# CUDA Programming Model

Xing Zeng, Dongyue Mou

---

## Outline

- Introduction
- Motivation
- Programming Model
- Memory Model
- CUDA API
- Example
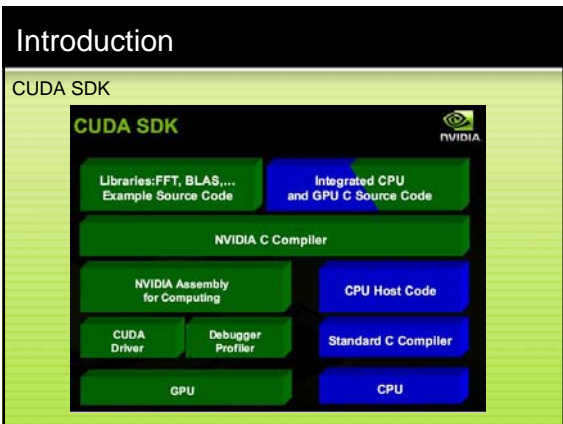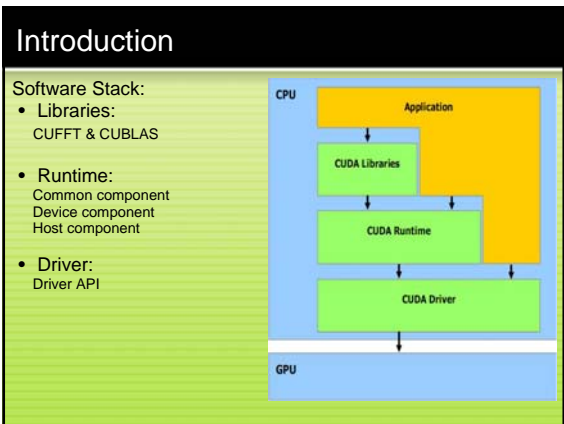- Pro & Contra
- Trend

---

## Outline

- **Introduction**
- Motivation
- Programming Model
- Memory Model
- CUDA API
- Example
- Pro & Contra
- Trend

---

## Introduction

What is CUDA?
- **Compute Unified Device Architecture.**
- A powerful parallel programming model for issuing and managing computations on the GPU without mapping them to a graphics API.

- Heterogenous - mixed serial-parallel programming
- Scalable - hierarchical thread execution model
- Accessible - minimal but expressive changes to C

---

## Introduction

Software Stack:
- Libraries:
  CUFFT & CUBLAS

- Runtime:
  Common component
  Device component
  Host component

- Driver:
  Driver API



---

## Introduction

CUDA SDK

## Outline

- Introduction
- **Motivation**
- Programming Model
- Memory Model
- CUDA API
- Example
- Pro & Contra
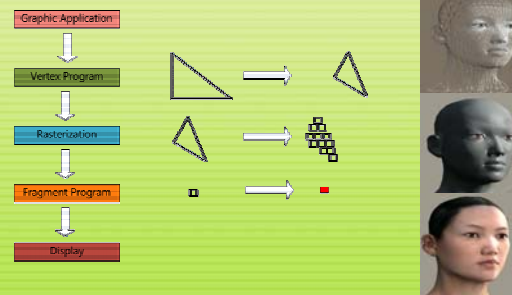- Trend

## Motivation

GPU Programming Model
GPGPU Programming Model
CUDA Programming Model

## Motivation

GPU Programming Model
GPGPU Programming Model
CUDA Programming Model

## Motivation

GPU Programming Model for Graphics



Graphic Application
Vertex Program
Rasterization
Fragment Program
Display

## Motivation

GPU Programming Model
GPGPU Programming Model
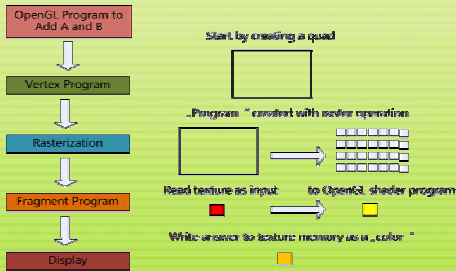CUDA Programming Model

## Motivation

GPGPU Programming Model
Trick the GPU into general-purpose
computing by casting problem as graphics

- Turn data into images ("texture maps")
- Turn algorithms into image synthesis ("rending passes")

Drawback:
- Tough learning curve
- potentially high overhead of graphics API
- highly constrained memory layout & access model
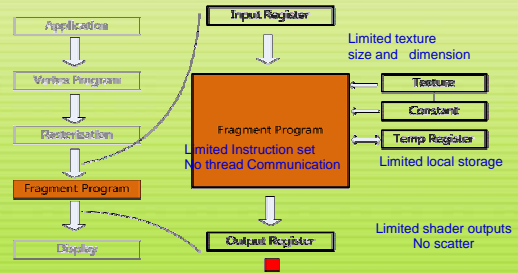- Need for many passes drives up bandwidth consumption

## Motivation

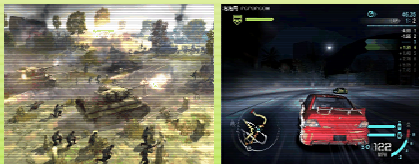GPGPU Programming to do A + B



## Motivation

What's wrong with GPGPU 1

APIs are specific to Graphics



Limited texture size and dimension

Limited Instruction set
No thread Communication

Limited local storage

Limited shader outputs
No scatter

## Motivation

What's wrong with GPGPU 2



## Motivation

GPU Programming Model
GPGPU Programming Model
CUDA Programming Model

## Outline

- Introduction
- Motivation
- **Programming Model**
- Memory Model
- CUDA API
- Example
- Pro & Contra
- Trend

## Programming Model
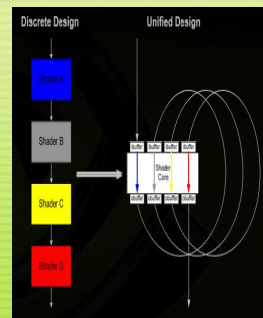
CUDA: Unified Design
Advantage:

HW: fully generally data-parallel arch-tecture.
- General thread launch
- Global load-store
- Parallel data cache
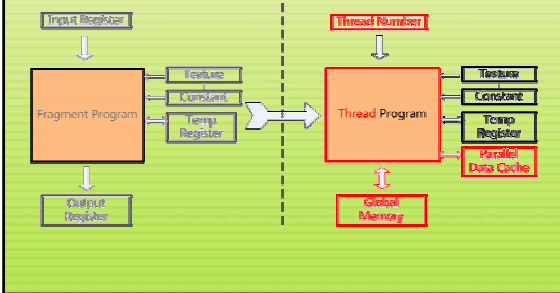- Scalar architecture
- Integers, bit operation

SW: program the GPU in C
- Scalable data parallel execuation/ memory model
- C with minimal yet powerful extensions

## Motivation

**From GPGPU to CUDA Programming Model**
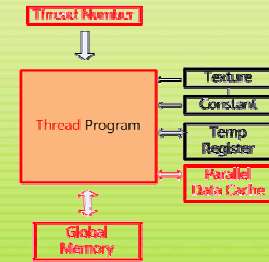


---

## Programming Model

**Feature 1:**
- Thread not pixel
- Full Integer and Bit Instructions
- No limits on branching, looping
- 1D, 2D, 3D threadID allocation

**Feature 2:**
- Fully general load/store to GPU memory
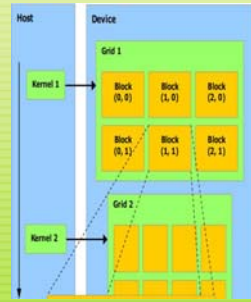- Untyped, not fixed texture types
- Pointer support

**Feature 3:**
- Dedicated on-chip memory
- Shared between threads for inter-threads communication
- Explicitly managed
- As fast as registers



---

## Programming Model

**Important Concepts:**
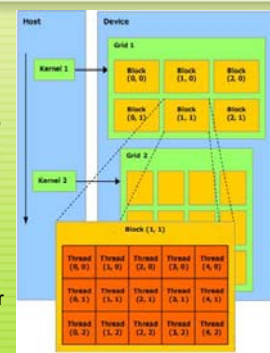- Device: GPU, viewed as a co-processor.
- Host: CPU
- Kernel: data-parallel, computed-intensive positions of application running on the device.



---

## Programming Model

**Important Concepts:**
- Thread: basic execution unit
- Thread block:
  A batch of thread. Threads in a block cooperate together, efficiently share data.
  Thread/block have unique id
- Grid:
  A batch of thread block. that excuate same kernel. Threads in different block in the same grid cannot directly communicate with each other



---

## Programming Model

Simple example ( Matrx addition ):

cpu c program:               cuda program:

```
void addVector (float *a, float *b,          __global__ void addVector (float *a, float *b,
              float *c, int N)                              float *c)
{                                            {
   int i, index;                                int i = threadIdx.x + blockDim.x*blockIdx.x;
   for (i = 0; i<N; i++) {                       c[i] = a[i] + b[i];
      c[index] = a[index] + b[index];         }
   }
}

void main()                                  void main()
{                                            {
   ....                                         // allocation & transfer data to GPU
   addVector(a, b, x, N);                        //Curuate on N/256 blocks, of 256 threads each
   ....                                         addVector <<< N/256, 256 >> ( d_A, d_B, d_C);
}                                               ....
                                             }
```

---

## Programming Model
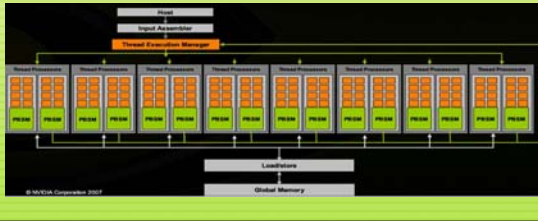
Hardware implementation:

A set of SIMD Multiprocessors with On-Chip shared memory
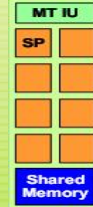
## Programming Model

G80 Example:
- 16 Multiprocessors, 128 Thread Processors
- Up to 12,288 parallel threads active
- Per-block shared memory accelerates processing.



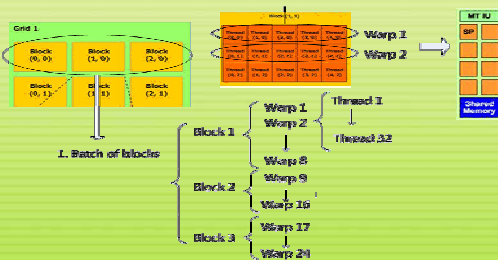## Programming Model

Streaming Multiprocessor (SM)
- Processing elements
  - 8 scalar thread processors
  - 32 GFLOPS peak at 1.35GHz
  - 8192 32-bit registers (32KB)
  - usual ops: float, int, branch...

- Hardware multithreading
  - up to 8 blocks (3 active) residents at once
  - up to 768 active threads in total

- 16KB on-chip memory
  - supports thread communication
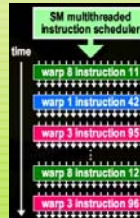  - shared amongst threads of a block



## Programming Model

Execution Model:



## Programming Model

Single Instruction Multiple Thread (SIMT) Execution:

- Groups of 32 threads formed into warps
  - always executing same instruction
  - share instruction fetch/dispatch
  - some become inactive when code path diverges
  - hardware automatically handles divergence

- Warps are primitive unit of scheduling
  - pick 1 of 24 warps for each instruction slot.
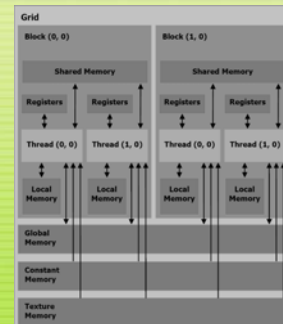  - all warps from all active blocks are time-sliced



## Outline

- Introduction
- Motivation
- Programming Model
- **Memory Model**
- CUDA API
- Example
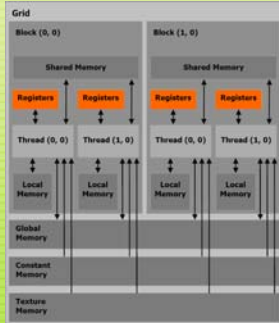- Pro & Contra
- Trend

## Memory Model

There are 6 Memory Types :

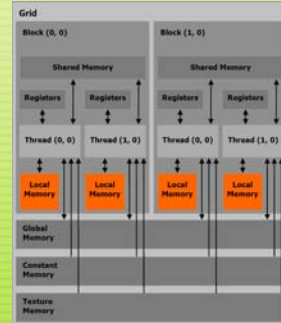## Memory Model

There are 6 Memory Types :

- **Registers**
  - on chip
  - fast access
  - per thread
  - limited amount
  - 32 bit



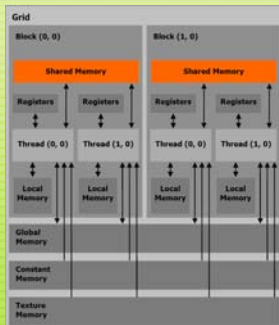## Memory Model

There are 6 Memory Types :

- Registers
- **Local Memory**
  - in DRAM
  - slow
  - non-cached
  - per thread
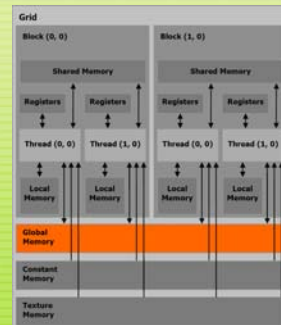  - relative large



## Memory Model

There are 6 Memory Types :

- Registers
- Local Memory
- **Shared Memory**
  - on chip
  - fast access
  - per block
  - 16 KByte
  - synchronize between threads



## Memory Model

There are 6 Memory Types :

- Registers
- Local Memory
- Shared Memory
- **Global Memory**
  - in DRAM
  - slow
  - non-cached
  - per grid
  - communicate between grids



## Memory Model

There are 6 Memory Types :

- Registers
- Local Memory
- Shared Memory
- Global Memory
- **Constant Memory**
  - in DRAM
  - cached
  - per grid
  - read-only



## Memory Model

There are 6 Memory Types :

- Registers
- Local Memory
- Shared Memory
- Global Memory
- Constant Memory
- **Texture Memory**
  - in DRAM
  - cached
  - per grid
  - read-only

## Memory Model

- Registers
- Shared Memory
  - on chip

- Local Memory
- Global Memory
- Constant Memory
- Texture Memory
  - in Device Memory



## Memory Model

- Global Memory
- Constant Memory
- Texture Memory
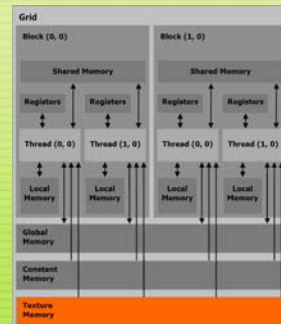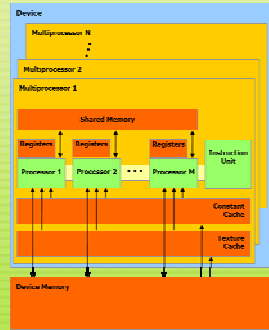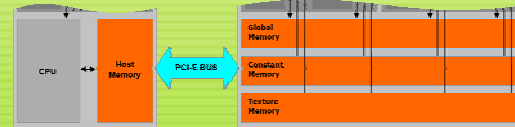  - managed by host code
  - persistent across kernels



## Outline

- Introduction
- Motivation
- Programming Model
- Memory Model
- **CUDA API**
- Example
- Pro & Contra
- Trend

## CUDA API

CUDA API provides a easily path for users to write programs for GPU device .

It consists of:

- A minimal set of extensions to C/C++
  - type qualifiers
  - call-syntax
  - build-in variables

- A runtime library to support the execution
  - host component
  - device component
  - common component

## CUDA API

CUDA C/C++ Extensions:
- New function type qualifiers

```
__host__   void HostFunc(...);   //executable on host
__global__ void KernelFunc(...); //callable from host
__device__ void DeviceFunc(...); //callable from device only
```
  - Restrictions for device code (`__global__` / `__device__`)
    - no recursive call
    - no static variable
    - no function pointer
    - `__global__` function is asynchronous invoked
    - `__global__` function must have `void` return type

## CUDA API

CUDA C/C++ Extensions:
- New variable type qualifiers

```
__device__ int GlobalVar; //in global memory, lifetime of app
__const__  int ConstVar;  //in constant memory, lifetime of app
__shared__ int SharedVar; //in shared memory, lifetime of blocks
```
  - Restrictions
    - no external usage
    - only file scope
    - no combination with `struct` or `union`
    - no initialization for `__shared__`

## CUDA API

CUDA C/C++ Extensions:
- New syntax to invoke the device code

```
KernelFunc<<< Dg, Db, Ns, S >>>(...);
```
  - ○ Dg: dimension of grid
  - ○ Db: dimension of block
  - ○ Ns: optional, shared memory for external variables
  - ○ S : optional, associated stream

- New build-in variables for indexing the threads
  - ○ gridDim: dimension of the whole grid
  - ○ blockIdx: index of the current block
  - ○ blockDim: dimension of each block in the grid
  - ○ threadIdx: index of the current thread

## CUDA API

CUDA Runtime Library:
- Common component
  - ○ Vector/Texture Types
  - ○ Mathematical/Time Functions

- Device component
  - ○ Mathematical/Time/Texture Functions
  - ○ Synchronization Function
    - ○ __syncthreads()
  - ○ Type Conversion/Casting Functions
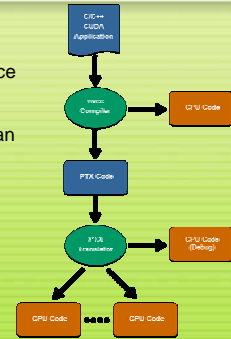
## CUDA API

CUDA Runtime Library:
- Host component

  - ○ Structure
    - ▪ Driver API
    - ▪ Runtime API

  - ○ Functions
    - ▪ Device, Context, Memory, Module, Texture management
    - ▪ Execution control
    - ▪ Interoperability with OpenGL and Direct3D

## CUDA API

The CUDA source file uses `.cu` as extension. It contains host and device source codes.

The CUDA Compiler Driver `nvcc` can compile it and generate CPU/PTX binary code.
(PTX: Parallel Thread Execution, a device independent VM code)

PTX code may be further translated for special GPU-Arch.



## Outline

- Introduction
- Motivation
- Programming Model
- Memory Model
- CUDA API
- **Example**
- Pro & Contra
- Trend

## Programming Model

Simple example ( Matrx addition ):
cpu c program:          cuda program:

```
void addVector (float *a, float *b,          __global__ void addVector (float *a, float *b,
             float *c, int N)                                float *c)
{                                            {
  int i, index;                                int i = threadIdx.x + blockDim.x*blockIdx.x;
  for (i = 0; i<N; i++) {                       c[i] = a[i] + b[i];
    c[index] = a[index] + b[index];          }
  }
}

                                             void main()
void main()                                  {
{
  ...                                          // allocation & transfer data to GPU
  addVector(a, b, c, N);                       //Compute on N/256 blocks of 256 threads each
  ...                                          addVector <<< N/256, 256 >> ( d_A, d_B, d_C);
}                                              ...
                                             }
```

Device code

Host code

## Outline

- Introduction
- Motivation
- Programming Model
- Memory Model
- CUDA API
- Example
- **Pro & Contra**
- Trend

## Pro & Contra

CUDA allows
- massive parallel computing
- with a relative low price
- high integrated solution
- personal supercomputing
- ecofriendly production
- easy to learn

## Pro & Contra

Problem ......
- slightly low precision
- limited support for IEEE-754
- no recursive function call
- hard to use for irregular join/fork logic
- no concurrency between jobs

## Outline

- Introduction
- Motivation
- Programming Model
- Memory Model
- CUDA API
- Example
- Pro & Contra
- **Trend**

## Trend

- More cores on-chip
- Better support for float point
- Flexiber configuration & control/data flow
- Lower price
- Support higher level programming language

## References

[1] CUDA Programming Guide, nVidia Corp.
[2] The CUDA Compiler Driver, nVidia Corp.
[3] Parallel Thread Execution, nVidia Corp.
[4] CUDA: A Heterogeneous Parallel Programming Model for Manycore Computing, ASPLOS 2008, gpgpu.org