported
- Double types are not supported, but reserved
- 3D Image writes are not supported
- Some restrictions are addressed through extensions

Optional Extensions



- Extensions are optional features exposed through OpenCL
- The OpenCL working group has already approved many extensions that are supported by the OpenCL specification:
 - Double precision floating-point types (Section 9.3)
 - Built-in functions to support doubles
 - Atomic functions (Section 9.5, 9.6, 9.7)
 - 3D Image writes (Section 9.8)
 - Byte addressable stores (write to pointers with types < 32bits) (Section 9.9)
 - Built-in functions to support half types (Section 9.10)

Performance Overview



OpenCL is about performance

- Standard to make use of the massive computing power of parallel processors like GPUs
- But, performance is generally not portable across devices:
 - There are multiple ways of implementing a given algorithm in OpenCL. Each can have vastly different performance characteristics for a given compute device!
- Achieving good performance on GPUs requires a basic understanding of GPU architecture

Heterogeneous Computing

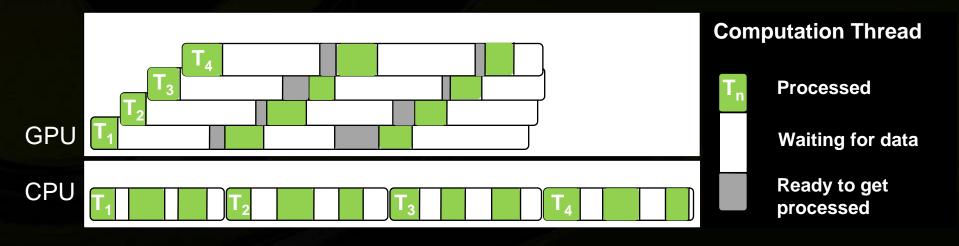


- Host + multiple devices = heterogeneous platform
- Distribute workload to:
 - Assign to each processor the type of work it does best
 CPU = serial, GPU = parallel
 - Keep all processors busy at all times
 - Minimize data transfers between processors or hide them by overlapping them with kernel execution
 - Overlapping requires data allocated with CL_MEM_ALLOC_HOST_PTR

GPU Computing: Highly Multithreaded



- GPU compute unit "hides" instruction and memory latency with computation
 - Switches from stalled threads to other threads at no cost (lightweight GPU threads)
 - Needs enough concurrent threads to hide latency
 - Radically different strategy than CPU core where memory latency is "reduced" via big caches



GPU Computing: Highly Multithreaded

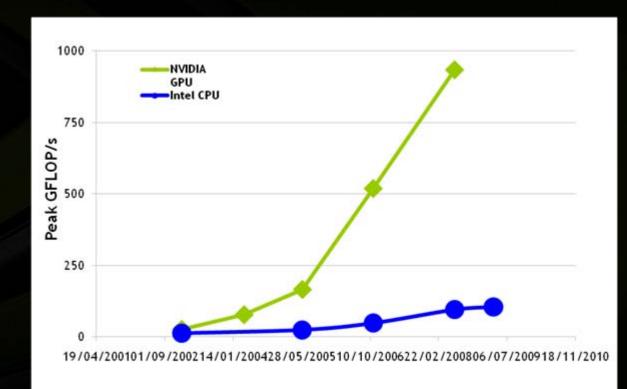


- Latency hiding is only possible if there is other work that can be done in parallel
- Therefore, kernels must be launched with hundreds of work-items per compute unit for good performance
 - Minimal work-group size of 64; higher is usually better (typically 1.2 to 1.5 speedup)
 - Number of work-groups is typically 100 or more

GPU Computing: High Arithmetic Intensity



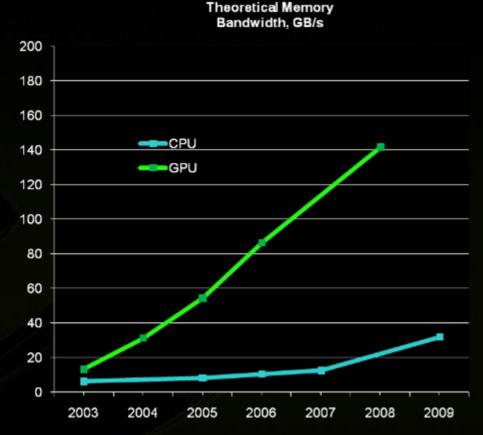
 GPU devotes many more transistors than CPU to arithmetic units => high arithmetic intensity



GPU Computing: High Memory Bandwidth



 GPUs offer high memory bandwidth, so applications can take advantage of high arithmetic intensity and achieve high arithmetic throughput



CUDA Memory Optimization



- Memory bandwidth will increase at a slower rate than arithmetic intensity in future processor architectures
- So, maximizing memory throughput is even more critical going forward
- Two important memory bandwidth optimizations:
 - Ensure global memory accesses are coalesced
 Up to an order of magnitude speedup!
 - Replace global memory accesses by shared memory accesses whenever possible

CUDA = SIMT Architecture



Same Instruction Multiple Threads

- Threads running on a compute unit are partitioned into groups of 32 threads (warps) in which all threads execute the same instruction simultaneously
- Minimize divergent branching within a warp
 - Different code paths within a warp get serialized
- Remove barrier calls when only threads within same warp need to communicate
 - Threads within a warp are inherently synchronized

CUDA = Scalar Architecture



- Use vector types for convenience, not performance
- Generally want more work-items rather than large vectors per work-item

Maximize Instruction Throughput



- Favor high-throughput instructions
- Use native_*() math functions whenever speed is more important than precision
- Use -cl-mad-enable compiler option
 - Enables use of FMADs, which can lead to large performance gains
- Investigate using the -cl-fast-relaxed-math compiler option
 - Enables many aggressive compiler optimizations

OpenCL Visual Profiler



• Analyze GPU HW performance signals, kernel occupancy, instruction throughput, and more

- Highly configurable tables and graphical views
- Save/load profiler sessions or export to CSV for later analysis
- Compare results visually across multiple sessions to see improvements
- Supported on Windows and Linux
- Included in the CUDA Toolkit

th Japan Spine																15
0020	Nertex 2	Chelona, fix.		chi.].miher	in fee											
i mater page 11 (acies)	Testers	Mathee	011m	(FO Test	Game	ption	ment	This Aread	appent pertext	-	-	-	manual		Angeland	. minim
1,14	a seal.	Individual	100	-	1.0	1		181	2			-		10		1008
- CR	4 1943		184	18.88						4	4					
ů.	1 1040	instants.		100.0	10			225						1.	1.	1.1
	4 20283		1.01	048						4	1					
	1 2744	industa.		841	1.0	2		100				1815	10	11	1.1	1.1
	14 1.579	humanitapi		100.0	1.2	1		102	8			8	8	15	11	1941
	1.1004		1.00	1147						41						
	4.104	hatirate.		188	10	1		398						101	30	198
	1. 1004.		100	10.000							4.					
	2 1612	Instants.	101	36747	10			104				194	14	4		
	1 100	Institution		101.0	14	<u> </u>	-	-				-			-	-
	2 180		161	1147	1	1	1	1	1.11	+	1		1	1	1	1
	1 1201	Instante.		-	10	1		-		611	1					
	1.100		101	1418	17		17	-			a (-			-	-
	5.188	befields.		100	aut .			128	8	-	-	182	10	-	10	001
	10.00	Instinutop		23	10	1	<u></u>	10	8			48	<u>_</u>	10		1004
	1 326		int	100	10	*		-		4	3		<u> </u>	1.0	~	1.0004
								-		<u>.</u>	1	-				100
	3 250	instants.	100	1291	14	<u>.</u>	*	28		1		303	K.,	361	81	20200
	3.000	instants.		301	12	-		-		1. ·	4	-		264	*	-
									8				-			1000
	1,200	instantiast		100	La .	1	*	-	<u></u>				-	535	21	
	CuteProfiler			_		_			_		_	_	_		_	- 5
file groffe	Options Winds	w Bilp														
00	8															
Settate		-	-	-	-	Session	22	-	-	-	-	-	-	10747		
	-							0.0000000	2.2							* 1
98-1 10 1	Profiler catpa	t wa	ve Summar	NO101												
200, 400 - 1 547, 400 - 1 577, 500 - 7					GRUtime											
10.1				1.1	0	322%	6.479	 33 	9.70%	32.94N	15	17%	19.40%	22	64%	2018
96 - 2 96 - 2 96 - 2									1.0			+ · · ·			-	-
ne -1 Ne -1 Ne -1			ster main	2.000				_								
			sgerjinalej stran 1 la	hat he		-	-									-
			stran 1 la	hat he							2					
n: -1 20-1 20-1			stran 1 la	fast he rt main encopy							2					
			strun_[]a ra sanax_g	fast ha nt main encopy (d main						3	P					
		al spera b	strum [lia mi samax g sx.main bi straw [up	fast, las ret_main encopy (d_main r_t_main ret_main						3	P					
96 - 1 196 - 1 196 - 1		al spera b	strun <u>ila</u> mi sanax g sa minihi stori jug stori jug	fast, hw rd, main emcapy (d, main r, na rh rd, main sp, main						3	r					
92 4 193 - 1 193 - 1		al spera b	strun []a na sanax g ex man in stran []a stas sass sass	fast, hw rd, main encopy (d, main con rd) rd, main (d, main (d, main						3	2					
92 4 193 - 1 193 - 1		nijageraji S	strun_[]a sanax_g sa_main_hi straw_[_up_ stass sovip_g stool_g	fast, hw rit main encopy (d, main v, na rib rit, main (d, main (d, main							2					
w 4 89 4 89 1		ndjageracju S	strun_[]a m sanax.g sz.msin_hi stran_[.u stran_[.u stran_] stran_j stran_j b[_msin_hi	fast, hw ret_main emcapy (d_main v_main p_main (d_main (d_main (d_main v_m_mb							2					
w 4 89 - 4 89 - 1		nijagemaju S nijagemajo K	strun_lia ma samax_g ex_main_in stran_l_up stran_l_up stran_g stran_g pid-main_in pid-telper_l	fast, las rt man encapy (d_man rt man rt, man rt, man (d_man (d_man rt, man rt, man												
m 4 89 4 89 9		ndjageracju S	strun_lia ma samax_g ex_main_in stran_l_up stran_l_up stran_g stran_g pid-main_in pid-telper_l	fact, fair rd_main encopy (d_main v_ma_rb) rd_main p_main (d_main v_ma_rb) rd_main v_ma_rb) rd_main v_ma_rb) rd_main v_ma_rb)									10.107		744	
m 4 80 4 80 1		nijagemaju S nijagemajo K	strun_lia ma samax_g ex_main_in stran_l_up stran_l_up stran_g stran_g pid-main_in pid-telper_l	hat, ha nt_main encapy (d_main capyha nt_main (d_main (d_main (d_main (d_main (d_main (d_main (d_main (d_main (d_main (d_main)) (d_main (d_main)) (d_main) (0 (2) 52+	5375	±411		2.57%	12.54%	18.	Ú%.	19476	22	ber	aús

OpenCL Information and Resources



- NVIDIA OpenCL Web Page:
 - http://www.nvidia.com/object/cuda_opencl.html
- NVIDIA OpenCL Forum:
 - http://forums.nvidia.com/index.php?showforum=134
- NVIDIA driver, profiler, code samples for Windows and Linux:
 - https://nvdeveloper.nvidia.com/object/get-opencl.html
- Khronos (current specification):
 - http://www.khronos.org/registry/cl/specs/opencl-1.0.48.pdf
- Khronos OpenCL Forum:
 - http://www.khronos.org/message_boards/viewforum.php?f=28