

Programmierung paralleler und verteilter Systeme

GPU Computing with OpenCL

Frank Feinbube

Operating Systems and Middleware
Prof. Dr. Andreas Polze

Agenda / Quicklinks

2

- [Motivation](#)
- [History of GPU Computing](#)
- [Programming Model](#)
- [Examples](#)
- [Development Support](#)
- [Hardware Characteristics](#)
- [Performance Tuning](#)
- [Further Readings](#)



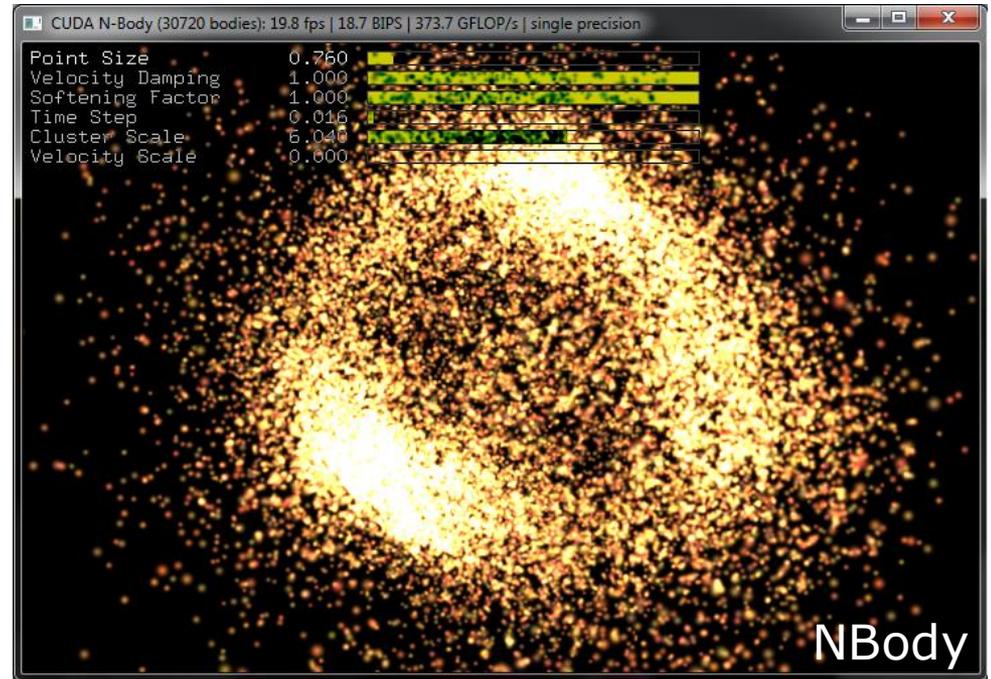
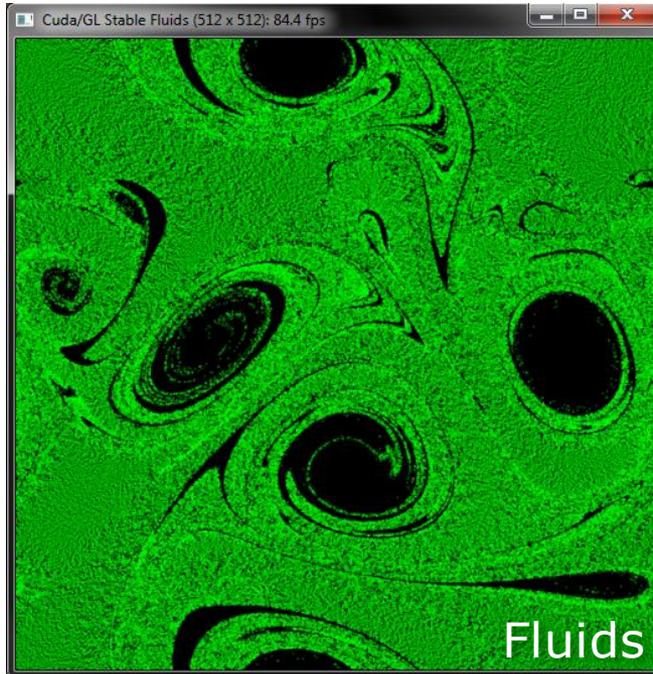
Intel Xeon Phi



NVIDIA Tesla

The Power of GPU Compute Devices

3



```
Using CUDA device [0]: GeForce GTX 275
Sorting 1048576 32-bit unsigned int keys and values
radixSort, Throughput = 74.6231 MElements/s, Time = 0.01405 s, Size = 1048576 elements, NumDevsUsed = 1, Workgroup = 256
PASSED
```

RadixSort

Wide Varsity of Application Domains

4

Research
Medical
Video and Photo
Energy
Finance
Military

<p>Statistical constraints on binary black hole inspi...</p> <p>50 x</p>	<p>BMC Bioinformatics</p>	<p>Acceleration of the Smith-Waterman Algorithm using...</p> <p>5 x</p>	<p>Figure 1. Model of the graphics pipeline.</p> <p>Graphic processors to speed-up simulations for the...</p> <p>420 x</p>
<p>Cmatch: Fast Exact String Matching on the GPU</p> <p>35 x</p>	<p>Quantitative Risk Analysis and Algorithmic Trading</p> <p>50 x</p>	<p>Harvesting graphics power for...</p> <p>150 x</p>	<p>Graphic-Card Cluster for Astrophysics (GraCCA)</p> <p>250 x</p>
<p>Quantum Chemistry Two-Electron Integral Evolution</p> <p>130 x</p>	<p>Accelerating Statistical Static Timing Analysis</p> <p>260 x</p>	<p>Distributed Password Recovery</p> <p>50 x</p>	<p>Folding@home</p> <p>100 x</p>

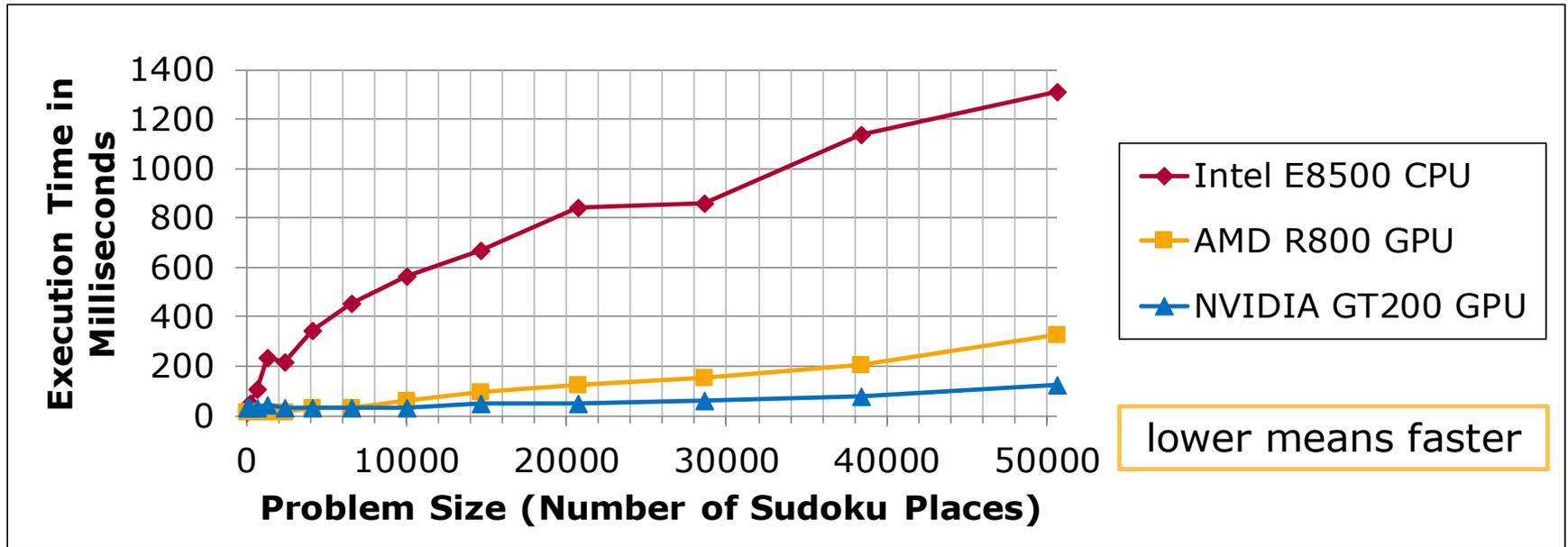
http://www.nvidia.com/object/cuda_apps_flash_new.html
http://www.nvidia.com/object/tesla_testimonials.html

Why GPU Compute Devices?

Short Term View: Cheap Performance

5

Performance

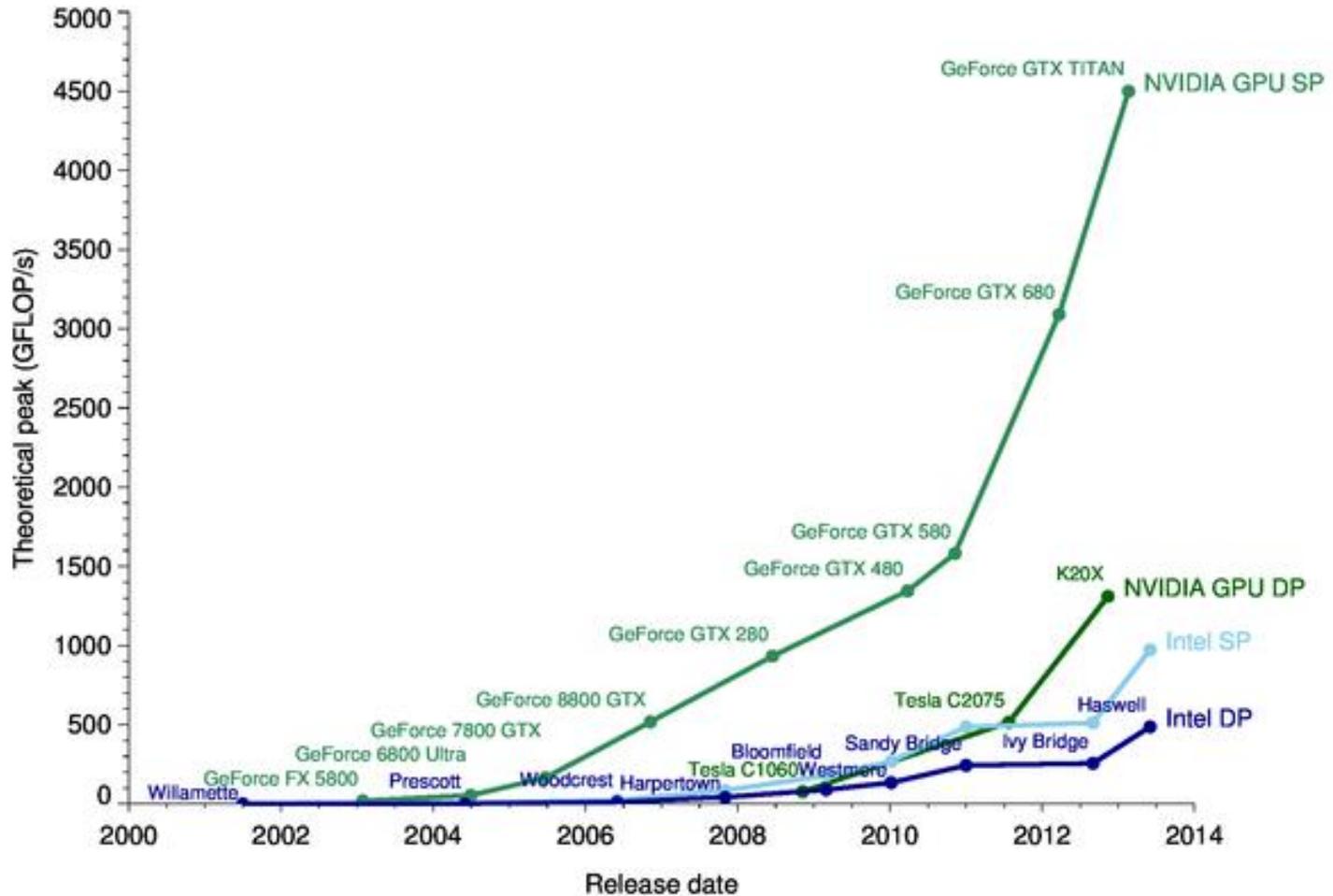


Energy / Price

- Cheap to buy and to maintain
- GFLOPS per watt: Fermi 1,5 / Kepler 5 / Maxwell 15

Why GPU Compute Devices? Middle Term View: More Performance

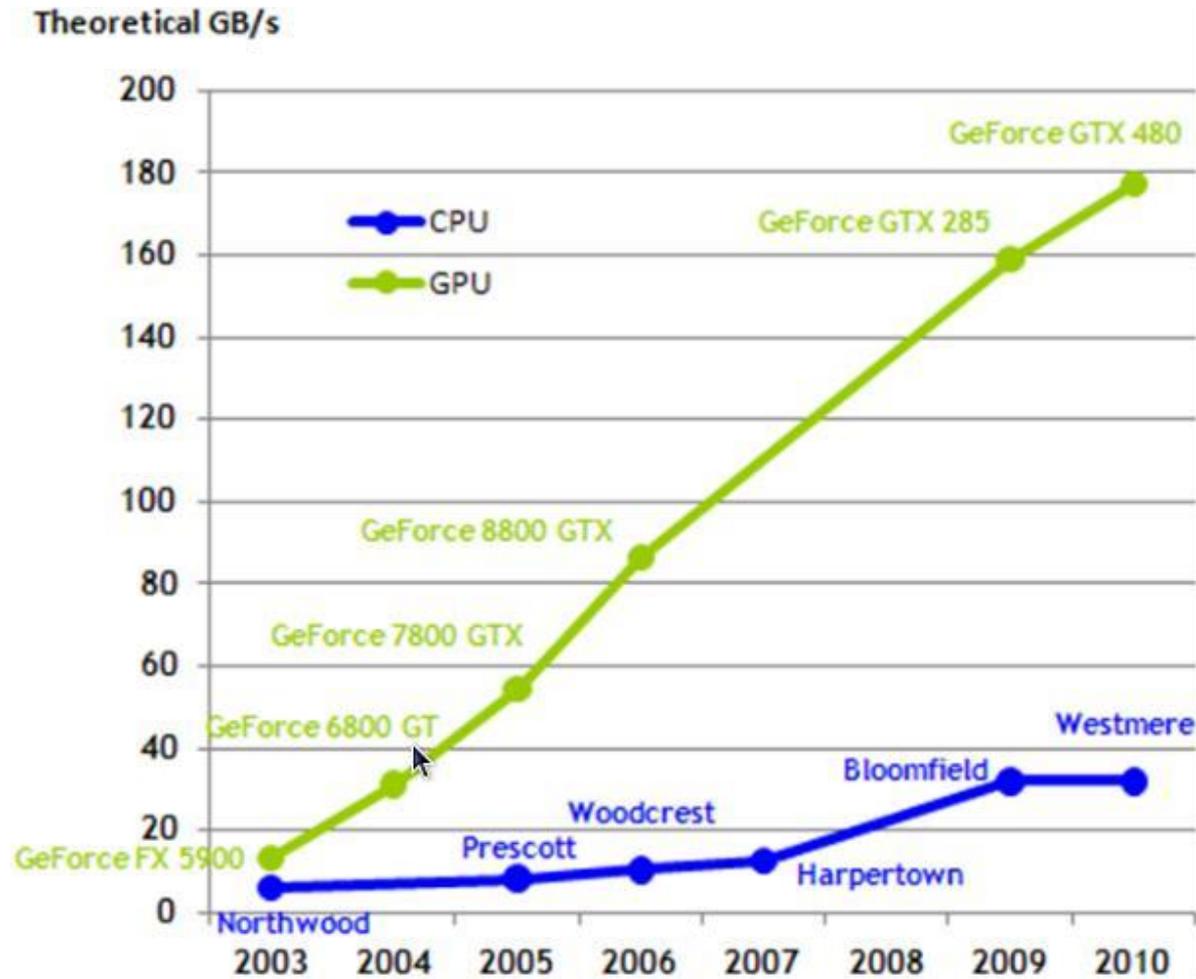
6



Why GPU Compute Devices?

Middle Term View: More Performance

7



Why GPU Compute Devices?

Long Term View: Hybrid Computing

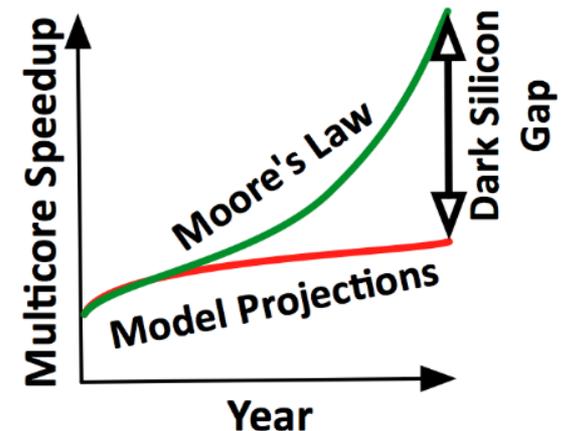
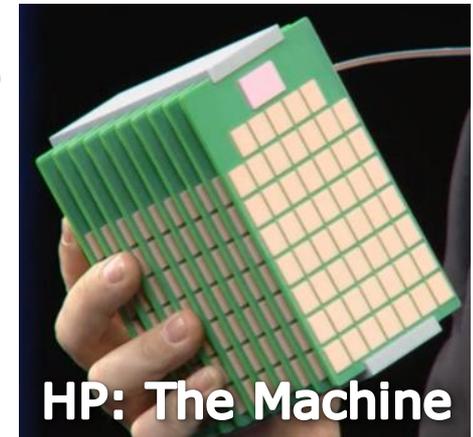
8

Dealing with massively multi-core:

- New architectures are evaluated (Intel, IBM, HP,...)
- Accelerators that accompany common general purpose CPUs (Hybrid Systems)

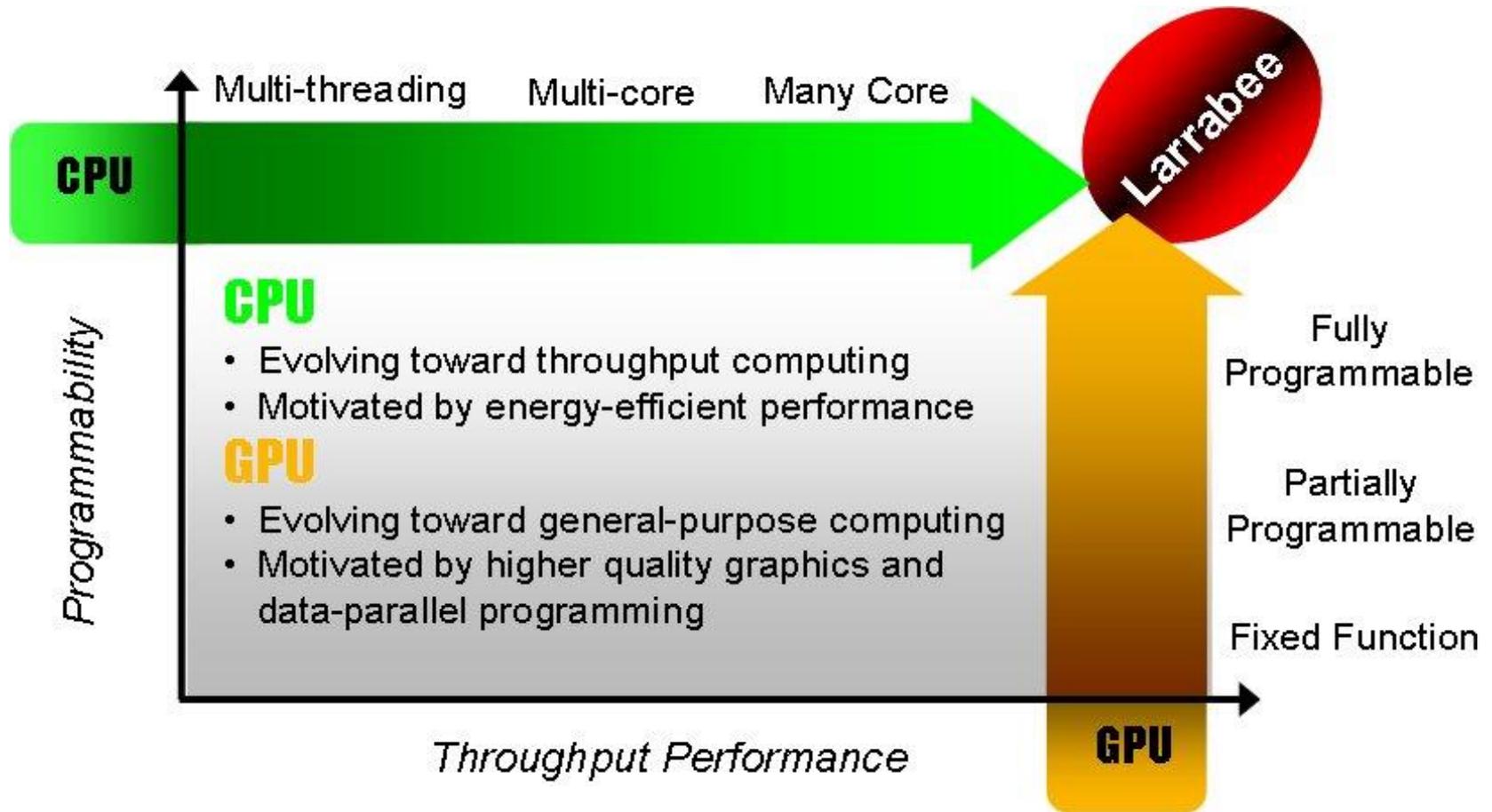
Hybrid Systems

- **GPU Compute Devices:**
High Performance Computing (3 of top 5 supercomputers are GPU-based!),
Business Servers, Home/Desktop Computers, Mobile and Embedded Systems
- **Special-Purpose Accelerators:**
(de)compression, XML parsing,
(en|de)cryption, regular expression matching



History of GPU Computing

9



History of GPU Computing

10

Fixed Function Graphic Pipelines

- 1980s-1990s; configurable, not programmable; first APIs (DirectX, OpenGL); Vertex Processing

Programmable Real-Time Graphics

- Since 2001: APIs for Vertex Shading, Pixel Shading and access to texture; DirectX9

Unified Graphics and Computing Processors

- 2006: NVIDIAs G80; unified processors arrays; three programmable shading stages; DirectX10

General Purpose GPU (GPGPU)

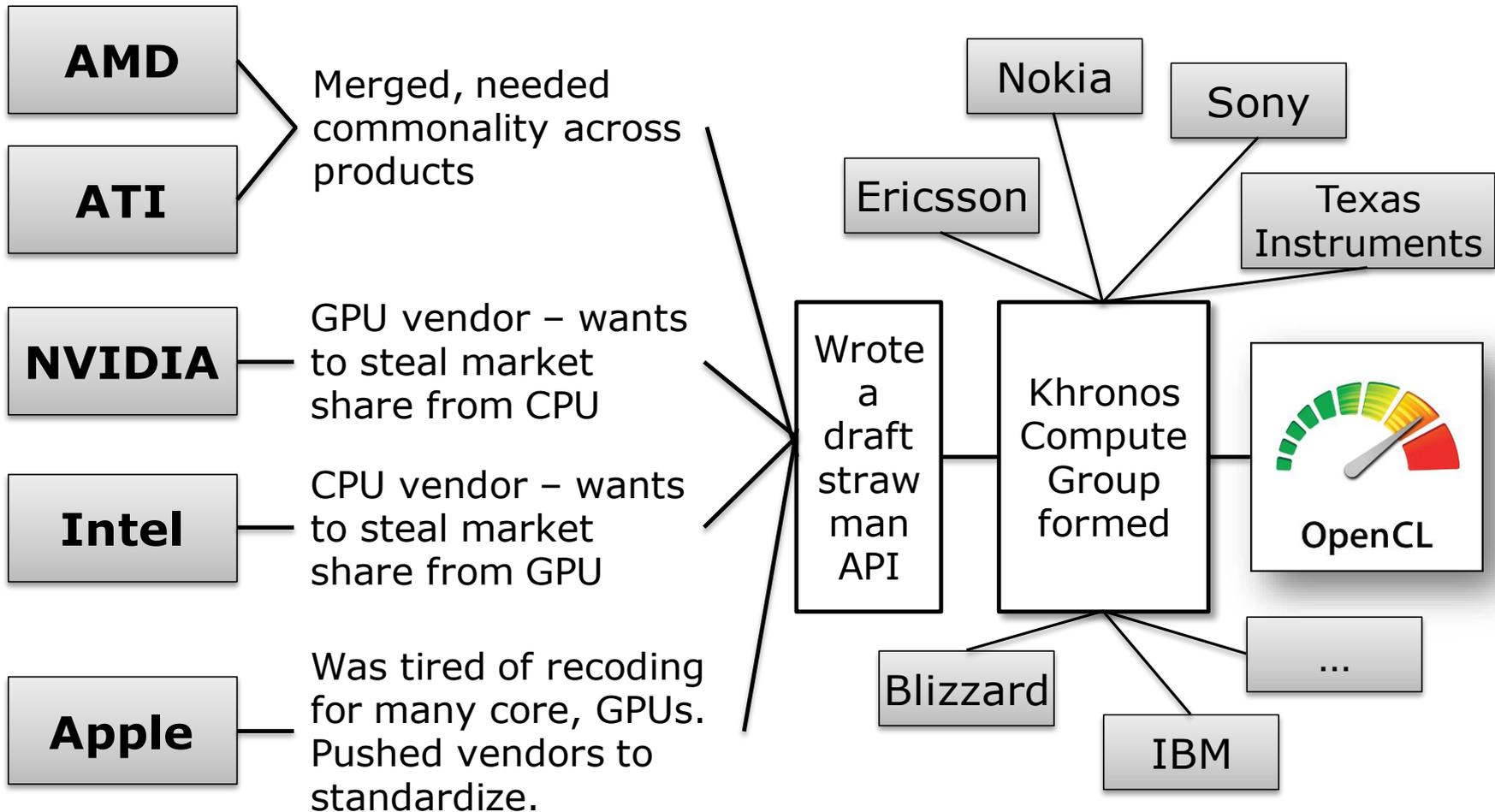
- compute problem as native graphic operations; algorithms as shaders; data in textures

GPU Computing

- Programming CUDA; shaders programmable; load and store instructions; barriers; atomics

Open Compute Language (OpenCL)

11



[5]

Open Compute Language (OpenCL)

12

- Hardware vendors, system OEMs, middleware vendors, application developers
- OpenCL became an important standard “on release” by virtue of the market coverage of the companies behind it.
- OpenCL implementations already exist for AMD, NVIDIA, Intel, IBM, ARM, ...

- Use all computational resources in system
 - Program GPUs, CPUs, and other processors as peers
 - Efficient C-based parallel programming model
 - Abstract the specifics of underlying hardware
- Abstraction is **low-level, high-performance but device-portable**
 - Approachable – but primarily targeted at expert developers
 - Ecosystem foundation – no middleware or “convenience” functions
- Implementable on a range of embedded, desktop, and server systems
 - HPC, desktop, and handheld profiles in one specification

Programming Models

13

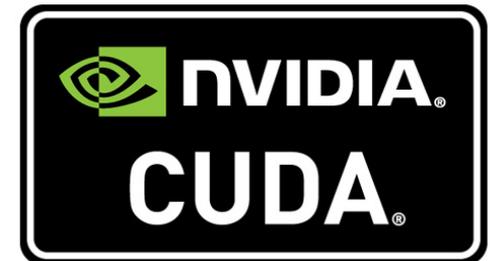
~~AMD: ATI Stream SDK~~

- Today using OpenCL



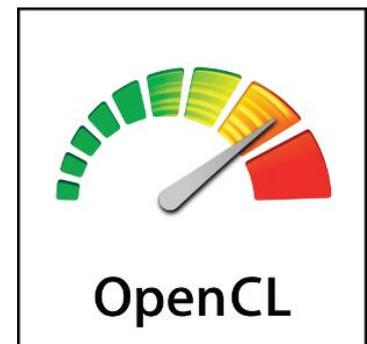
NVIDIA: Common Unified Device Architecture

- CUDA C/C++ compiler, libraries, runtime
- Mature: literature, examples, tool, development support



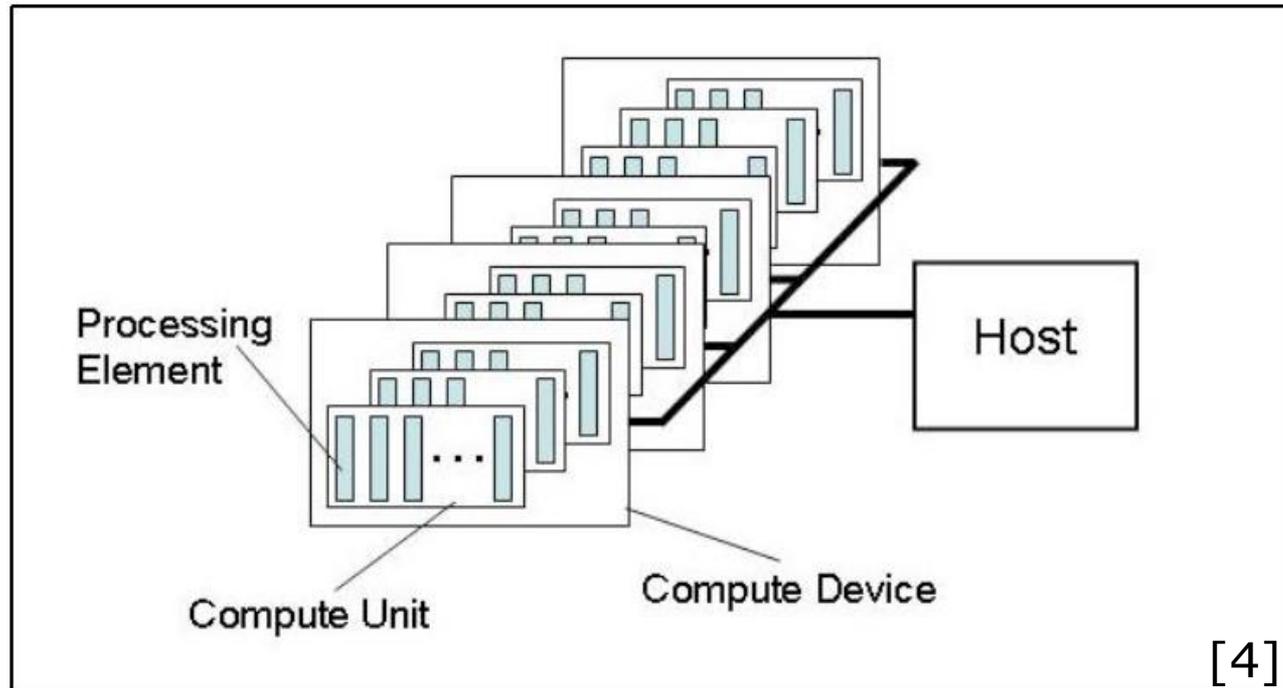
Khronos Group: OpenCL

Open standard for portable, parallel programming of heterogeneous parallel computing CPUs, GPUs, and other processors



OpenCL Platform Model

14



- OpenCL exposes CPUs, GPUs, and other Accelerators as “devices”
- Each “device” contains one or more “compute units”, i.e. cores, SMs,...
- Each “compute unit” contains one or more SIMD “processing elements”

The BIG idea behind OpenCL

15

OpenCL execution model ... execute a kernel at each point in a problem domain.

Traditional loops

```
void
trad_mul(int n,
         const float *a,
         const float *b,
         float *c)
{
    int i;
    for (i=0; i<n; i++)
        c[i] = a[i] * b[i];
}
```



Data Parallel OpenCL

```
kernel void
dp_mul(global const float *a,
       global const float *b,
       global float *c)
{
    int id = get_global_id(0);
    c[id] = a[id] * b[id];
} // execute over "n" work-items
```

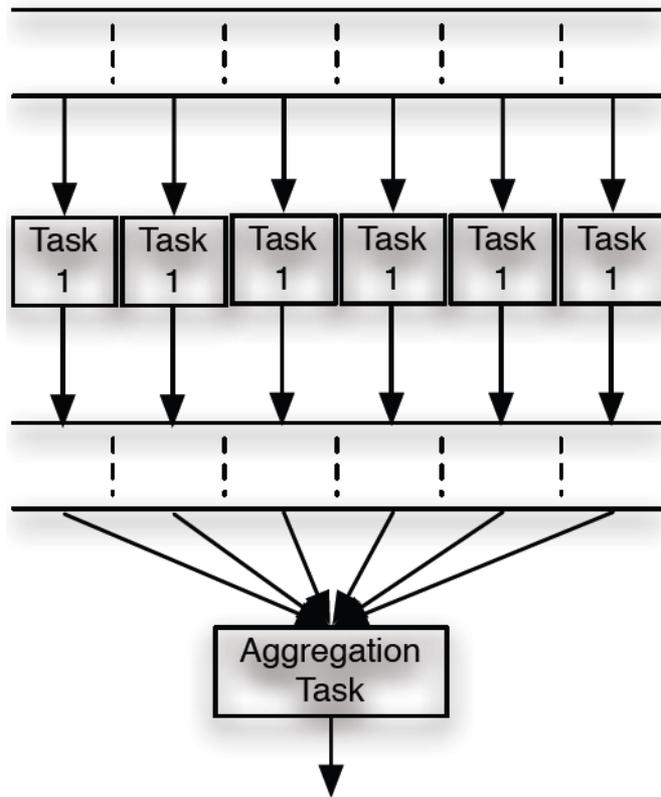
E.g., process a 1024 x 1024 image with one kernel invocation per pixel or $1024 \times 1024 = 1,048,576$ kernel executions

[5]

Data Parallelism versus Task Parallelism

16

Data Parallelism

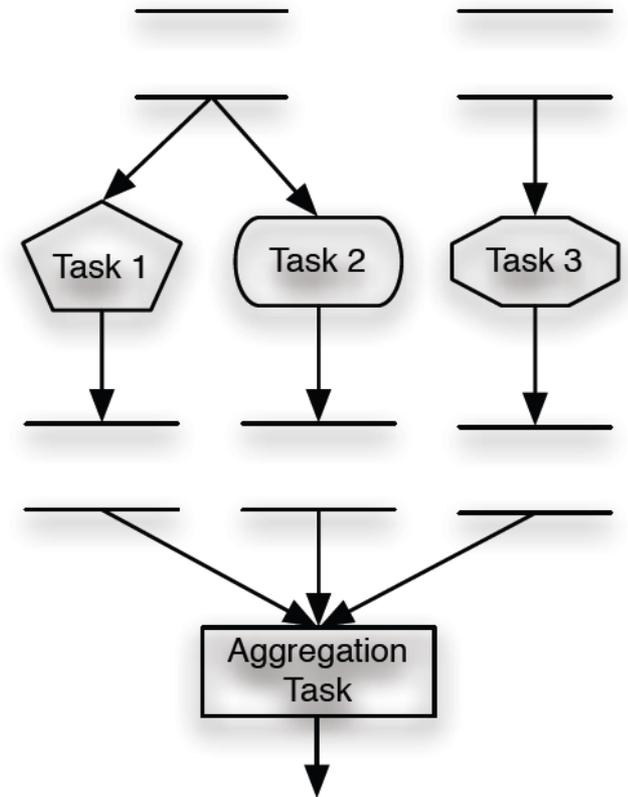


Task Parallelism

Input Data

Parallel Processing

Result Data

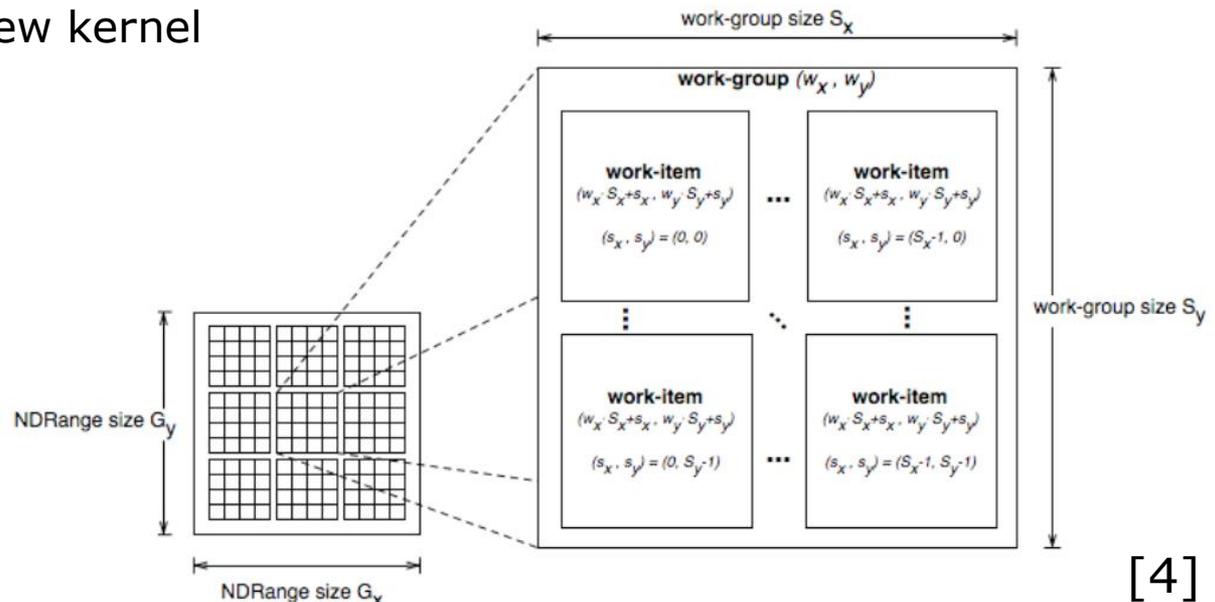


OpenCL is designed for SIMD / SPMD approaches

OpenCL Execution Model

17

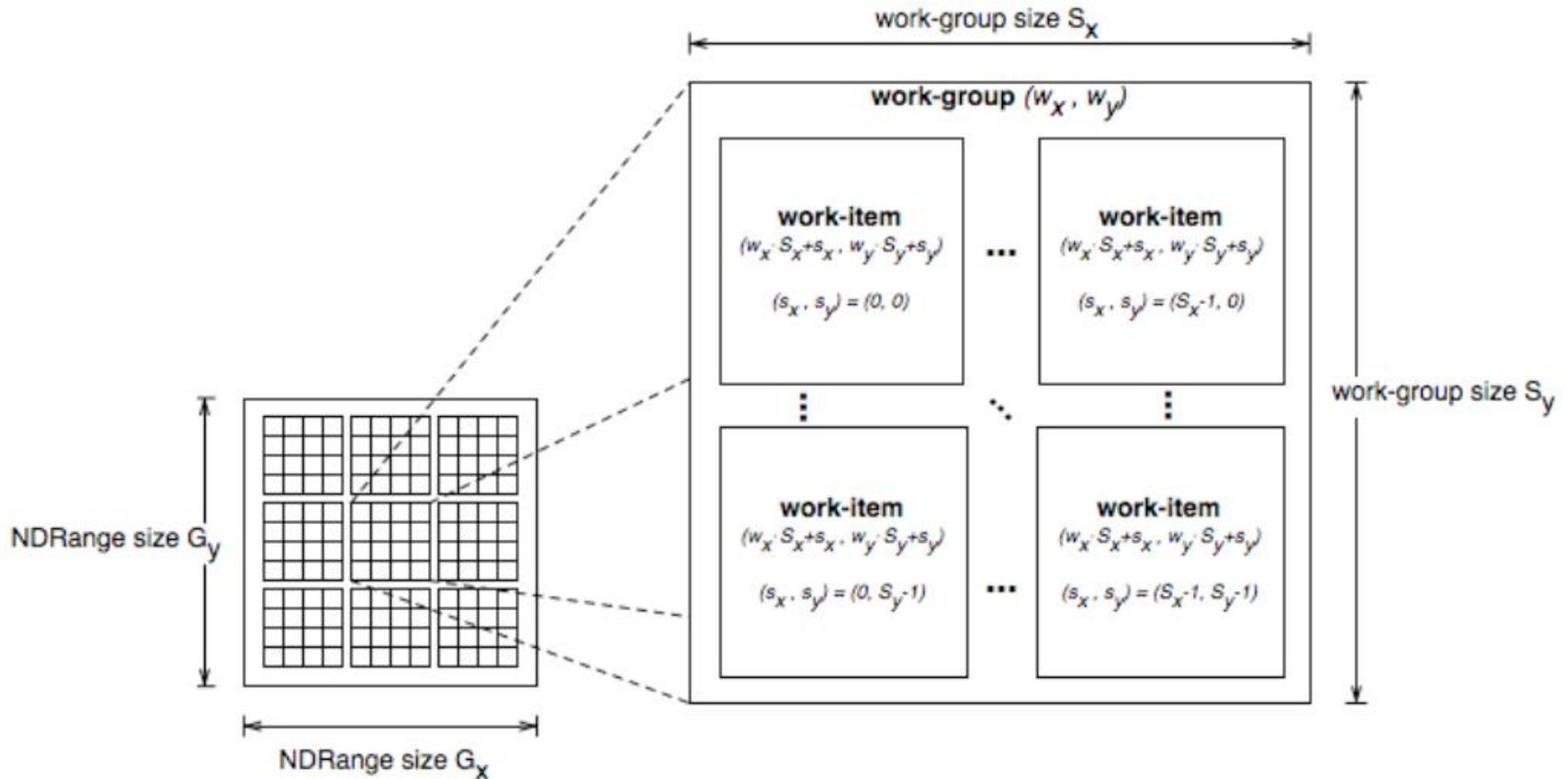
- Parallel work is submitted to devices by launching kernels
- Kernels run over global dimension index ranges (NDRange), broken up into “work groups”, and “work items”
- Work items executing within the same work group can synchronize with each other with barriers or memory fences
- Work items in different work groups can't sync with each other, except by launching a new kernel



[4]

OpenCL Execution Model

18



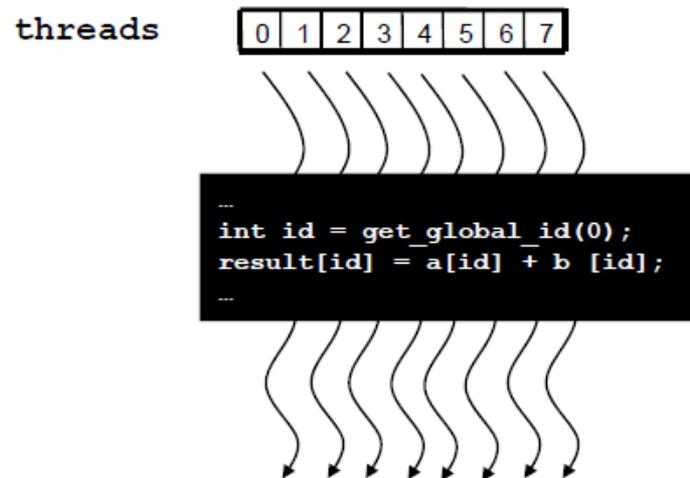
An example of an NDRange index space showing work-items, their global IDs and their mapping onto the pair of work-group and local IDs. [4]

OpenCL Execution Model

19

An OpenCL kernel is executed by an array of work items.

- All work items run the same code (SPMD)
- Each work item has an index that it uses to compute memory addresses and make control decisions



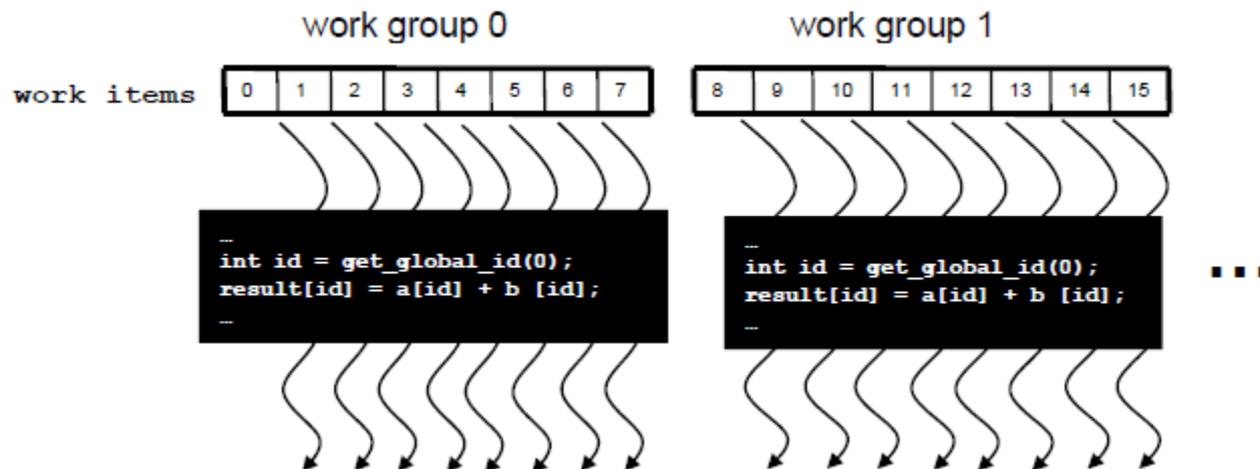
[1]

Work Groups: Scalable Cooperation

20

Divide monolithic work item array into work groups

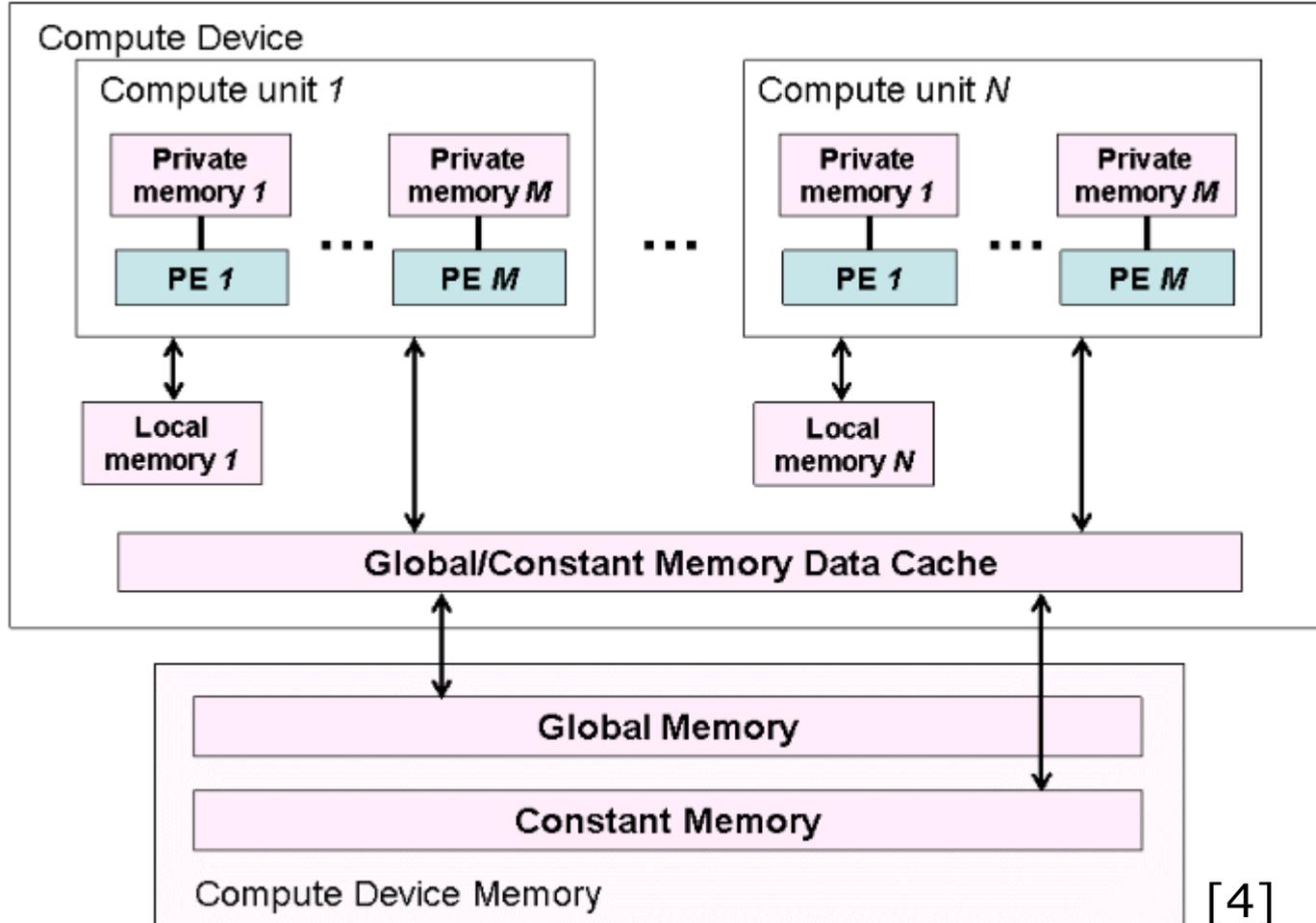
- Work items within a work group cooperate via **shared memory, atomic operations** and **barrier synchronization**
- Work items in different work groups cannot cooperate



[1]

OpenCL Memory Architecture

21



[4]

Private
Per work-item

Local
Shared within
a workgroup

**Global/
Constant**
Visible to
all workgroups

Host Memory
On the CPU

OpenCL Memory Architecture

22

- Memory management is explicit: you must move data from host → global → local... and back

Memory Type	Keyword	Description / Characteristics
Global Memory	<code>__global</code>	Shared by all work items; read/write; may be cached (modern GPU), else slow; huge
Private Memory	<code>__private</code>	For local variables; per work item; may be mapped onto global memory (Arrays on GPU)
Local Memory	<code>__local</code>	Shared between workitems of a work group; may be mapped onto global memory (not GPU), else fast; small
Constant Memory	<code>__constant</code>	Read-only, cached; add. special kind for GPUs: texture memory

OpenCL Work Item Code

23

A subset of ISO C99 - without some C99 features

- headers, function pointers, recursion, variable length arrays, and bit fields

A superset of ISO C99 with additions for

- Work-items and workgroups
- Vector types (2,4,8,16): endian safe, aligned at vector length
- Image types mapped to texture memory
- Synchronization
- Address space qualifiers

Also includes a large set of built-in functions for image manipulation, work-item manipulation, specialized math routines, vectors, etc. [5]

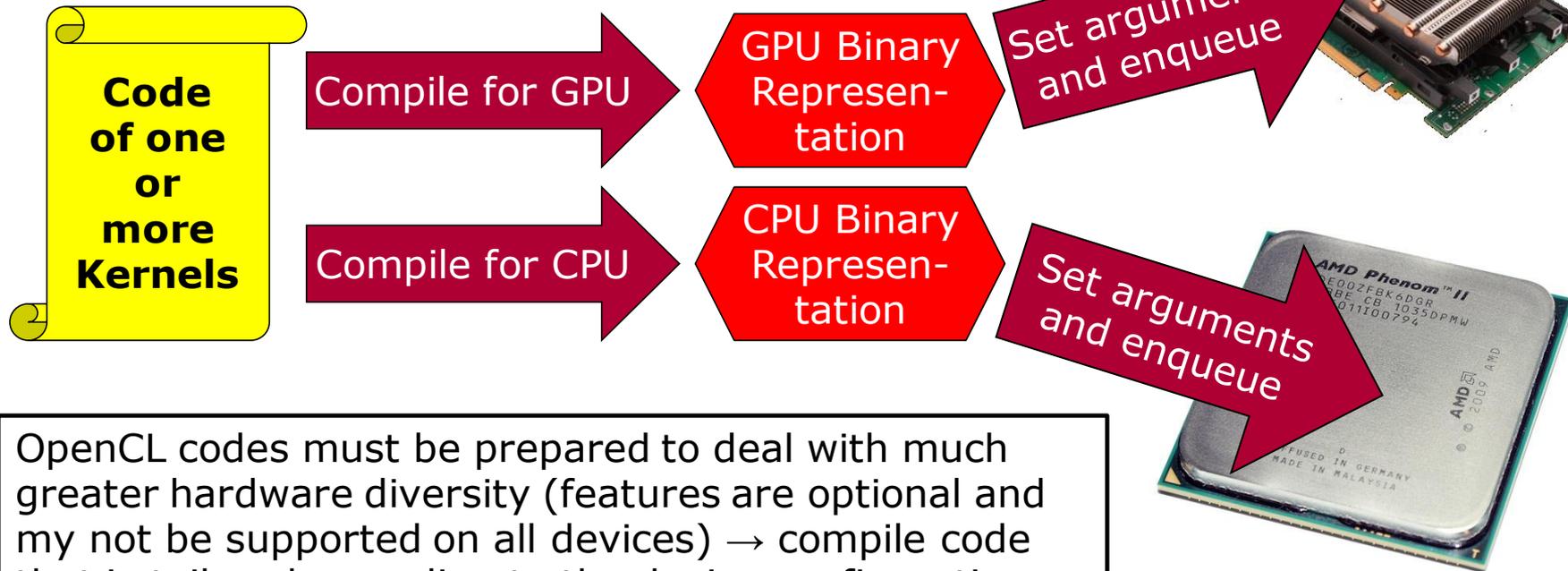
Building and Executing OpenCL Code

24

Kernel

Program

Device



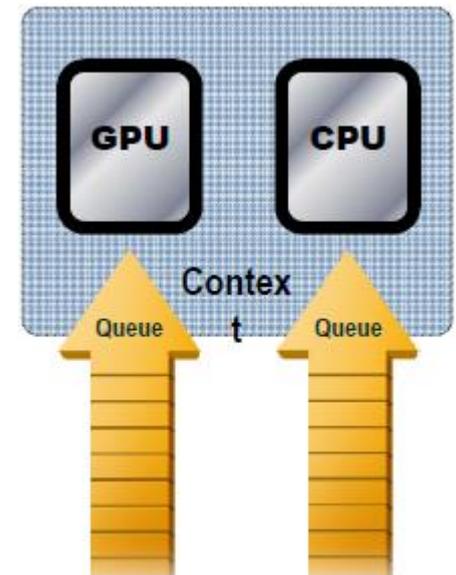
OpenCL codes must be prepared to deal with much greater hardware diversity (features are optional and may not be supported on all devices) → compile code that is tailored according to the device configuration

OpenCL Execution Model

25

An OpenCL application runs on a host which submits work to the compute devices. Kernels are executed in contexts defined and manipulated by the host.

- **Work item:** the basic unit of work on an OpenCL device.
- **Kernel:** the code for a work item. Basically a C function
- **Program:** Collection of kernels and other functions (Analogous to a dynamic library)
- **Context:** The environment within which work-items executes ... includes devices and their memories and command queues.
- **Queue:** used to manage a device. (copy memory, start work item, ...) In-order vs. out-of-order execution

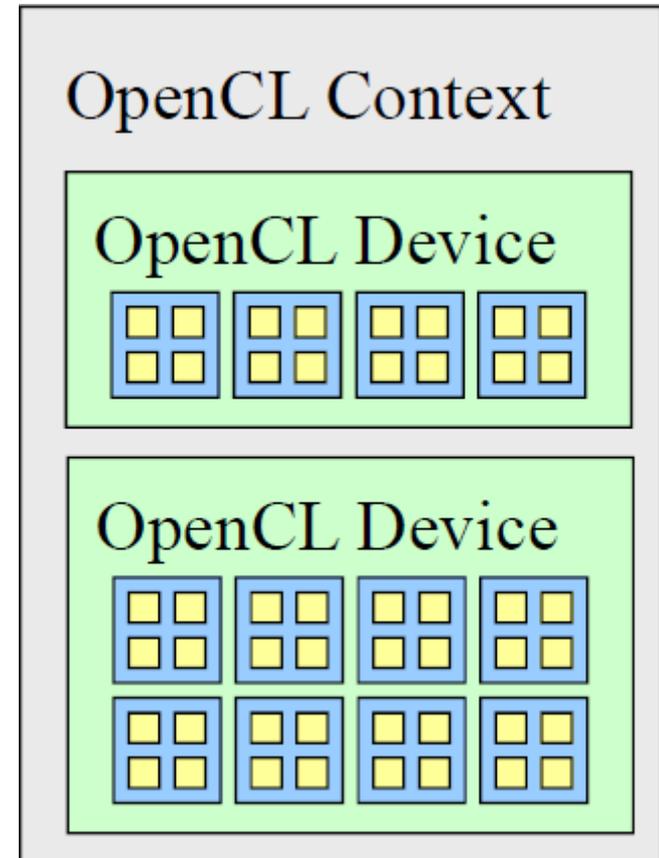


[5]

OpenCL Context

26

- Contains one or more devices
- OpenCL memory objects are associated with a context, not a specific device
- `clCreateBuffer()` is the main data object allocation function
 - error if an allocation is too large for any device in the context
- Each device needs its own work queue(s)
- Memory transfers are associated with a command queue (thus a specific device)

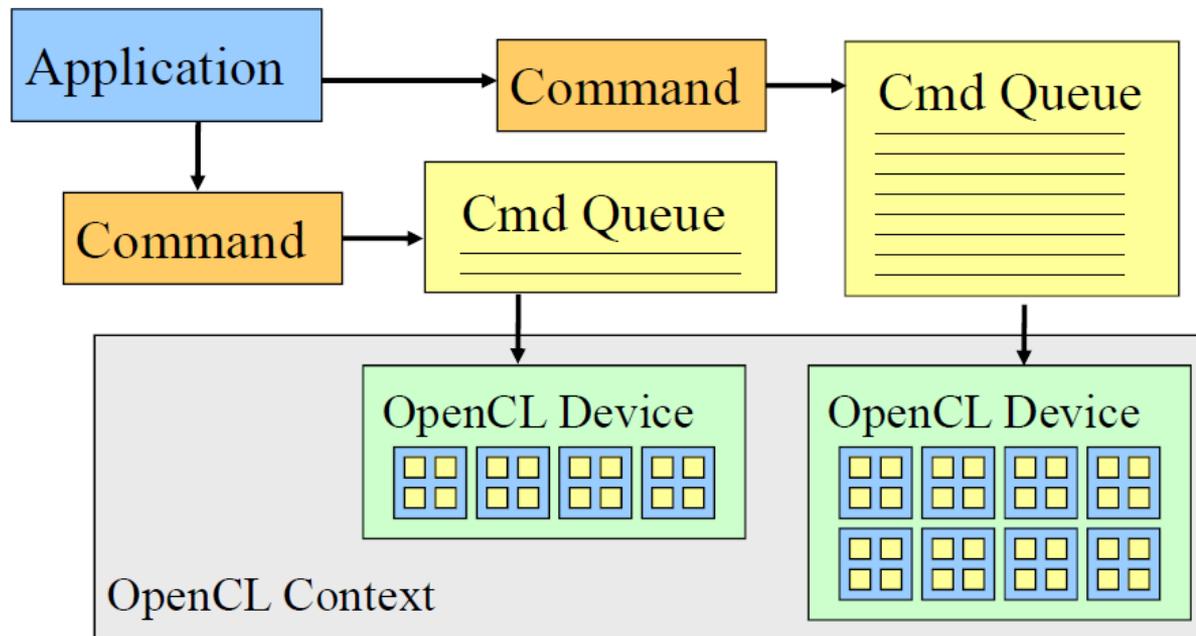


[1]

OpenCL Device Command Execution

27

- Command-queue - coordinates execution of kernels
 - Kernel execution commands
 - Memory commands: transfer or mapping of memory object data
 - Synchronization commands: constrains the order of commands



[1]

Vector Addition: Kernel

28

- Kernel body is instantiated once for each work item; each getting an unique index

```

__kernel void vec_add (__global const float *a,
                      __global const float *b,
                      __global float *c)
{
    int gid = get_global_id(0);
    c[gid] = a[gid] + b[gid];
}

```

» Code that actually executes on target devices

[5]

Vector Addition: Host Program

```

// create the OpenCL context on a GPU device
cl_context = clCreateContextFromType(0,
    CL_DEVICE_TYPE_GPU, NULL, NULL, NULL);

// get the list of GPU devices associated with
// context
clGetContextInfo(context, CL_CONTEXT_DEVICES, 0,
    NULL, &cb);
devices = malloc(cb);
clGetContextInfo(context, CL_CONTEXT_DEVICES, cb,
    devices, NULL);

// create a command-queue
cmd_queue = clCreateCommandQueue(context,
    devices[0], 0, NULL);

// allocate the buffer memory objects
memobjs[0] = clCreateBuffer(context,
    CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
    sizeof(cl_float)*n, srcA,
    NULL);
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_ONLY
    | CL_MEM_COPY_HOST_PTR, sizeof(cl_float)*n, srcB,
    NULL);
memobjs[2] =
    clCreateBuffer(context, CL_MEM_WRITE_ONLY,
        sizeof(cl_float)*n,
    NULL,
    NULL);

// create the program
program = clCreateProgramWithSource(context, 1,
    &program_source, NULL, NULL);

// build the program
err = clBuildProgram(program, 0, NULL, NULL, NULL,
    NULL);

// create the kernel
kernel = clCreateKernel(program, "vec_add", NULL);

// set the args values
err = clSetKernelArg(kernel, 0, (void *) &memobjs[0],
    sizeof(cl_mem));
err |= clSetKernelArg(kernel, 1, (void *) &memobjs[1],
    sizeof(cl_mem));
err |= clSetKernelArg(kernel, 2, (void *) &memobjs[2],
    sizeof(cl_mem));

// set work-item dimensions
global_work_size[0] = n;

// execute kernel
err = clEnqueueNDRangeKernel(cmd_queue, kernel, 1,
    NULL, global_work_size, NULL, 0, NULL, NULL);

// read output array
err = clEnqueueReadBuffer(cmd_queue, memobjs[2],
    CL_TRUE, 0, n*sizeof(cl_float), dst, 0, NULL, NULL);

```

Vector Addition: Host Program

Define platform and queues

```
devices[0], 0, NULL);
```

```
// allocate the buffer memory objects
```

Define Memory objects

```
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_ONLY,
    | CL_MEM_COPY_HOST_PTR, sizeof(cl_float)*n, SrcB,
    NULL);
memobjs[2] =
    clCreateBuffer(context, CL_MEM_WRITE_ONLY,
```

Create the program

Build the program

Create and setup kernel

```
// set the args values
err = clSetKernelArg(kernel, 0, (void *) &memobjs[0],
    sizeof(cl_mem));
err |= clSetKernelArg(kernel, 1, (void *) &memobjs[1],
    sizeof(cl_mem));
err |= clSetKernelArg(kernel, 2, (void *) &memobjs[2],
    sizeof(cl_mem));
```

```
// set work-item dimensions
```

```
global_work_size[0] = 1;
```

Execute the kernel

```
// execute
```

```
err = clEnqueueNDRangeKernel(cmd_queue, kernel, 1,
    NULL, global_work_size, NULL, -1, NULL, NULL);
```

```
// read
```

```
err = clReadBuffer(cmd_queue, memobjs[2],
    0, TRUE,
```

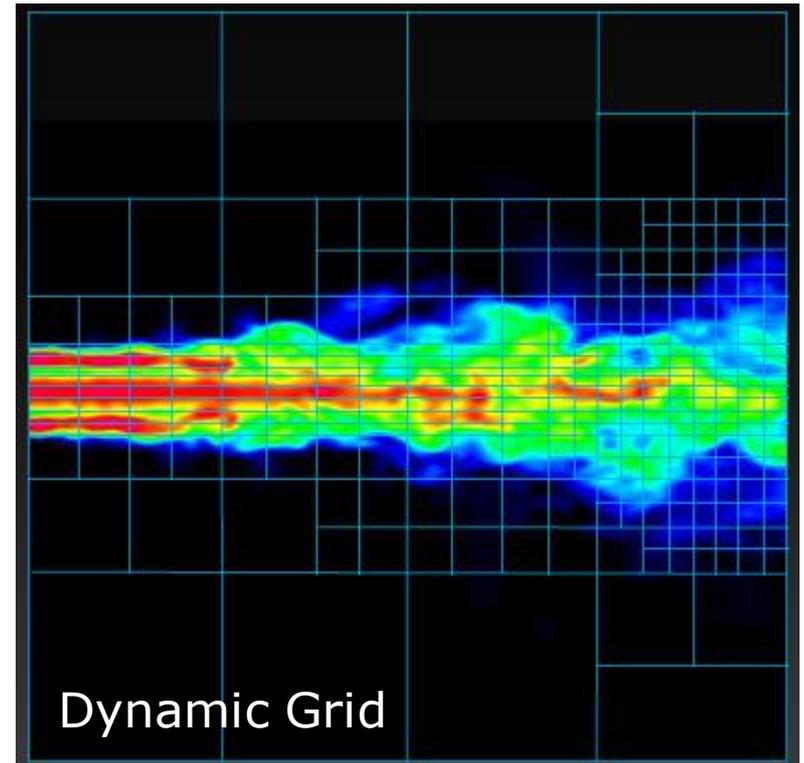
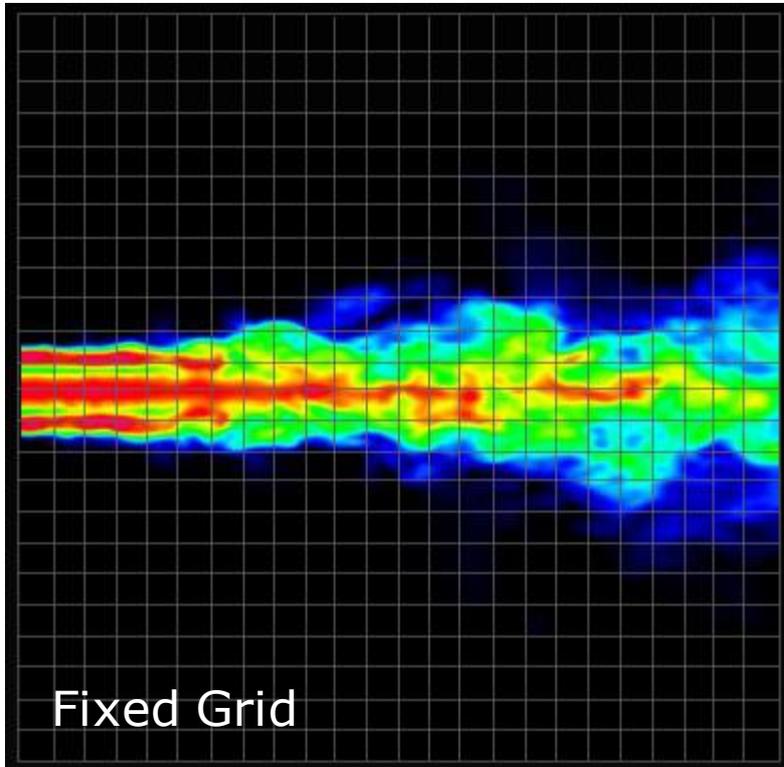
„standard“ overhead
for an OpenCL program

OpenCL "Hello Device"

Java+OpenCL "Fractal Explorer"

Dynamic Parallelism: The Idea

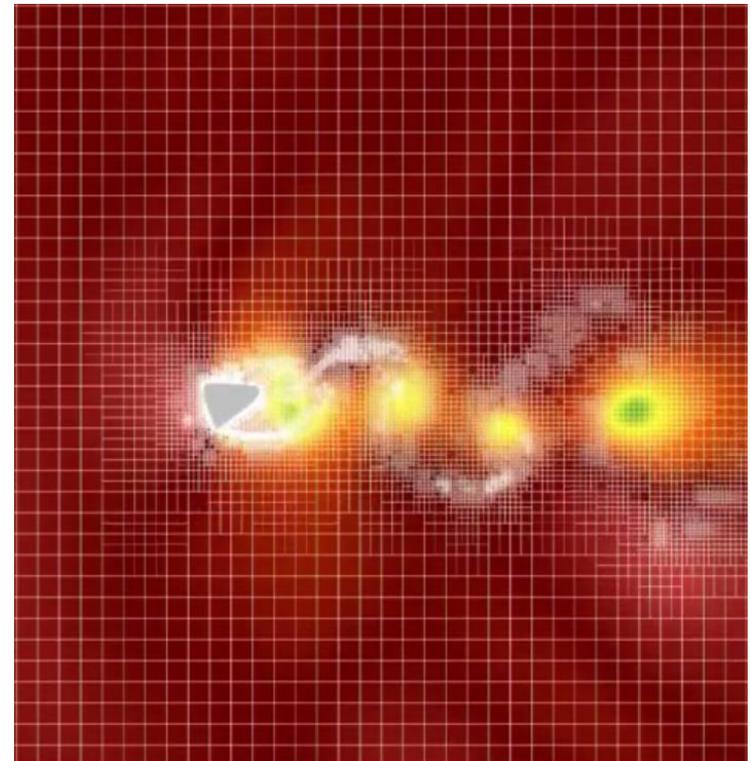
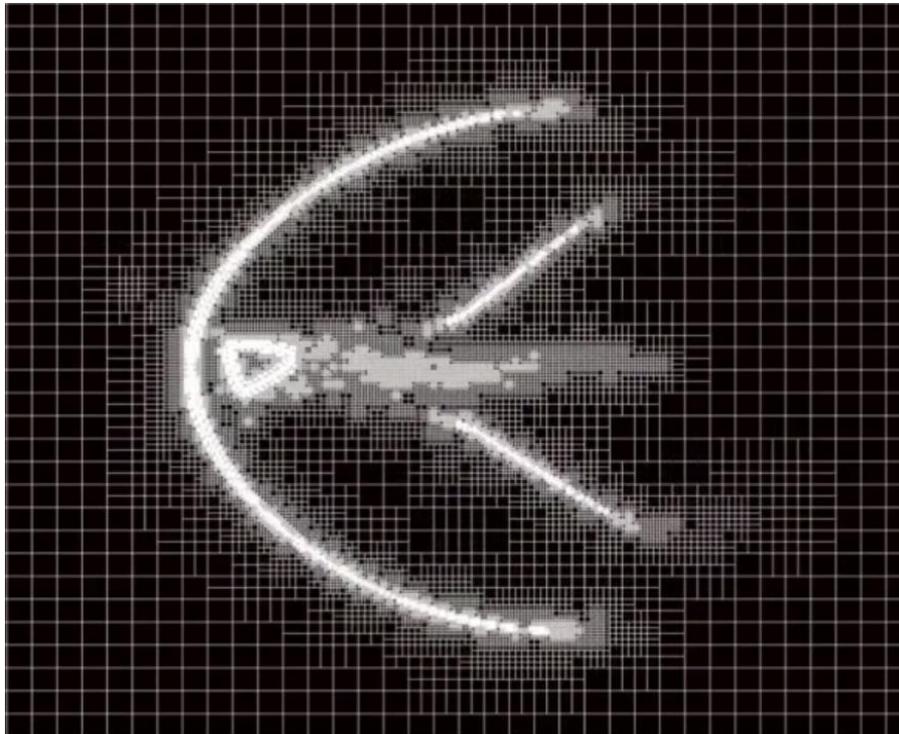
32



[10]

Dynamic Parallelism in Action: SpaceX

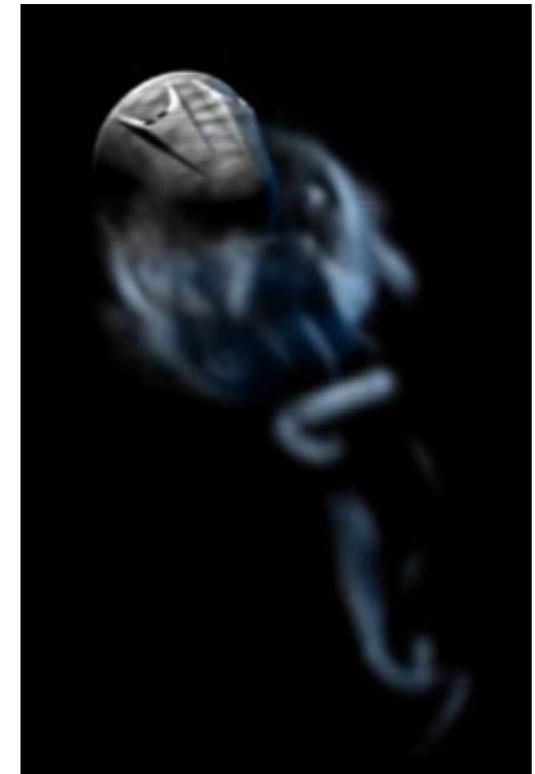
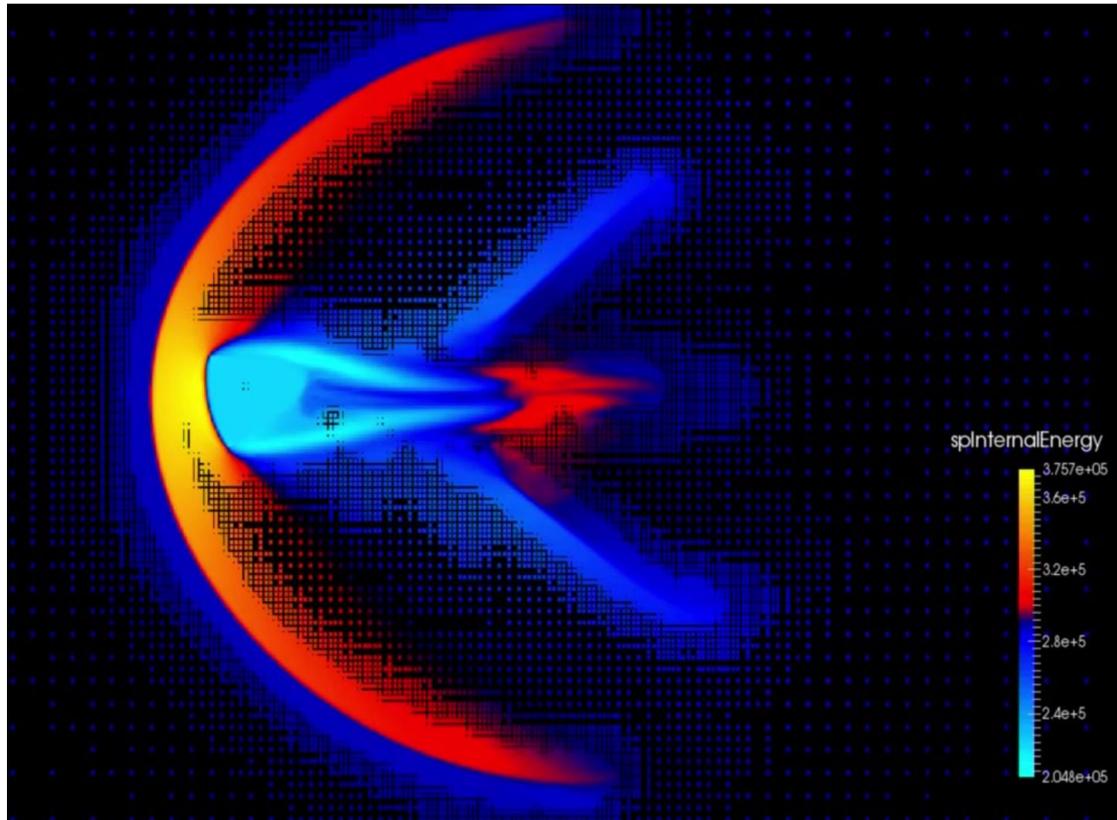
33



<http://on-demand.gputechconf.com/gtc/2015/video/S5398.html>

Dynamic Parallelism in Action: SpaceX

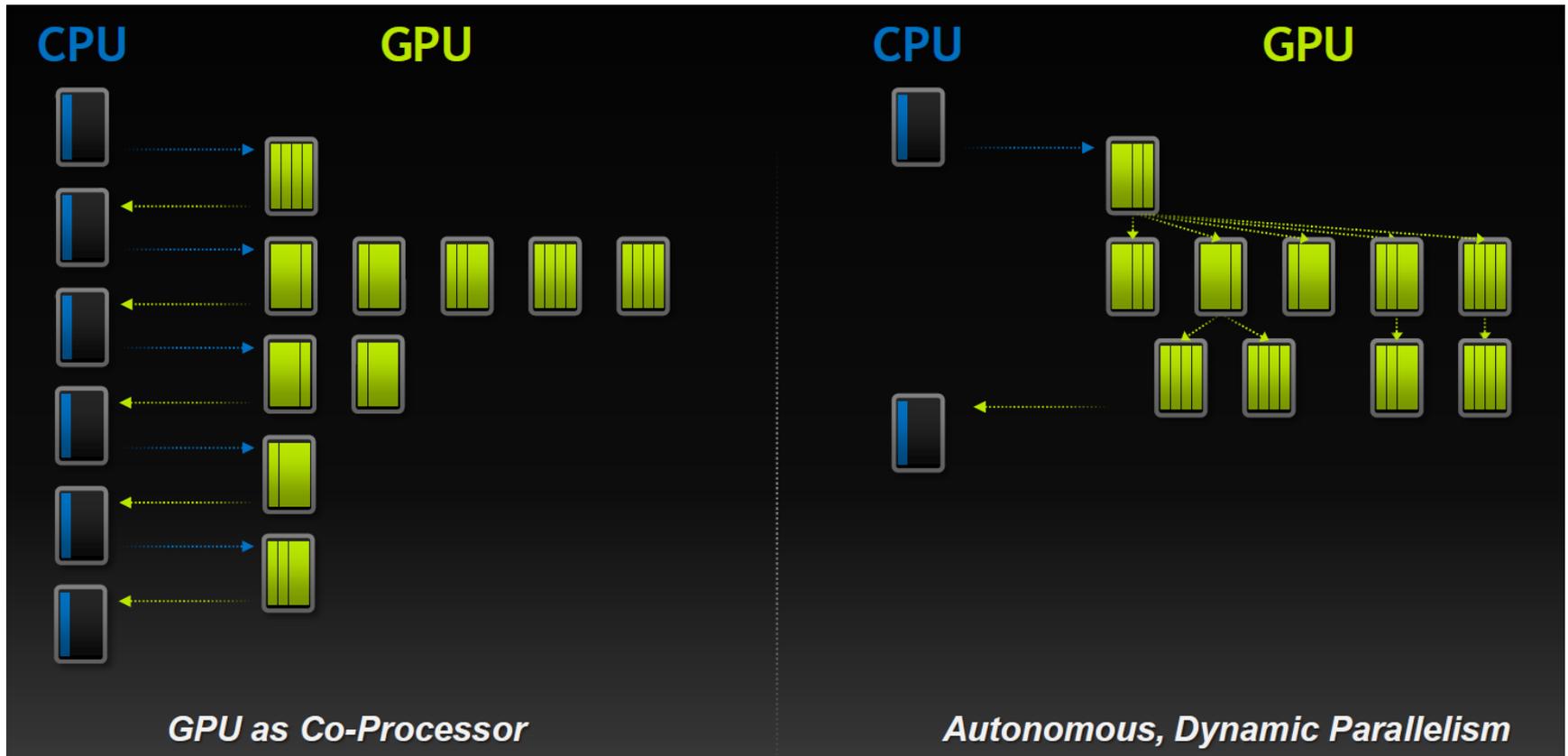
34



<http://on-demand.gputechconf.com/gtc/2015/video/S5398.html>

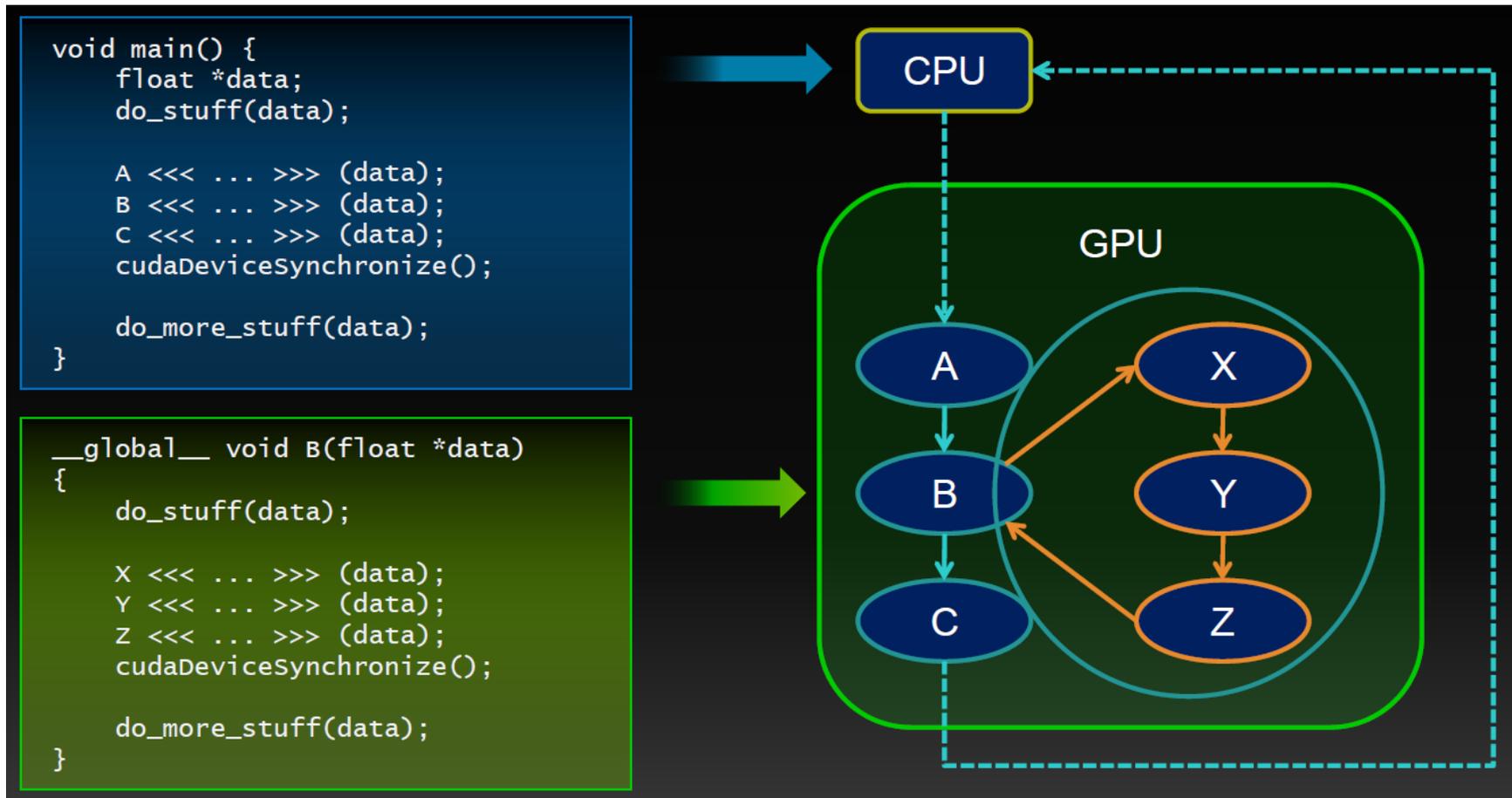
Dynamic Parallelism: Task Distribution

35



Dynamic Parallelism: Familiar Syntax

36



Dynamic Parallelism: Code Example

37

- CUDA Runtime syntax & semantics
- Launch is per-thread
- Sync includes all launches by any thread in the block
- *cudaDeviceSynchronize()* does not imply *syncthreads*
- Asynchronous launches only
(note bug in program, here!)

Code Example

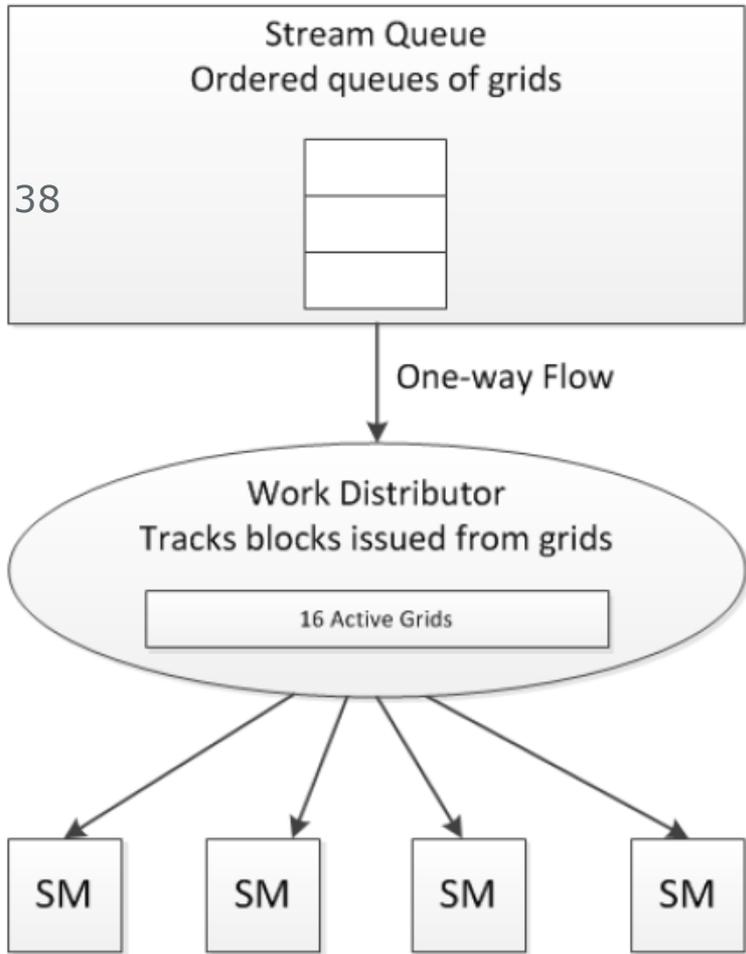
```

__device__ float buf[1024];
__global__ void dynamic(float *data)
{
    int tid = threadIdx.x;
    if(tid % 2)
        buf[tid/2] = data[tid]+data[tid+1];
    __syncthreads();

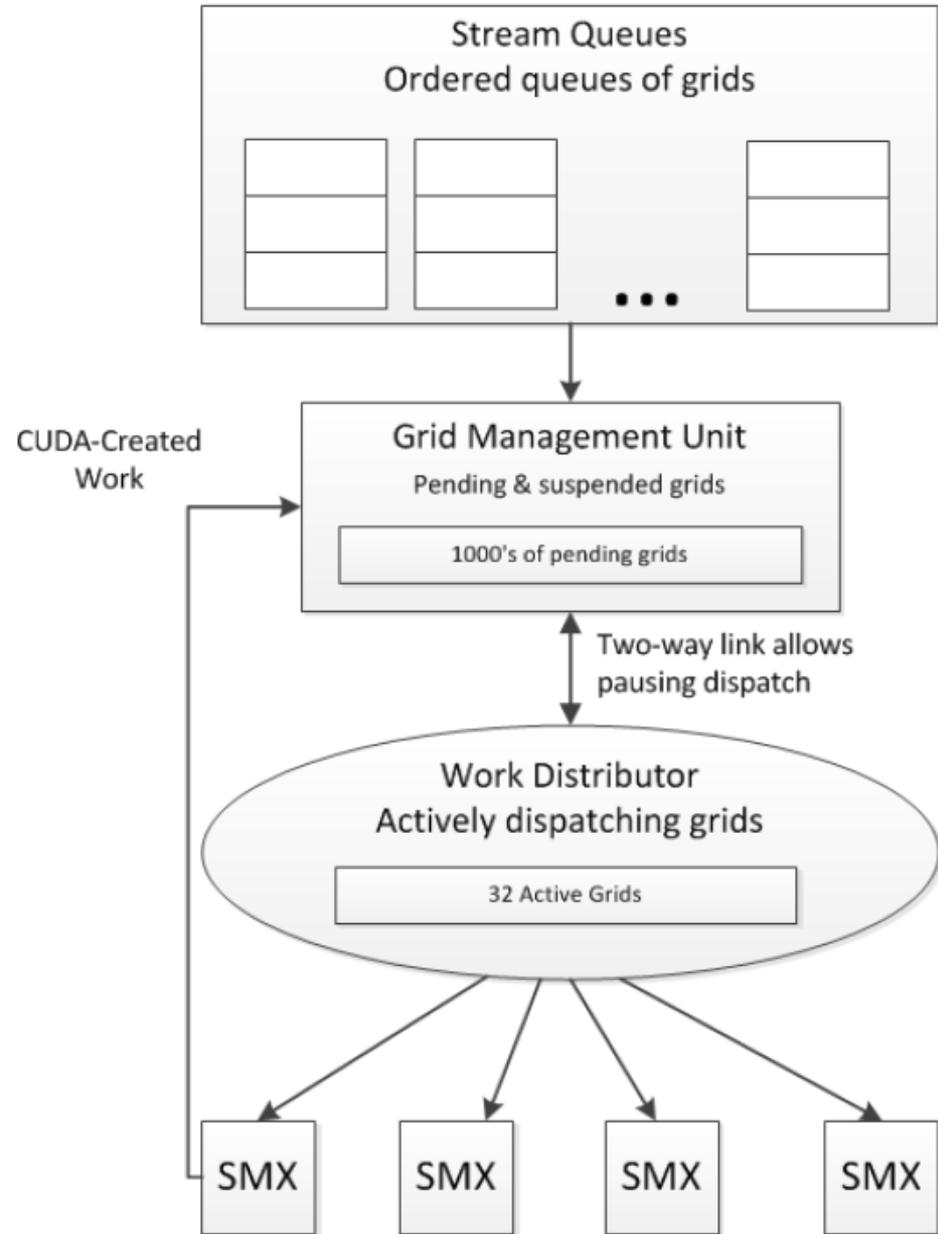
    if(tid == 0) {
        launch<<< 128, 256 >>>(buf);
        cudaDeviceSynchronize();
    }
    __syncthreads();

    cudaMemcpyAsync(data, buf, 1024);
    cudaDeviceSynchronize();
}
                
```

Fermi Workflow



Kepler Workflow



Development Support

39

Software development kits: NVIDIA and AMD; Windows and Linux

Special libraries: AMD Core Math Library, BLAS and FFT libraries by NVIDIA, OpenNL for numerics and CULA for linear algebra; NVIDIA Performance Primitives library: a sink for common GPU accelerated algorithms

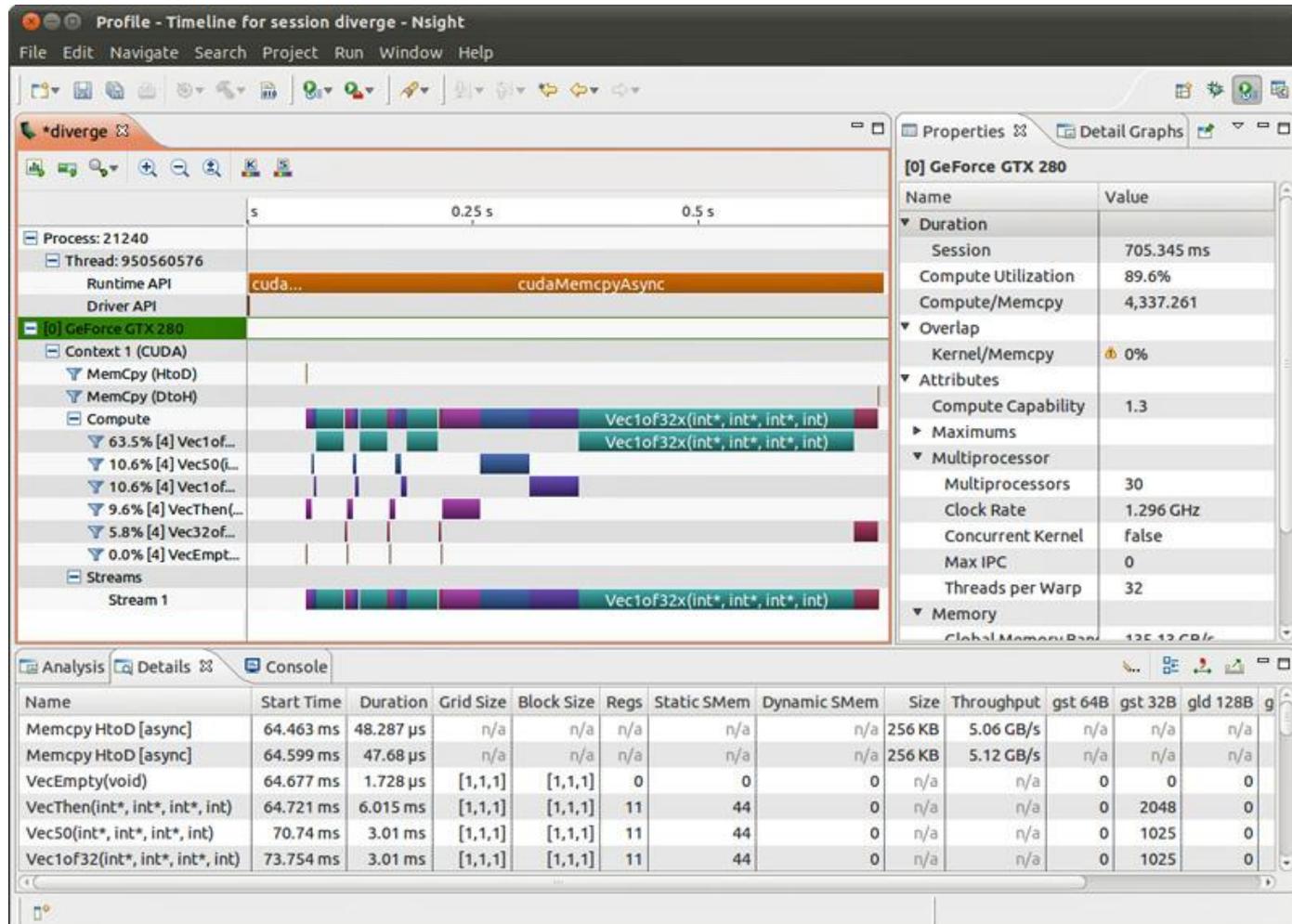
Profiling and debugging tools:

- NVIDIAs Parallel Nsight for Microsoft Visual Studio
- AMDs ATI Stream Profiler
- AMDs Stream KernelAnalyzer:
displays GPU assembler code, detects execution bottlenecks
- gDEDebugger (platform-independent)

Big knowledge bases with tutorials, examples, articles, show cases, and developer forums

Nsight

40



Towards new Platforms

41

WebCL [Draft] <http://www.khronos.org/webcl/>



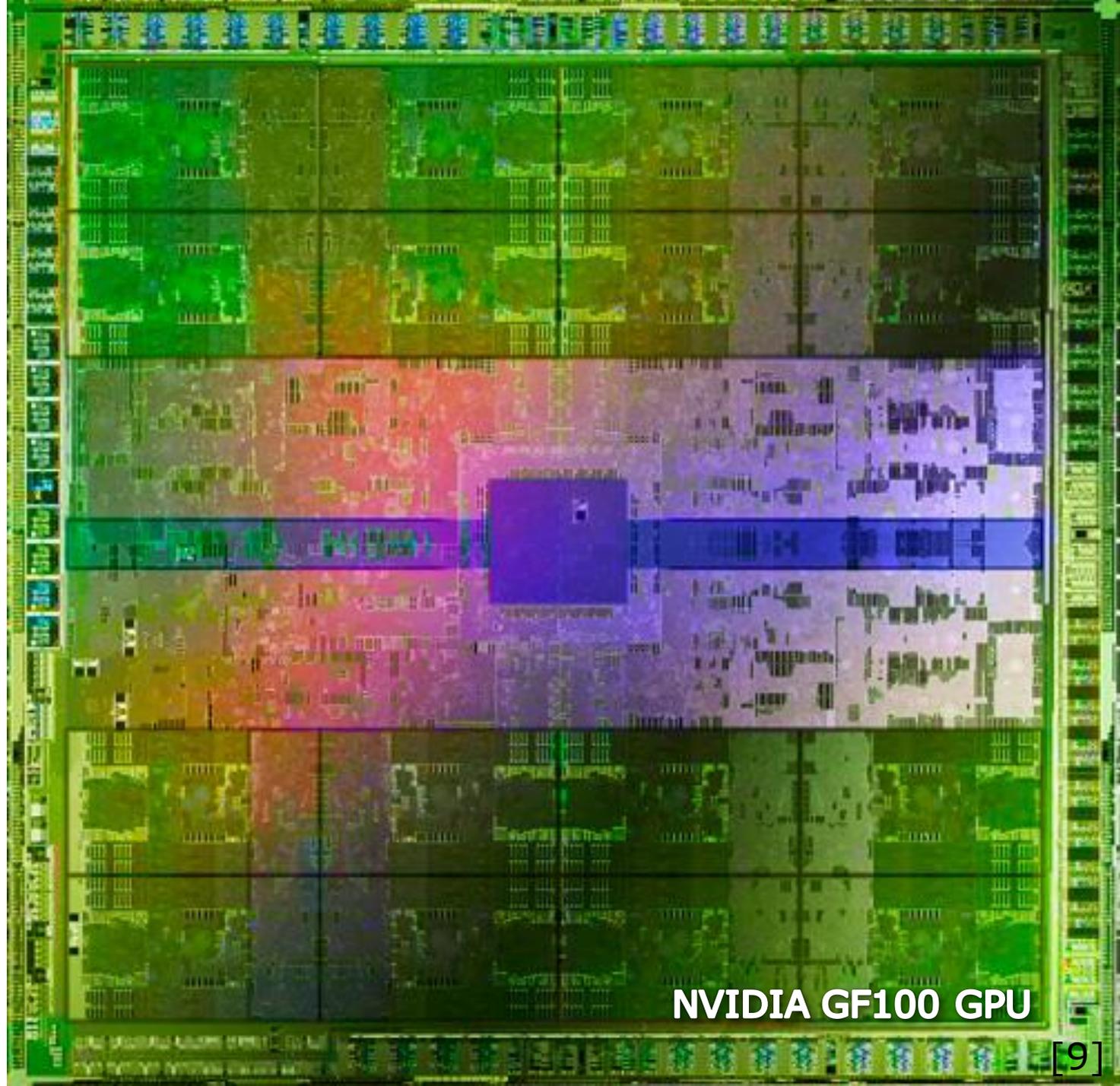
- JavaScript binding to OpenCL
- Heterogeneous Parallel Computing (CPUs + GPU) within Web Browsers
- Enables compute intense programs like physics engines, video editing...
- Currently only available with add-ons (Node.js, Firefox, WebKit)

Android installable client driver extension (ICD)

- Enables OpenCL implementations to be discovered and loaded as a shared object on Android systems.

PGI for multicore ARM processors

GPU Hardware in Detail



NVIDIA GF100 GPU

GF100

Host Interface

GigaThread Engine



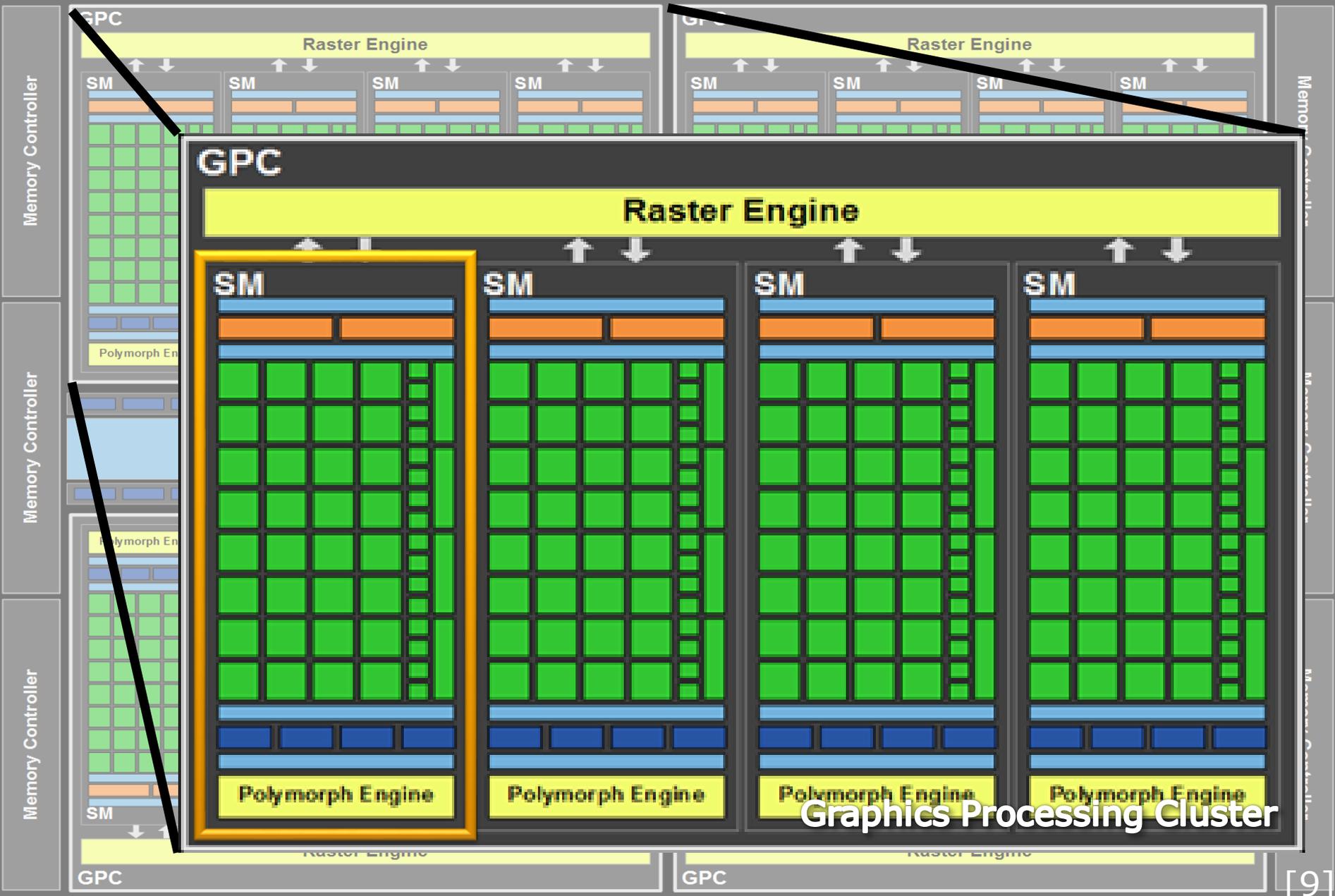
Graphics Processing Cluster

L2 Cache

GF100

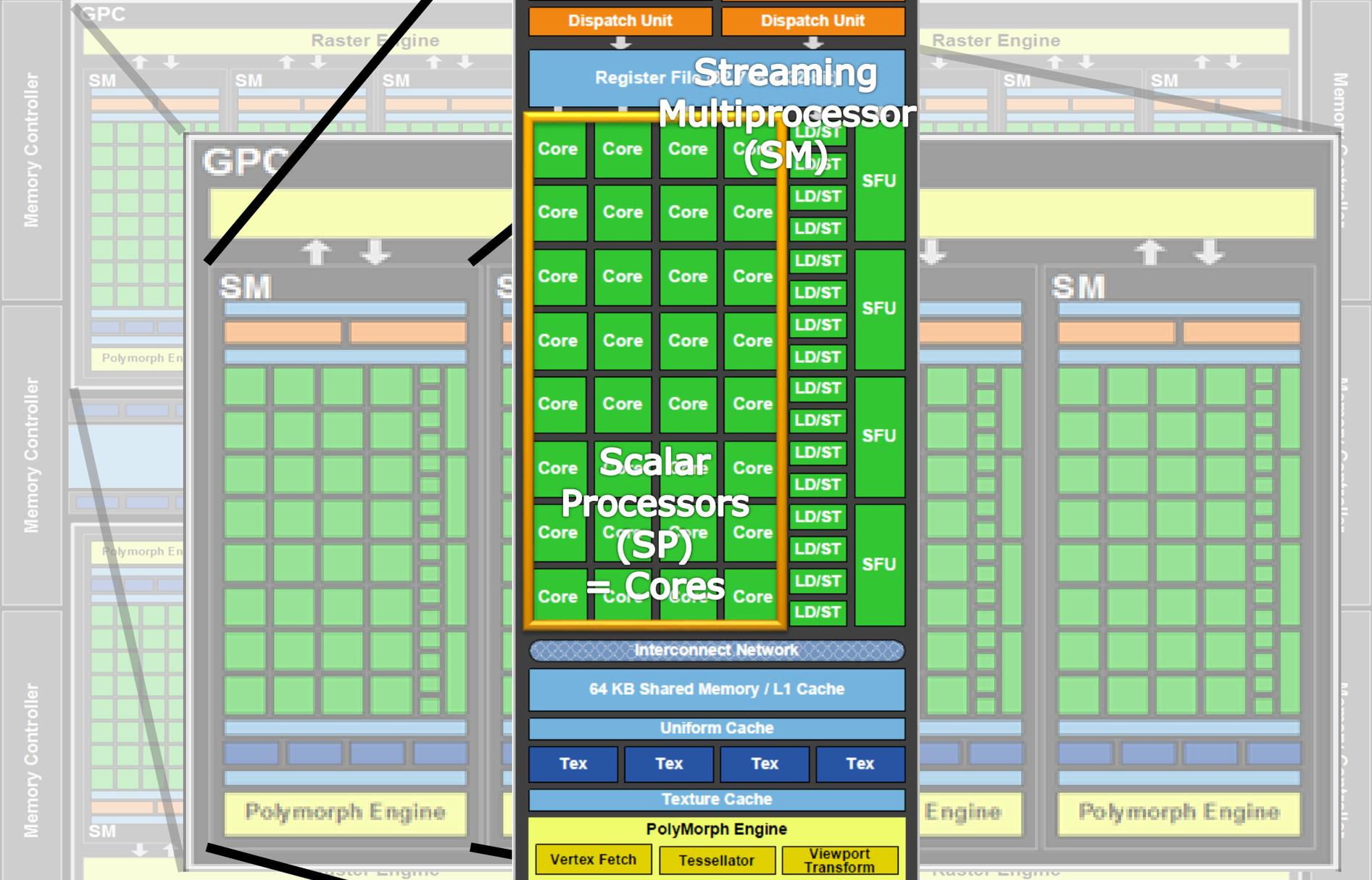
Host Interface

GigaThread Engine

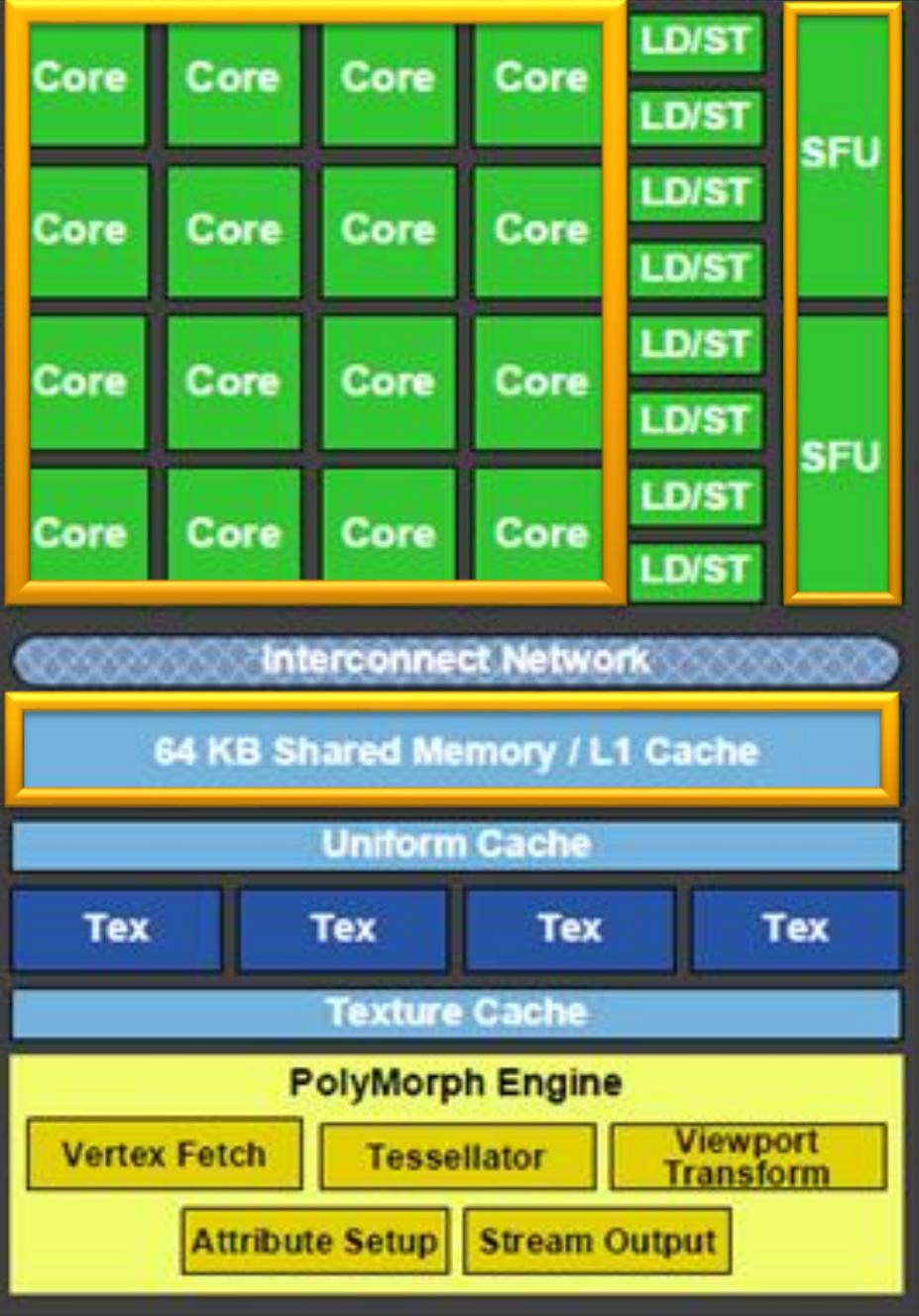
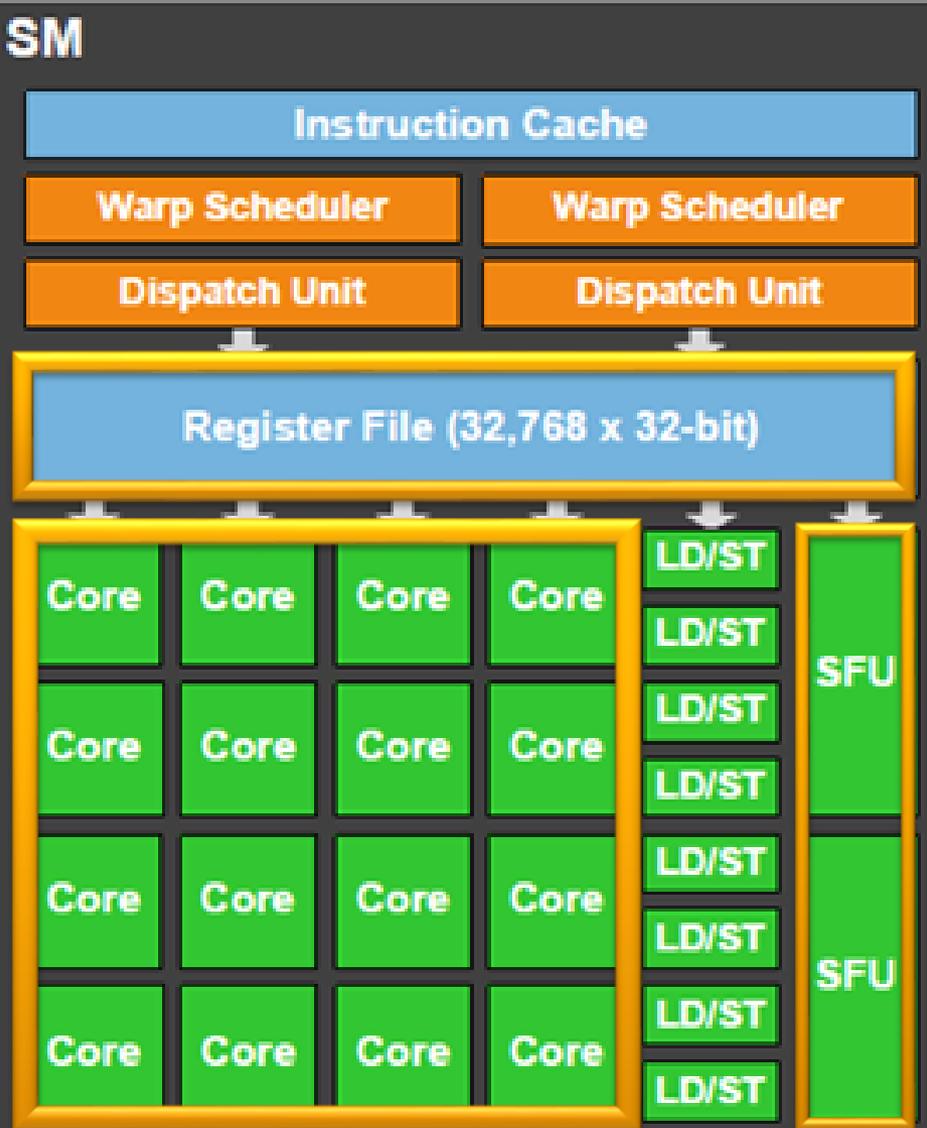


Graphics Processing Cluster

GF100



GF100



GT200 – early architecture

48

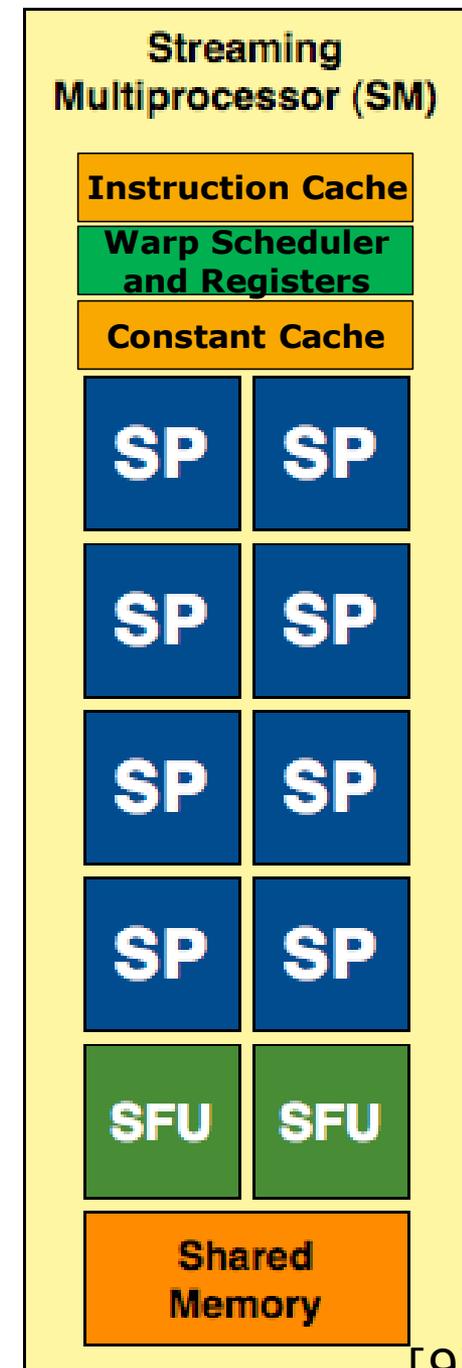
Simpler architecture, but same principles

Several Work Groups reside on one SM

- Amount depends on available resources (Shared Memory (=Local Memory in OpenCL), Registers)
- More Work Groups → better latency hiding
 - Latencies occur for memory accesses, pipelined floating-point arithmetic and branch instructions

Thread execution in “Warps” (called “wavefronts” on AMD)

- Native execution size (32 Threads for NVIDIA)
- Zero-Overhead Thread Scheduling: If one warp stalls (accesses memory) next warp is selected for execution



Warp Execution Example

49

Application creates 200.000 „Tasks“

→ Global Work Group Size: 200.000 Work Items

Programmer decides to use a Local Work Group Size of 100 Work Items

→ Number of Work Groups: 2.000 Work Groups

One Work Item requires 10 registers and 20 byte of Shared Memory; a SM has 16 KB of Shared Memory and 16.384 registers

→ Number of Work Items per SM:

$$\text{Max}(16.384/10, 16\text{KB}/20\text{B}) = 819 \text{ Work Items}$$

→ Number of Work Groups per SM:

$$819/100 = 8 \text{ Work Groups per SM}$$

Even if 7 Work Groups are waiting for memory, 1 can be executed.

Warp Execution Example

50

Each of the Work Groups contains 100 Work Items; the Warp Size (native execution size of a SM) is 32

- Number of Threads Executed in parallel: 32 Threads
- Number of „Rounds“ to execute a Work Group: $100/32 = 4$
- Threads running in the first 3 rounds: 32 Threads
- Threads running in the last round: $100 - 32 * 4 = 4$ Threads

If one of the threads accesses memory: whole warp stalls

If one of the threads follows a differing execution path: it is executed in an additional separate round

OpenCL Platforms

51

AMD APP SDK (supports OpenCL CPU and accelerated processing unit Devices)	X86 + SSE2 (or higher) compatible CPUs 64-bit & 32-bit ; ^[77] Linux 2.6 PC, Windows Vista/7 PC	AMD Fusion E-350, E-240, C-50, C-30 with HD 6310/HD 6250	AMD Radeon /Mobility HD 6800, HD 5x00 series GPU, iGPU HD 6310/HD 6250	ATI FirePro Vx800 series GPU
Intel SDK for OpenCL Applications 2013 ^[78] (supports Intel Core processors and Intel HD Graphics 4000/2500)	Intel CPUs with SSE 4.1, SSE 4.2 or AVX support. ^{[79][80]} Microsoft Windows , Linux	Intel Core i7, i5, i3 ; 2nd Generation Intel Core i7/5/3, 3rd Generation Intel Core Processors with Intel HD Graphics 4000/2500	Intel Core 2 Solo, Duo Quad, Extreme	Intel Xeon 7x00, 5x00, 3x00 (Core based)
IBM Servers with OpenCL Development Kit for Linux on Power running on Power VSX ^{[81][82]}	IBM Power 755 (PERCS), 750	IBM BladeCenter PS70x Express	IBM BladeCenter JS2x, JS43	IBM BladeCenter QS22
IBM OpenCL Common Runtime (OCR) ^[83]	X86 + SSE2 (or higher) compatible CPUs 64-bit & 32-bit; ^[84] Linux 2.6 PC	AMD Fusion , Nvidia Ion and Intel Core i7, i5, i3; 2nd Generation Intel Core i7/5/3	AMD Radeon, Nvidia GeForce and Intel Core 2 Solo, Duo, Quad, Extreme	ATI FirePro, Nvidia Quadro and Intel Xeon 7x00, 5x00, 3x00 (Core based)
Nvidia OpenCL Driver and Tools ^[85]	Nvidia Tesla C/D/S	Nvidia GeForce GTS/GT/GTX	Nvidia Ion	Nvidia Quadro FX/NVX/Plex

Compute Capability by version

52

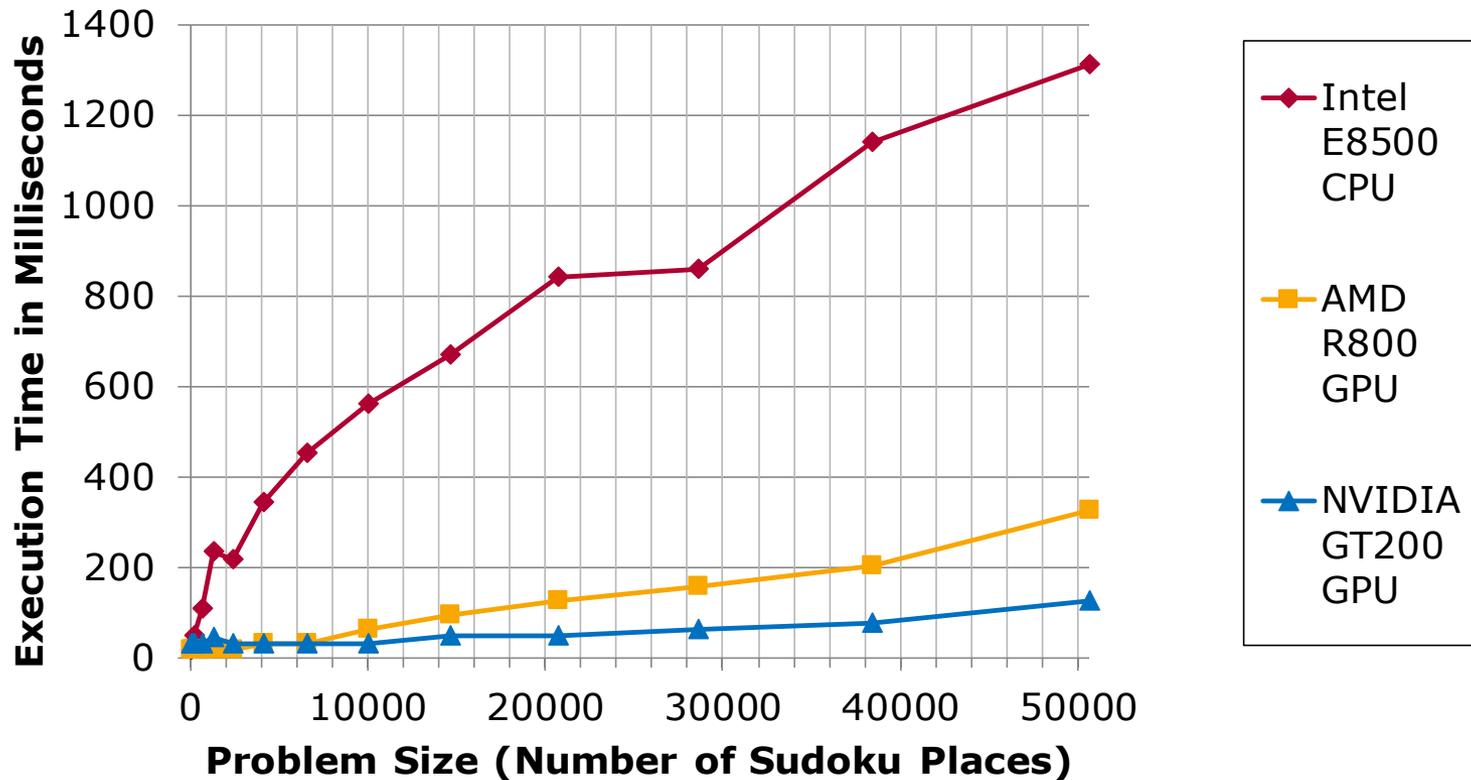
	1.0	1.1	1.2	1.3	2.x	3.0	Since 3.5
double precision floating point operations	No			Yes			
caches	No			Yes			
max # concurrent kernels	1			8		Dynamic Parallelism	
max # threads per block	512			1024			
max # Warps per MP	24	32		48	64		
max # Threads per MP	768	1024		1536	2048		
register count (32 bit)	8192	16384		32768	65536		
max shared mem per MP	16KB			16/48KB	48-112KB		
# shared memory banks	16			32			

Plus: varying amounts of cores, global memory sizes, bandwidth, clock speeds (core, memory), bus width, memory access penalties ...

The Power of GPU Computing

53

big performance gains for small problem sizes

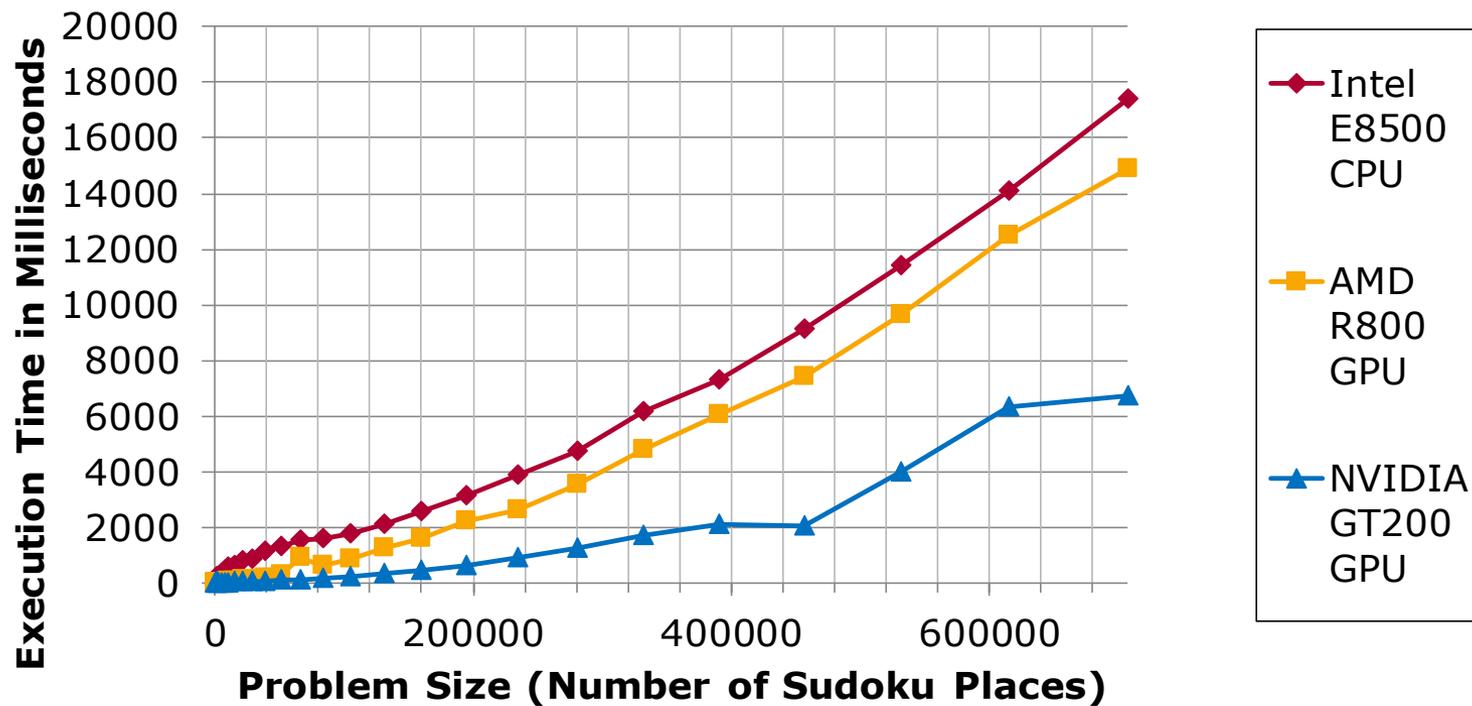


* less is better

The Power of GPU Computing

54

small/moderate performance gains for large problem sizes
 → further optimizations needed



* less is better

Best Practices for Performance Tuning

55

Algorithm Design

- Asynchronous, Recompute, Simple

Memory Transfer

- Chaining, Overlap Transfer & Compute

Control Flow

- Divergent Branching, Predication

Memory Types

- Local Memory as Cache, rare resource

Memory Access

- Coalescing, Bank Conflicts

Sizing

- Execution Size, Evaluation

Instructions

- Shifting, Fused Multiply, Vector Types

Precision

- Native Math Functions, Build Options

Divergent Branching and Predication

56

Divergent Branching

- Flow control instruction (`if`, `switch`, `do`, `for`, `while`) can result in different execution paths
- Data parallel execution → varying execution paths will be serialized
- Threads converge back to same execution path after completion

Branch Predication

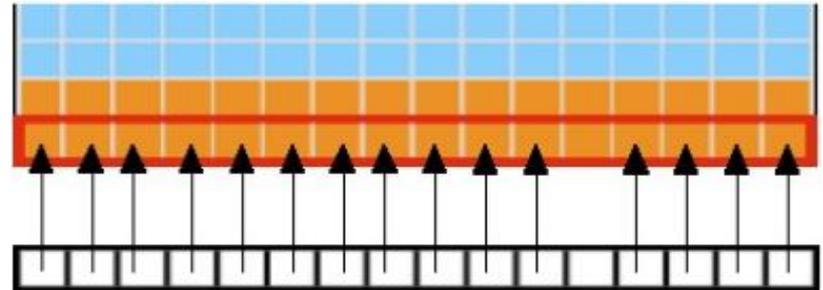
- Instructions are associated with a per-thread condition code (predicate)
 - All instructions are scheduled for execution
 - Predicate true: executed normally
 - Predicate false: do not write results, do not evaluate addresses, do not read operands
- Compiler may use branch predication for `if` or `switch` statements
- Unroll loops yourself (or use `#pragma unroll` for NVIDIA)

Coalesced Memory Accesses

57

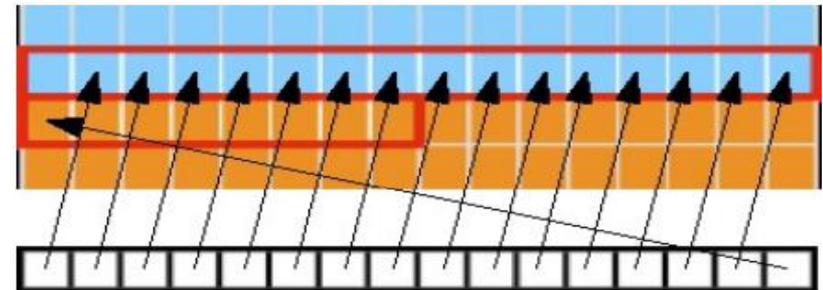
Simple Access Pattern

- Can be fetched in a single 64-byte transaction (red rectangle)
- Could also be permuted *



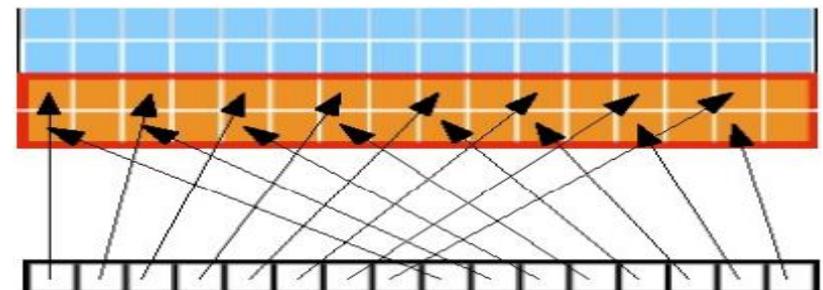
Sequential but Misaligned Access

- Fall into single 128-byte segment: single 128-byte transaction, else: 64-byte transaction + 32-byte transaction *



Strided Accesses

- Depending on stride from 1 (here) up to 16 transactions *



* 16 transactions with compute capability 1.1

Use Caching: Local, Texture, Constant

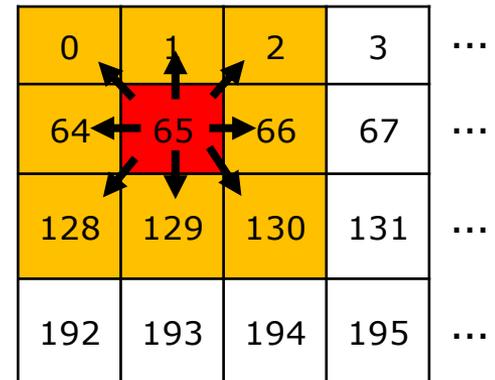
58

Local Memory

- Memory latency roughly 100x lower than global memory latency
- Small, no coalescing problems, prone to memory bank conflicts

Texture Memory

- 2-dimensionally cached, read-only
- Can be used to avoid uncoalesced loads from global memory
- Used with the image data type



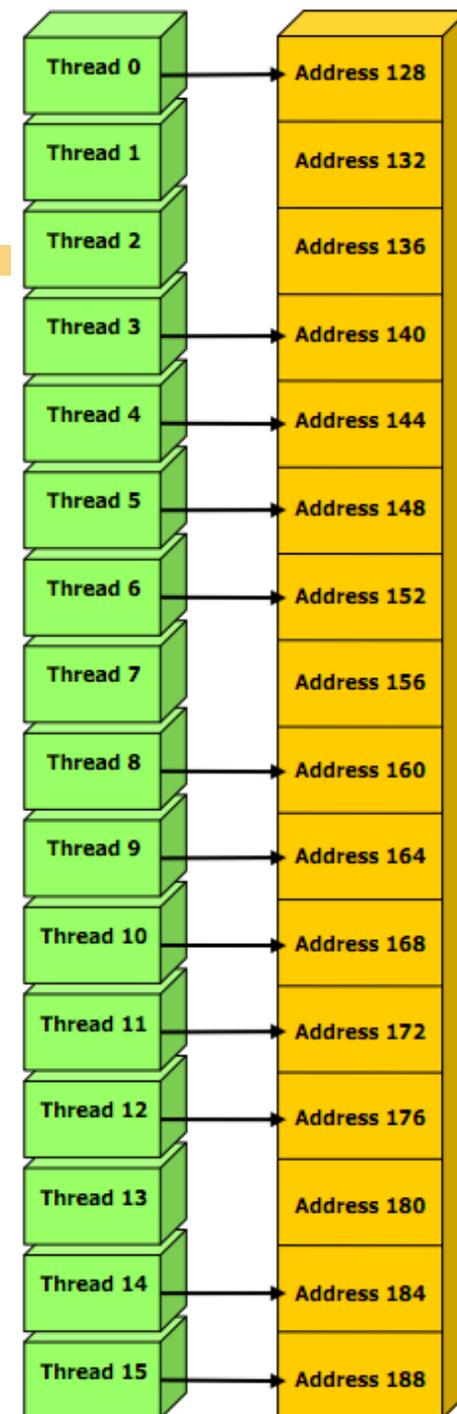
Constant Memory

- Lineary cached, read-only, 64 KB
- as fast as reading from a register for the same address
- Can be used for big lists of input arguments

Memory Bank Conflicts

59

- Access to (Shared) Memory is implemented via hardware memory banks
- If a thread accesses a memory address this is handled by the responsible memory bank
- Simple Access Patterns like this one are fetched in a single transaction



Memory Bank Conflicts

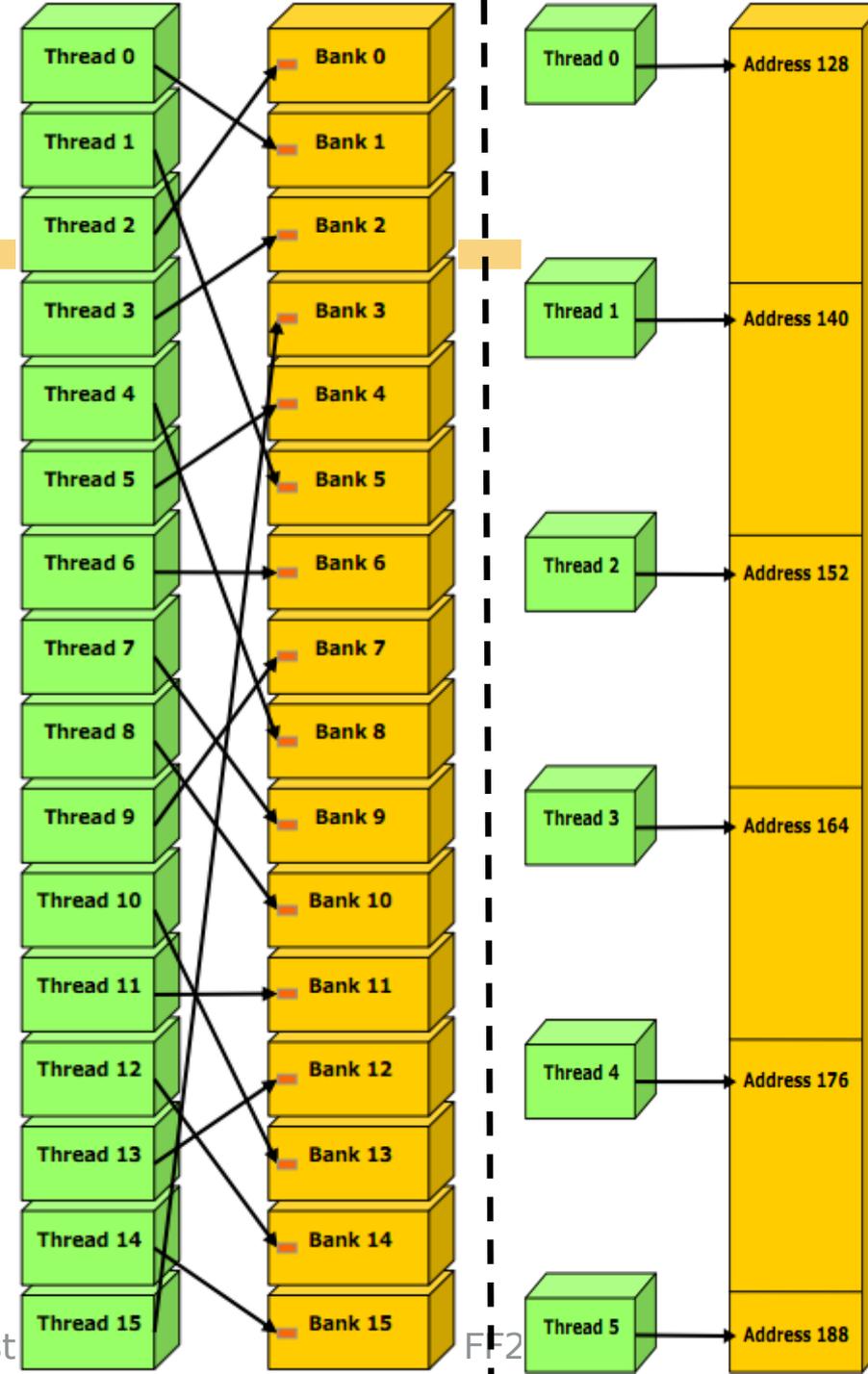
60

Permuted Memory Access (left)

- Still one transaction on cards with compute capability ≥ 1.2 ; otherwise 16 transactions are required

Strided Memory Access (right)

- Still one transaction on cards with compute capability ≥ 1.2 ; otherwise 16 transactions are required

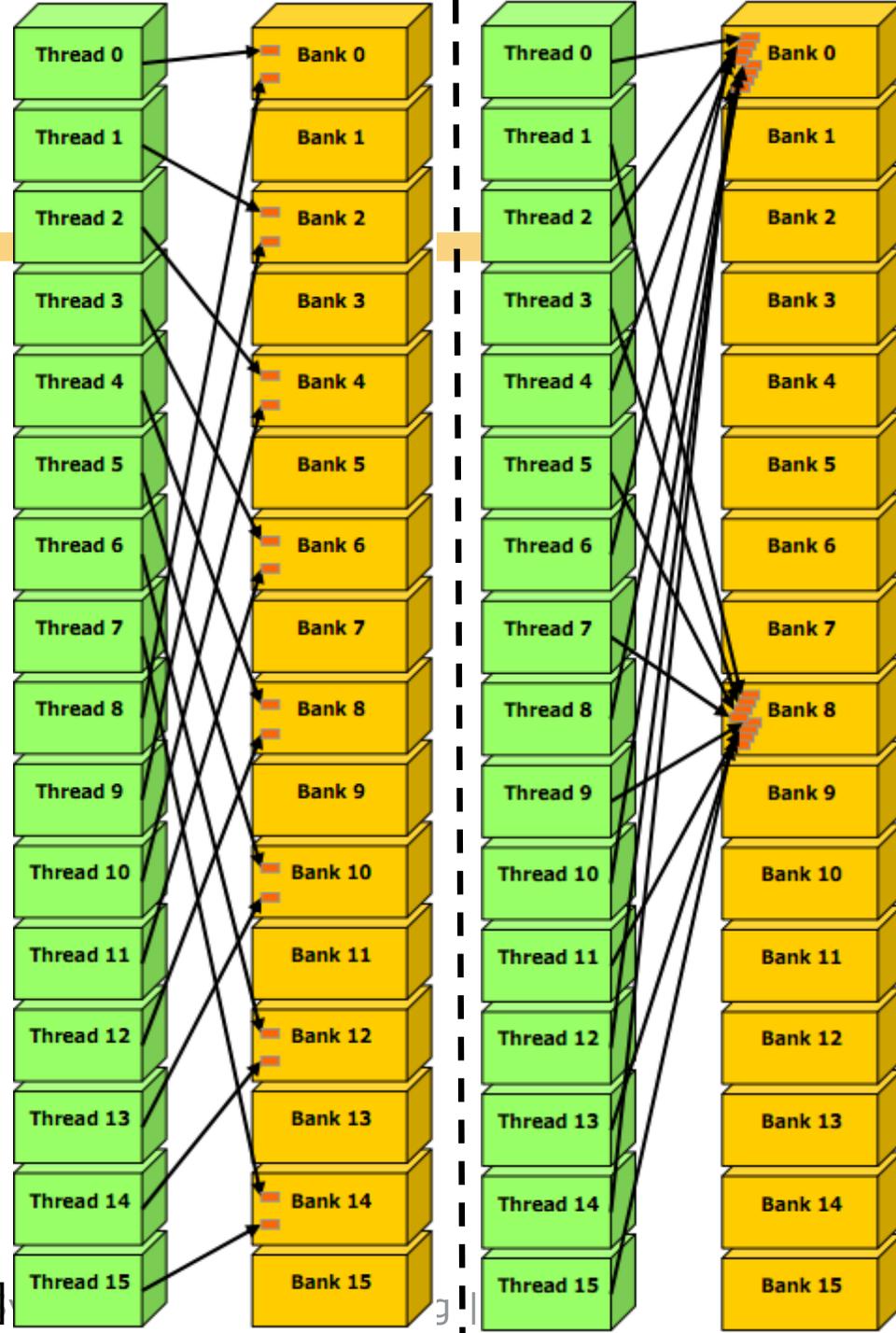


Memory Bank Conflicts

61

Bank conflicts

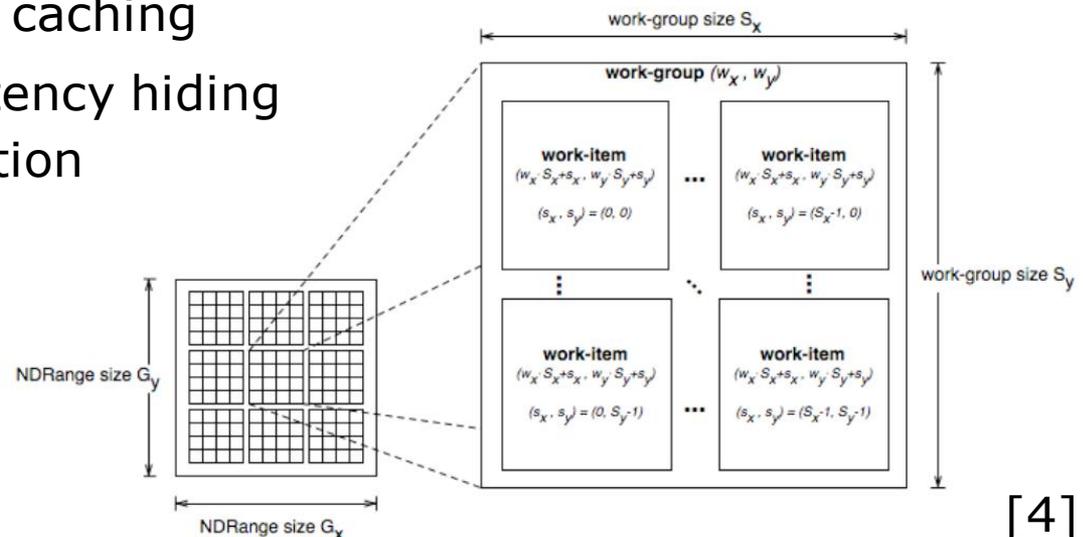
- Left figure: 2 bank conflicts → resulting bandwidth is $\frac{1}{2}$ of the original bandwidth
- Right figure: 8 bank conflicts → resulting bandwidth is $\frac{1}{8}$ of the original bandwidth



Sizing: What is the right execution layout?

62

- Local work item count should be a multiple of native execution size (NVIDIA 32, AMD 64), but not to big
- Number of work groups should be multiple of the number of multiprocessors (hundreds or thousands of work groups)
- Can be configured in 1-, 2- or 3-dimensional layout: consider access patterns and caching
- Balance between latency hiding and resource utilization
- Experimenting is required!



[4]

Instructions and Precision

63

- Single precision floats provide best performance
- Use shift operations to avoid expensive division and modulo calculations
- Special compiler flags
- AMD has native vector type implementation; NVIDIA is scalar
- Use the native math library whenever speed trumps precision

Functions	Throughput
single-precision floating-point add, multiply, and multiply-add	8 operations per clock cycle
single-precision reciprocal, reciprocal square root, and native_logf(x)	2 operations per clock cycle
native_sin, native_cos, native_exp	1 operation per clock cycle

Further Readings

64

<http://www.dcl.hpi.uni-potsdam.de/research/gpureadings/>

- [1] Kirk, D. B. & Hwu, W. W. ***Programming Massively Parallel Processors: A Hands-on Approach***. 1 ed. Morgan Kaufmann.
- [2] Herlihy, M. & Shavit, N. *The Art of Multiprocessor Programming*.
- [3] Sanders, J. & Kandrot, E. ***CUDA by Example: An Introduction to General-Purpose GPU Programming*** . 1 ed. Addison-Wesley Professional.
- [4] Munshi, A. (ed.). The OpenCL Specification - v1.1. The Khronos Group Inc.
- [5] **Mattson, T. The Future of Many Core Computing: Software for many core processors.**
- [6] NVIDIA. NVIDIA OpenCL Best Practices Guide
- [7] Rob Farber. CUDA, Supercomputing for the Masses. Dr. Dobb's
- [8] NVIDIA. OpenCL Programming for the CUDA Architecture
- [9] Ryan Smith, NVIDIA's GeForce GTX 480 and GTX 470: 6 Months Late, Was It Worth the Wait?
- [10] Stephen Jones. Introduction to Dynamic Parallelism