



Computational Science on Many-Core Architectures

360.252

Karl Rupp



Institute for Microelectronics
Vienna University of Technology
<http://www.iue.tuwien.ac.at>



Zoom Channel 95028746244
Wednesday, October 14, 2020

Introducing Myself

Current Positions

- Postdoctoral Researcher at I μ E
- Co-Founder and Managing Director at BrickXter GmbH

Professional Interests

- Efficient computation on modern hardware
- Semiconductor device simulation
- Circuit design
- Making technology useful for “the average Joe”

Sideline Activities

- PETSc developer (<https://www.mcs.anl.gov/petsc/>)
- ViennaCL developer (<http://viennacl.sourceforge.net/>)

Subjects

- Amdahl's Law
- FLOPs, Bandwidth, and Latency
- Performance Modeling
- Graphics Processing Units (SIMT processing, thread block synchronization)
- Programming Models (Annotation-driven such as OpenMP, native such as CUDA)
- Field Programmable Gate Arrays
- Emerging Many-Core Architectures

Course Objectives

Main Objective

Maximize students' useful knowledge on using many-core architectures within the available time

Modalities

- 11 lectures (45 - 60 minutes each)
- 10 exercises (DIY-approach)
- Slides and exercise material:
<https://owncloud.tuwien.ac.at/index.php/s/xjvjXDl077CC4zv>

Outcome

- Hands-on experience
- You will create some of the fastest GPU kernels in the world

Inspiration

Tell me and I forget.
Teach me and I remember.
Involve me and I learn.

Benjamin Franklin

Related CSE Courses

Term 1

- 360.242 - Numerical Simulation and Scientific Computing I
- 101.826 - Numerical Computation

Term 2

- 184.726 - Advanced Multiprocessor Programming
- 101.773 - Numerical Methods for PDEs

Part 1: Hands-On Exercises

- approx. 100 points over 10 exercises (excl. bonus points)
- approx. 40 percent of overall grade
- minimum of 50 percent of total points

Part 2: Oral exam

- oral exam (most likely virtual because of COVID-19)
- approx. 60 percent of overall grade
- “Fail” on oral exam means “Fail” on course

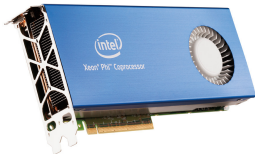
Introduction

Many-Core Architectures

- High FLOP/Watt ratio
- High memory bandwidth
- (Usually) Attached via PCI-Express



AMD FirePro W9100
320 GB/sec



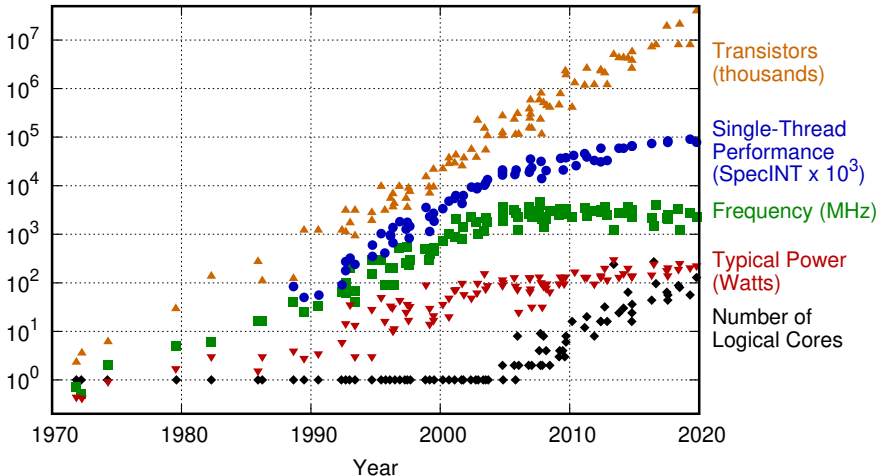
INTEL Xeon Phi
320 (220?) GB/sec



NVIDIA Tesla K20
250 (208) GB/sec

Introduction

48 Years of Microprocessor Trend Data

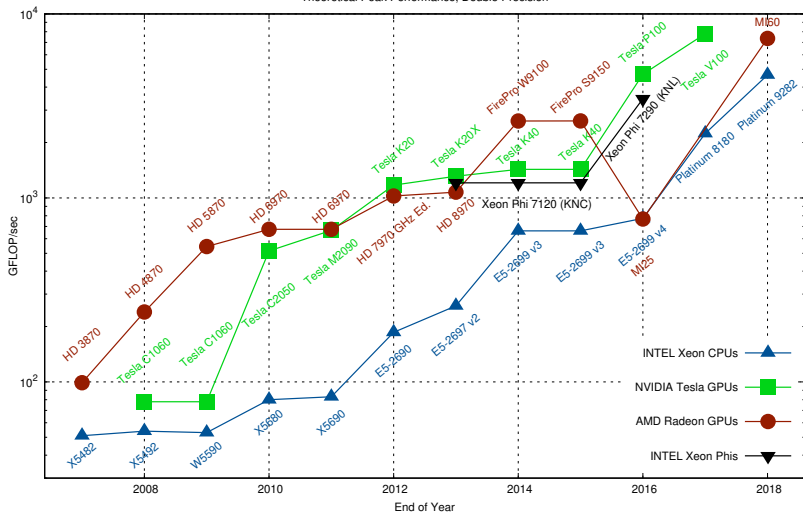


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2019 by K. Rupp

Introduction

Theoretical Peak Performance

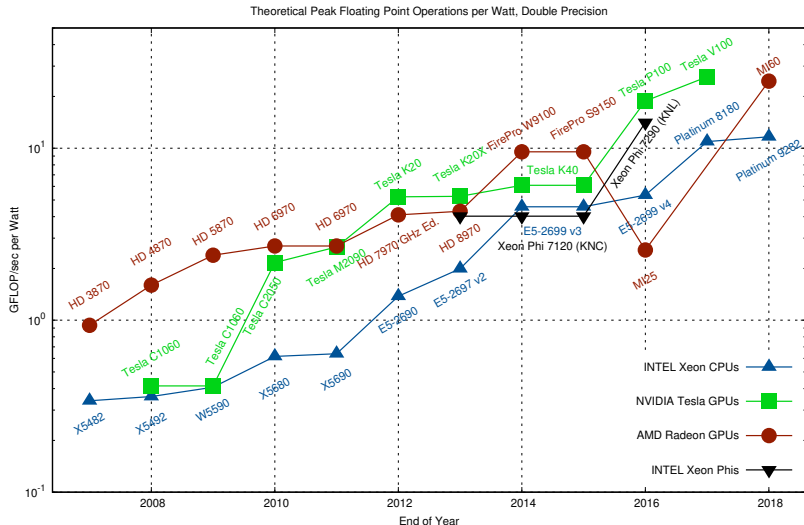
Theoretical Peak Performance, Double Precision



<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Introduction

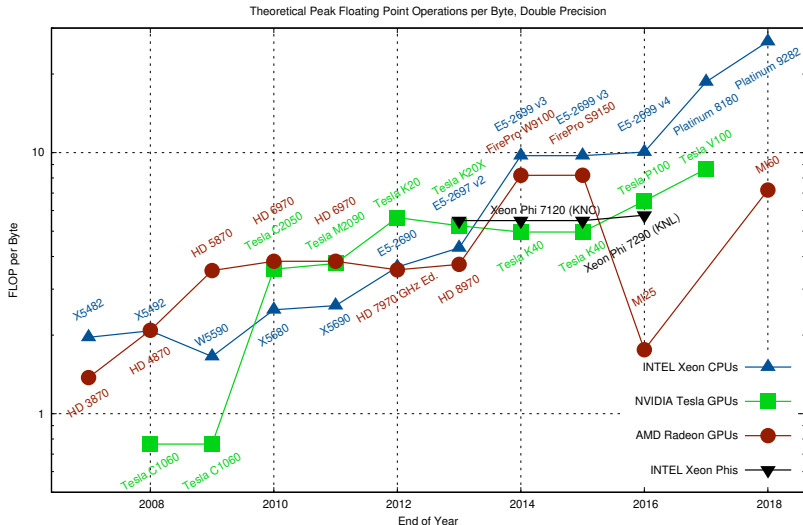
Theoretical Peak Performance per Watt



<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Introduction

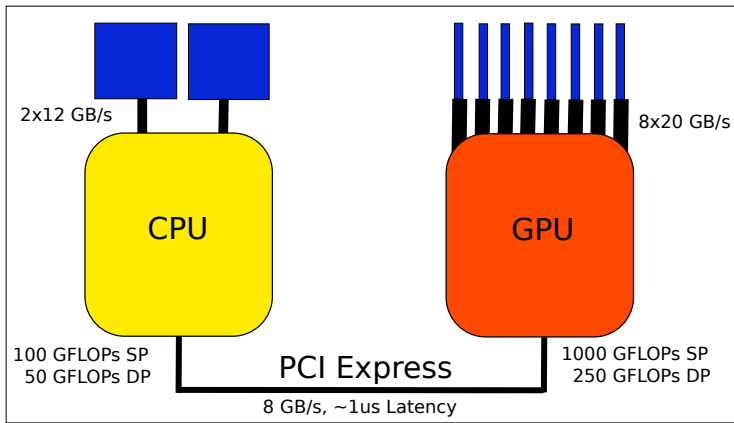
Theoretical Peak Performance (FLOPs) per Byte of Memory Bandwidth



<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

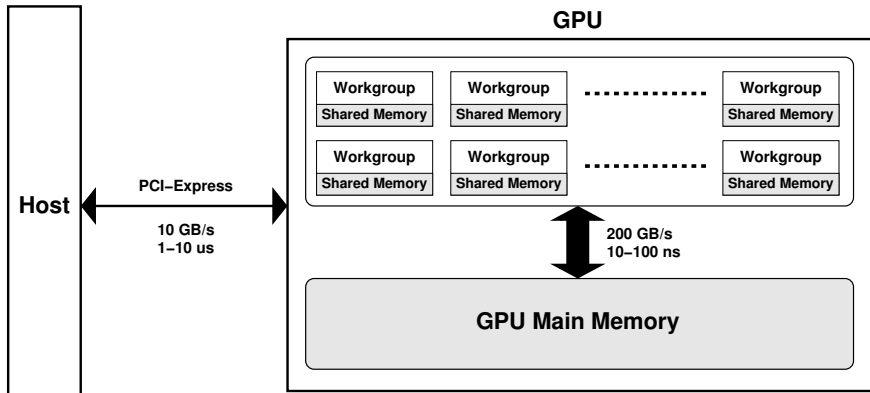
GPU Overview

Computing Architecture Schematic



- Good for large FLOP-intensive tasks, high memory bandwidth
- PCI-Express can be a bottleneck
- \gg 10-fold speedups (usually) not backed by hardware

GPU Overview



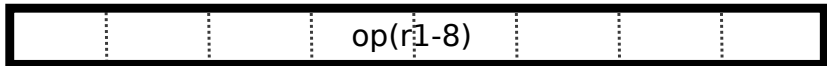
Details

- Workgroups consist of 32-64 hardware threads
- Up to 24 hardware workgroups
- Shared memory small: approx. 32-64 KB

GPU Overview

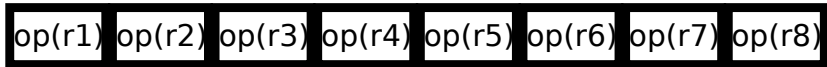
Reminder: AVX

- One instruction for all elements of a vector register



Single Instruction Multiple Threads (SIMT)

- One instruction for all threads in workgroup
- Each thread has separate registers
- Efficient if all threads execute the same instruction

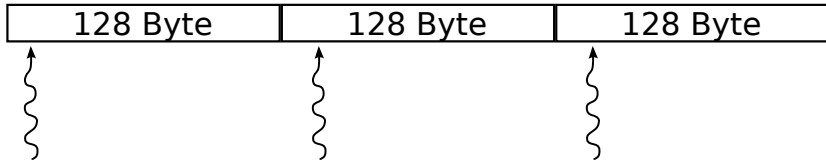


GDDR5

- Optimized for throughput
- Channel width: multiple of 32 bits
- High bus width: 256 bits, 384 bits

Structured Memory Access

- Memory controllers use 32/64/128 byte transactions
- Partial transactions degrade effective bandwidth



GPU Overview

Host-Device Communication

- PCI-Express v2: 8 GB/sec max
- PCI-Express v3: 16 GB/sec max
- Latency: about 10 μ s

