



INTRODUCTION TO PARALLEL AND GPU COMPUTING

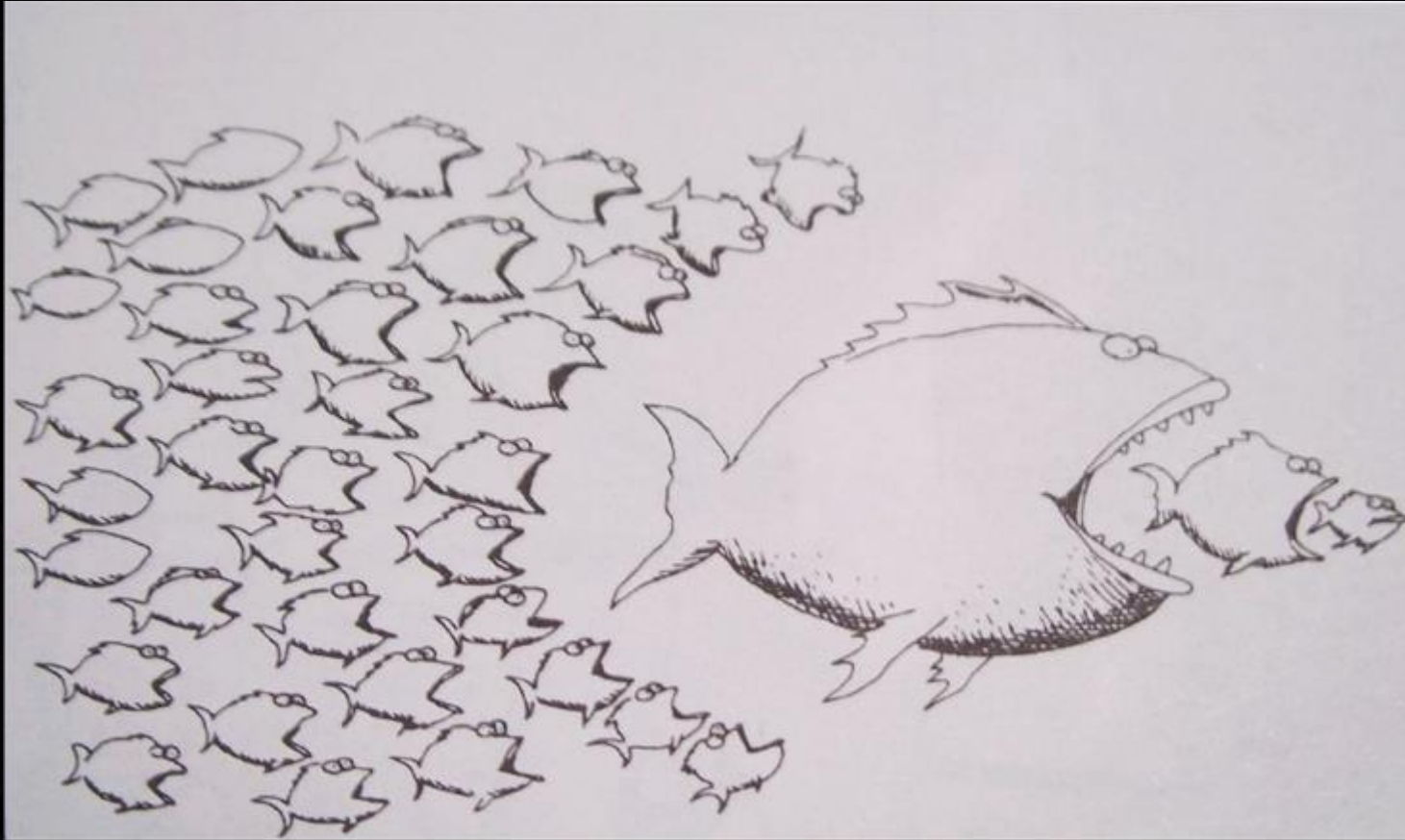
Carlo Nardone

Sr. Solution Architect, NVIDIA EMEA

AGENDA

- 1 Intro
- 2 Processors Trends
- 3 Multiprocessors
- 4 GPU Computing

PARALLEL COMPUTING, ILLUSTRATED

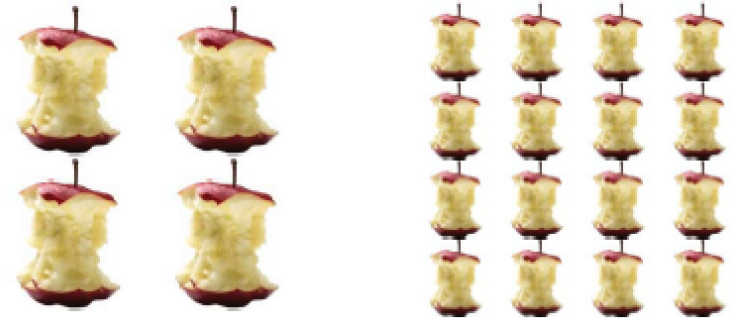


FROM MULTICORE TO MANYCORE

Past



Present



Future



PARALLEL COMPUTING IS OLD



Cray 1 (1976)

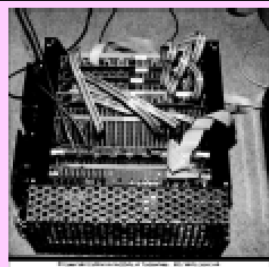


Cray 2 (1985)



Cray C-90 (1991)

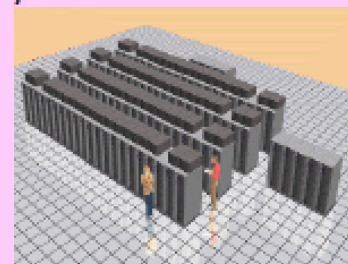
Vector Computers



Cosmic cube (1983)



Paragon (1993)



ASCI Red (1997)

Massively Parallel Processors (MPP)



Clusters (late 80's)

Cluster Computers

Linux PC Clusters (~1995)

Source: Tim Mattson
Late 70's

Late 80's

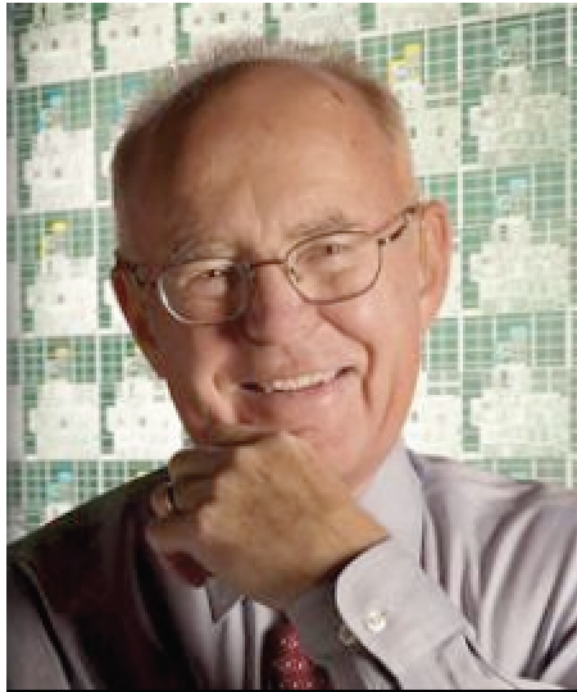
Late 90's

15

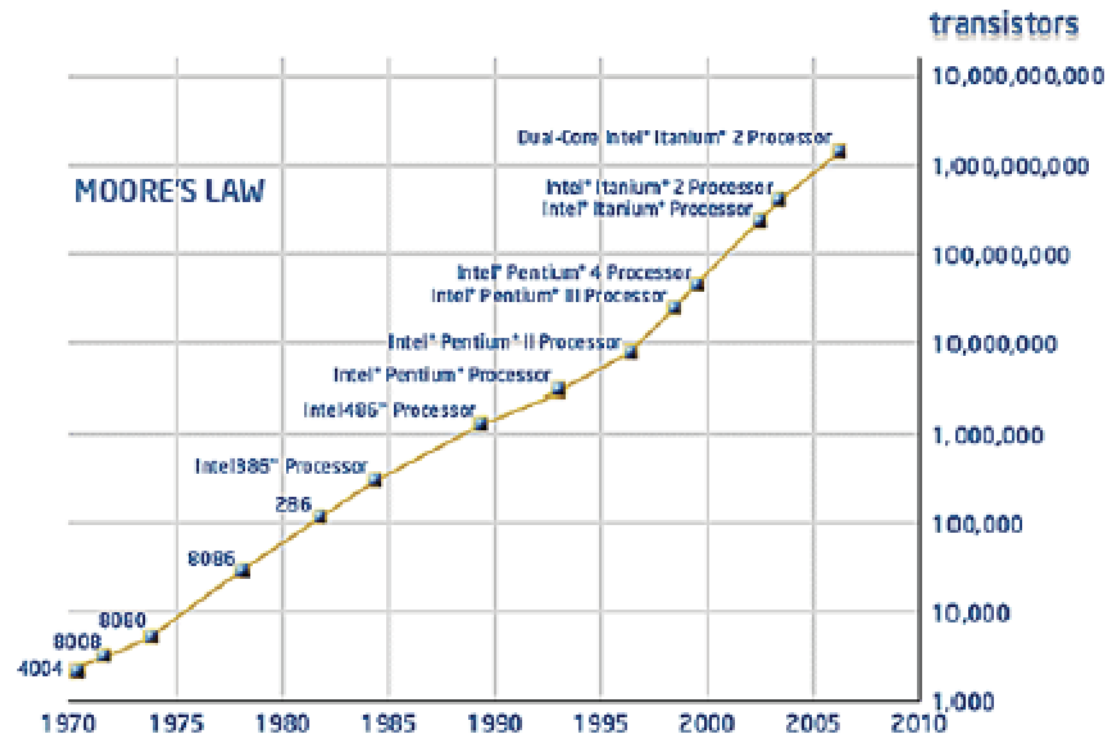
AGENDA

- 1 Intro
- 2 Processors Trends
- 3 Multiprocessors
- 4 GPU Computing

MOORE'S LAW



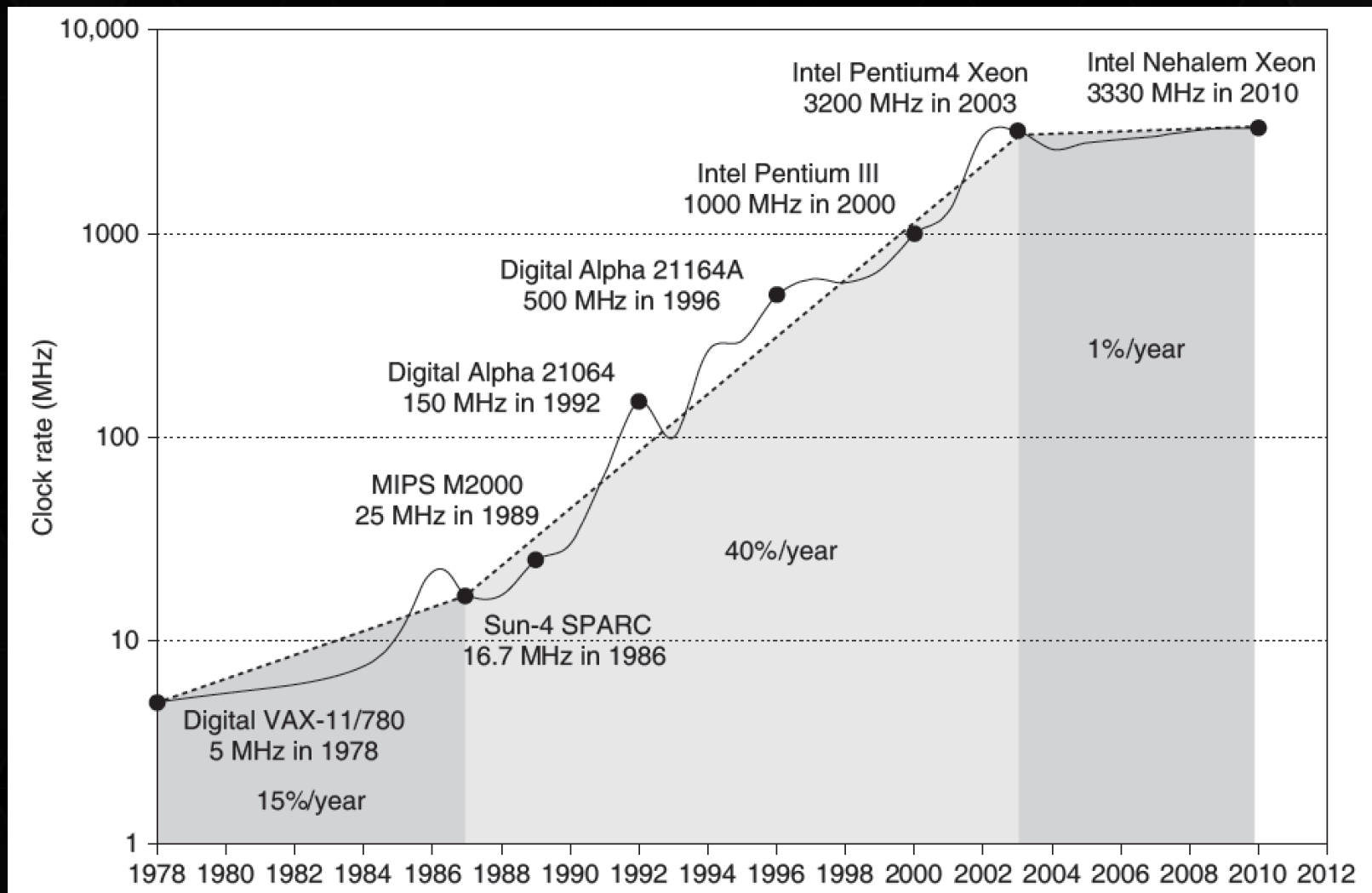
Gordon Moore (Intel co-founder) predicted in 1965 that transistor density of semiconductor chips would double every 18 months.



Moore's Law

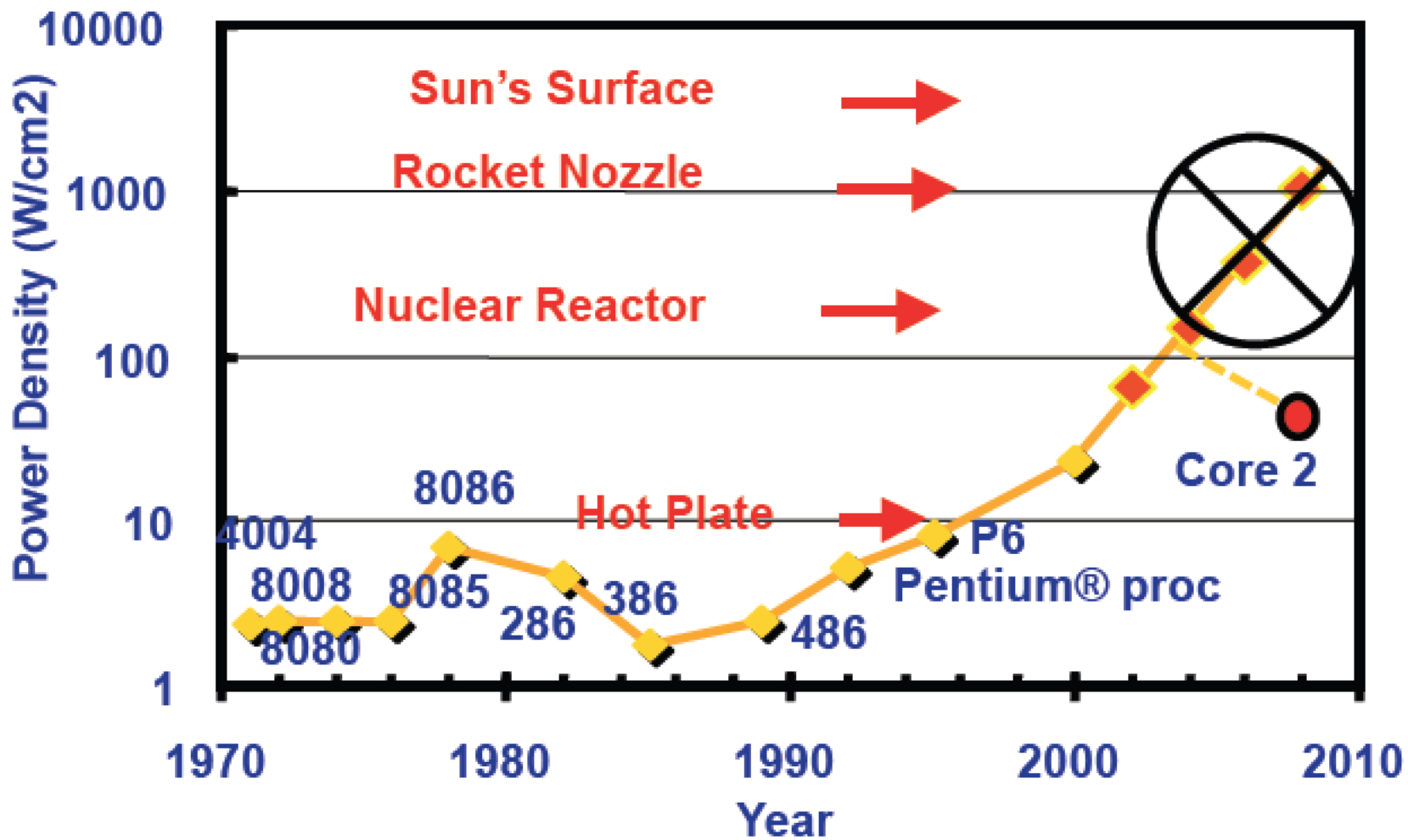
Source: U.Delaware CISC879 / J.Dongarra

MICROPROCESSOR FREQUENCY TREND



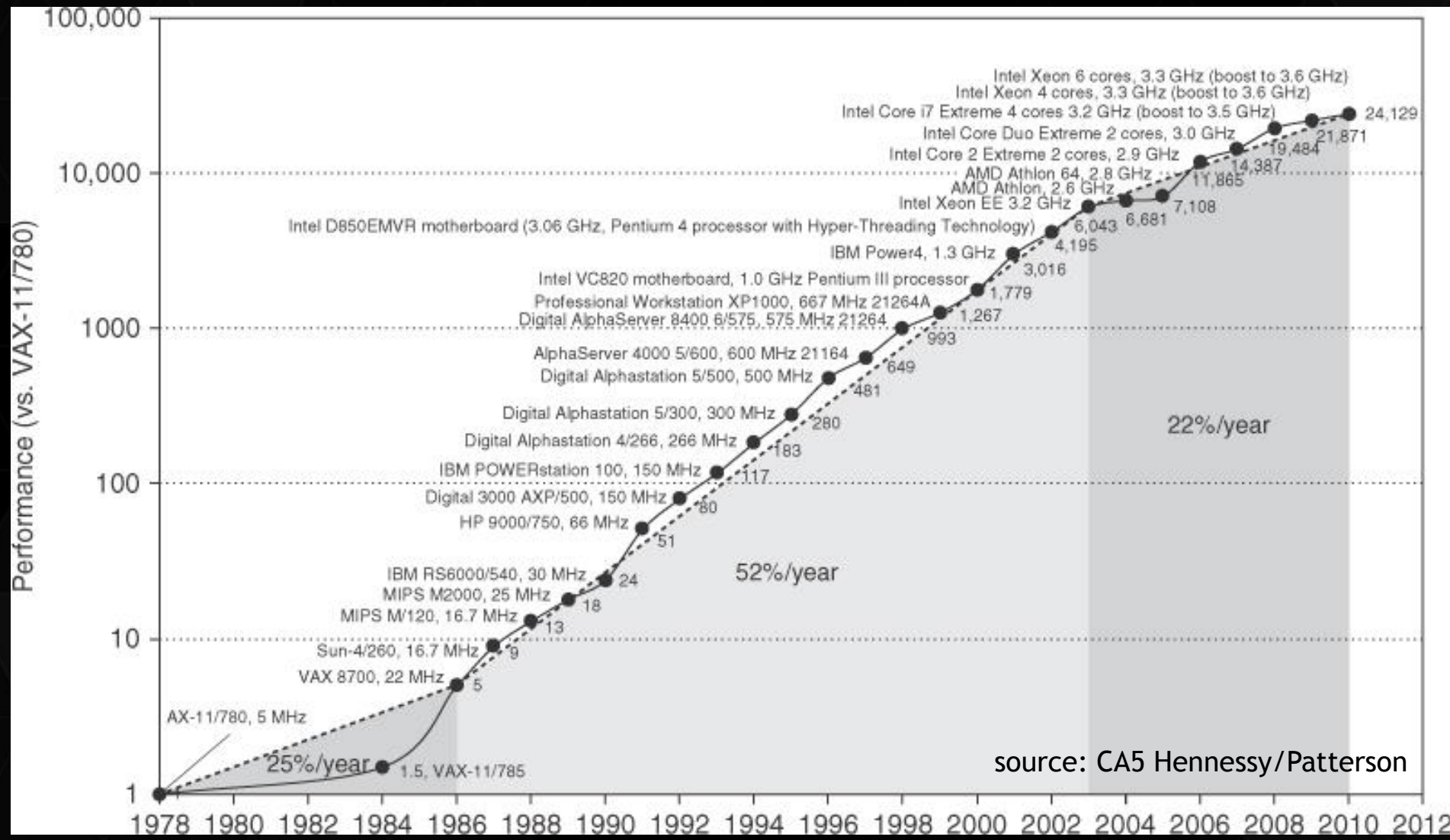
source: CA5
Hennessy/Patterson

THE POWER WALL

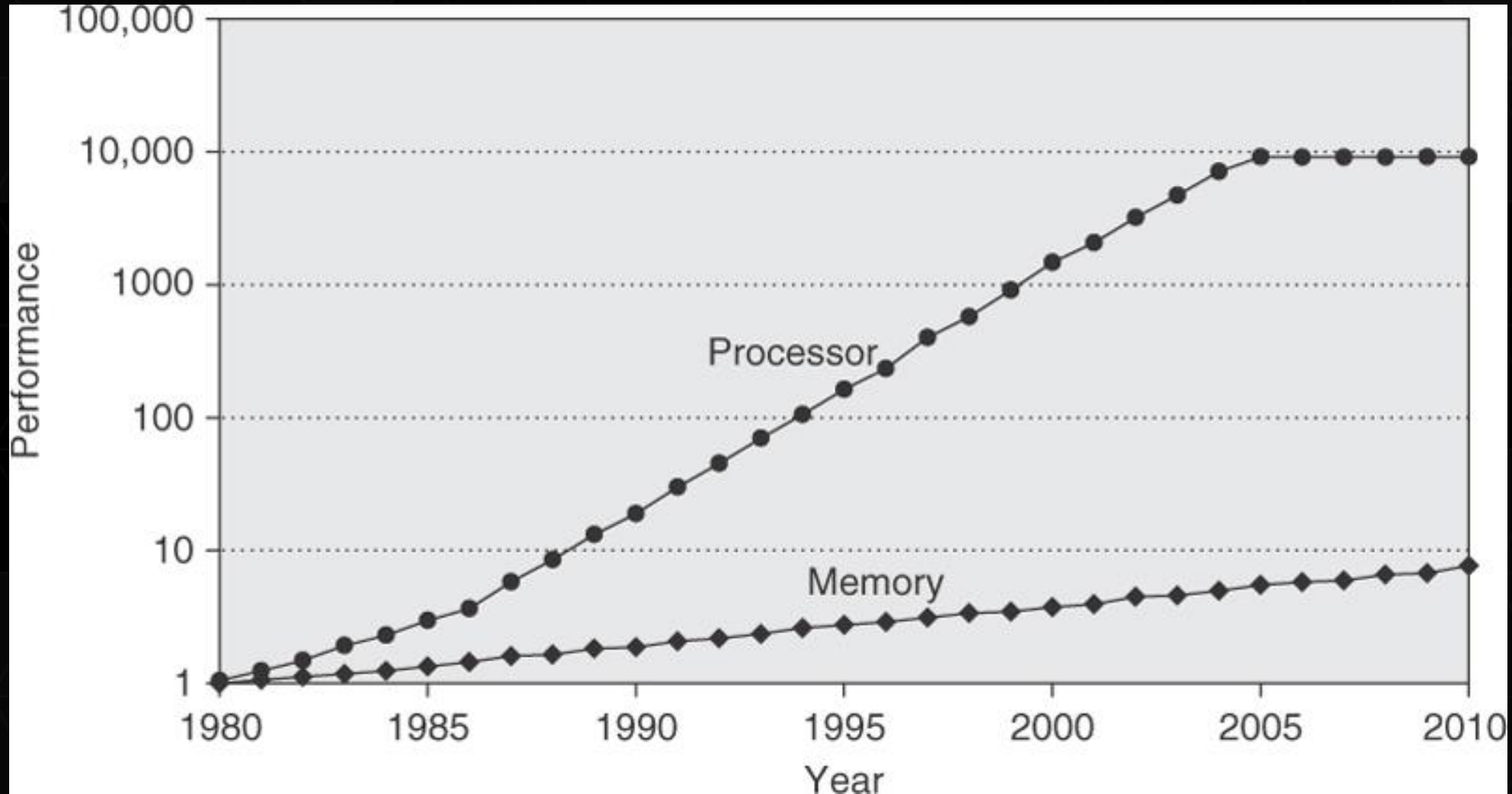


Source: S. Borkar (Intel)

MICROPROCESSORS PERFORMANCE TREND

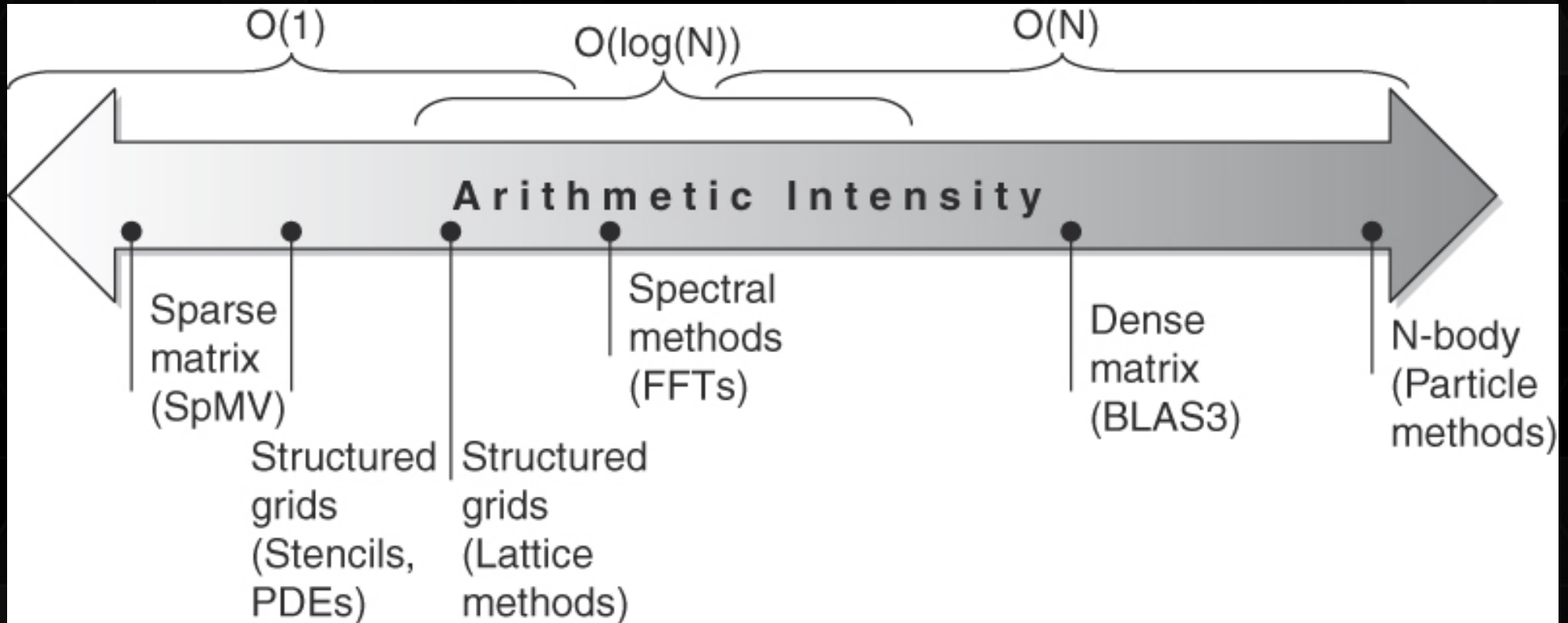


THE PROCESSOR-MEMORY GAP



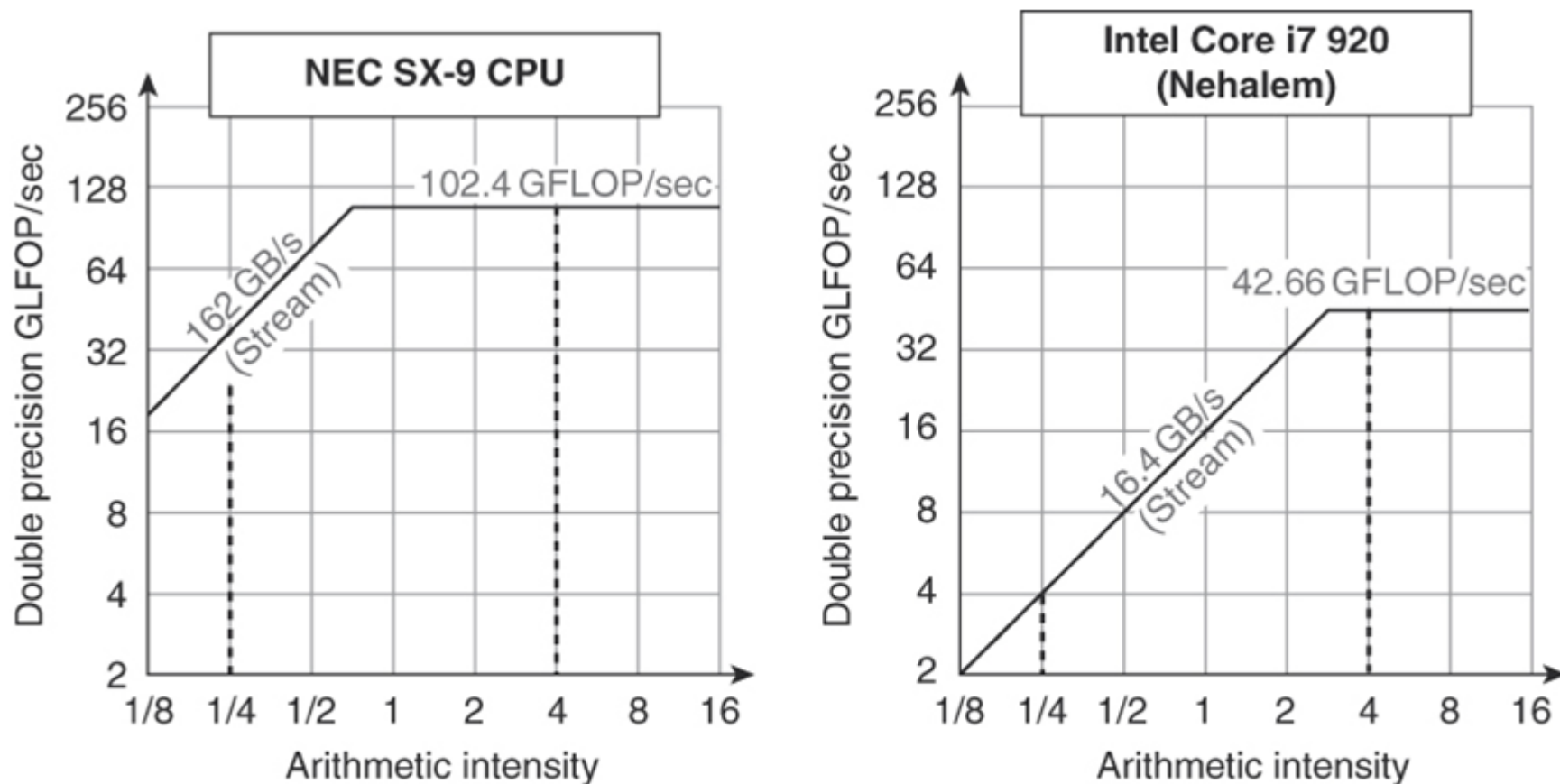
source: CA5
Hennessy &
Patterson

ARITHMETIC INTENSITY



Arithmetic intensity, specified as the number of floating-point operations to run the program divided by the number of bytes accessed in main memory [Williams et al. 2009]. Some kernels have an arithmetic intensity that scales with problem size, such as dense matrix, but there are many kernels with arithmetic intensities independent of problem size. Source: CA5 Hennessy & Patterson

BANDWIDTH AND THE ROOFLINE MODEL



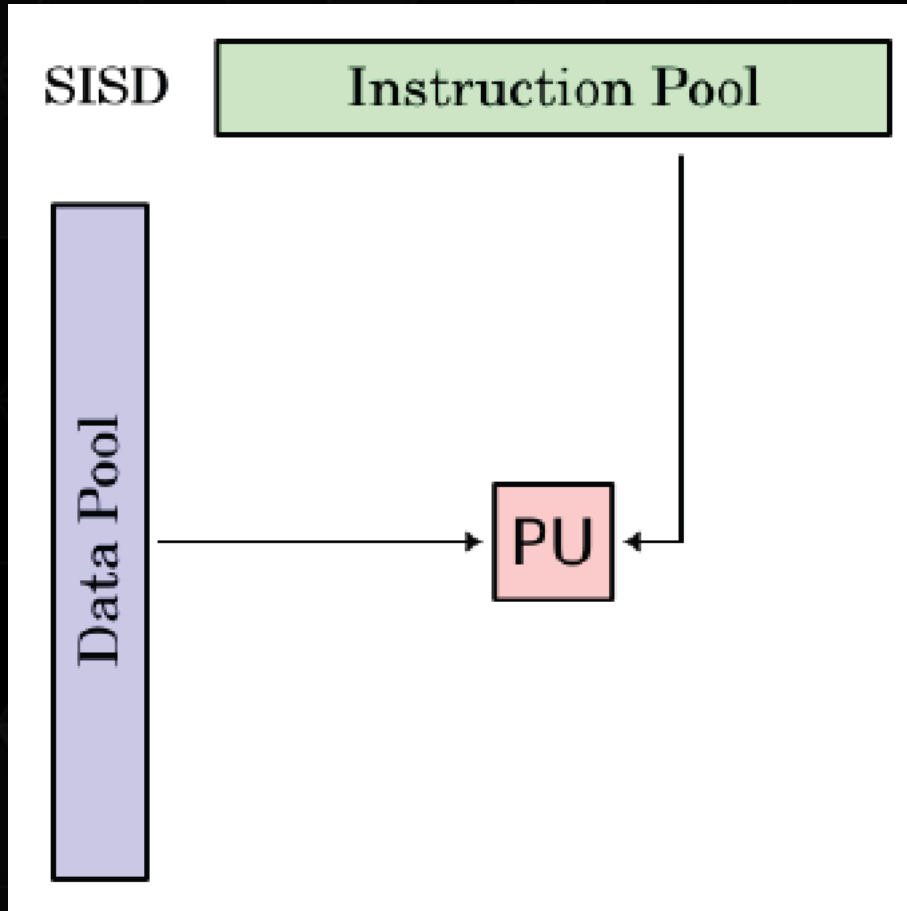
Source: HPCA5

Figure 4.11 Roofline model for one NEC SX-9 vector processor on the left and the Intel Core i7 920 multicore computer with SIMD Extensions on the right [Williams et al. 2009]. This Roofline is for unit-stride memory accesses and double-precision floating-point performance. .

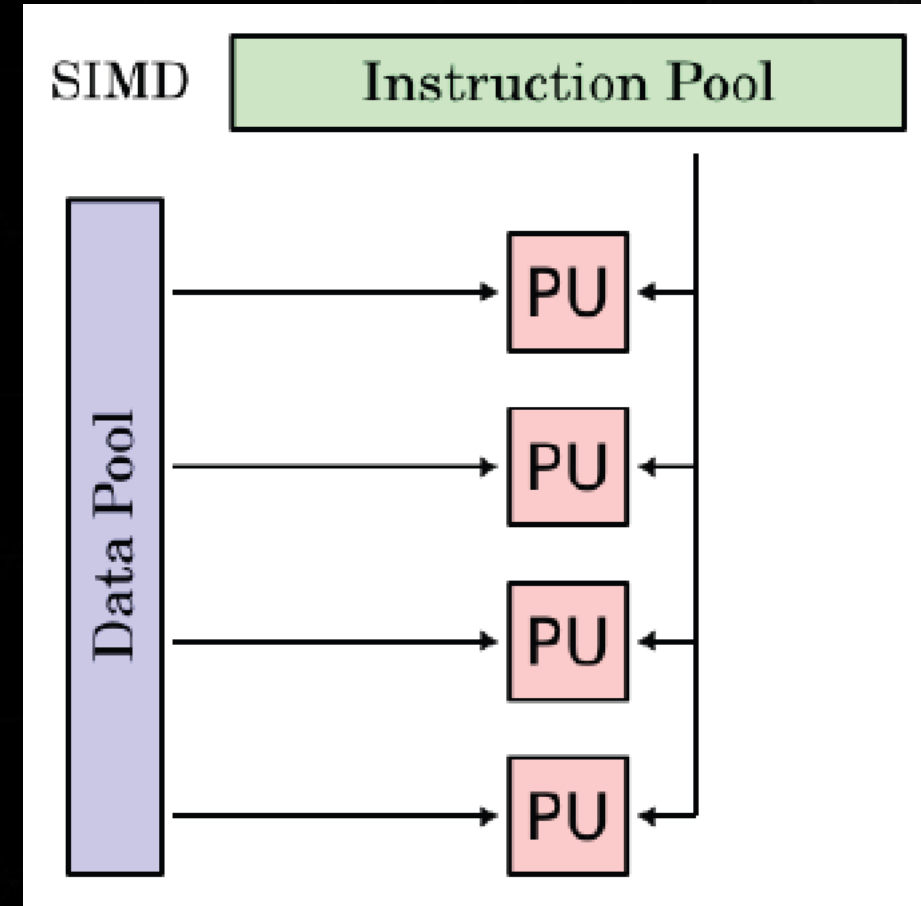
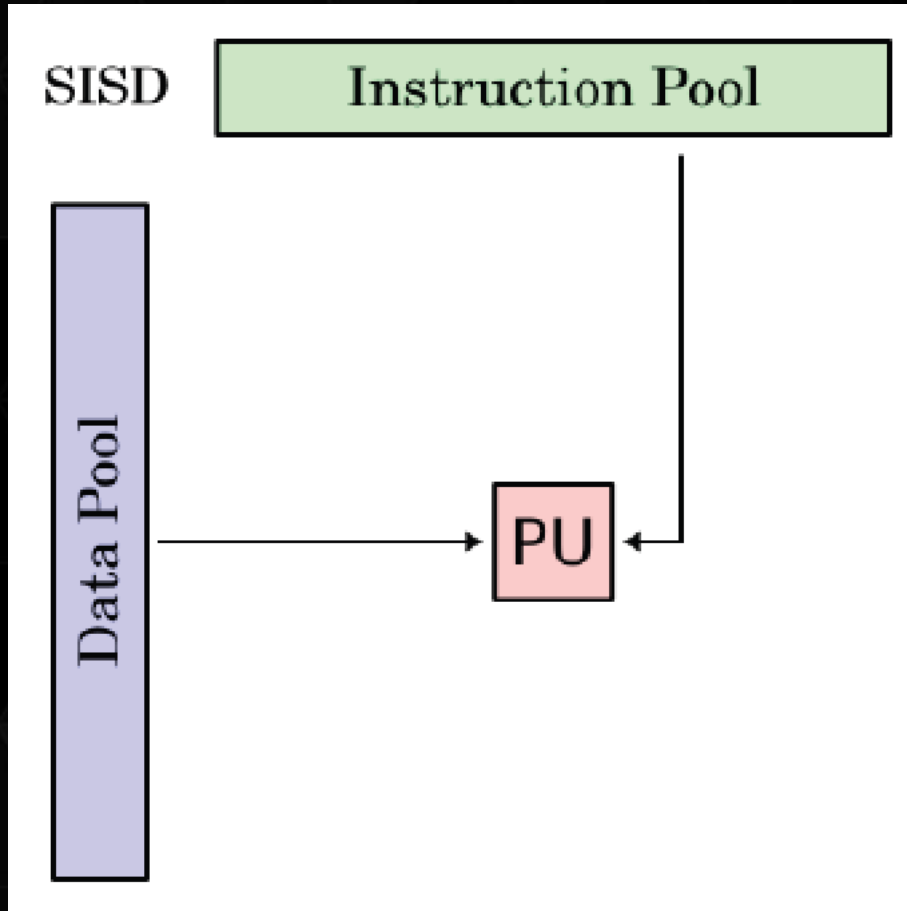
AGENDA

- 1 Intro
- 2 Processors Trends
- 3 **Multiprocessors**
- 4 GPU Computing

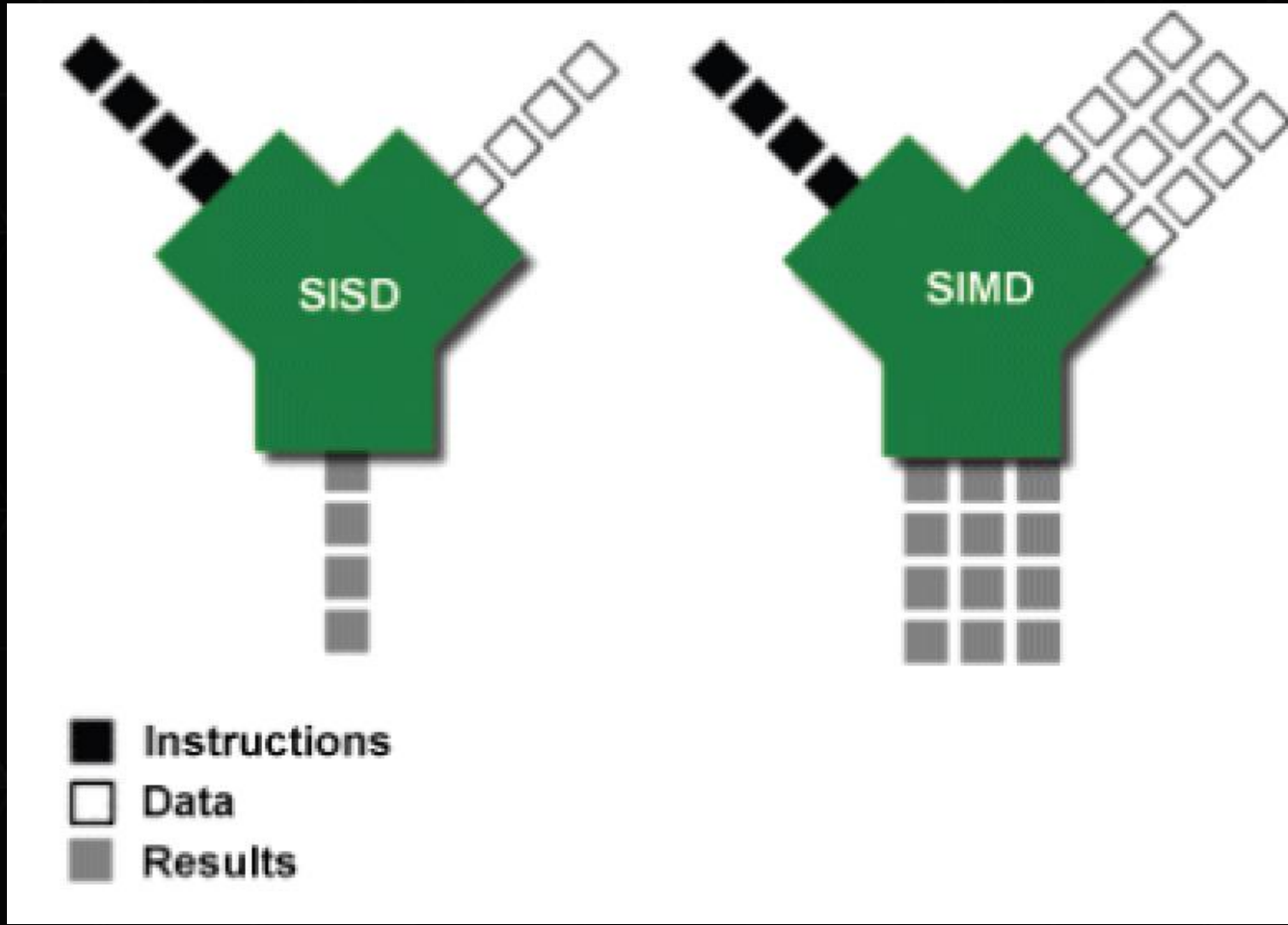
FLYNN'S TAXONOMY: SISD



FLYNN'S TAXONOMY: SIMD

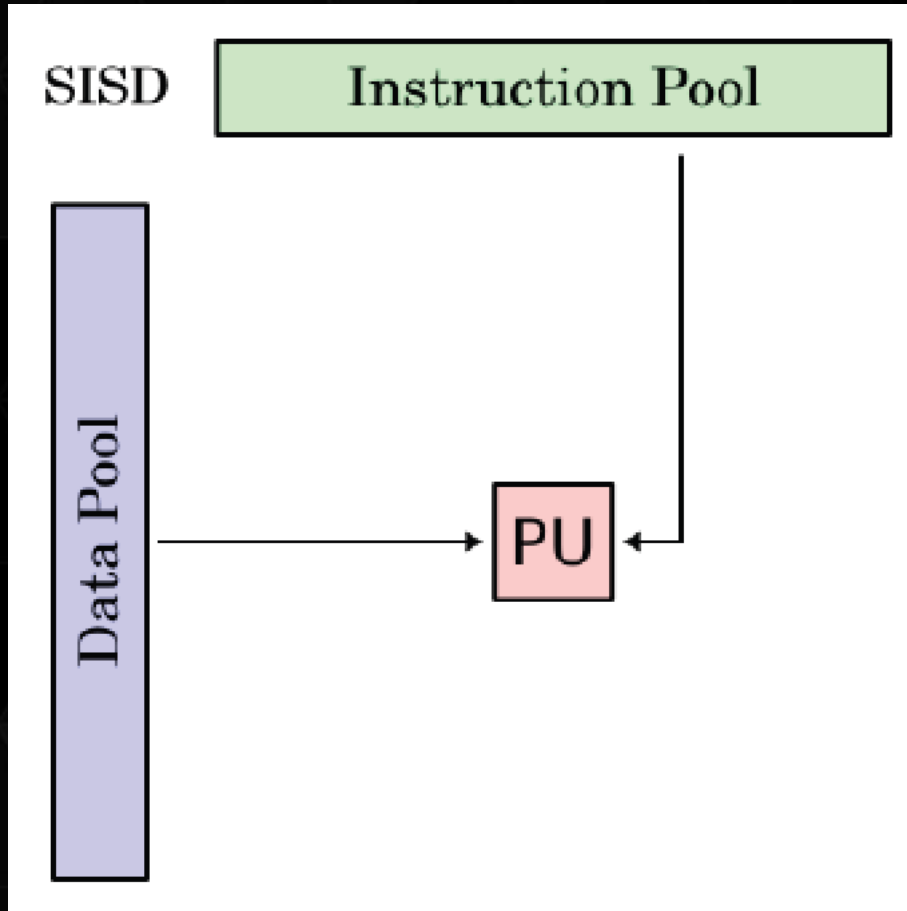


FLYNN'S TAXONOMY: SIMD (2)

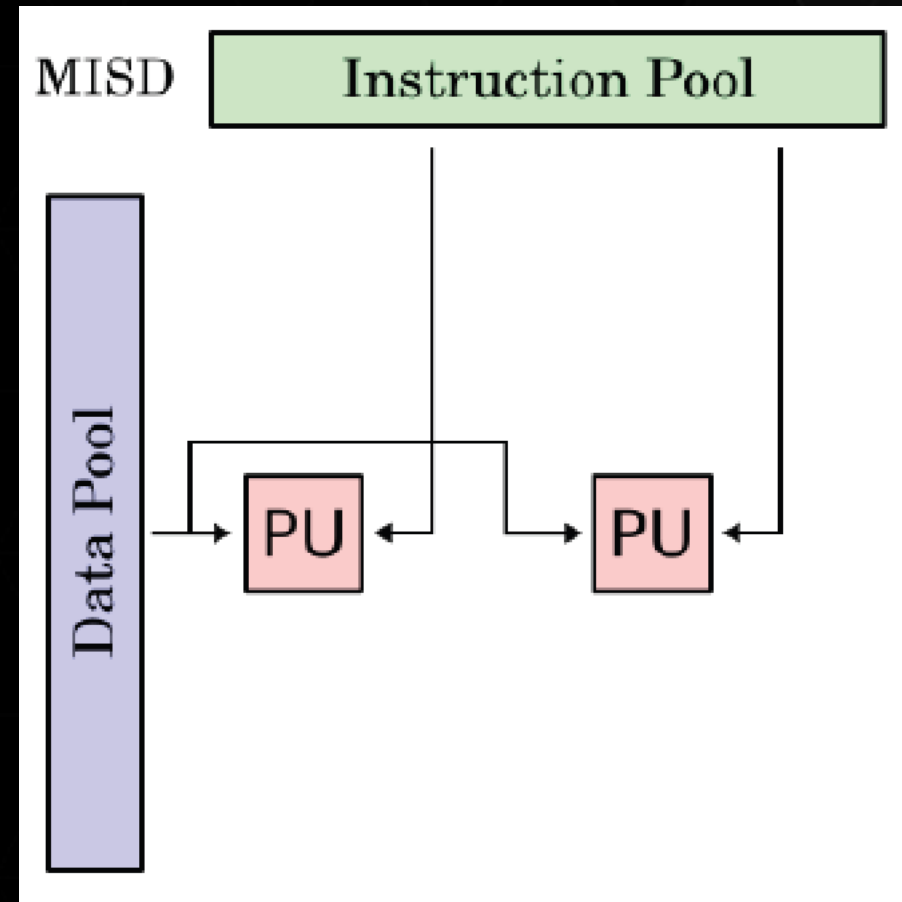
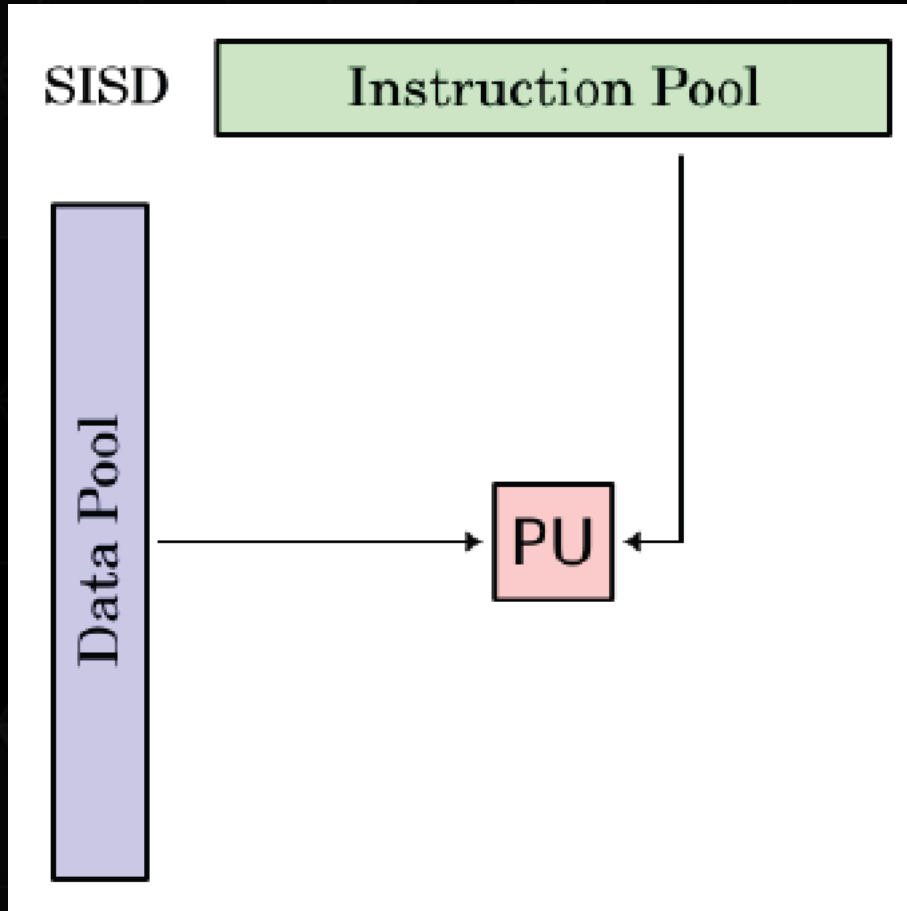


source: CA5
Hennessy &
Patterson

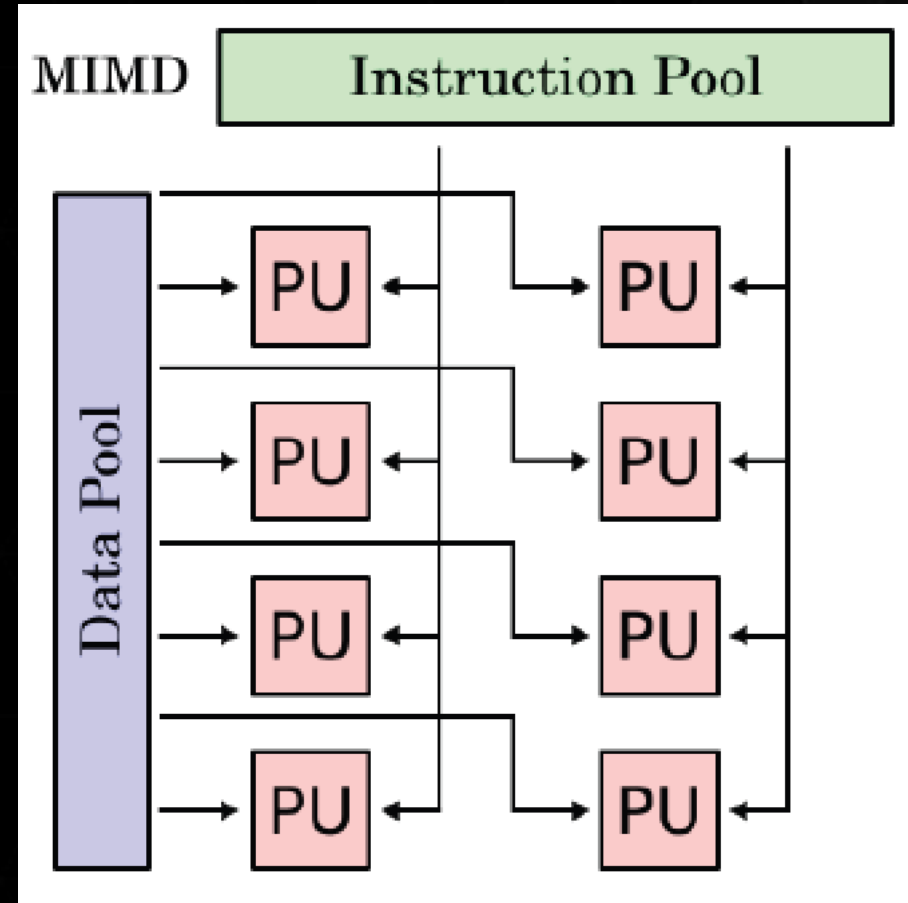
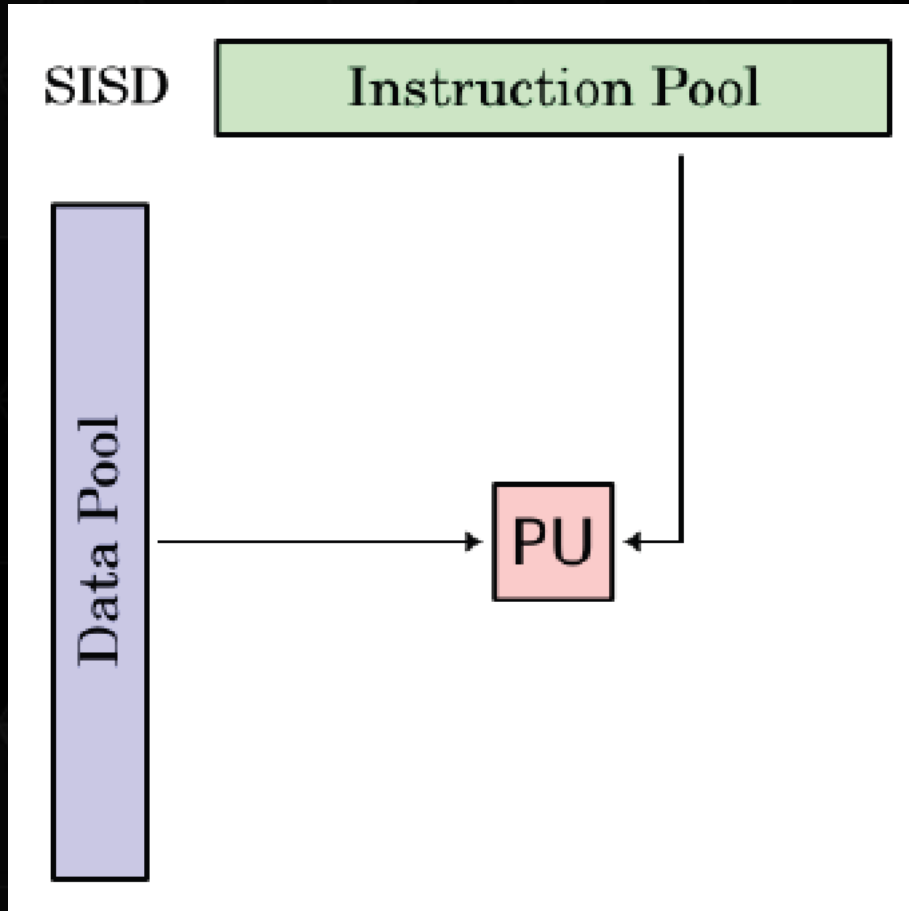
FLYNN'S TAXONOMY



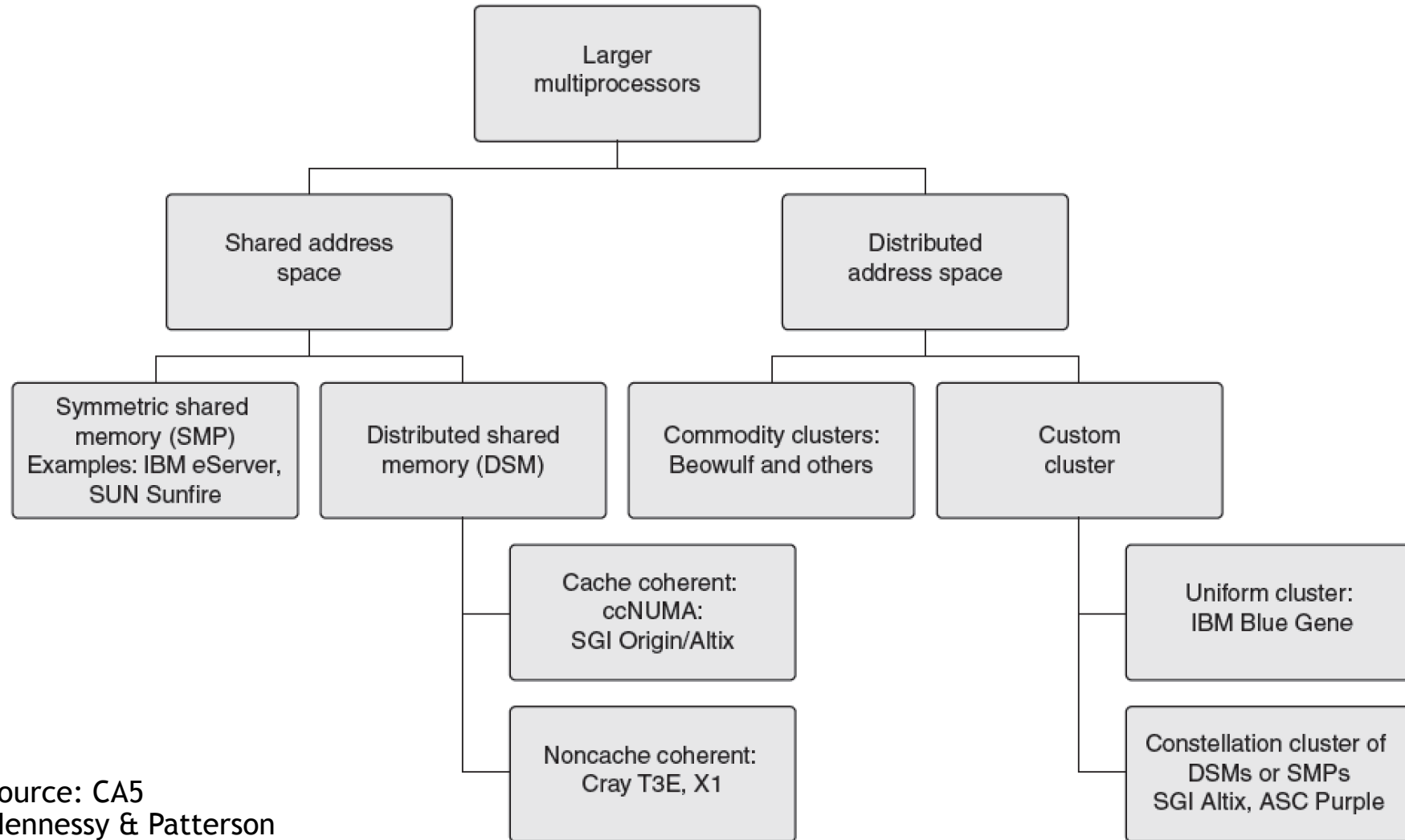
FLYNN'S TAXONOMY: MISD (?)



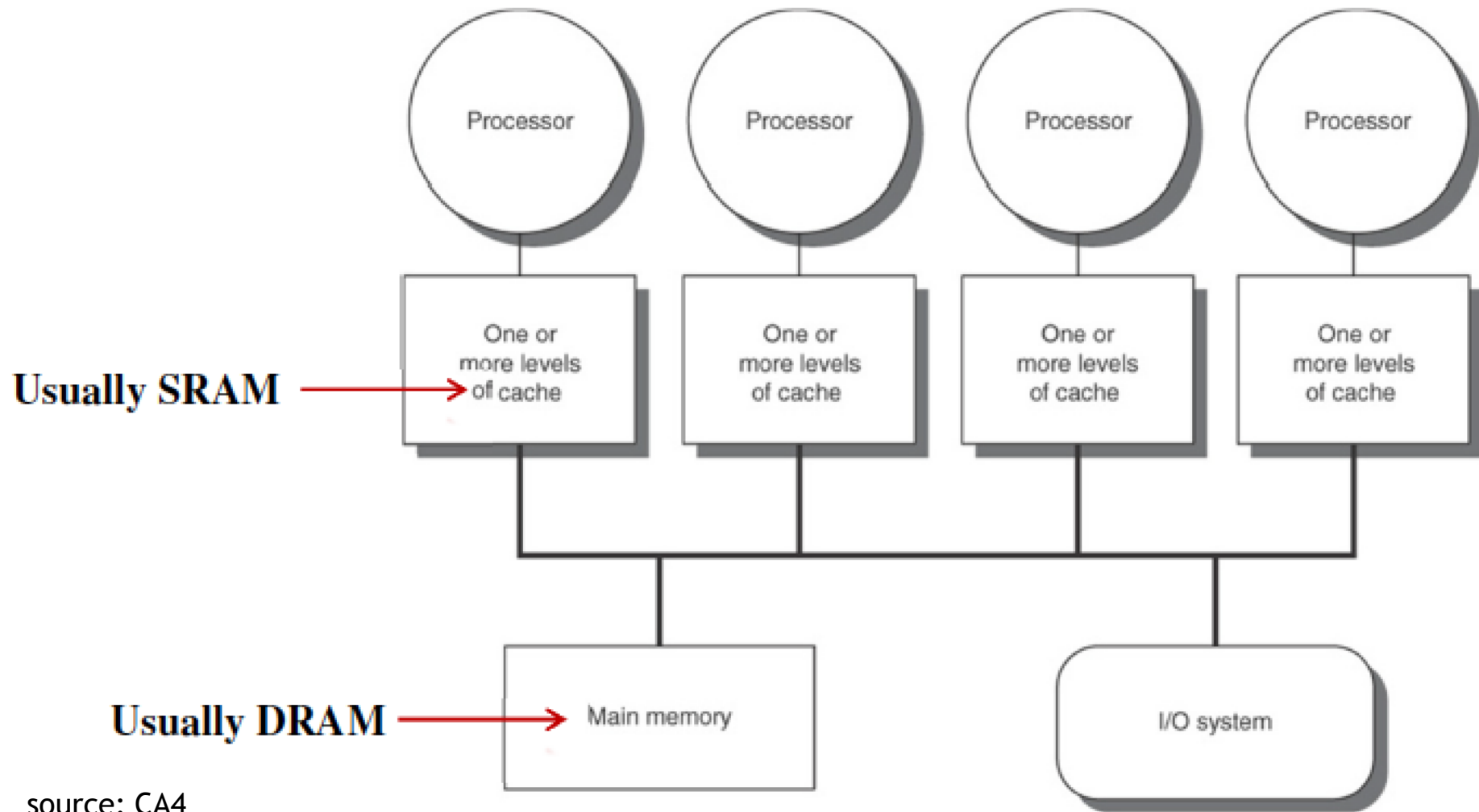
FLYNN'S TAXONOMY: MIMD



MULTIPROCESSORS

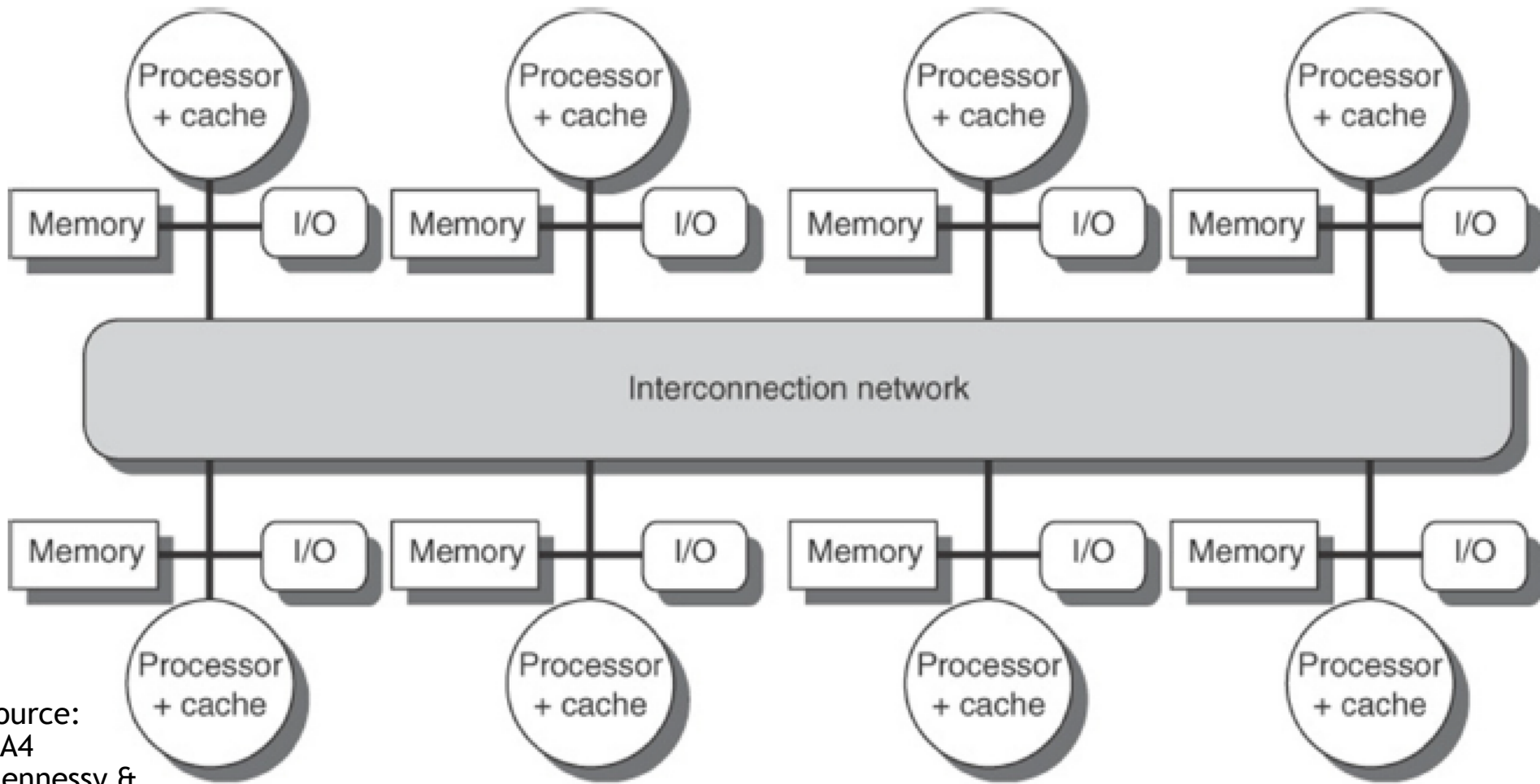


MULTIPROCESSOR SMP CLASS



source: CA4
Hennessy & Patterson

MULTIPROCESSOR DISTRIBUTED MEM CLASS

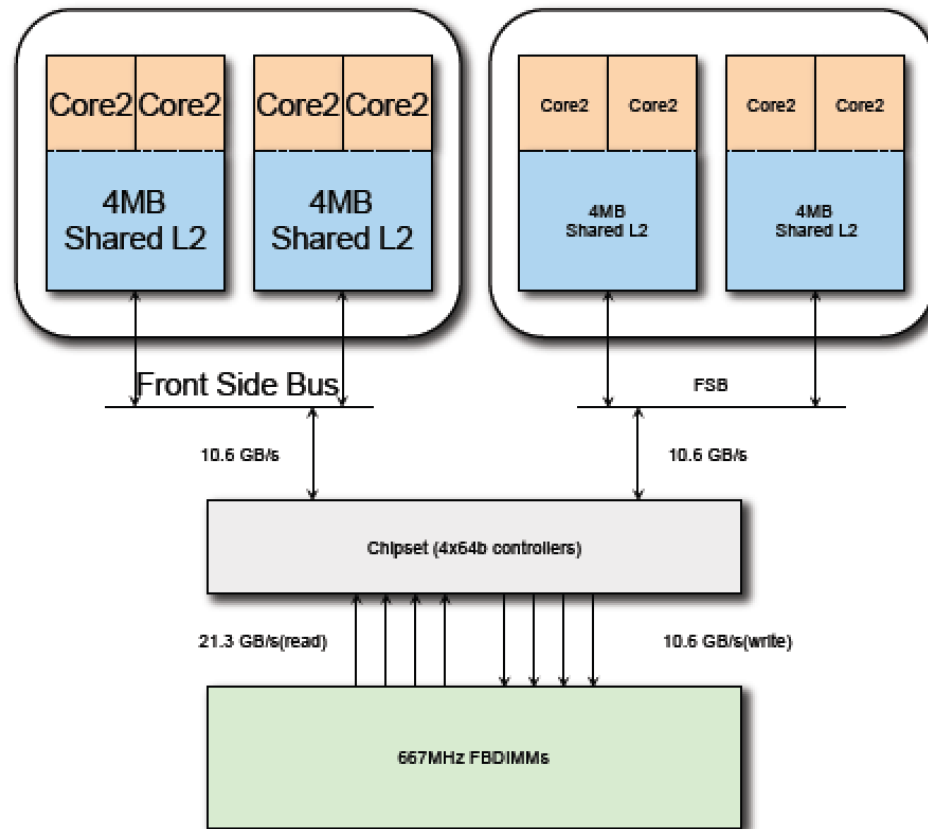


source:
CA4
Hennessy &
Patterson

© 2007 Elsevier, Inc. All rights reserved.

X86 ARCHITECTURE

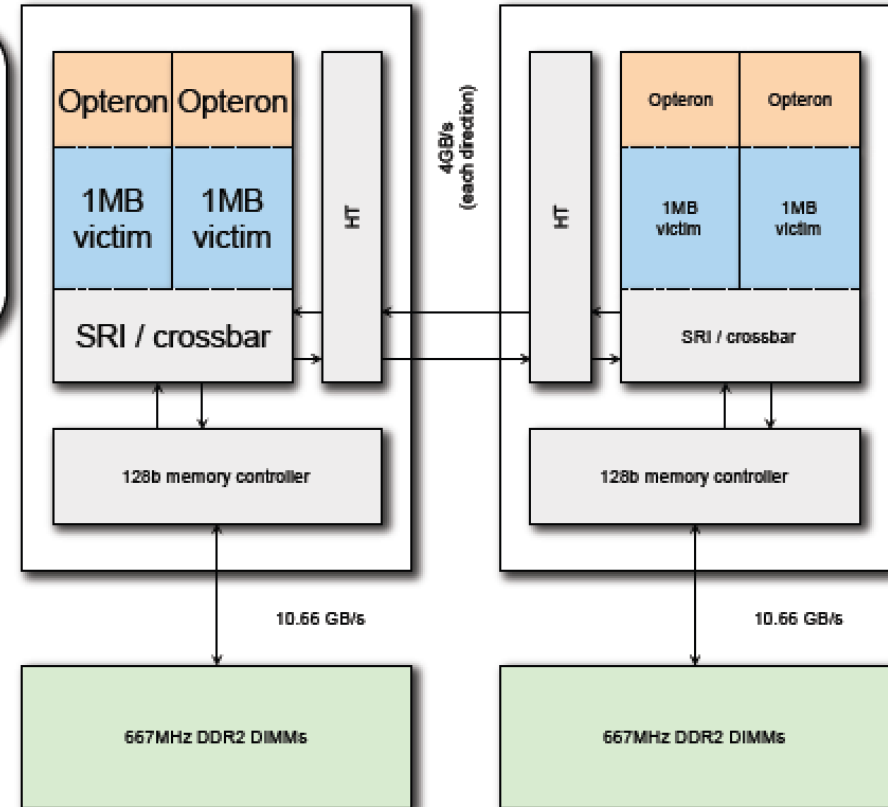
Intel Clovertown



Uniform Memory Access

Intel Nehalem

AMD Opteron



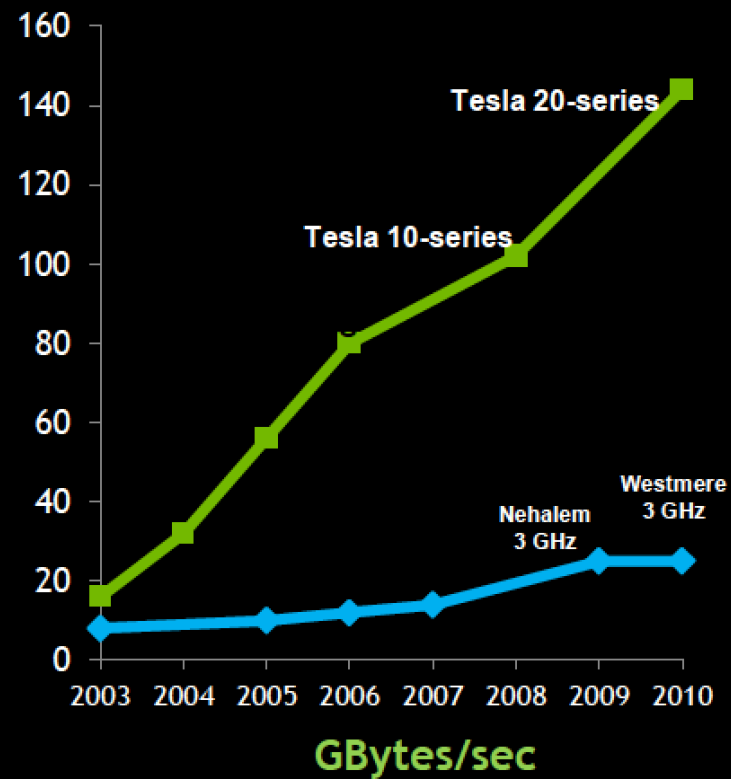
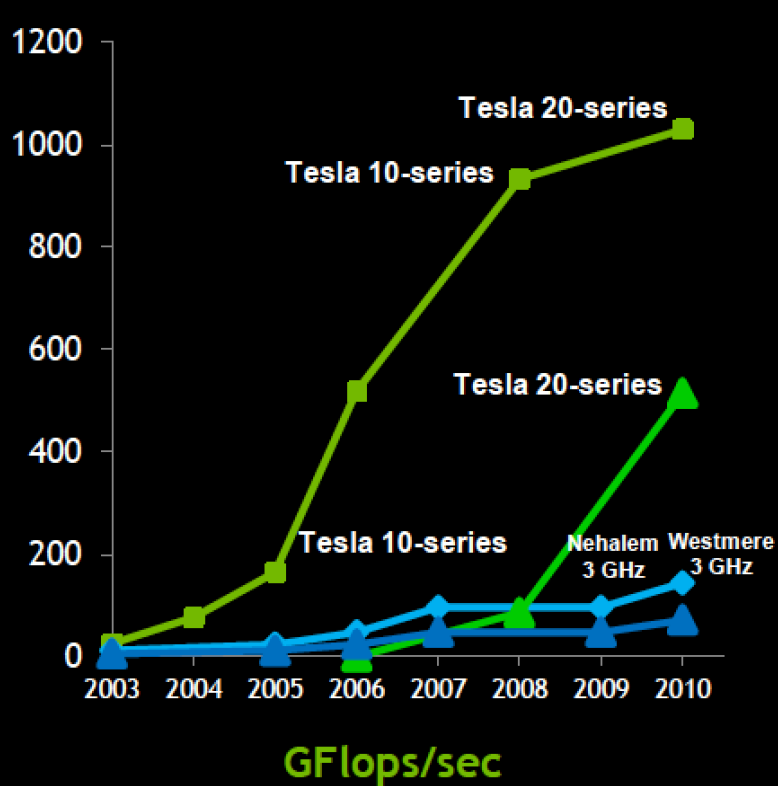
Non-uniform Memory Access

Adapted from Sam Williams, John Shalf, LBL/NERSC et al.

AGENDA

- 1 Intro
- 2 Processors Trends
- 3 Multiprocessors
- 4 **GPU Computing**

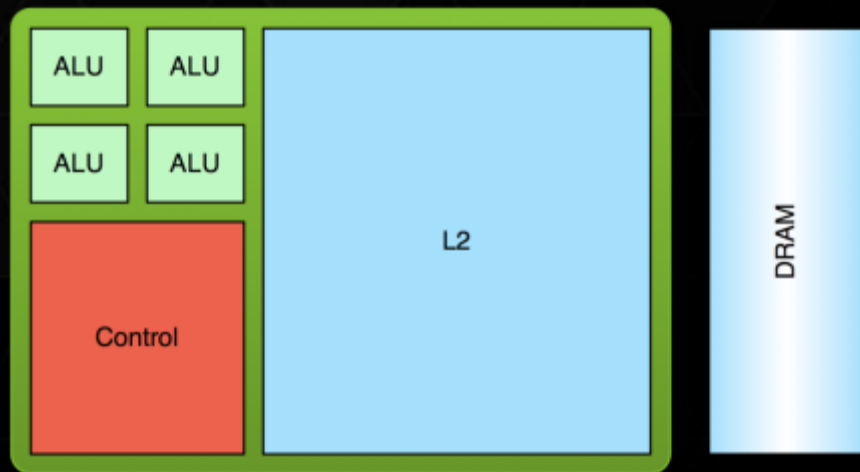
Why GPU Computing?



■ Single Precision: NVIDIA GPU ◆ Single Precision: x86 CPU
 ■ Double Precision: NVIDIA GPU ◆ Double Precision: x86 CPU

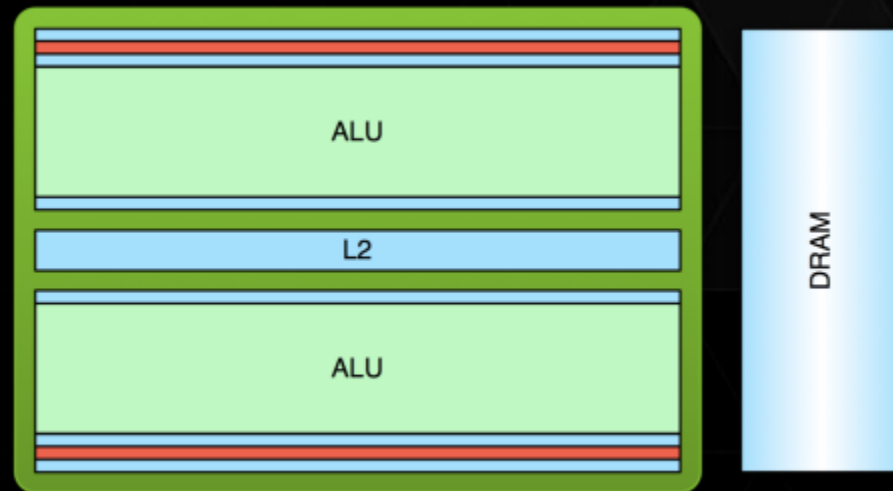
■ NVIDIA GPU ◆ X86 CPU
 ECC off

LOW LATENCY OR HIGH THROUGHPUT?



CPU

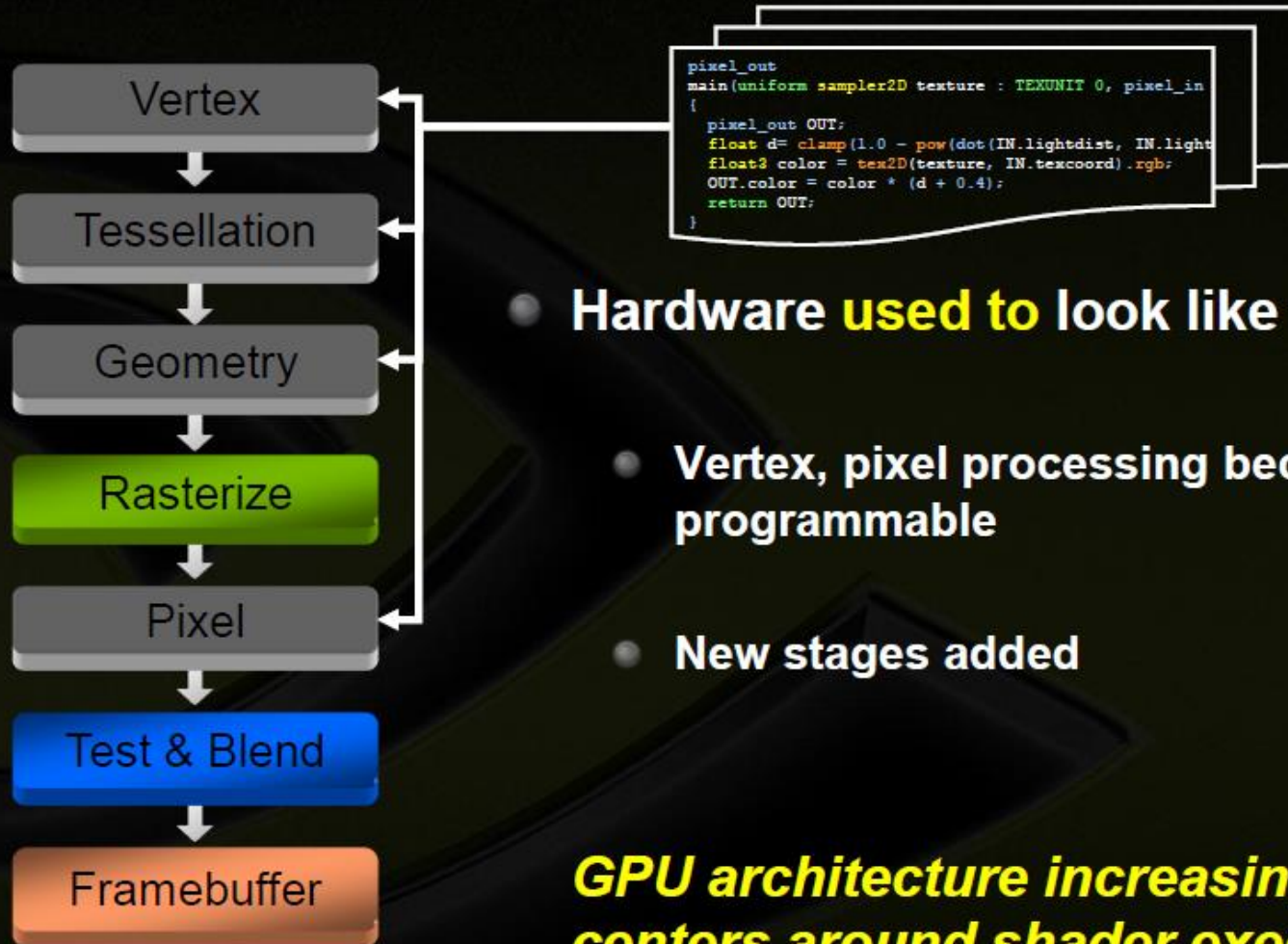
- Optimised for low-latency access to cached data sets
- Control logic for out-of-order and speculative execution



GPU

- Optimised for data-parallel, throughput computation
- Architecture tolerant of memory latency
- More transistors dedicated to computation

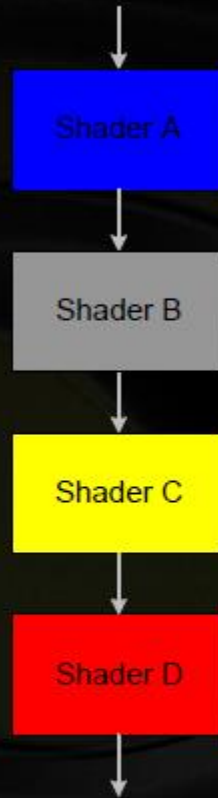
The Graphics Pipeline



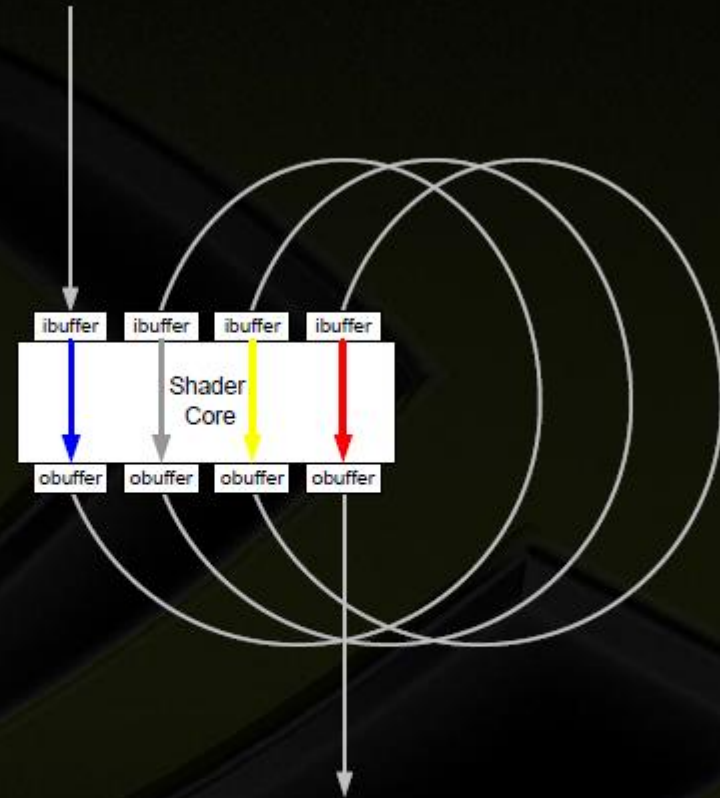
GPU architecture increasingly centers around shader execution

GPU UNIFIED DESIGN

Discrete Design

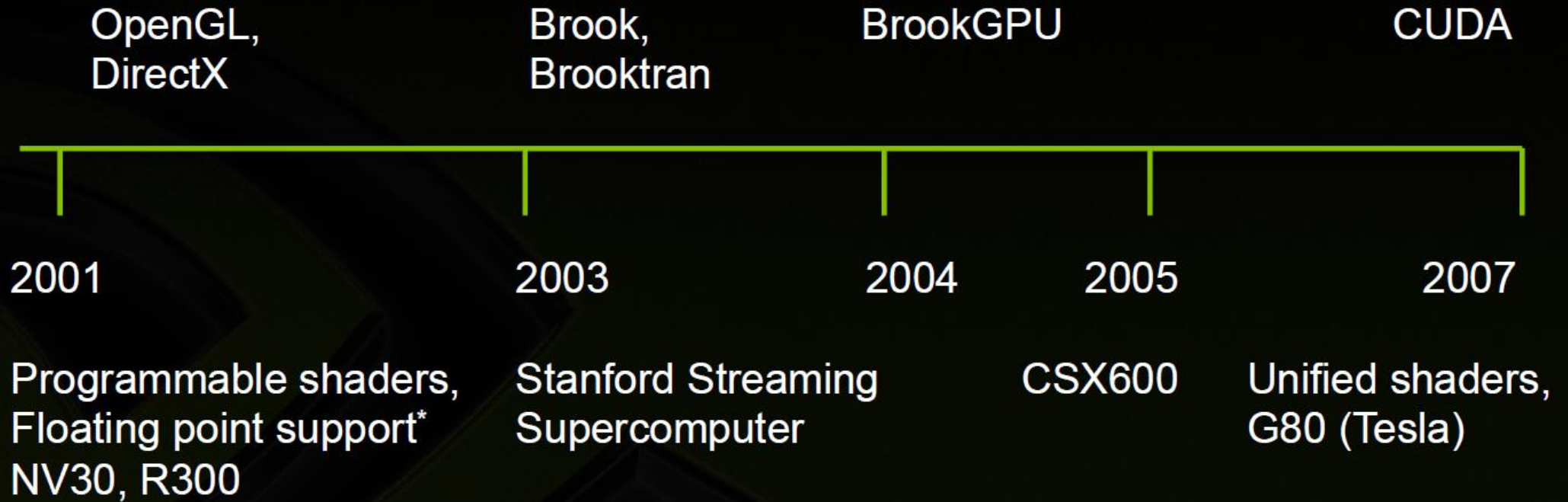


Unified Design



Vertex shaders, pixel shaders, etc. become *threads* running different programs on a flexible core

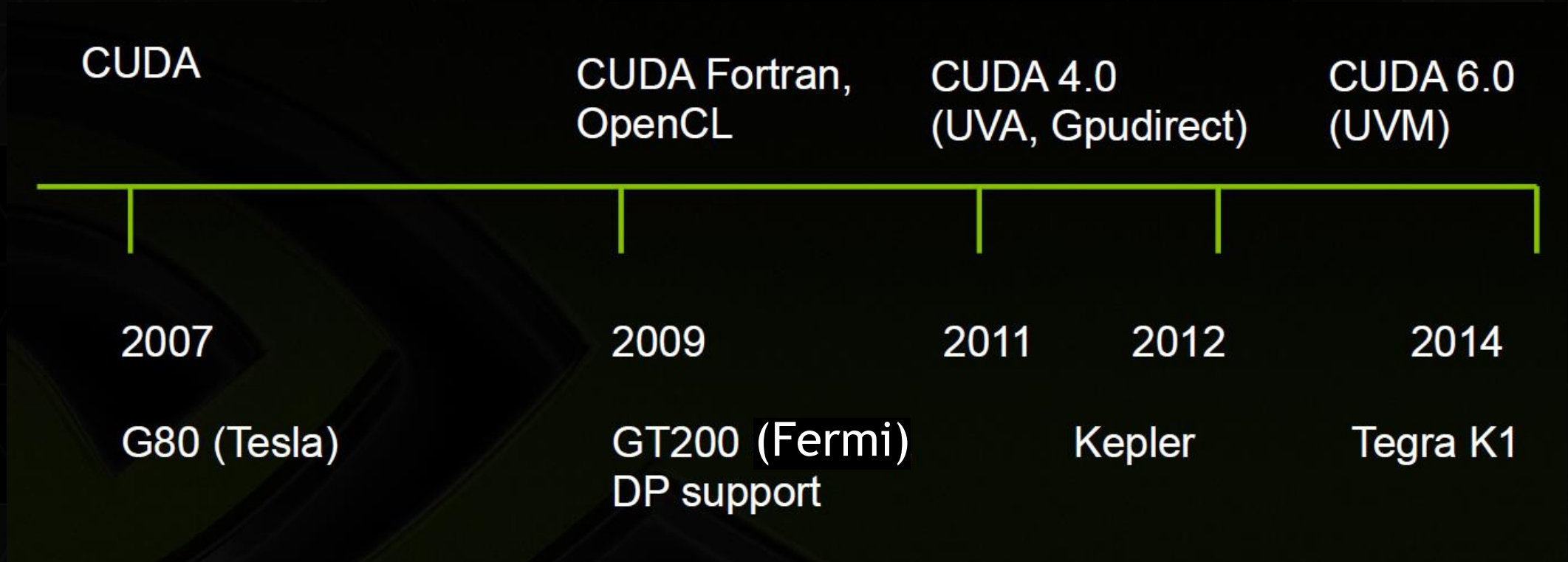
GPGPU COMPUTING, B.C. ERA



WHY GPU ARE ATTRACTIVE FOR HPC?

- Massive multithreaded manycore chips
- High flops count (SP and later DP)
- High memory bandwidth, (later) including ECC
- Lots of programming languages
- Lots of programming tools
- Widely used in Computational Physics

GPGPU COMPUTING, A.C. ERA



GROWTH IN GPU COMPUTING

2008

150,000
CUDA Downloads



27
CUDA Apps



60
Universities
Teaching



4,000
Academic
Papers



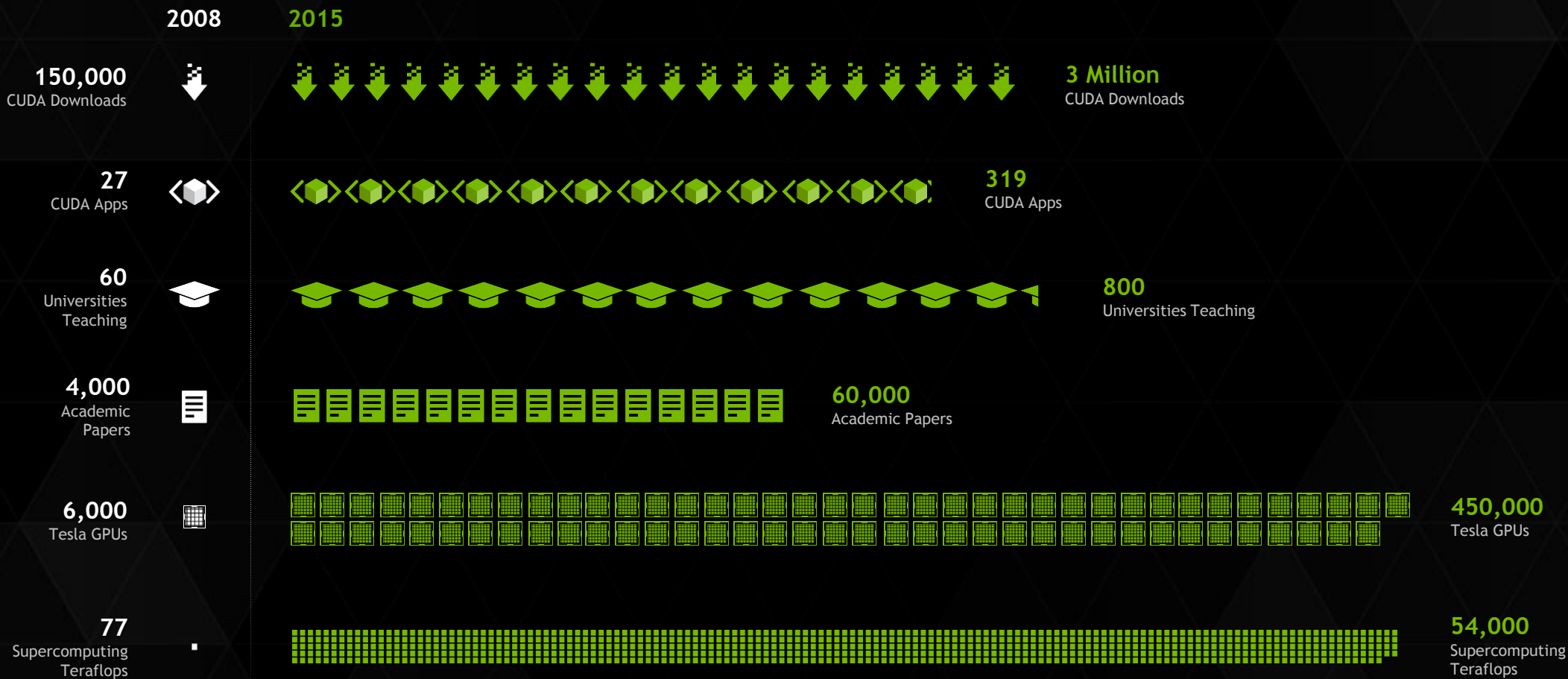
6,000
Tesla GPUs



77
Supercomputing
Teraflops



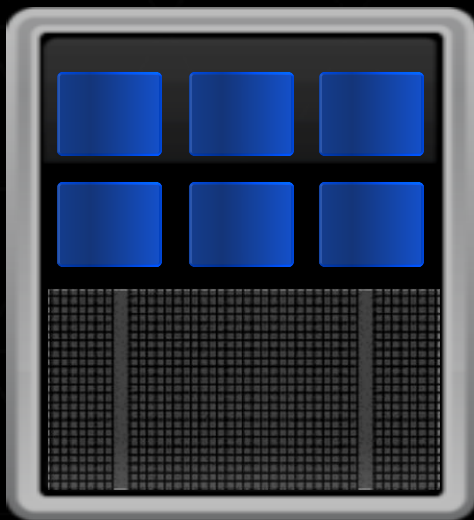
GROWTH IN GPU COMPUTING



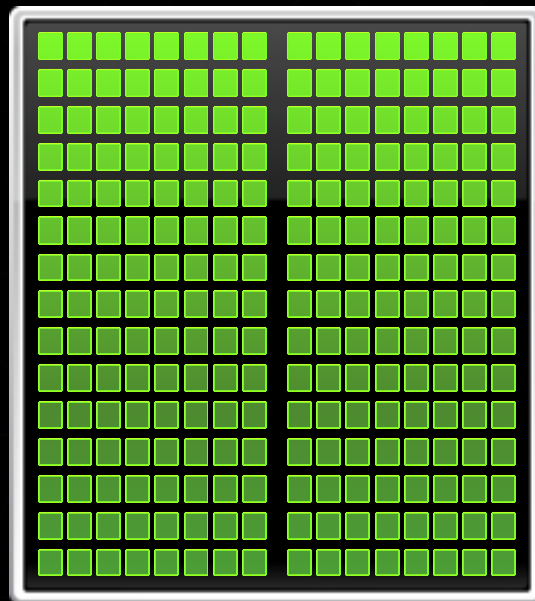
ACCELERATED, HETEROGENEOUS COMPUTING

10x Performance & 5x Energy Efficiency for HPC

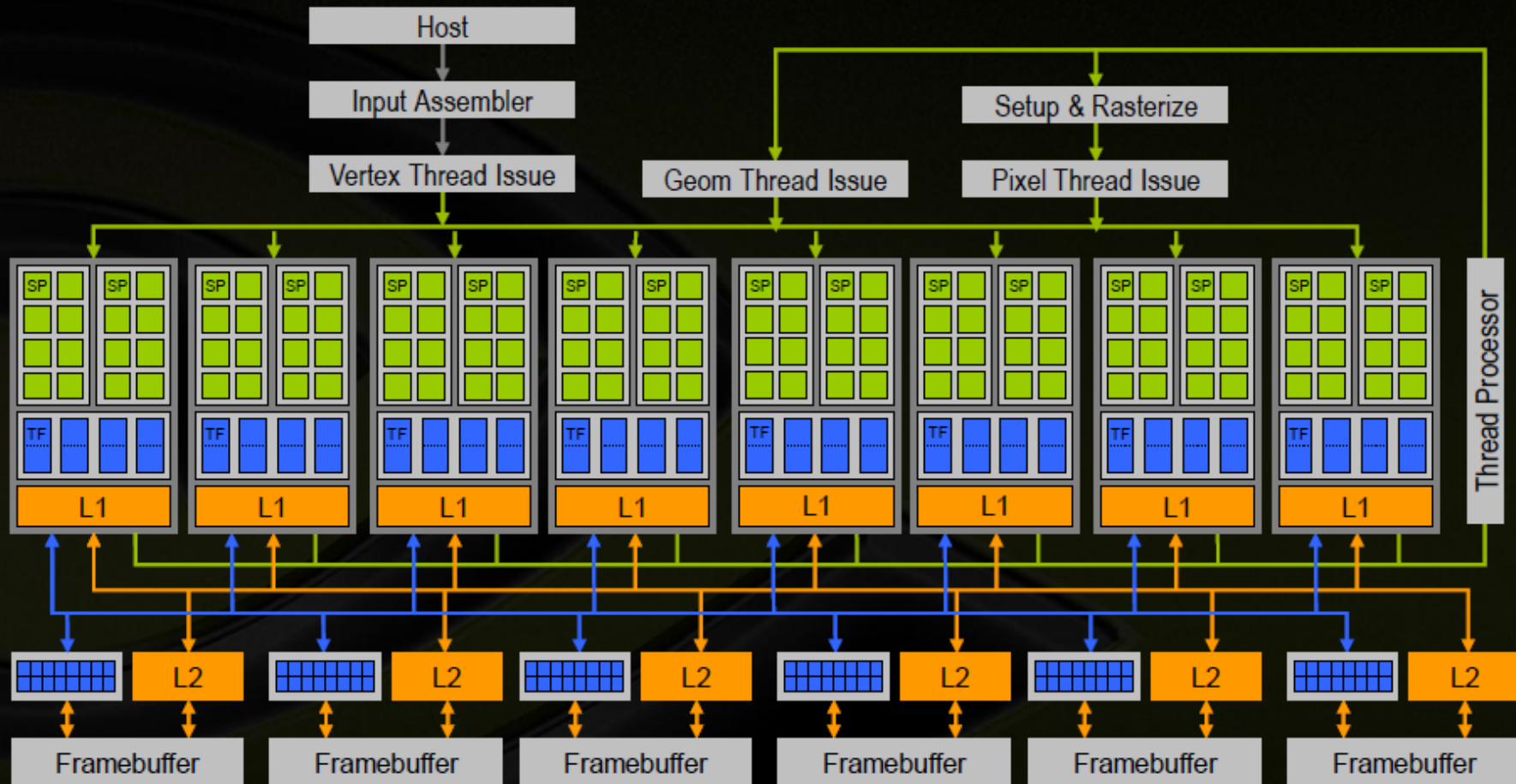
CPU
Optimized for
Serial Tasks



GPU Accelerator
Optimized for
Parallel Tasks



GEFORCE 8: 1ST MODERN GPU ARCH



NVIDIA KEPLER GK110 PROCESSOR

SMX



3x Performance per Watt

Hyper-Q

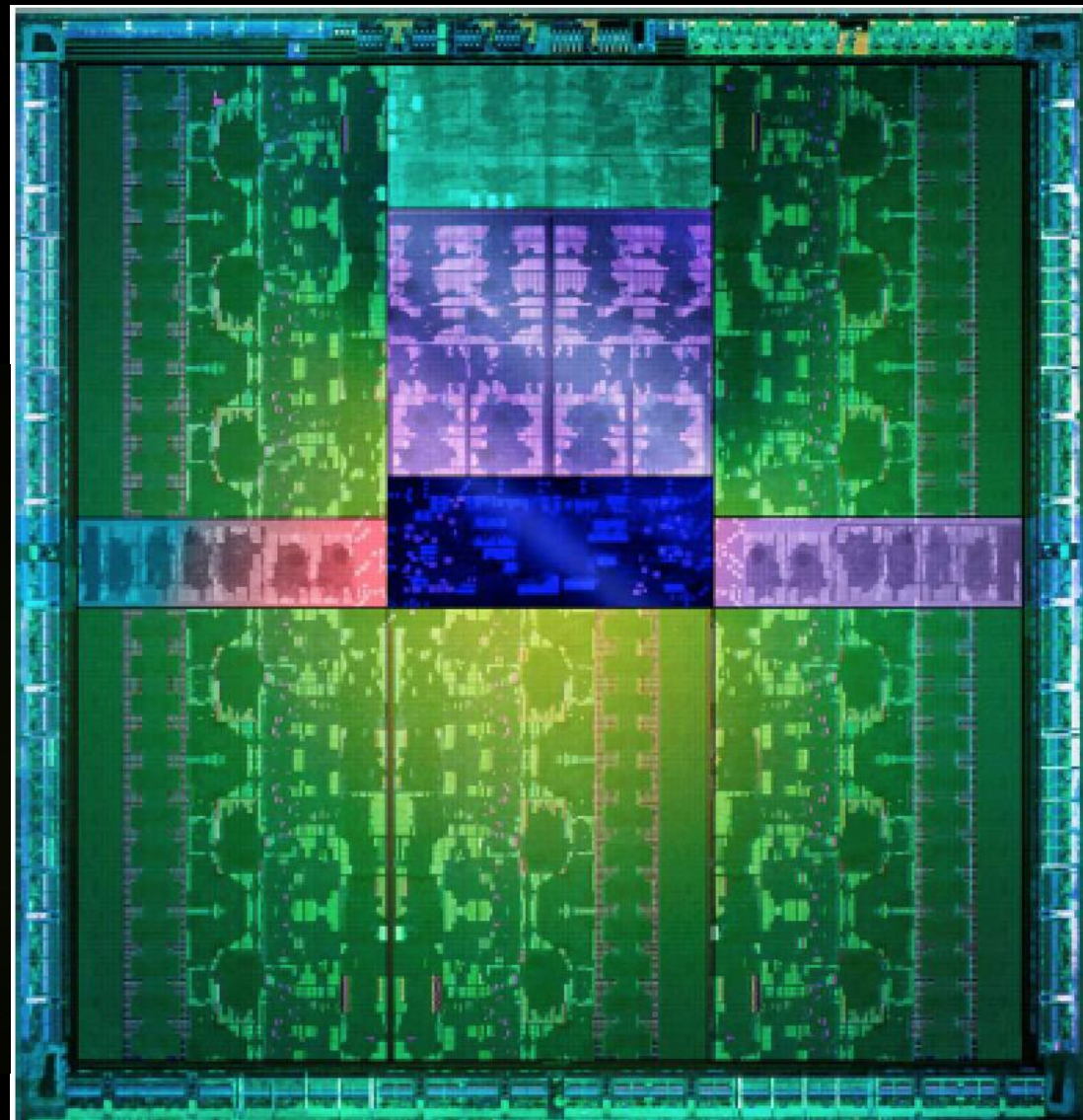


Easy Speed-up for Legacy MPI
Apps

Dynamic
Parallelism



Parallel Programming Made Easier
than Ever



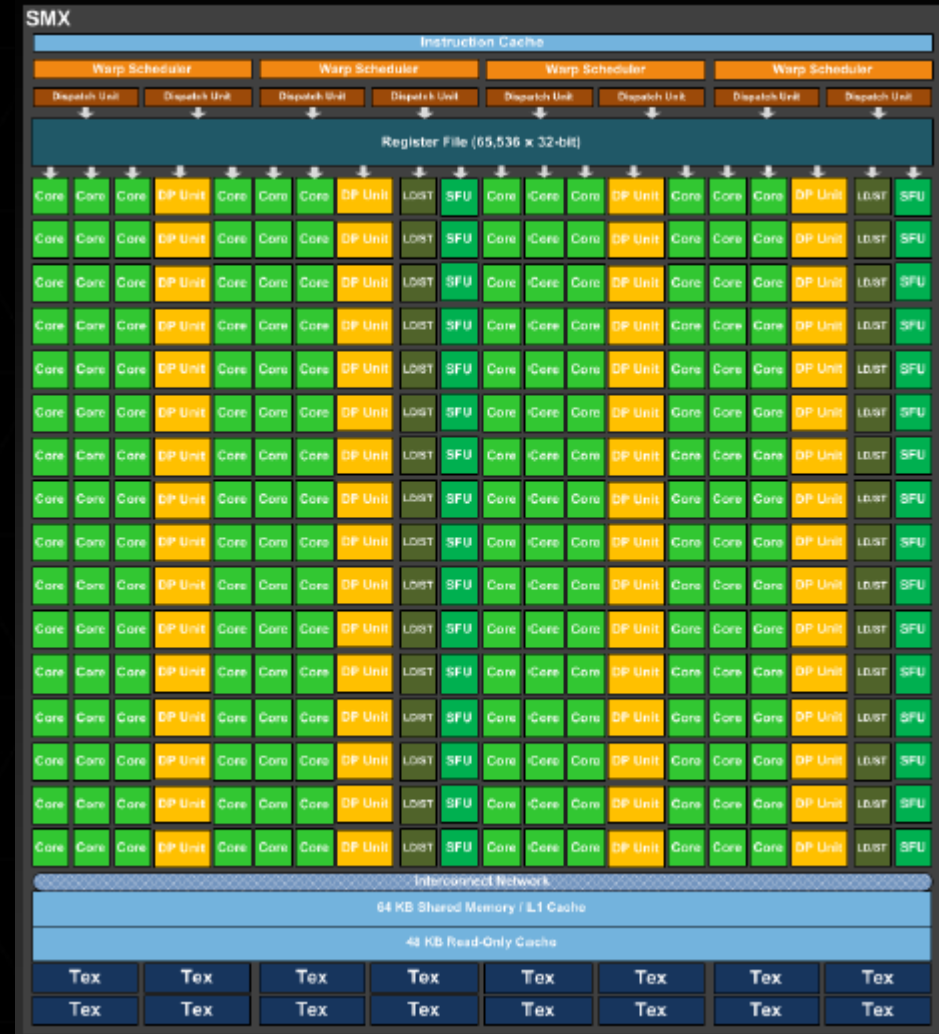
Kepler GK110 Block Diagram

Architecture

- 7.1B Transistors
- 15 SMX units
- > 1 TFLOP FP64
- 1.5 MB L2 Cache
- 384-bit GDDR5



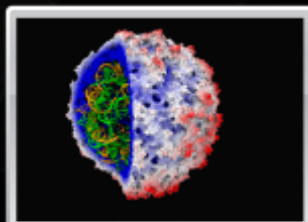
KEPLER SMX DIAGRAM





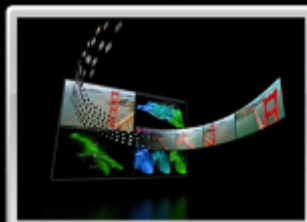
146X

Medical Imaging
U of Utah



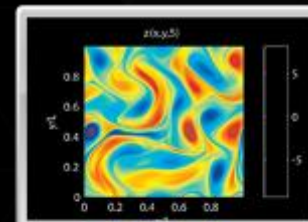
36X

Molecular Dynamics
U of Illinois, Urbana



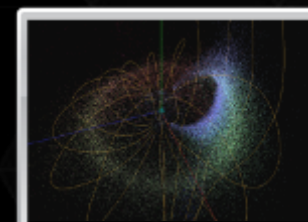
18X

Video Transcoding
Elemental Tech



50X

Matlab Computing
AccelerEyes



100X

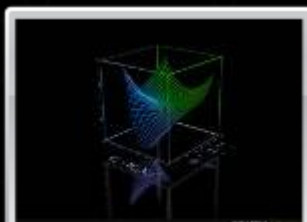
Astrophysics
RIKEN

GPUs Accelerate Science



149X

Financial Simulation
Oxford



47X

Linear Algebra
Universidad Jaime



20X

3D Ultrasound
Techniscan



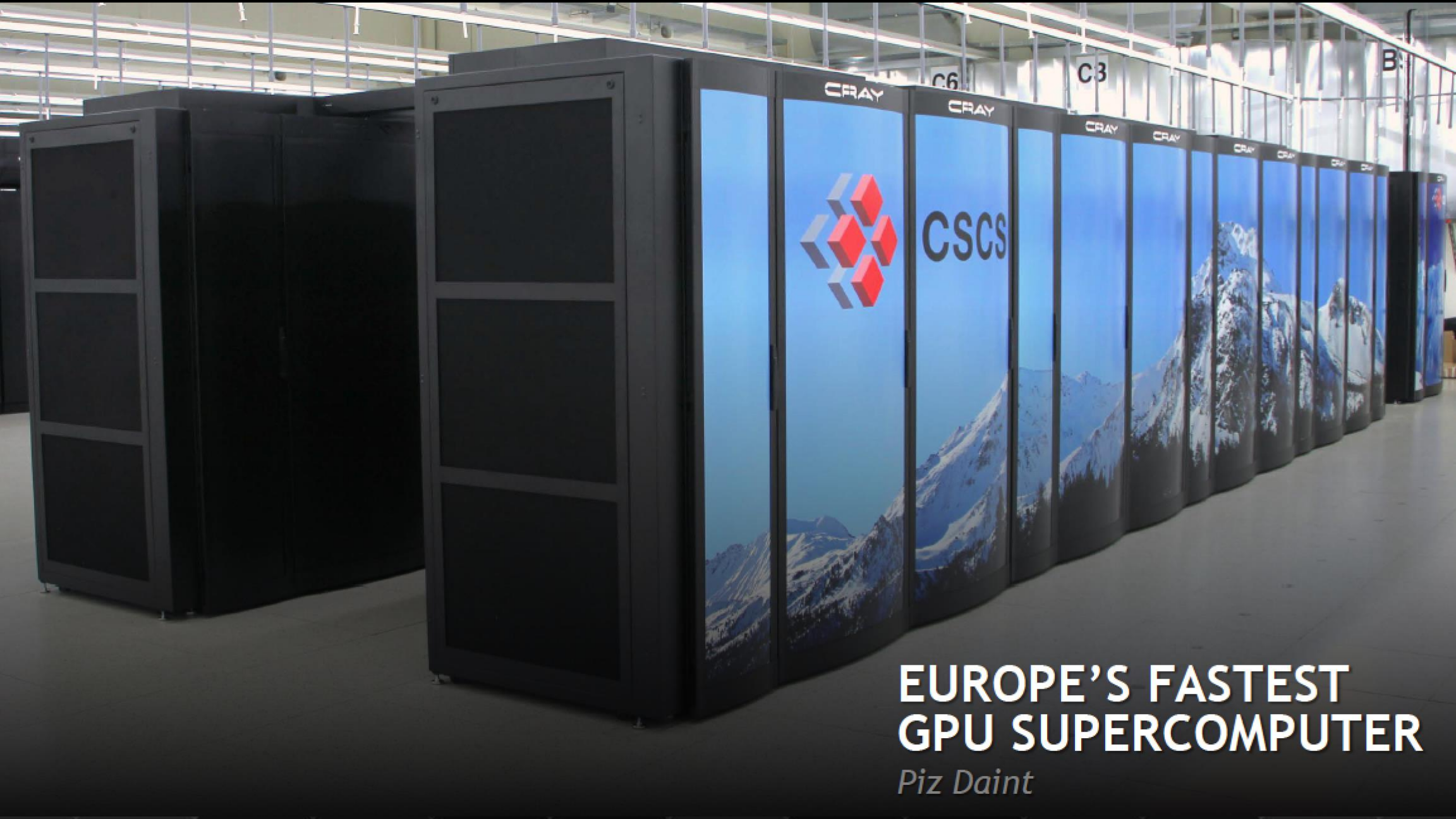
130X

Quantum Chemistry
U of Illinois, Urbana



30X

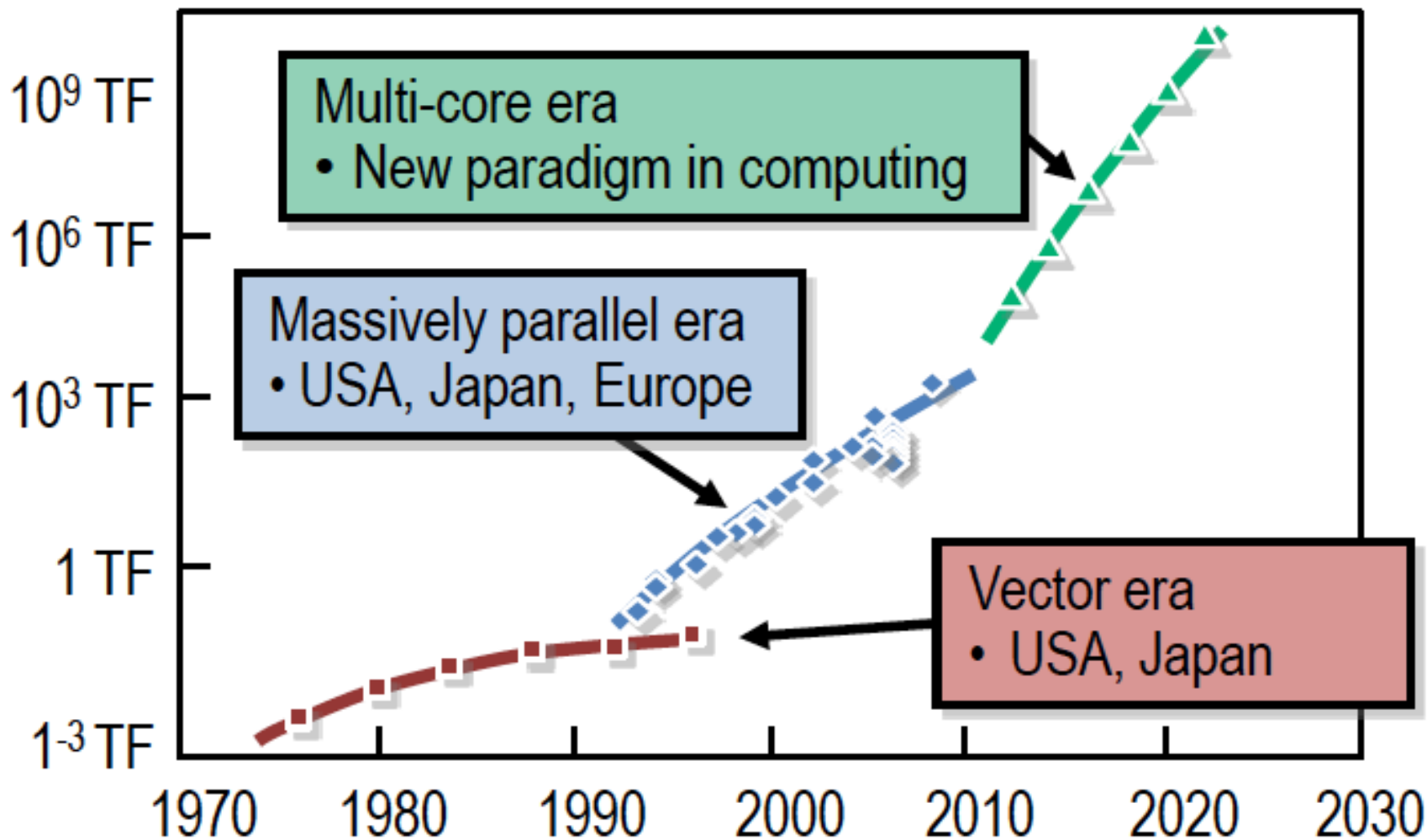
Gene Sequencing
U of Maryland



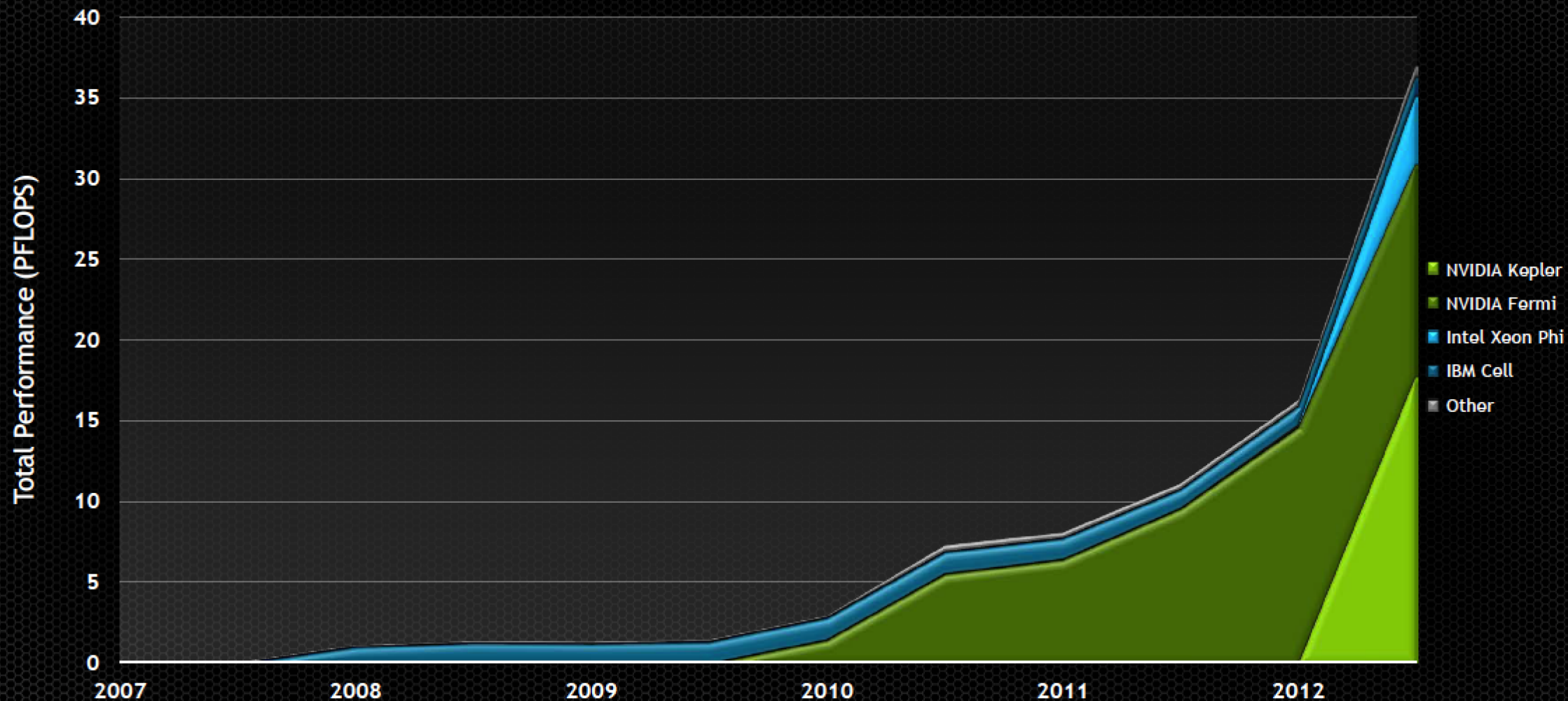
EUROPE'S FASTEST GPU SUPERCOMPUTER

Piz Daint

SUPERCOMPUTING ERA

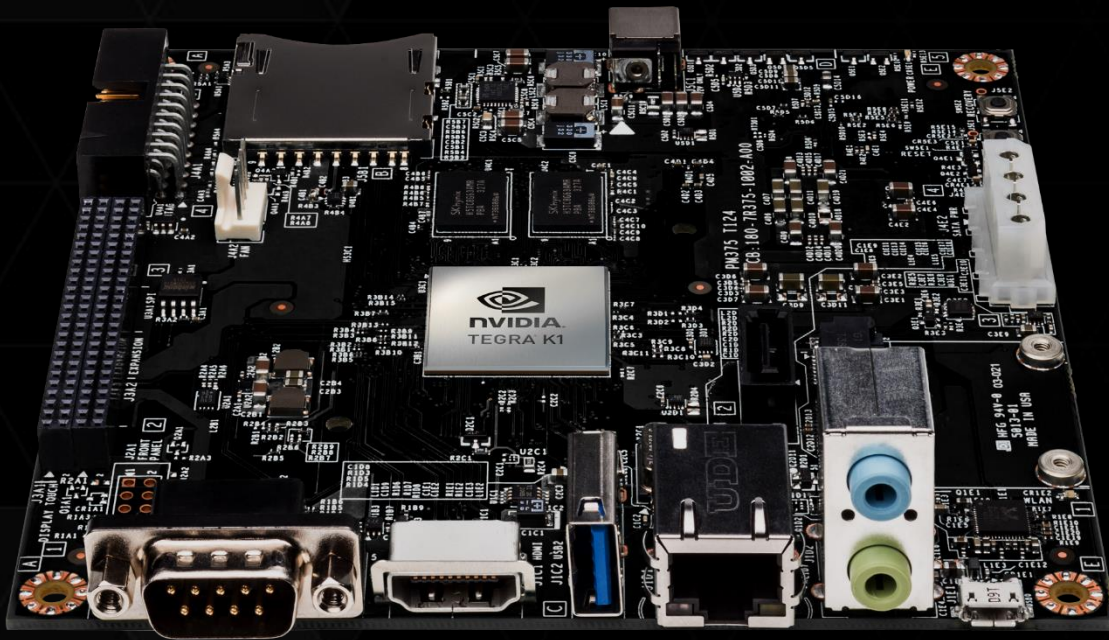


Top500: Performance from Accelerators



JETSON TK1

THE WORLD'S 1st EMBEDDED SUPERCOMPUTER



Development Platform for Embedded
Computer Vision, Robotics, Medical

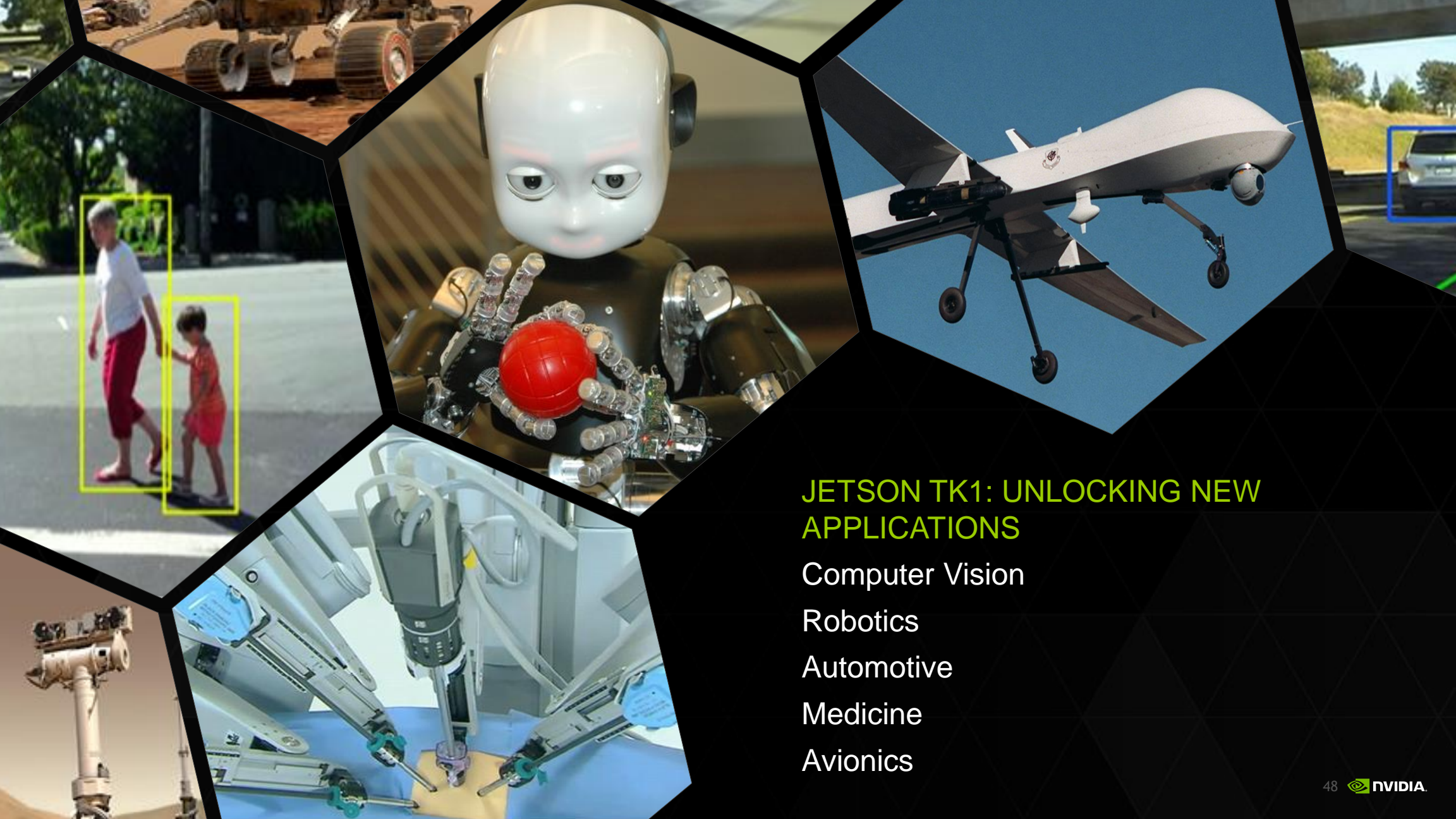
Tegra K1 SoC

Quad core A15 + Kepler GPU

192 CUDA Enabled cores

326 Gflops @ 5 Watt

\$192



JETSON TK1: UNLOCKING NEW APPLICATIONS

Computer Vision

Robotics

Automotive

Medicine

Avionics

US TO BUILD TWO FLAGSHIP SUPERCOMPUTERS



SUMMIT

150-300 PFLOPS
Peak Performance



**Lawrence Livermore
National Laboratory**

SIERRA

> 100 PFLOPS
Peak Performance

IBM POWER9 CPU + NVIDIA Volta GPU

NVLink High Speed Interconnect

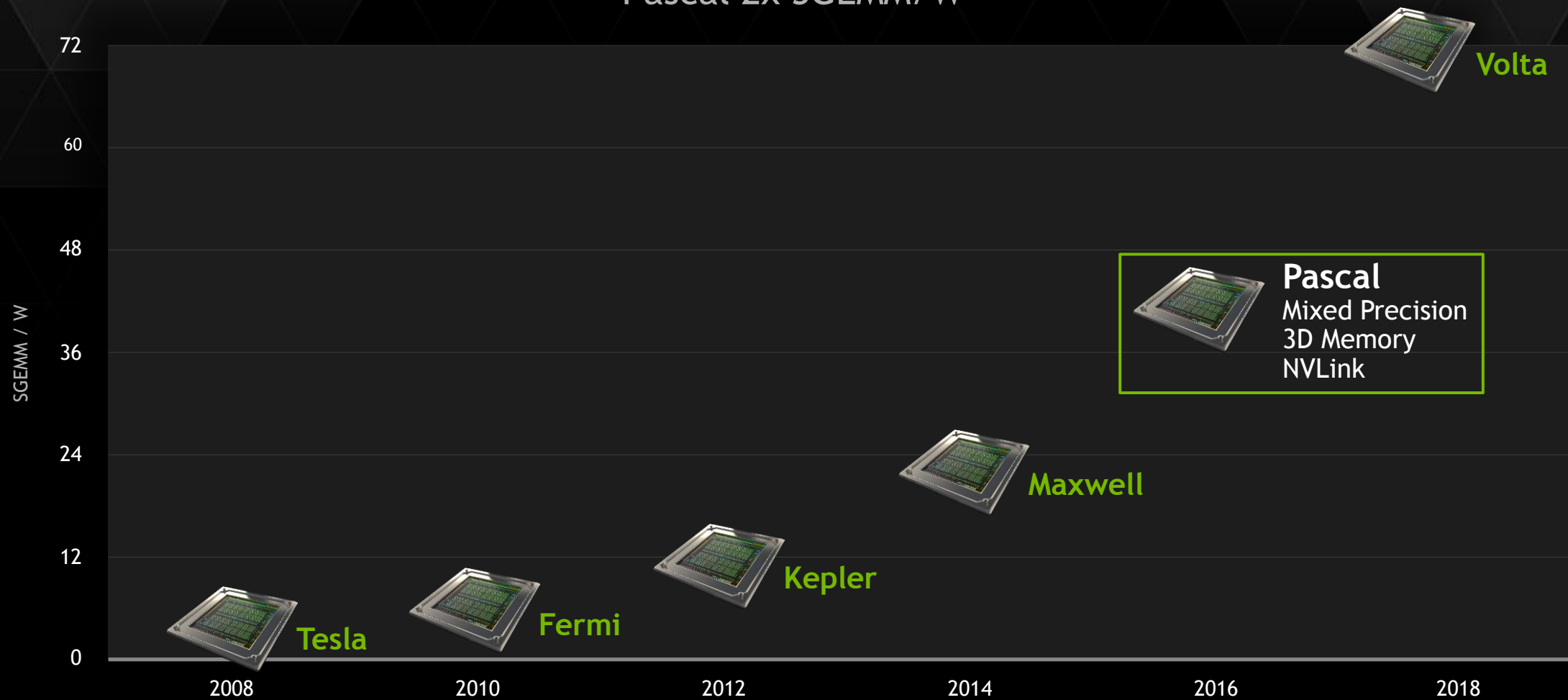
>40 TFLOPS per Node, >3,400 Nodes

2017

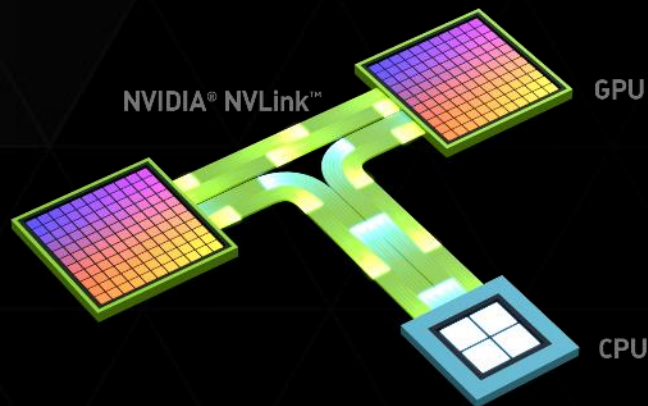
Major Step Forward on the Path to Exascale

GPU ROADMAP

Pascal 2x SGEMM/W

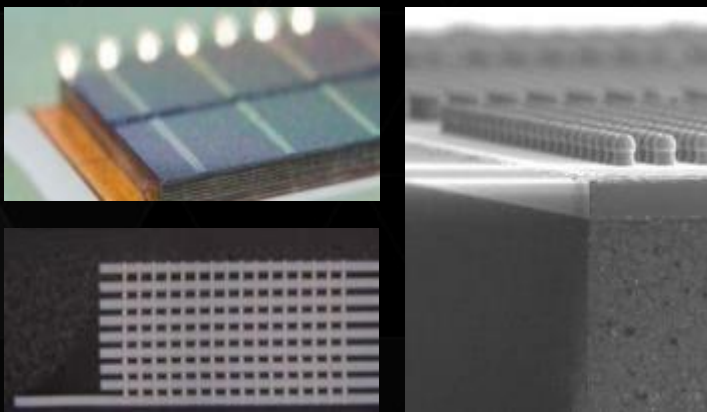


PASCAL GPU FEATURING NVLINK AND STACKED MEMORY



NVLINK

- GPU high speed interconnect
- 80-200 GB/s



3D Stacked Memory

- 4x Higher Bandwidth (~1 TB/s)
- 3x Larger Capacity
- 4x More Energy Efficient per bit

THANK YOU



nvidia®

cnardone@nvidia.com

+39 335 5828197