

INTRODUCTION TO GPU COMPUTING

Ciprian Mihai Mitu, Mihai Niculescu
Institute for Space Sciences
Romania

Contents

- Parallel Computing
- What is CUDA?
- Why CUDA?
- CUDA Programming Model
- CUDA Execution Model
- Conclusions

Parallel computing

- Definition: solving complex problems using more processors.
- Michael J.Flynn – 1966 classification:
 - Number of programs
 - Number of sets of data

Parallel Computing Classification



What is CUDA?

- **Compute Unified Device Architecture - NVIDIA**
 - General purpose parallel computing for GPU
 - Requirements:
 - ISA(Instruction Set Architecture)
 - NV driver
 - CUDA Toolkit
 - (CUDA SDK)

Cuda Toolkit

- nvcc C compiler
- CUFFT + CUBLAS libraries
- Cuda Profiler
- Nv gdb
- CUDA Runtime
- Programming manual

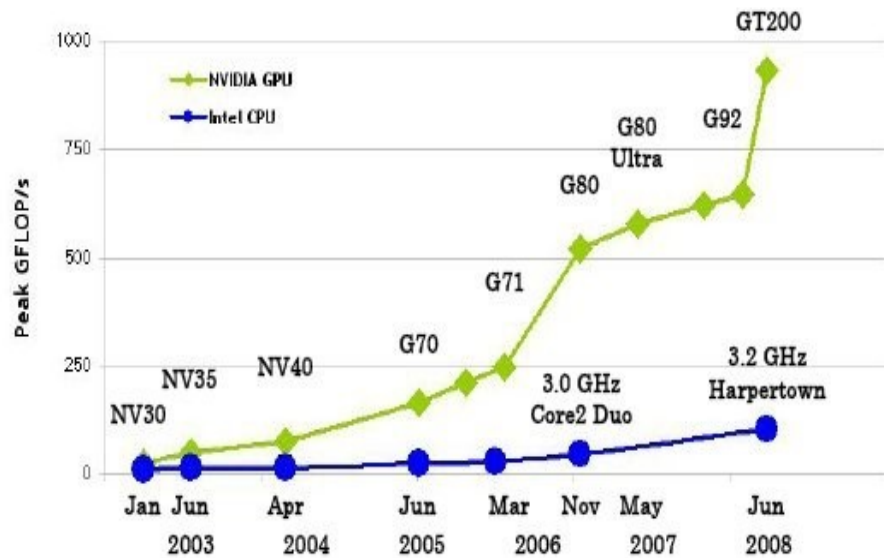
Cuda SDK

- Lots of examples
 - Mandelbrot
 - Nbody
 - Fluids
 - ...
- Dev easy

Why CUDA?

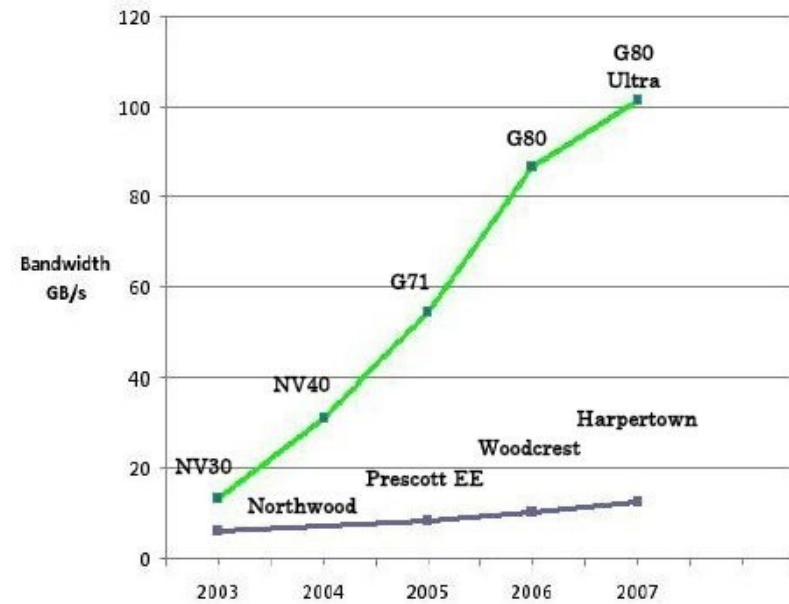
- Price ~ \$ 200
- Performance - speedups > 20x
- Language Bindings C/C++ and FORTRAN
- Generality: diverse parallel algorithms
- Scalability

GPU vs CPU

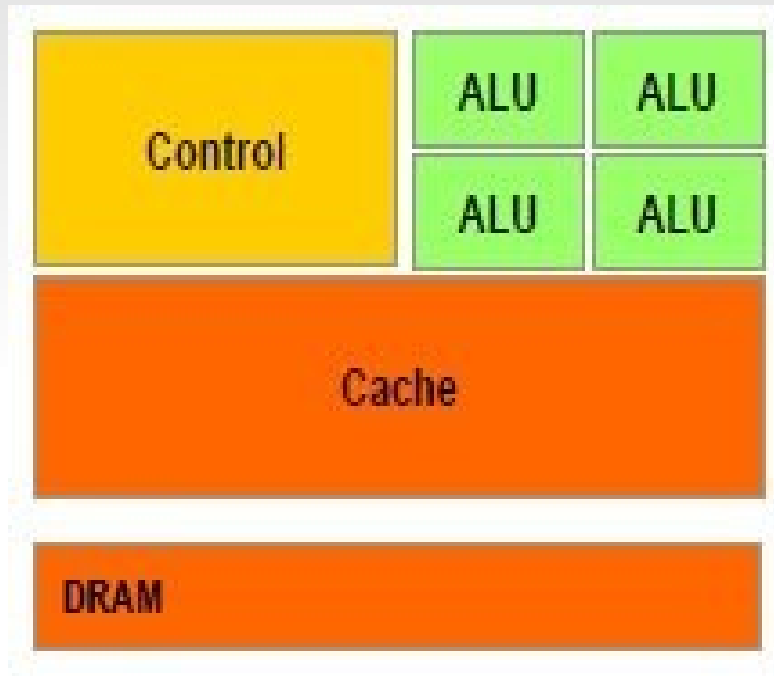


GT200 = GeForce GTX 280	G71 = GeForce 7900 GTX	NV35 = GeForce FX 5950 Ultra
G92 = GeForce 9800 GTX	G70 = GeForce 7800 GTX	NV30 = GeForce FX 5800
G80 = GeForce 8800 GTX	NV40 = GeForce 6800 Ultra	

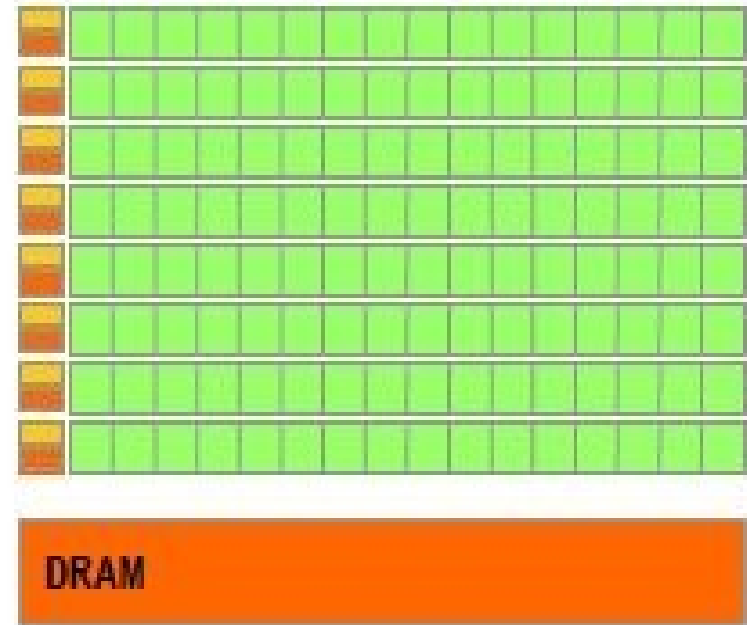
GT200 = GeForce GTX 280	G71 = GeForce 7900 GTX	NV35 = GeForce FX 5950 Ultra
G92 = GeForce 9800 GTX	G70 = GeForce 7800 GTX	NV30 = GeForce FX 5800
G80 = GeForce 8800 GTX	NV40 = GeForce 6800 Ultra	



Why is GPU so fast?



CPU

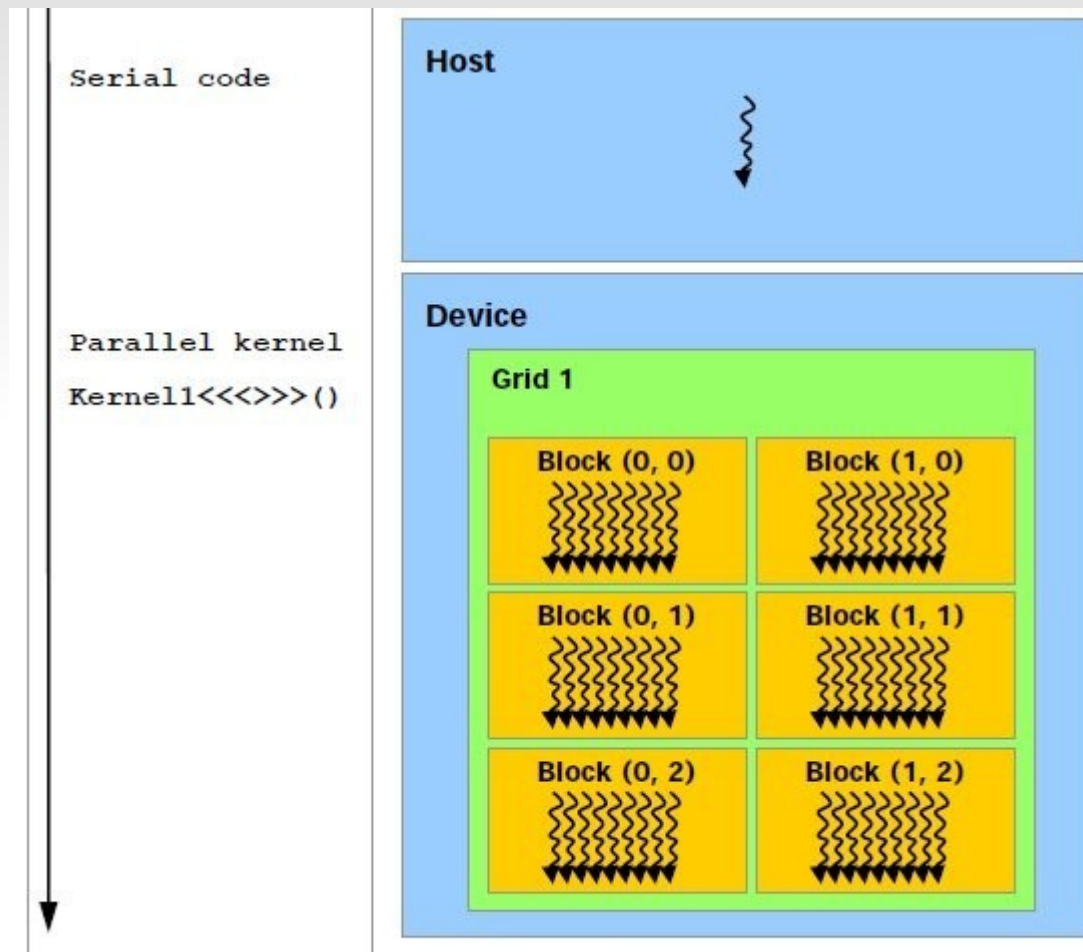


GPU

Programming Model

- GPU is a computing device:
 - coprocessor for the CPU (host)
 - own memory (DRAM)
 - many threads in parallel (>1000)
 - SPMD – single program, multiple data

Execution Model



Serial code executes on the host while parallel code executes on the device.

Performance Optimization

- Expose as much parallelism as possible
- Optimize memory usage to achieve maximum memory throughput
- Optimize instruction usage to achieve maximum instruction throughput
- Maximize occupancy to hide latency

Conclusions

- CUDA is a powerful parallel programming model
 - Heterogeneous
 - Scalable
 - Accessible

References

- http://www.nvidia.com/object/cuda_home_new.html
- <http://forums.nvidia.com>