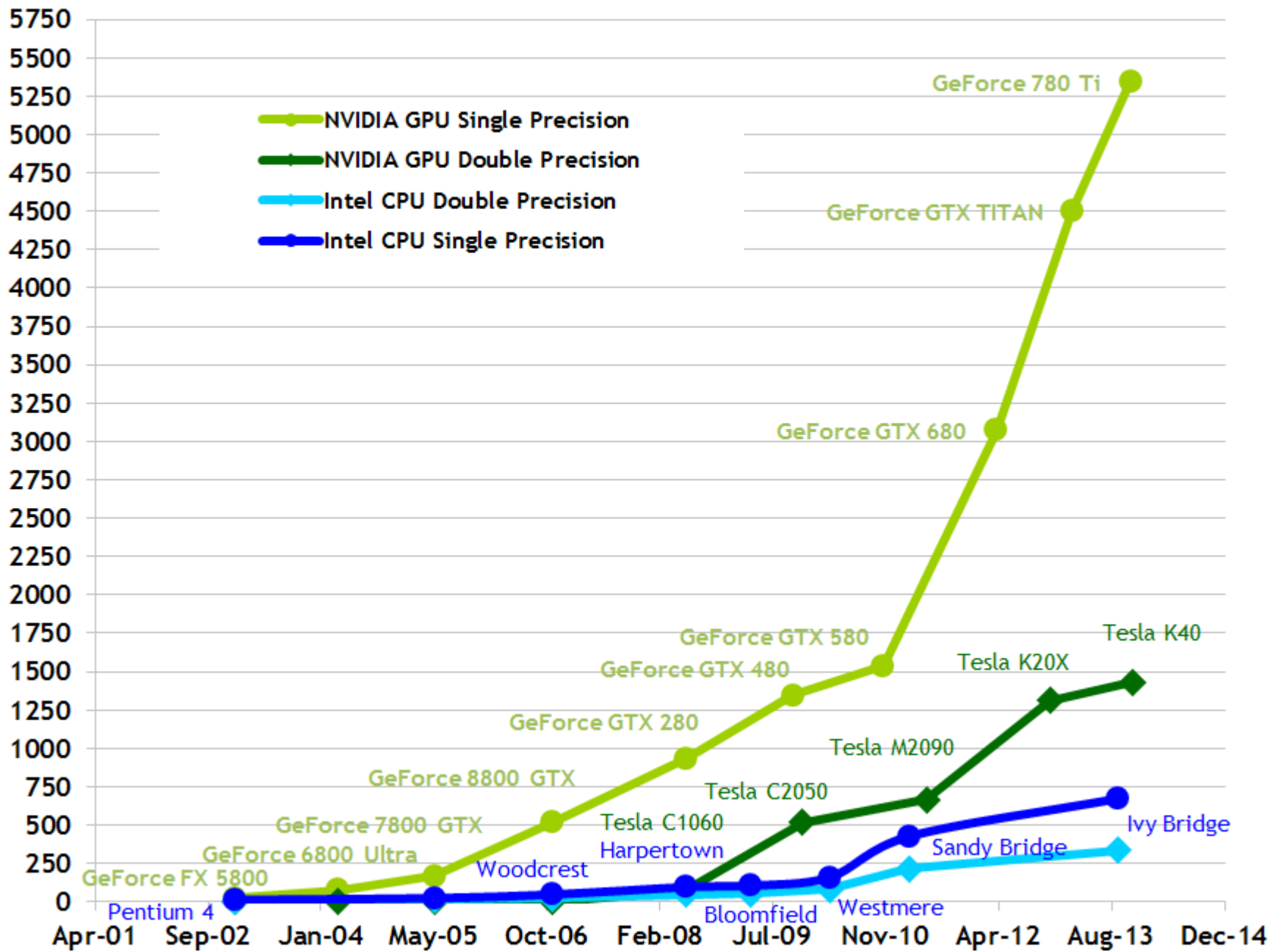
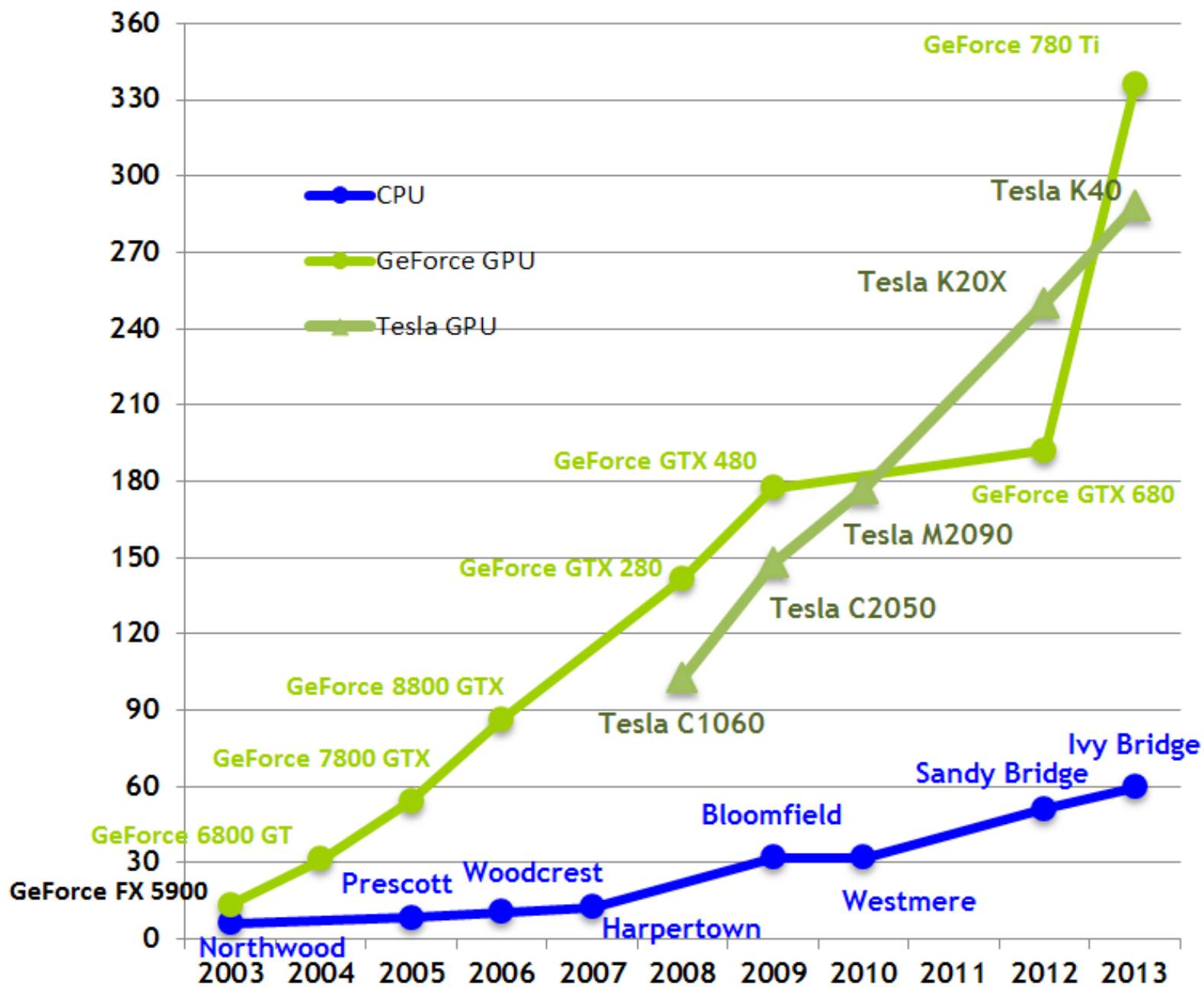


# NVIDIA Architecture

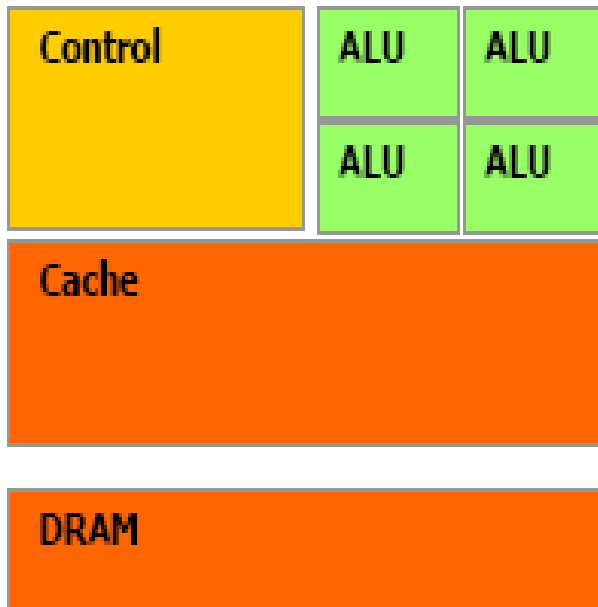
# Theoretical GFLOP/s



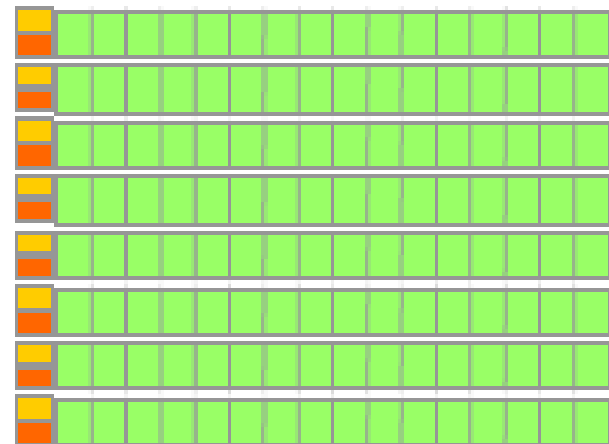
## Theoretical GB/s



# CPU v/s GPU



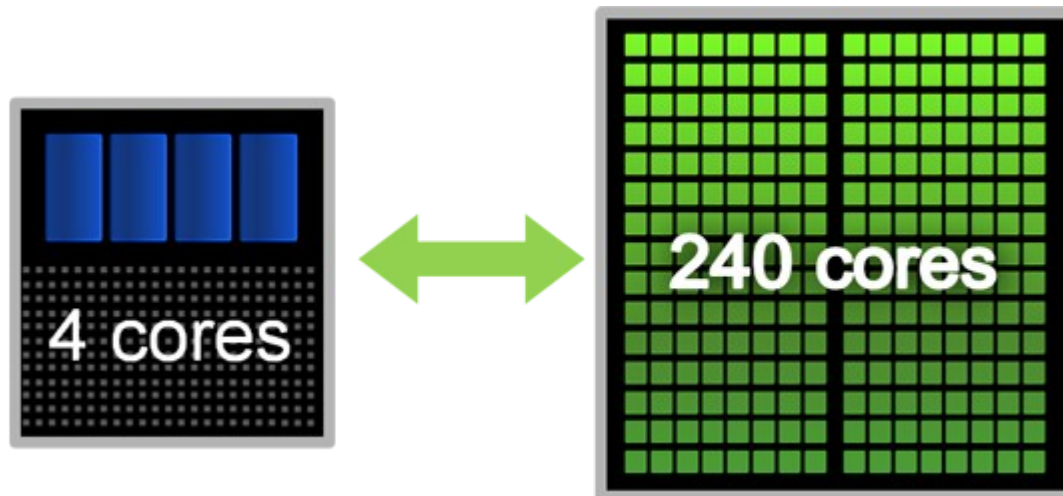
CPU



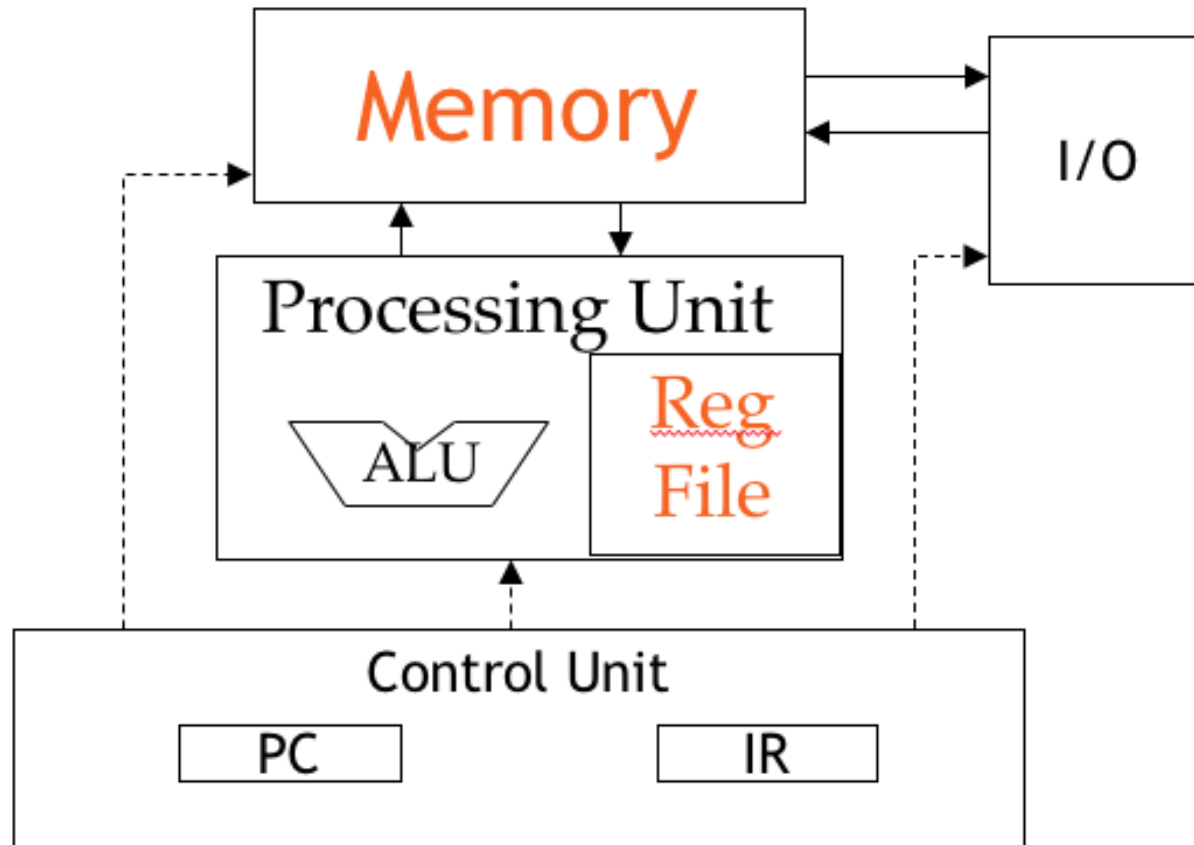
GPU

# GPU and CPU

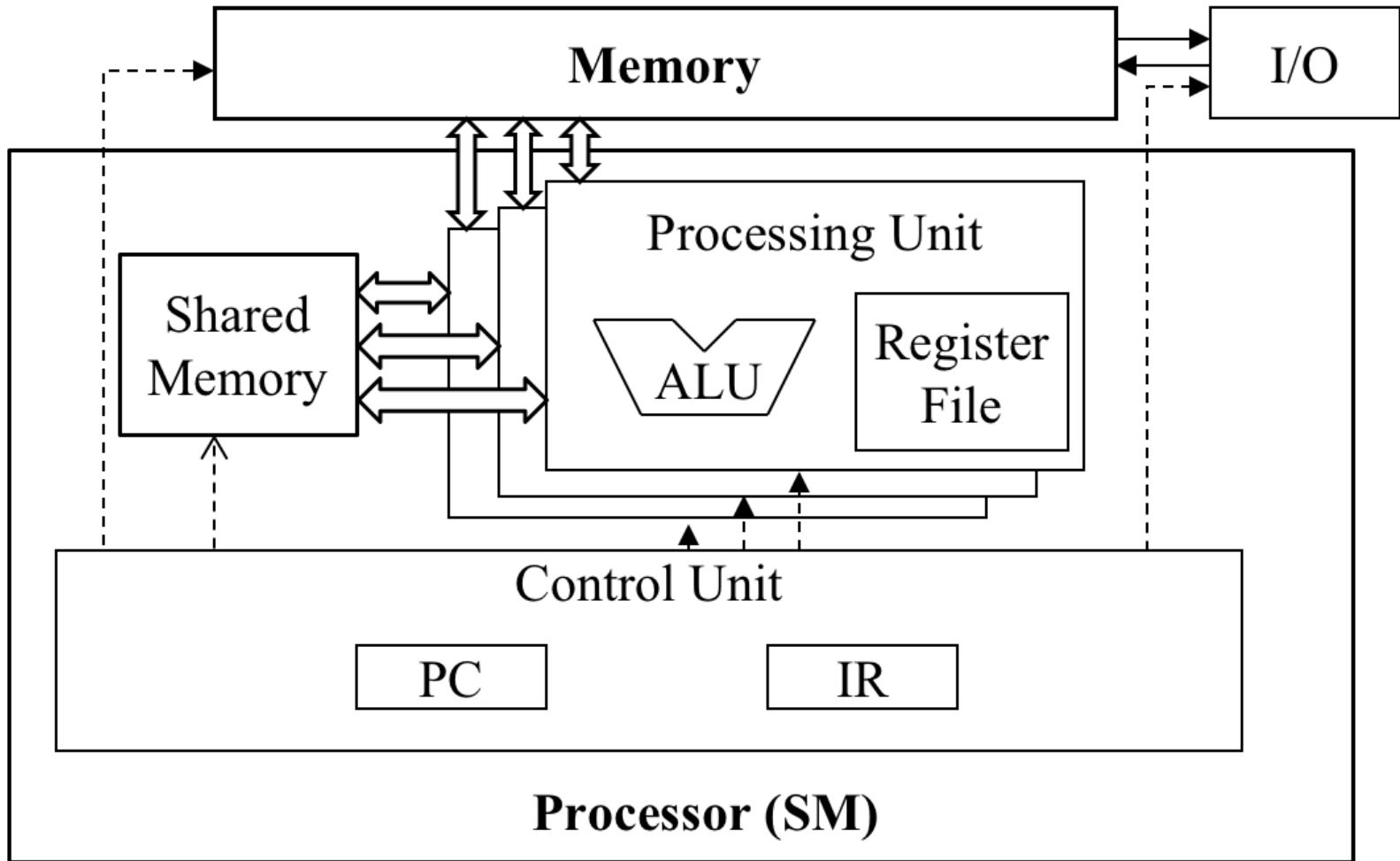
- Typically GPU and CPU coexist in a heterogeneous setting
- “Less” computationally intensive part runs on CPU (coarse-grained parallelism), and more intensive parts run on GPU (fine-grained parallelism)
- NVIDIA’s GPU architecture is called CUDA (Compute Unified Device Architecture) architecture, accompanied by CUDA programming model, and CUDA C language



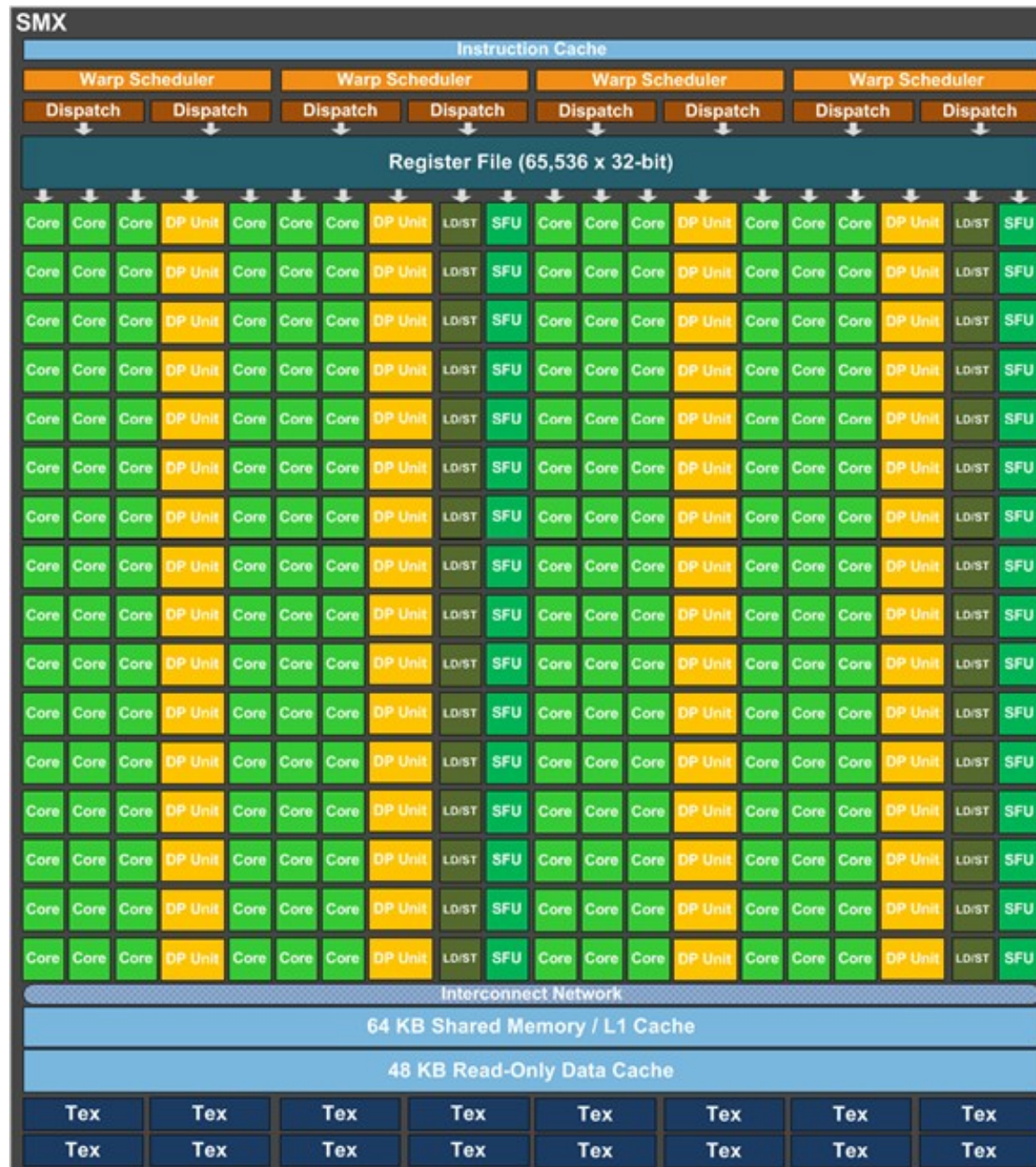
# The von Neumann Core



# GPU Core



# Onde os blocos são executados

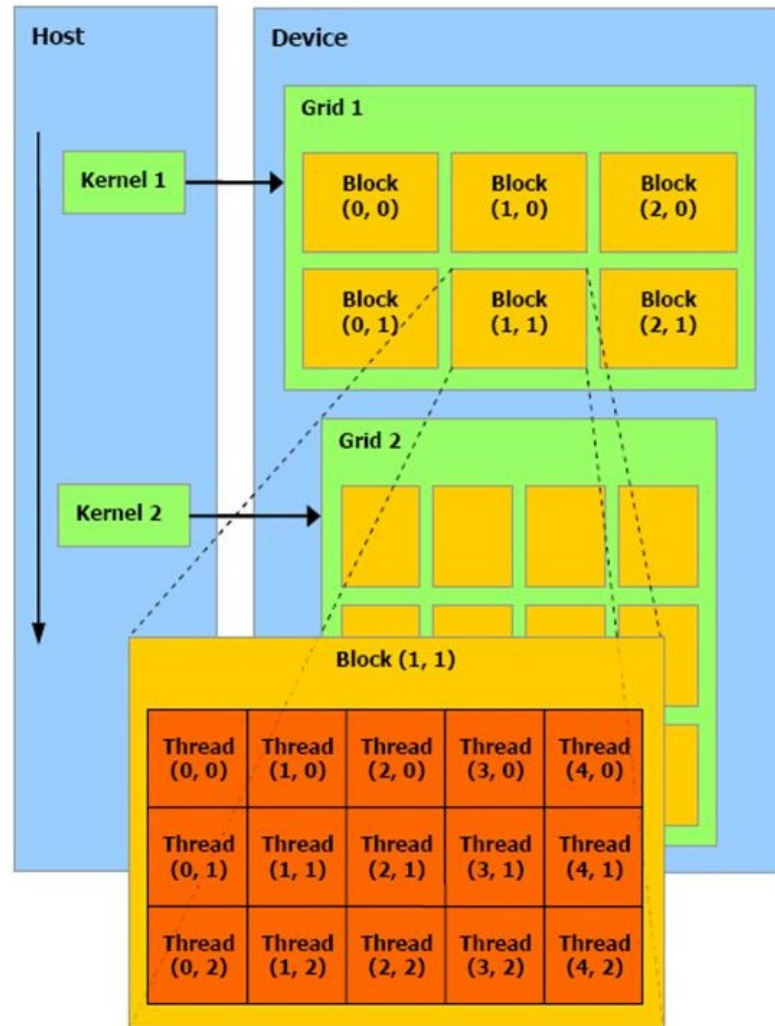




# Onde os blocos são executados



# Programming Model

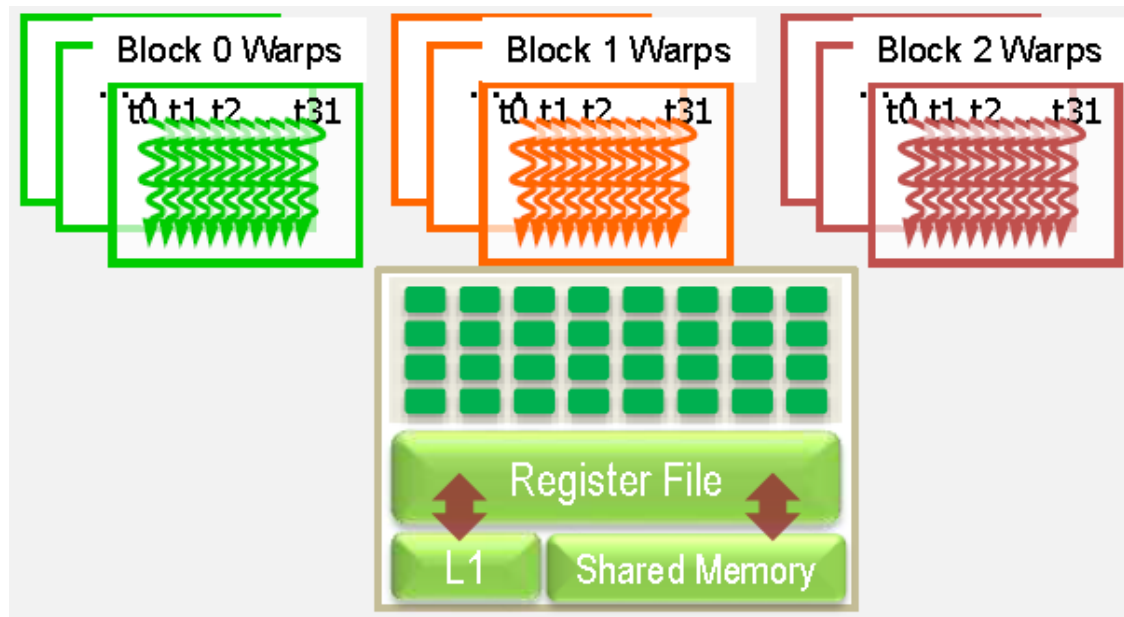


# Programming Model

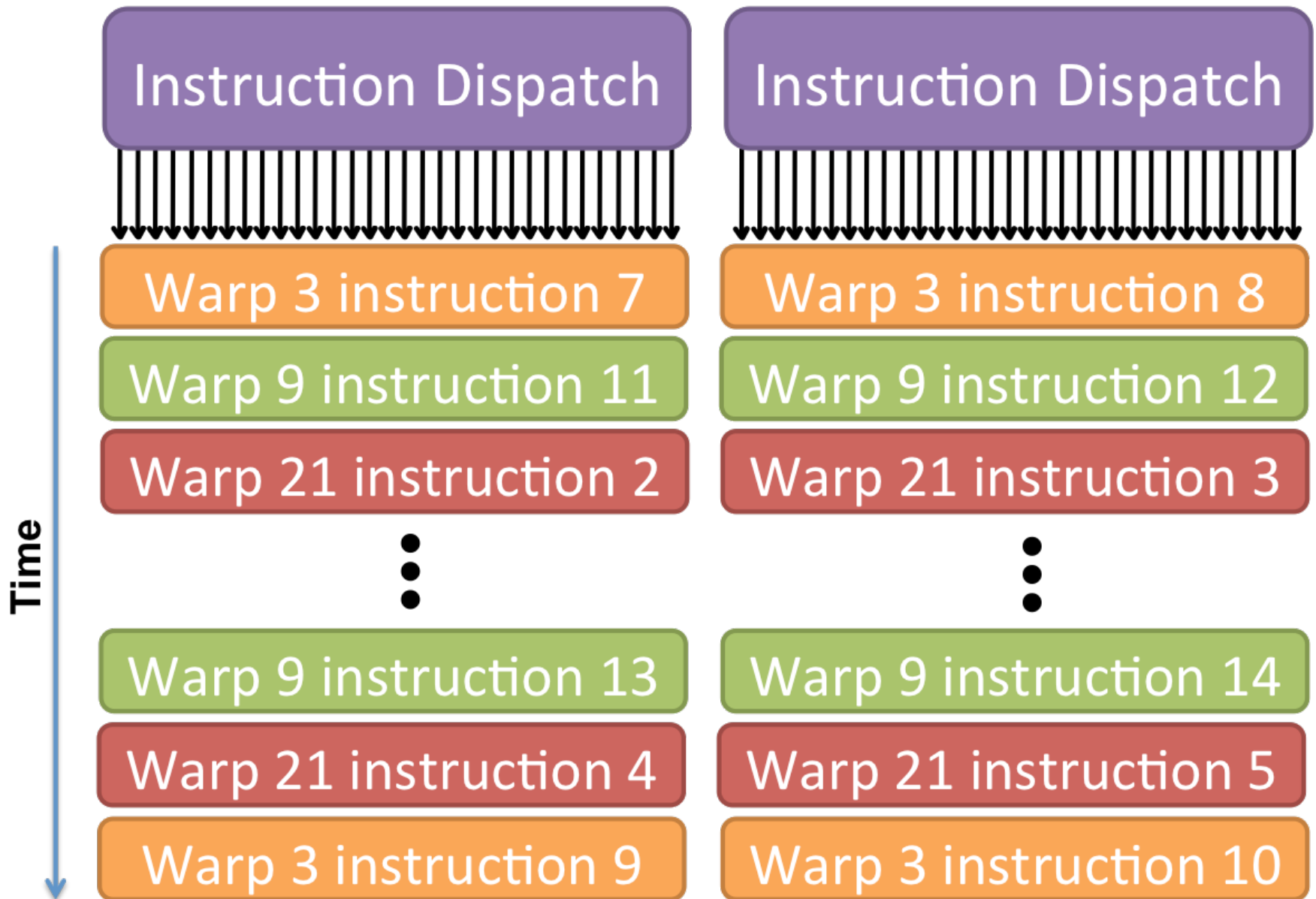
- All threads in a grid run the same kernel code (SPMD)
- Each thread has indexes that it uses to compute memory addresses and make control decisions.
  - $\text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$
- Threads run in groups of 32 called warps
  - Threads in a warp execute in SIMD
- Every thread in a warp executes the same instruction at a time

# Warp example

- If 3 blocks are assigned to an SM and each block has 256 threads, how many Warps are there in an SM?
  - Each Block is divided into  $256/32 = 8$  Warps
  - There are  $8 * 3 = 24$  Warps

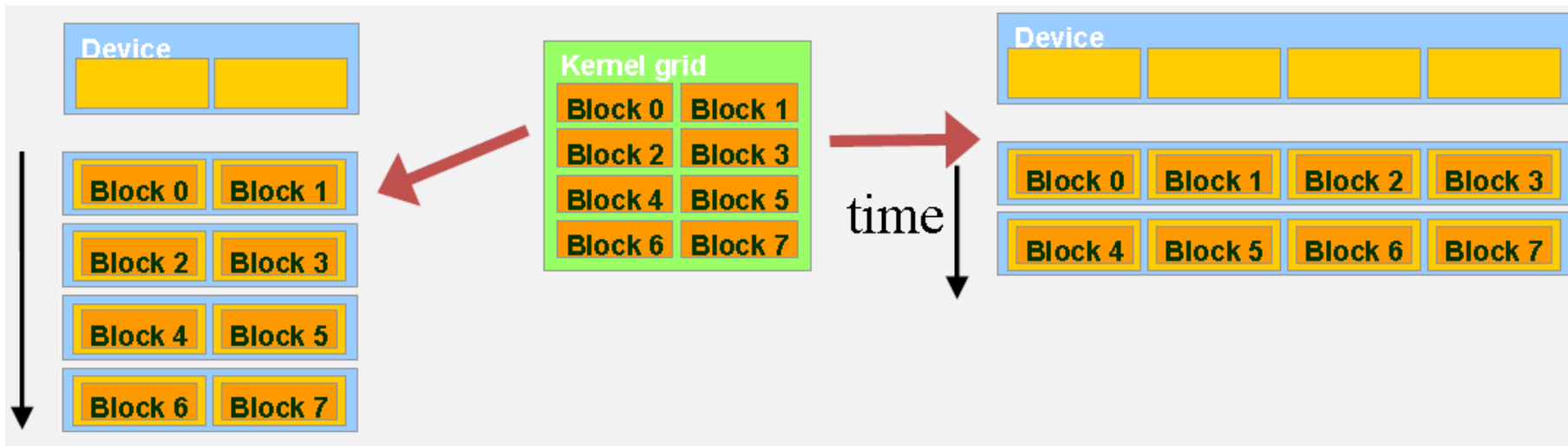


# Warp Scheduler



# Transparent Scalability

- Each block can execute in any order relative to others.
  - Hardware is free to assign blocks to any processor at any time
  - A kernel scales to any number of parallel processors

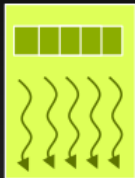


# Execution Model

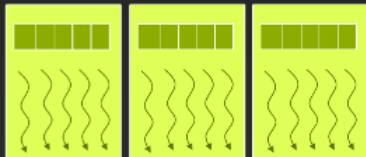
## Software



Thread



Thread Block

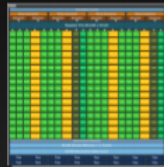


Grid

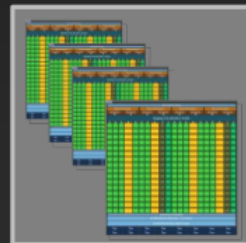
## Hardware



CUDA  
Core



Multiprocessor



Device

**Threads are executed by scalar CUDA Cores**

**Thread blocks are executed on multiprocessors**

**Thread blocks do not migrate**

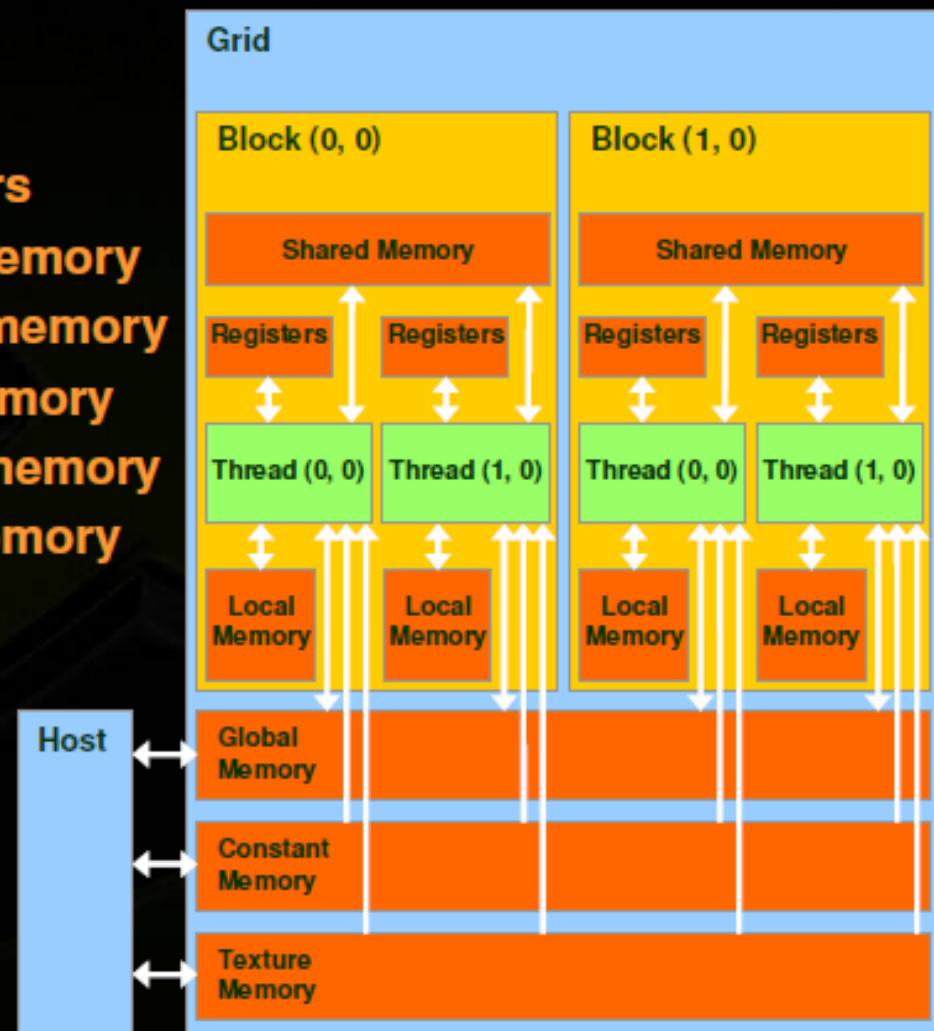
**Several concurrent thread blocks can reside on one multiprocessor - limited by multiprocessor resources (shared memory and register file)**

**A kernel is launched as a grid of thread blocks**



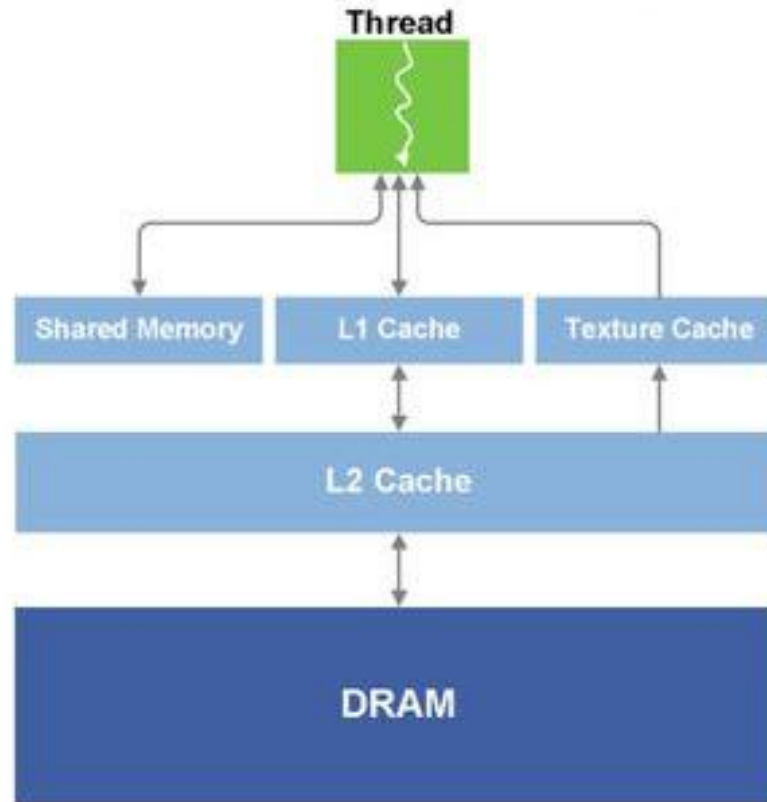
# CUDA Memory Spaces

- Each thread can:
  - Read/write per-thread **registers**
  - Read/write per-thread **local memory**
  - Read/write per-block **shared memory**
  - Read/write per-grid **global memory**
  - Read only per-grid **constant memory**
  - Read only per-grid **texture memory**
  
- The host can read/write **global, constant, and texture memory (stored in DRAM)**





# Thread Memory Access



# NVIDIA Kepler K40

- 2880 streaming processors/cores (SPs) organized as 15 streaming multiprocessors (SMs)
- Each SM contains 192 cores
- Memory size of the GPU system: 12 GB
- Clock speed of a core: 745 MHz

# P100

Tesla P100 accelerators have four 4-die HBM2 stacks, for a total of 16 GB of memory, and **720** GB/s peak bandwidth, which is 3 times higher than the Tesla M40 memory bandwidth



Tesla Products	Tesla K40	Tesla M40	Tesla P100
GPU	GK110 (Kepler)	GM200 (Maxwell)	GP100 (Pascal)
SMs	15	24	56
TPCs	15	24	28
FP32 CUDA Cores / SM	192	128	64
FP32 CUDA Cores / GPU	2880	3072	3584
FP64 CUDA Cores / SM	64	4	32
FP64 CUDA Cores / GPU	960	96	1792
Base Clock	745 MHz	948 MHz	1328 MHz
GPU Boost Clock	810/875 MHz	1114 MHz	1480 MHz
Compute Performance - FP32	5.04 TFLOPS	6.82 TFLOPS	10.6 TFLOPS
Compute Performance - FP64	1.68 TFLOPS	0.21 TFLOPS	5.3 TFLOPS
Texture Units	240	192	224
Memory Interface	384-bit GDDR5	384-bit GDDR5	4096-bit HBM2
Memory Size	Up to 12 GB	Up to 24 GB	16 GB
L2 Cache Size	1536 KB	3072 KB	4096 KB
Register File Size / SM	256 KB	256 KB	256 KB
Register File Size / GPU	3840 KB	6144 KB	14336 KB
TDP	235 Watts	250 Watts	300 Watts
Transistors	7.1 billion	8 billion	15.3 billion
GPU Die Size	551 mm <sup>2</sup>	601 mm <sup>2</sup>	610 mm <sup>2</sup>
Manufacturing Process	28-nm	28-nm	16-nm