General-purpose programming on GPU Exploiting multi-GPU systems

> Eugenio Rustico rustico@dmi.unict.it

D.M.I. - Università di Catania

Updated: May 20, 2011







- 一一司

Università di Catania

Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

Introduction	CUDA contexts		CUDA 4
Overview			

1 Introduction

2 CUDA contexts

3 Domain subdivision

4 Problems

5 CUDA 4

<ロト < 聞 > < 置 > < 置 > 一 置 の < @</p>

Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU Università di Catania

Introduction	CUDA contexts		CUDA 4
		4	► 5000

Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU Università di Catania

Introduction	CUDA contexts		CUDA 4

With less than 1.000, it is possible to setup a 1 TERAFLOPS workstation...



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

Introduction	CUDA contexts		CUDA 4

With less than 1.000, it is possible to setup a 1 TERAFLOPS workstation...

... is it really possible to achieve such (theoretical) peak power?



Introduction	CUDA contexts	Domain subdivision	Problems	CUDA 4



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

Introduction	CODA CONTEXES		C0D/4 4

To exploit a multi-GPU system we have to think of the problem at **two** levels of parallelism:



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

introduction	CODA contexts	Domain subdivision	CODA 4

To exploit a multi-GPU system we have to think of the problem at **two** levels of parallelism:

• A grid of parallel sequences of operations (threads)

introduction	CODA contexts	Domain subdivision	CODA 4

To exploit a multi-GPU system we have to think of the problem at **two** levels of parallelism:

- A grid of parallel sequences of operations (threads)
- A set of such grids

introduction	CODA CONTEXES		C0DA 4

To exploit a multi-GPU system we have to think of the problem at **two** levels of parallelism:

- A grid of parallel sequences of operations (threads)
- A set of such grids

It is absolutely not trivial unless different threads are *completely* independent each other

	CUDA contexts			CUDA 4
Within the f created.	irst call to a CUDA	runtime routine a <i>cont</i>	<i>ext</i> is implicitly	



Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

Università di Catania

		CUDA contexts	Domain subdivision		CUDA 4
Within the first call to a CUDA runtime routing a contact is implicitly					

Within the first call to a CUDA runtime routine a *context* is implicitly created.

A CUDA context is a sort of *session* of CUDA operations associated to a thread. One context is associated to one GPU (n-to-1), and can operate only with the associated GPU.

Introduction	CUDA contexts	Domain subdivision	Problems	CUDA 4

Within the first call to a CUDA runtime routine a *context* is implicitly created.

A CUDA context is a sort of *session* of CUDA operations associated to a thread. One context is associated to one GPU (n-to-1), and can operate only with the associated GPU.

One GPU may have several contexts associated; a CPU thread may have only one active context, unless the low-level API is used to switch between contexts.

Introduction	CUDA contexts	Domain subdivision	Problems	CUDA 4

Within the first call to a CUDA runtime routine a *context* is implicitly created.

A CUDA context is a sort of *session* of CUDA operations associated to a thread. One context is associated to one GPU (n-to-1), and can operate only with the associated GPU.

One GPU may have several contexts associated; a CPU thread may have only one active context, unless the low-level API is used to switch between contexts.

We can see a context as a hidden set of settings and structures that are used by the CUDA runtime, when we call a CUDA function like a cudaMemCpy(), to answer the question: which device are we working on? With which settings?

CUDA contexts		CUDA 4



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

	CUDA contexts			CUDA 4
To be al	ole to use multiple Gl	PUs simultaneously, w	ve need one conte	xt per
GPU.	•	<b>3</b> .		·

With low-level API, we can explicitly switch between context (cuCtx...() methods)



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU Università di Catania

CUDA contexts		CUDA 4

- With low-level API, we can explicitly switch between context (cuCtx...() methods)
- With high-level API, we need one separate thread per GPU (multithreaded applications)

CUDA contexts		CUDA 4

- With low-level API, we can explicitly switch between context (cuCtx...() methods)
- With high-level API, we need one separate thread per GPU (multithreaded applications)

High-level APIs are way easier to use, but designing a robust multi-threaded application may be challenging. In general, we will need a thread-safe synchronization mechanism (e.g. signals, semaphores, barriers, etc.).

- With low-level API, we can explicitly switch between context (cuCtx...() methods)
- With high-level API, we need one separate thread per GPU (multithreaded applications)

High-level APIs are way easier to use, but designing a robust multi-threaded application may be challenging. In general, we will need a thread-safe synchronization mechanism (e.g. signals, semaphores, barriers, etc.).

It is a recommended practice to call cudaThreadExit() at the and of every host thread to explicitly clean up the context.

CUDA contexts	Domain subdivision	CUDA 4

There are two general approaches:



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

CUDA contexts	Domain subdivision	CUDA 4

There are two general approaches:

Divide in the data domain (split input, assign to GPUs, manage overlapping)

	CUDA contexts	Domain subdivision	CUDA 4

There are two general approaches:

- Divide in the data domain (split input, assign to GPUs, manage overlapping)
- 2 Divide in the computation domain (pass all input data across all GPUs in pipeline)

	CUDA contexts	Domain subdivision	CUDA 4

There are two general approaches:

- Divide in the data domain (split input, assign to GPUs, manage overlapping)
- 2 Divide in the computation domain (pass all input data across all GPUs in pipeline)

Approach n.1 is in most cases the simplest and most efficient one.

CUDA contexts	Domain subdivision	CUDA 4

How should we divide an image?



Università di Catania

Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

CUDA contexts	Domain subdivision	CUDA 4

### How should we divide an image?

GPU 0	GPULO	CPU 1
GPU 1	GFU U	GFUI
GPU 2		CPU 3
GPU 3	GPU Z	GrU 5

Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

- 🔹 🚍

Image: Image:

CUDA contexts	Domain subdivision	CUDA 4

## How should we divide an image?

GPU 0	GPU 0	GPU 1
GPU 1	GFU U	GPUI
GPU 2		
GPU 3	GPU Z	GPU 5

## Why first is preferable? When it is not?



Second has a better perimeter/area ratio (i.e. divergences and special cases), however: sparse border accesses (is memory per row?), two neighbors for every GPU (imagine 2x3!)

Image: Image:

CUDA contexts	Domain subdivision	CUDA 4

One GPU may be more powerful than another, or may have to deal also with visualization (e.g. an attached monitor), or may be in a PCI slot slower than other, or may simply have more operations or less coalesced memory requirements.



	CUDA contexts	Domain subdivision	CUDA 4

One GPU may be more powerful than another, or may have to deal also with visualization (e.g. an attached monitor), or may be in a PCI slot slower than other, or may simply have more operations or less coalesced memory requirements.

Dynamic *load balancing* may become fundamental to obtain a real speedup. Can you imagine a generic load balancing technique suitable for different applications?

Introduction	CUDA contexts	Domain subdivision	Problems	CUDA 4

One GPU may be more powerful than another, or may have to deal also with visualization (e.g. an attached monitor), or may be in a PCI slot slower than other, or may simply have more operations or less coalesced memory requirements.

Dynamic *load balancing* may become fundamental to obtain a real speedup. Can you imagine a generic load balancing technique suitable for different applications?

Answer: the only general technique is timing *a posteriori*. And it is very complex: requires some basic signal processing to get rid of oscillations and avoid local minima.

CUDA contexts	Problems	CUDA 4



Università di Catania

Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

Most applications present a form of data **interdependence**. Example: if we need to apply in sequence two convolution effects to an image, kernel nr. 2 requires the updated borders of kernel nr. 1 from other GPUs. Imagine a continuous timeline (streaming)...

Most applications present a form of data **interdependence**. Example: if we need to apply in sequence two convolution effects to an image, kernel nr. 2 requires the updated borders of kernel nr. 1 from other GPUs. Imagine a continuous timeline (streaming)...

Reality is even worse than this: **all transfers must pass through the host**, no direct GPU-GPU transfers up to CUDA 3.2.

Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

Most applications present a form of data **interdependence**. Example: if we need to apply in sequence two convolution effects to an image, kernel nr. 2 requires the updated borders of kernel nr. 1 from other GPUs. Imagine a continuous timeline (streaming)...

Reality is even worse than this: **all transfers must pass through the host**, no direct GPU-GPU transfers up to CUDA 3.2.

The higher the number of GPUs, the more memory transfers will request exclusive use of PCI bus. The more memory transfers, the less benefits of using multiple GPUs. Need **asynchronous operations** to cover latencies (next lecture).

CUDA contexts	Problems	CUDA 4

An example of expected vs. real speedup:



# of active cells



Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

CUDA contexts	Problems	CUDA 4

An example of expected vs. real speedup:



Kernel launch latency and saturation threshold cause

time(data/2) > time(data)/2

No matter how I split the input data: with K GPUs, in most cases, speedup will be < K.

CUDA contexts	Problems	CUDA 4

Examples:

- MAGFLOW (image-like)
- SPH (list-like)



Università di Catania

Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU



Università di Catania

Eugenio Rustico rustico@dmi.unict.it

General-purpose programming on GPU

 Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.



Eugenio Rustico rustico@dmi.unict.it General-purpose programming on GPU

- Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.
- Implicit context switch: just repeat cudaSetDevice() in high level API to switch device

- Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.
- Implicit context switch: just repeat cudaSetDevice() in high level API to switch device
- Device to device direct memory transfers, cudaMemCpy() bursts and direct access

- Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.
- Implicit context switch: just repeat cudaSetDevice() in high level API to switch device
- Device to device direct memory transfers, cudaMemCpy() bursts and direct access
- New libraries (NPP) and middlewares (Thrust)

- Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.
- Implicit context switch: just repeat cudaSetDevice() in high level API to switch device
- Device to device direct memory transfers, cudaMemCpy() bursts and direct access
- New libraries (NPP) and middlewares (Thrust)
- Dynamic memory allocation and virtual functions in device code

- Unified Virtual Addressing: on x86\_64 systems, a unique address space for CPU and GPUs.
- Implicit context switch: just repeat cudaSetDevice() in high level API to switch device
- Device to device direct memory transfers, cudaMemCpy() bursts and direct access
- New libraries (NPP) and middlewares (Thrust)
- Dynamic memory allocation and virtual functions in device code

New method to retrieve memory type by address and simpler cudaMemcpy() (cudaMemcpyDefault(): automatically determine transfer direction)

# Before NVIDIA GPUDirect<sup>™</sup> v2.0

## **Required Copy into Main Memory**



Image: Image:

# NVIDIA GPUDirect<sup>™</sup> v2.0:

## **Peer-to-Peer Communication**

## Direct Transfers between GPUs

