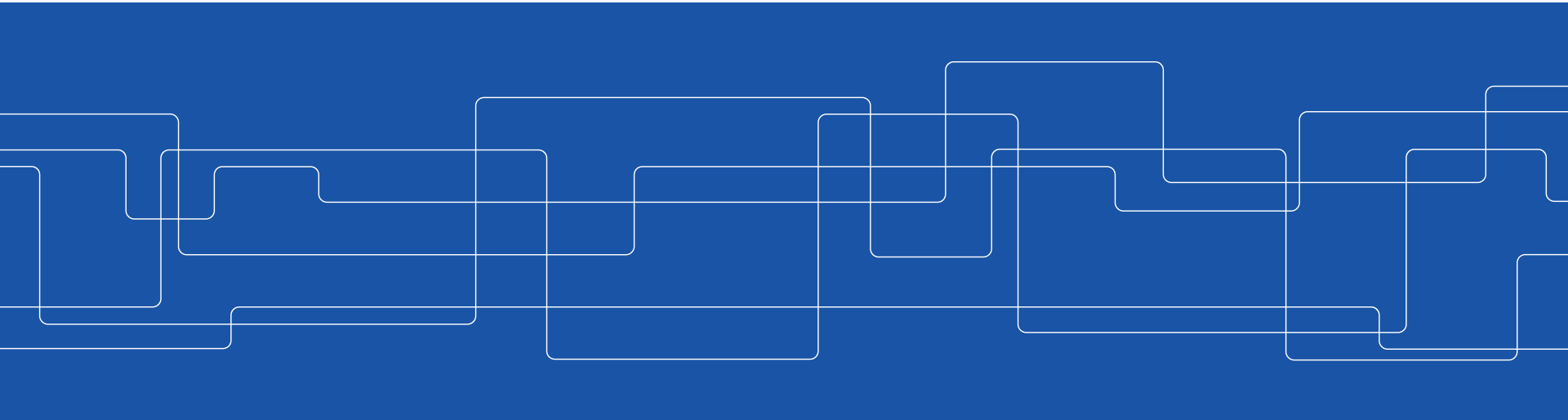




Introduction to GPUs

Stefano Markidis

KTH Royal Institute of Technology



GPUs

GPU = Graphical Processing Unit = specialized microcircuit to accelerate the creation and manipulation of images in video frame for display devices.

GPUs are used in game consoles, embedded systems (like systems on cars for automatic driving), computers and supercomputers.

- Since 2012, GPUs are the main workforce for training deep-learning networks

Some important GPU vendors: **NVIDIA**, **AMD**, ...



The Rise of GPUs in HPC

Top500

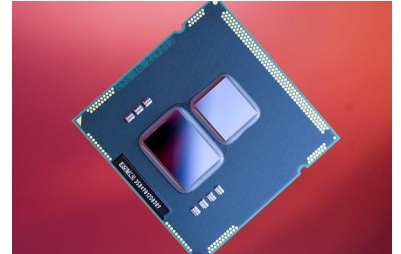
GPUs are a core technology in many world's fastest and most energy-efficient supercomputers

- Top500: #1, #3, #5 and #6 (fastest European SC)
- GPUs compete well in terms of **FLOPS/Watt**

Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM DOE/SC/Oak Ridge National Laboratory United States	2,282,544	122,300.0	187,659.3	8,806
2	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway , NRCPC National Supercomputing Center in Wuxi China	10,649,600	93,014.6	125,435.9	15,371
3	Sierra - IBM Power System S922LC, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband , IBM DOE/NNSA/LLNL United States	1,572,480	71,610.0	119,193.6	
4	Tianhe-2A - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000 , NUDT National Super Computer Center in Guangzhou China	4,981,760	61,444.5	100,678.7	18,482
5	AI Bridging Cloud Infrastructure (ABCI) - PRIMERGY CX2550 M4, Xeon Gold 6148 20C 2.4GHz, NVIDIA Tesla V100 SXM2, Infiniband EDR , Fujitsu National Institute of Advanced Industrial Science and Technology (AIST) Japan	391,680	19,880.0	32,576.6	1,649
6	Piz Daint - Cray XC50, Xeon E5-2690v3 12C 2.6GHz, Aries interconnect , NVIDIA Tesla P100 , Cray Inc. Swiss National Supercomputing Centre (CSCS) Switzerland	361,760	19,590.0	25,326.3	2,272

Where do you find GPUs ?

- **Integrated:** Every laptop has an integrated GPU built into its processor, i.e. Intel HD or Iris Graphics.
- **Dedicated:** A standalone GPU uses its own processor and memory. Most dedicated GPUs are removable. They require more power but also provide higher performance
 - **In HPC, we use dedicated GPUs**



Source: PC Authority



Source: bit-tech.net

Question: What is the main difference between the two?

GPU Design Motivation: Process Pixels in Parallel

Data parallel

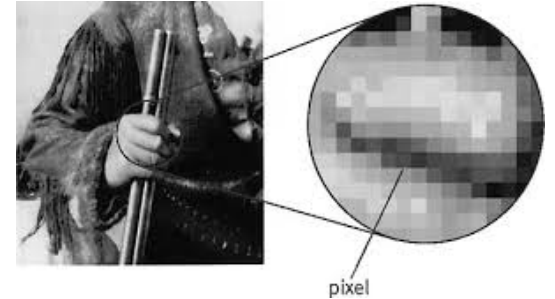
- In 1080i and 1080p videos, 1920 x 1080 pixels = 2M pixels per video frame → compute intensive
- Lots of parallelism at low clock speed → power efficient

Computation on each pixel is independent from computation on other pixels.

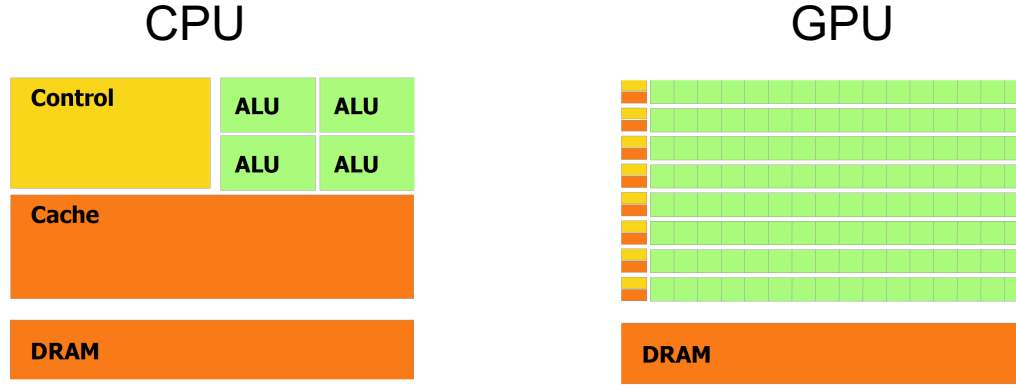
- No need for synchronization

Large data-locality = access to data is regular

- No need for large caches



What are the differences?



CPU has tens of massive cores, CPU excels at irregular control-intensive work

- Lots of hardware for control, fewer ALUs

GPU has thousands of small cores, GPU excels at regular math-intensive work

- Lots of ALUs, little hardware for control

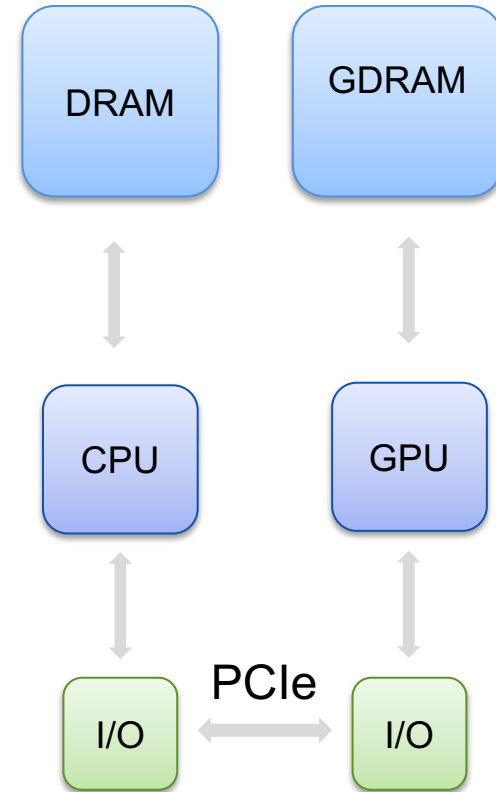
GPUs as Accelerators

GPU are simple, lower power and highly parallel

Problem: Still require OS, IO and scheduling

Solution: “Hybrid System”

- CPU provides management
- “Accelerators” (or co-processors) such as GPUs provide compute power



GPU Hardware Model

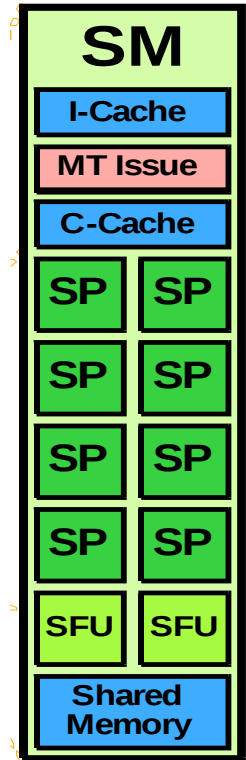
In order to program a GPU program, it is important to understand the Hardware Model.

The fundamental computing entity is

- **Streaming Processor (SP)** or **CUDA core**

A **Streaming Multiprocessor (SM)**:

- A collection of 8/32/192 CUDA Cores (depends on SM architecture)
- All CUDA cores in SM run the same instructions
- Has some fast cache shared memory
- Can synchronize

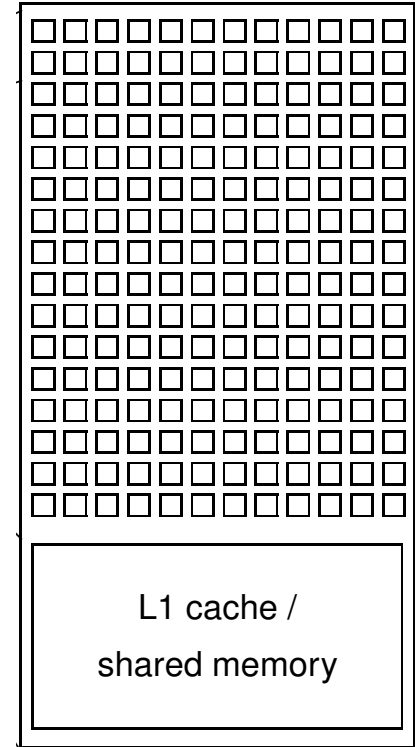




SMX = Next Generation SM

SM Architecture introduced in Kepler:

- Unified clock to save power
- **192 cores per SMX**





Today Lab: Which GPU?

You are going to ask for a **K420**

```
salloc --nodes=1 --gres=gpu:K420:1 -t  
00:05:00 -A ... - -reservation=...
```

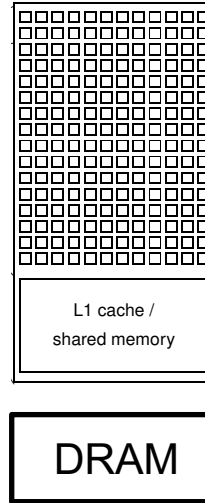
also available on Tegner: `gpu:K80`



K420

1 SMX:

- 192 cores!



Questions: how many cores per node on Beskow?



K80

2 GPUs:

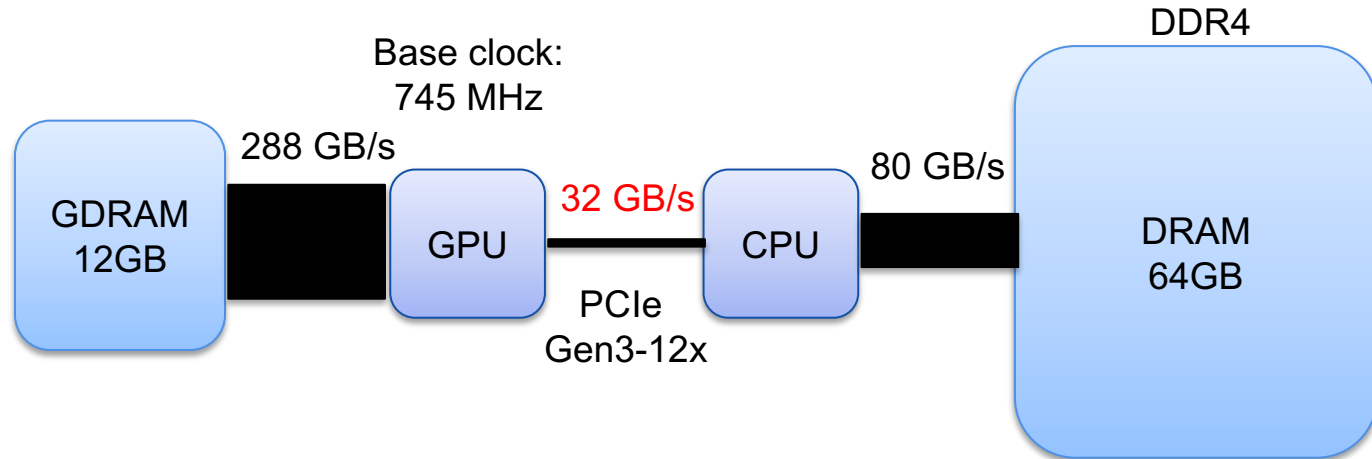
- 13 SMX per GPU

Questions: how many cores?

$$2 \times 13 \times 192 = \mathbf{4992!}$$

Weakness of GPU (but not for Volta... NVLink)

GPU is very fast (huge parallelism) but getting data from/to GPU is slow



NVIDIA TESLA K40 = the most common GPU on supercomputers in Nov. 2016 Top500 list



Is GPU good for my *non-graphics* application?

It depends on the application:

- **Compute-intensive applications** with little synchronization benefit the most from GPU:
 - Deep-learning network training 8×-10×, GROMACS 2×-3×, LAMMPS 2×-8×, QMCPack 3×.
- Irregular applications, such as sorting and constraint solvers, are faster on CPU.

General strategy when you work on your code: **take the computational-heavy part of your code and run it on GPU**



Low-Level Programming GPUs

- **OpenCL (Open Computing Language):** based on C, not only for GPUs but also for other “accelerators” (DSP, FPGA, ...)
- **CUDA (compute unified device architecture):** extension to C language. Only for **NVIDIA GPUs**.



High-Level Programming Interfaces

- **OpenMP**: compiler directives and library for accelerators
- **OpenACC**: compiler directives and library for NVIDIA GPUs

- **Thrust**: C++ template library resembling C++ STL.
- **OpenCV**: Computer vision library using GPU
- **CUDA-based libraries for math**: cuBLAS, cuFFT, cuDNN, ...

**Compiler
+ runtime
library**

**Libraries
atop
CUDA**



Let's move to CUDA now ...