# GPU Computing with NVIDIA CUDA

Introduction to parallel computing

Wageningen on 21.06.2019

heiko.loewe@dell.com
christian.schramm@dell.com
most images courtesy of NVIDIA

# Introduction to GPU computing with CUDA

A little bit of background

How could you increase the speed of a computing process?

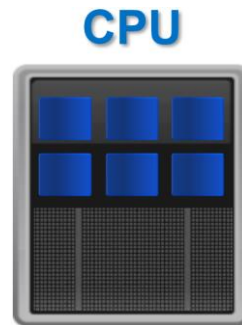-higher clock speed
-more work per clock cycle
-more processors



WE JUST NEED MORE POWER!

DELLEMC

# Introduction to GPU computing with CUDA

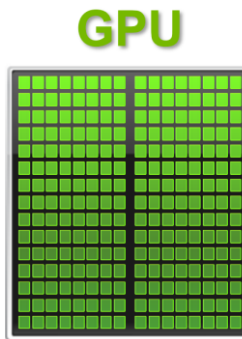A little bit of background

**Central Processing Unit (CPU)**
-consists of a few cores
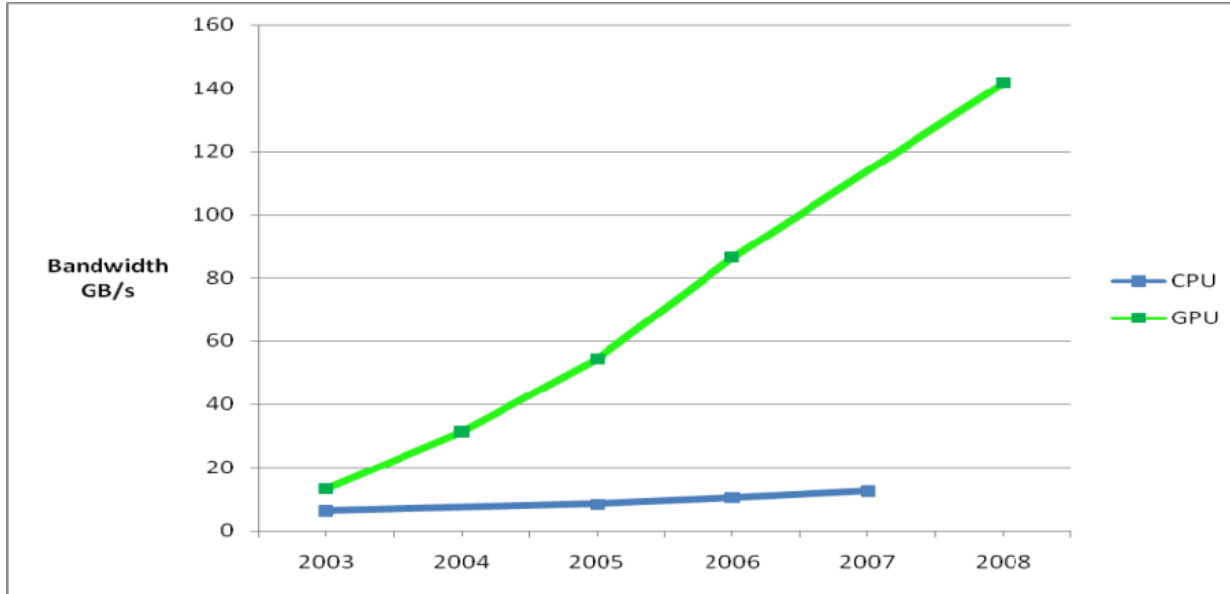-each one is powerful and optimized for **sequential** processing.

**Graphic Processing Unit (GPU)**
-consists of hundreds and thousands of smaller, less powerful cores
-the architecture is designed for handling multiple tasks **simultaneously**.

# CPU versus GPU

Supercomputing revolution



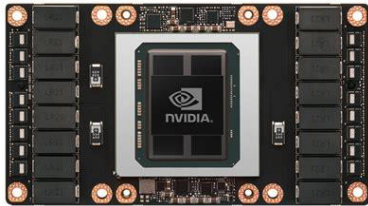| Model | Micro-architecture | Launch | Chips | Core clock (MHz) | Shaders | | | Memory | | | | | Processing power (GFLOPS)[a] | | | CUDA compute ability[b] | TDP (watts) | Notes, form_factor |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Cuda cores (total) | Base clock (MHz) | Max boost clock (MHz)[c] | Bus type | Bus width (bit) | Size (GB) | Clock (MT/s) | Bandwidth (GB/s) | Single precision (MAD+MUL) | Single precision (MAD or FMA) | Double precision (FMA) | | | |
| Units | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | MHz ⬍ | MHz ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | ⬍ | W ⬍ | |
| K80 GPU Accelerator[170] | | November 17, 2014 | 2× GK210 | N/A | 4992 | 560 | 875 | GDDR5 | 2× 384 | 2× 12 | 5000 | 2× 240 | No | 5591–8736 | 1864–2912 | 3.7 | 300 | Internal PCIe GPU (full-height, dual-slot) |
| T4 GPU Accelerator (PCIe card)[186][187] | Turing | September 12, 2018 | 1× TU104 | N/A | 2560 | 585 | 1590 | GDDR6 | 256 | 16 | Unknown | 320 | No | 8100 | Unknown | 7.5 | 70 | PCIe card |

# What is a GPU?

A little bit of background

3dfx Voodoo and NVIDIA GeForce 256

VGA interfaces

NVIDIA K80 (look familiar?)
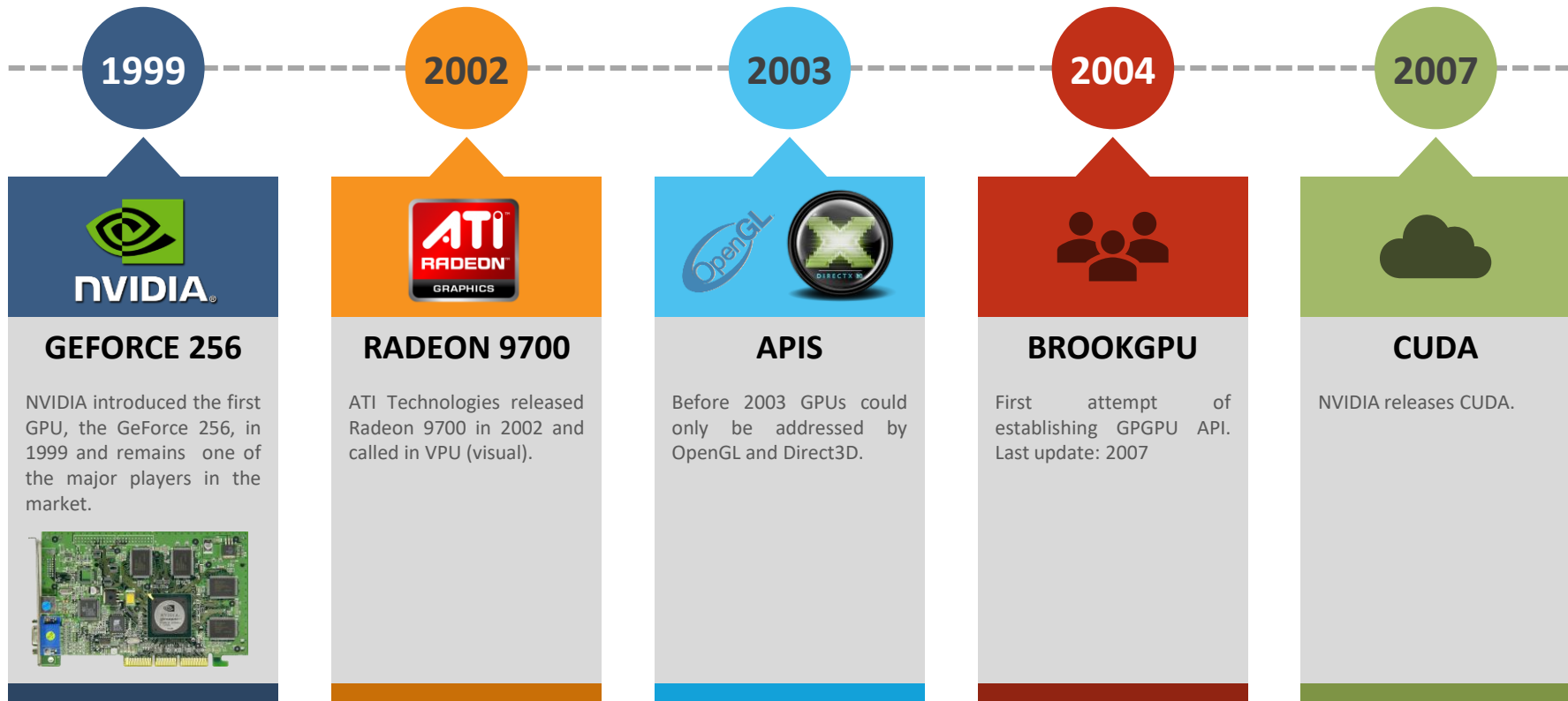
You cannot even attach a monitor.
…still PCI though

NVIDIA Pascal

This one does not even have PCI.
Clips right onto the main board.

**D≪LL**EMC

# What is a GPU?

A little bit of background

## 1999

### GEFORCE 256

NVIDIA introduced the first GPU, the GeForce 256, in 1999 and remains one of the major players in the market.

## 2002

### RADEON 9700

ATI Technologies released Radeon 9700 in 2002 and called in VPU (visual).

## 2003

### APIS

Before 2003 GPUs could only be addressed by OpenGL and Direct3D.

## 2004

### BROOKGPU

First attempt of establishing GPGPU API. Last update: 2007

## 2007

### CUDA

NVIDIA releases CUDA.

DELLEMC

# What is CUDA?

A little bit of background

**CUDA** stands for **C**ompute **U**nified **D**evice **A**rchitecture.
It is a parallel computing platform (using a GPU) and a programming model
(using code). CUDA is an extension of C and fully supports C++.

**Flynn's Taxonomy** introduced in 1966:
Single Instruction Single Datum (PC)
Single Instruction Multiple Data (GPU)
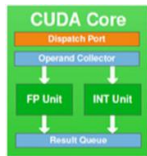Multiple Instruction Single Datum (Fault Tolerance)
Multiple Instruction Multiple Data (distributed systems, autonomous processors)

**D∕ELL**EMC

# What is a GPU?

Architecture

**CUDA Core:**
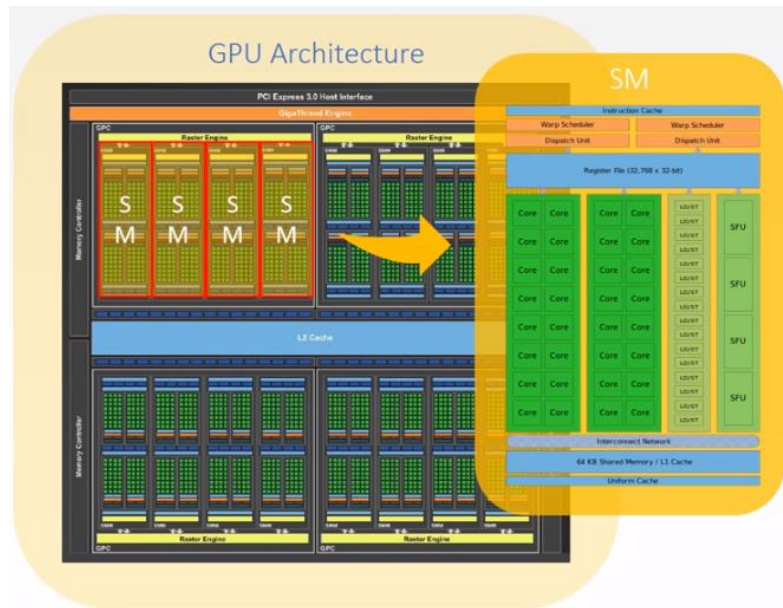-Smallest building block of a GPU.
-Executes computations ("threads")

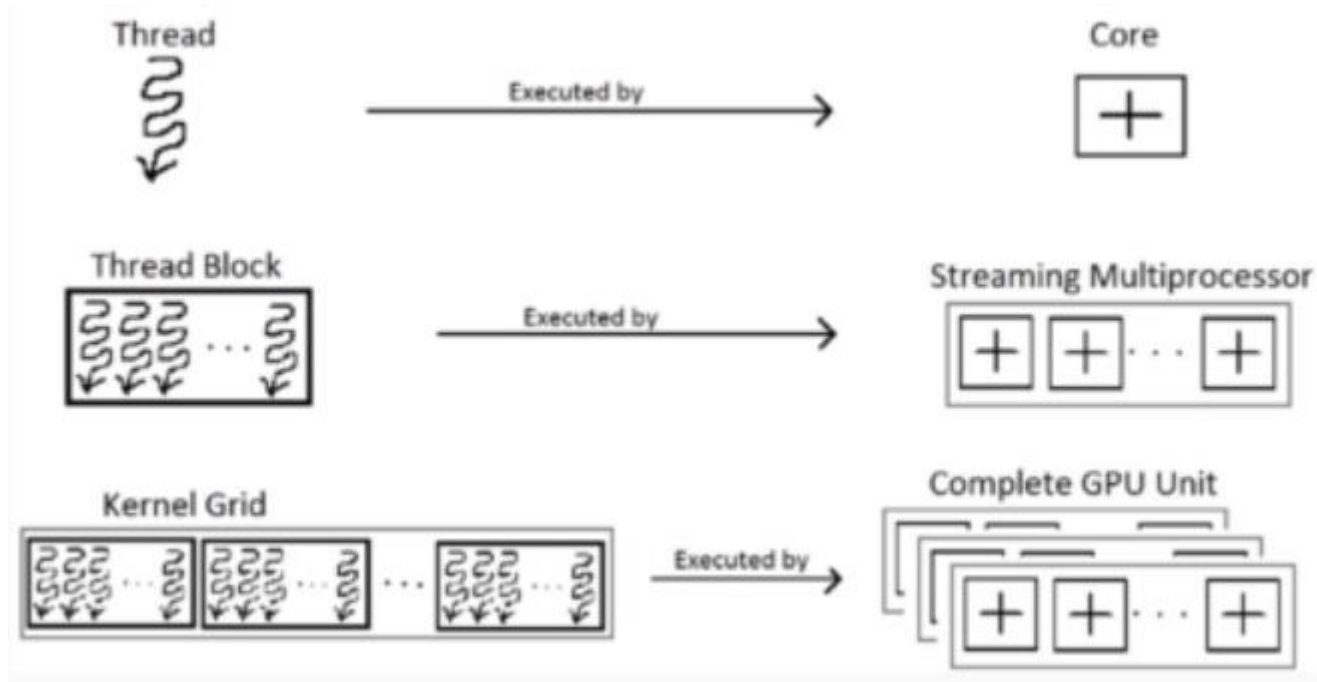**Stream Multiprocessor:**
-Collection of CUDA Cores including a Scheduler

**GPU:**
-Collection of SMs

**DELL**EMC

# What is a GPU?

Execution



© Copyright 2019 Dell Inc.

**D∅LL**EMC

# Heterogeneous Computing

A little bit of background

**Host:**
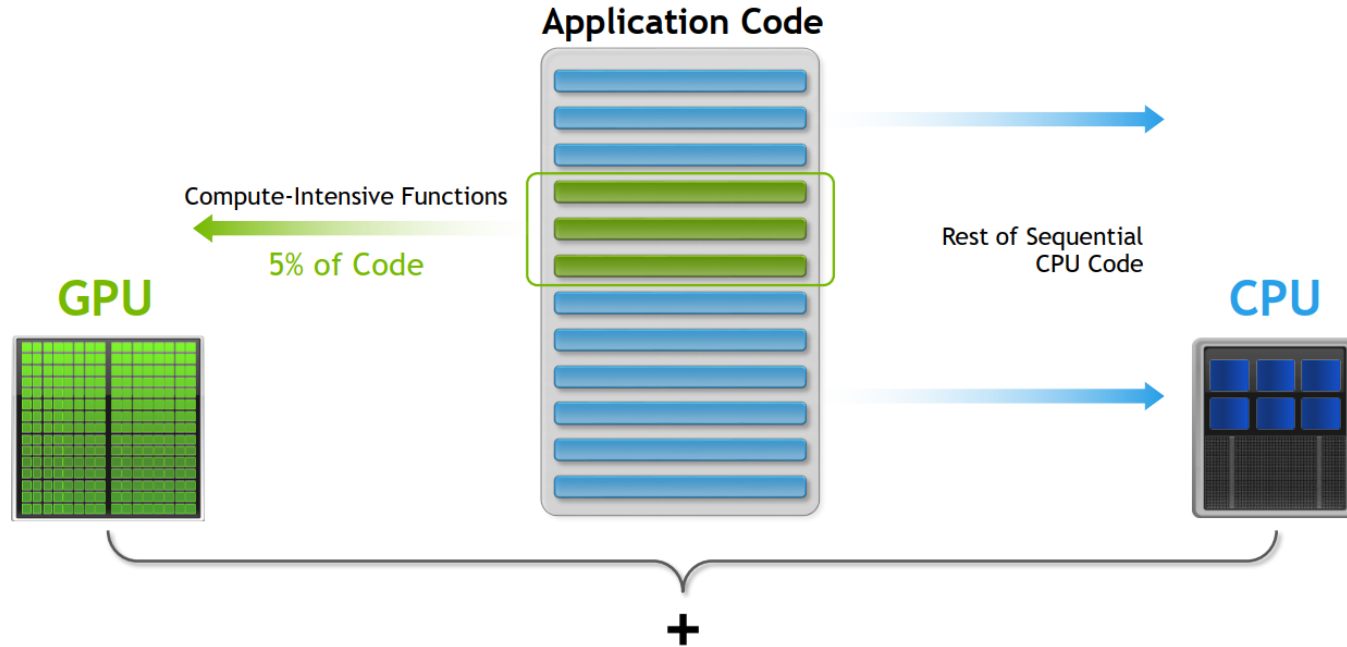The computer that has it own CPU and memory ("host memory")

**Device:**
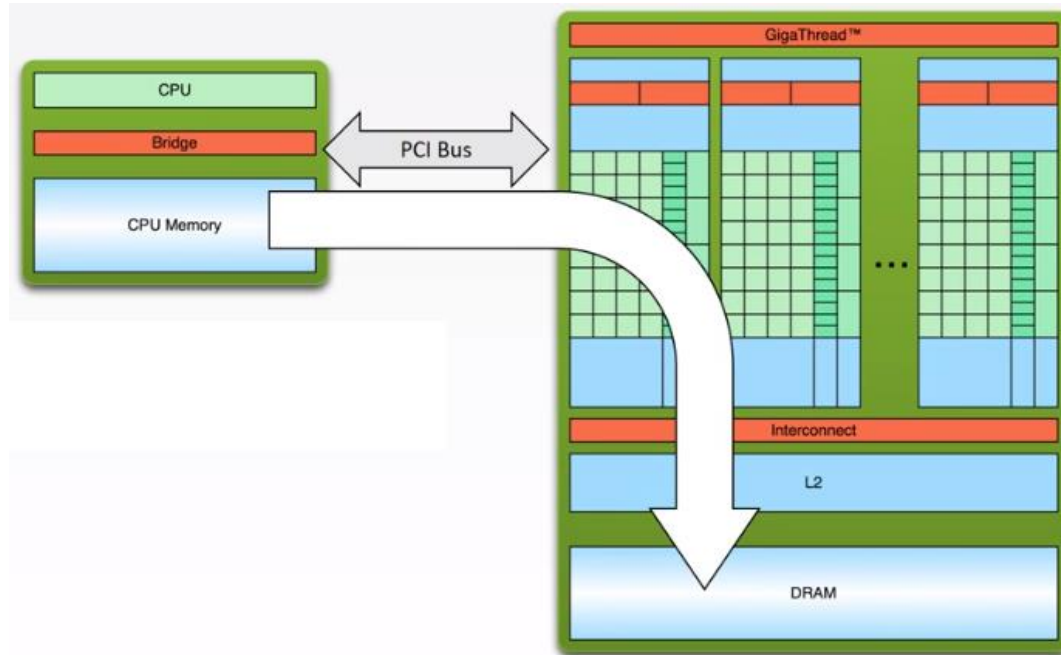GPU that has it own memory ("device memory")

NVIDIA K80

**D&LL**EMC

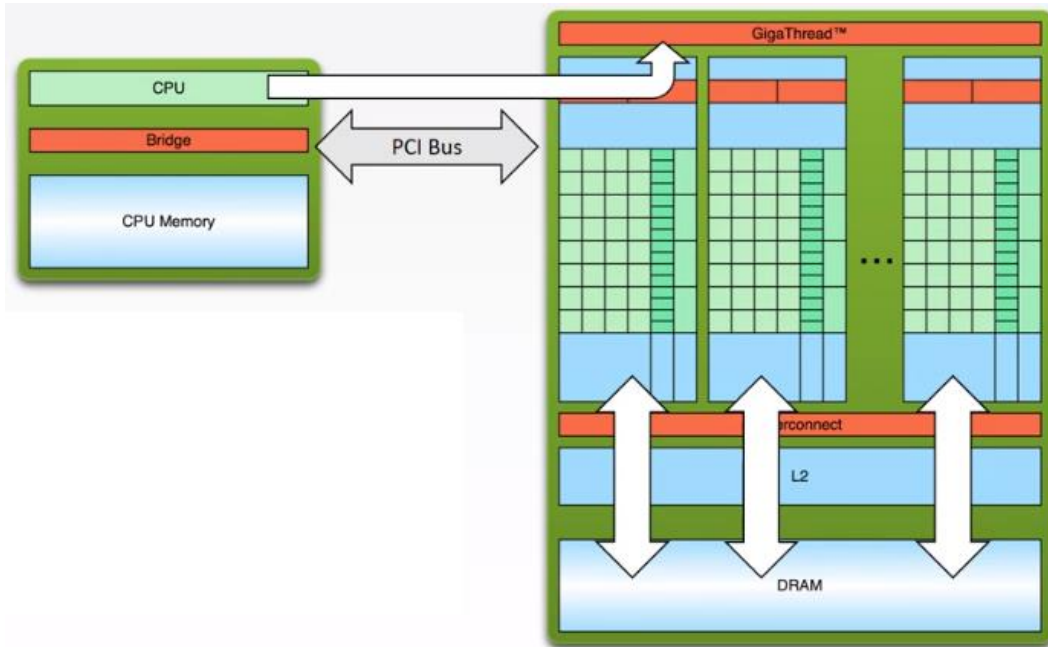# Putting these two together

How GPU acceleration works

# How GPU acceleration works

Copy input data from CPU memory to GPU memory
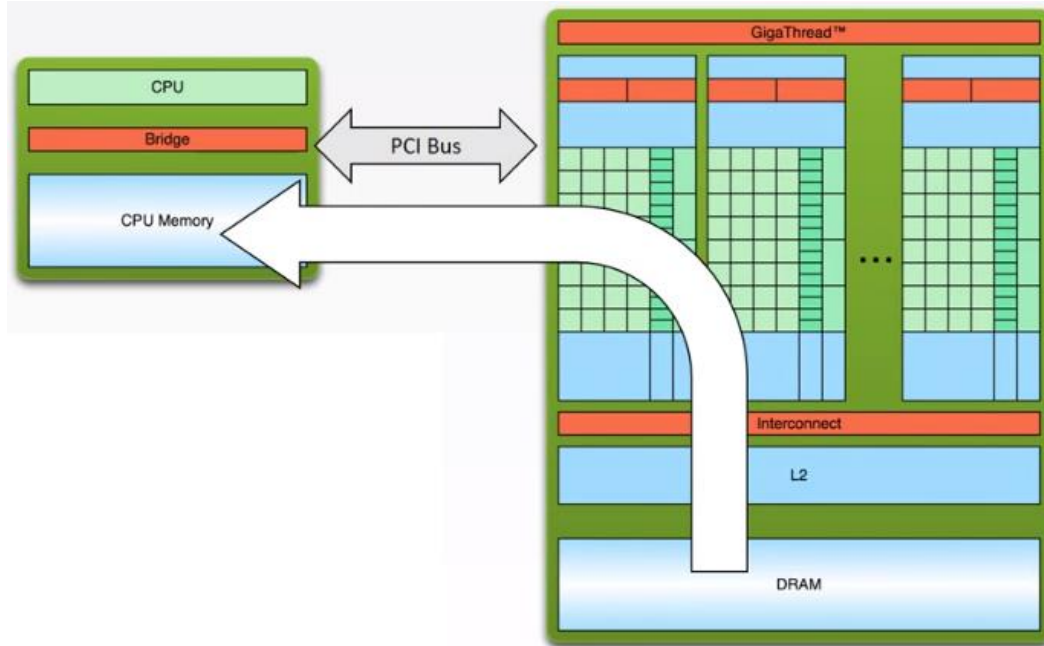
**DELL**EMC

# How GPU acceleration works

Load GPU program and execute, caching data on chip for performance

DELLEMC

# How GPU acceleration works

Copy results from GPU memory to CPU memory

**DELL**EMC

# CUDA Toolkits

…and over to Heiko

**DELL**EMC