

OpenCL

A brief introduction for NVIDIA CUDA programmers

Open Computing Language

A language extension (mainly C) for accessing heterogeneous computational devices (CPUs, GPUs, accelerator cards).

- Parallel code execution on single or multiple devices of many vendors (e.g. NVIDIA, AMD/ ATI, Intel, IBM, ARM, ...).
- Vendor neutral and royalty free.
- Specifications under review by Khronos OpenCL working group - developed with industry leaders.
- Profiles for desktop machines, server farms (MPI-integration) and even hand-held devices.
- Interoperable with OpenGL and Microsoft Direct 3D graphics APIs.

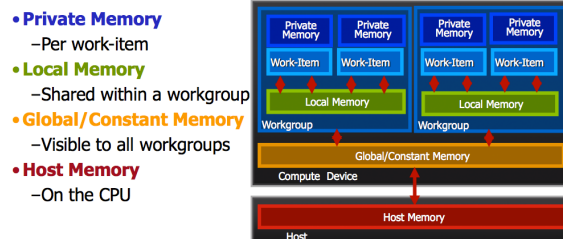
Design goals

An easy way to allow, with sparse new syntax and extensions, parallel computing (data- or task-parallel) with all available devices in a modern system and without writing new code for every type of device - or short:

One portable program uses all available resources.

Abstract hardware model

- *Platform model*: One host system controls multiple computing *devices* (e.g. CPUs or GPUs). Each device contains computing units (*workgroups*) consisting of processing elements (*work-items*).
- *Memory model*: An abstracted shared memory model (like in CUDA), that is independent of the based hardware.

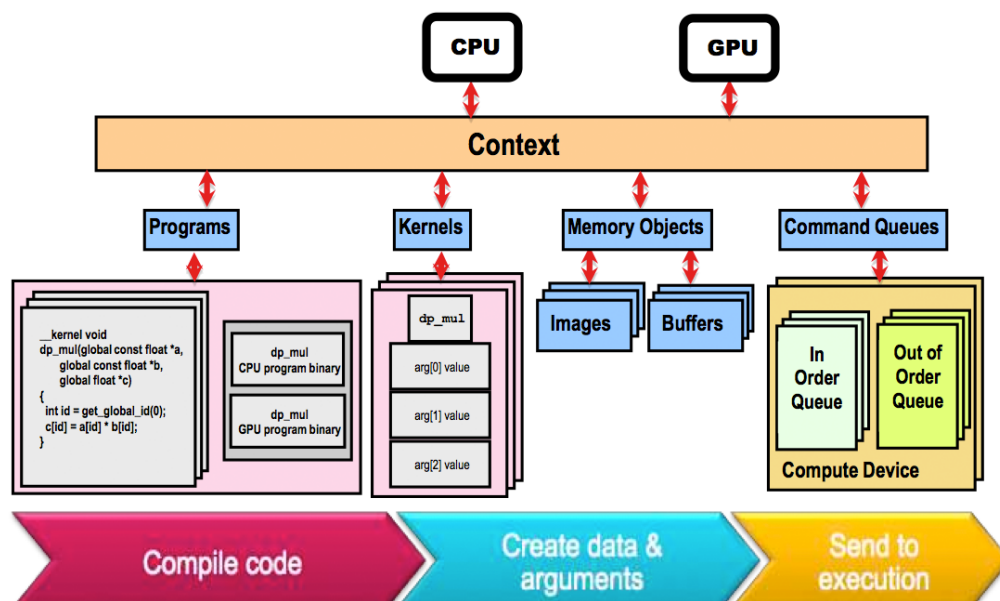


Execution model

The application runs on the host and submits work to the device(s). Basic elements needed to execute OpenCL code¹

- *Work-item*: Basic unit on a OpenCL device (e.g. one ALU/ core).
- *Command queue*: Queues Kernel executions and synchronizes if needed (in-order² or out-of-order³) - event handling.
- *Kernel*: The code to be executed on the devices - extended C function.
- *Program*: Collection of kernels and other functions - analogous to a dynamic library. It has to be compiled for every device architecture at runtime.
- *Context*: Environment to execute the work item in, includes devices, memories and command queues.

Workflow overview



¹ All elements but the work-items must be created by the programmer in the host application.

² In-order queues: One queue per device.

³ Out-of-order queues: Multiple queues per device

Benchmark

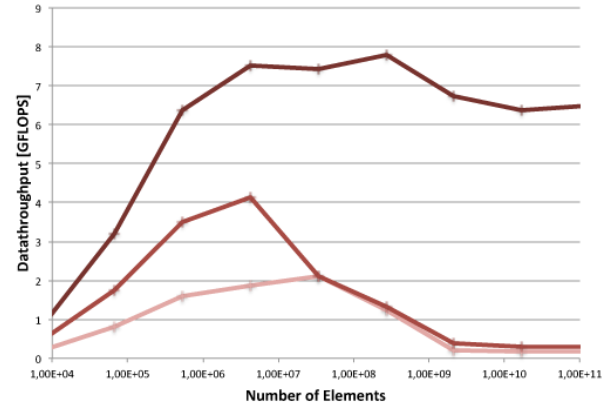
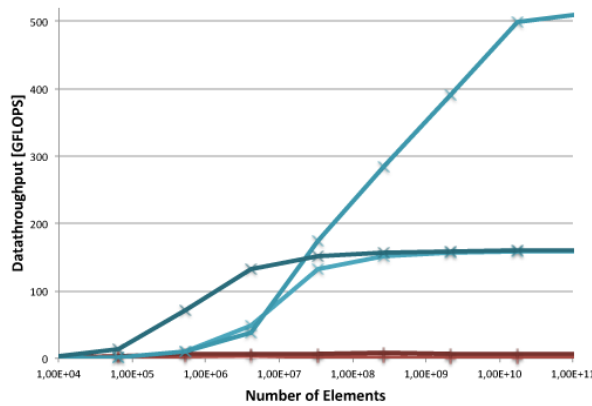
System: CPU: Intel Core i3 - 2 Cores (2 Threads each) @ 3.1 GHz - 3MB cache
RAM: Corsair - 8GB DDR3 (1333 MHz)
GPU: MSI - NVIDIA GeForce GTX 5600Ti - 384 cores @ 880MHz - 1024MB GGDR5 (4200MHz)

Software: Compiler: Microsoft Visual Studio 2008 Professional x86
OS: Microsoft Windows 7 Professional SP1 x64
APIs: OpenCL 1.2⁴, NVIDIA CUDA 4.0⁴, NVIDIA CUBLAS 2.0⁴, OpenMP 2.0⁵

Problem: $A^{2048 \times 8192} \cdot B^{4096 \times 2048} = C^{4096 \times 8196} \Rightarrow 1.374 \cdot 10^{11}$ elements to compute

Results:

		Execution time [s]	Throughput [GFLOPS/s]
CPU	Single threaded	818.3	0.16
	OpenMP	460.05	0.30
	OpenCL 1.2	21.76	6.32
GPU	NVIDIA CUDA 4.0	0.87	158.23
	OpenCL 1.2	0.86	159.72
	NVIDIA CUBLAS 2.0	0.24	580.90



Conclusion: CPU: OpenCL executes the code not on every core available but on every possible thread on a core (e.g. OpenMP doesn't). In addition OpenCL uses a very advanced event handling to queue the kernel execution.

GPU: OpenCL offers approximately the same computation speed as other GPGPU libraries such as NVIDIA CUDA or ATI Stream. The disadvantage to other libraries lies in the more complex workflow that has to be maintained and in the lag of delivered libraries for common tasks like CUBLAS as a fast BLAS⁶ implementation for NVIDIA CUDA.

Summary: OpenCL vs. NVIDIA CUDA

	OpenCL 1.2	NVIDIA CUDA 4.0
Free of charge	✓	✓
Vendor independent/ Multiple SDK-vendors	✓	
Open standard	✓	
Use different types of heterogeneous hardware	✓	
Cross-platform	✓	✓
Interoperability with graphic libraries	✓	✓
Included libraries for common tasks		✓
Peer-to-peer-communication between devices		✓
Global virtual addressing space		✓

Conclusion: OpenCL is slightly more complicated to code (see workflow) than NVIDIA CUDA but is not limited to using GPUs and can use devices from multiple vendors. NVIDIA CUDA is more enhanced and offers more features and well-engineered libraries such as CUBLAS.

References:

Khronos Group: <http://khronos.org/opencv/>
NVIDIA: <http://developer.nvidia.com/opencv/>
Intel: <http://software.intel.com/en-us/articles/vcs-source-tools-opencv-sdk/>
AMD/ ATI: <http://developer.amd.com/zones/OpenCLZone/>

Online resources:

Further information as well as the benchmark and example source code is available at <http://www.janbeneke.de/ParComp>.

⁴Part of the NVIDIA GPU Computing SDK 4.0.

⁵Part of Microsoft Visual Studio 2008 Professional x86.

⁶Basic Linear Algebra Subprograms: A de facto API standard for basic linear algebra operations.