

CME 213, ME 339—Spring 2021

Eric Darve, ICME, Stanford



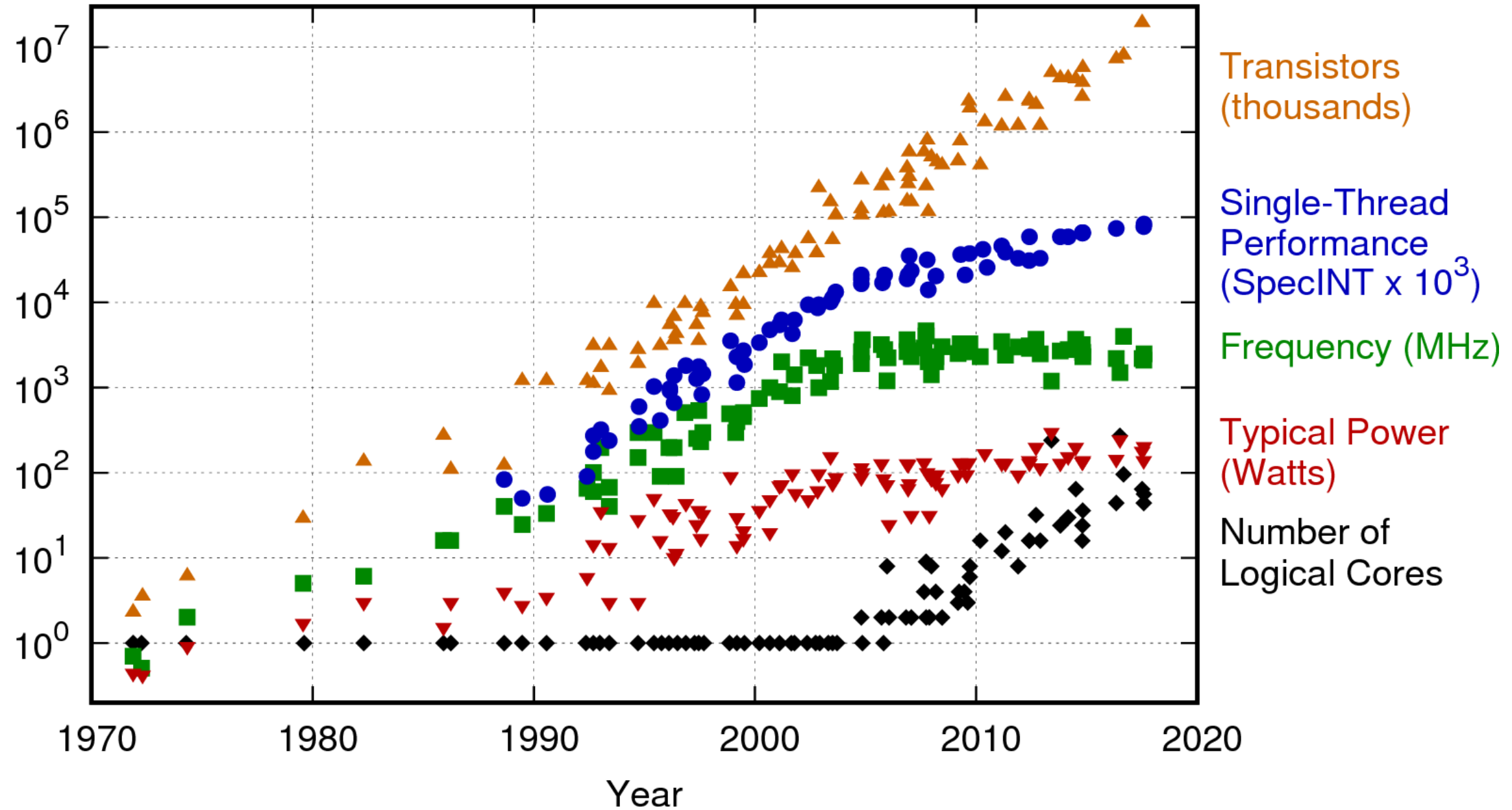
“If debugging is the process of removing bugs, then programming must be the process of putting them in.” (Edsger W. Dijkstra)

CME 213 so far:

- C++ threads
- OpenMP: for loop and task
- Sorting algorithms on shared memory

Onwards to GPU computing!

# 42 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 2010-2017 by K. Rupp

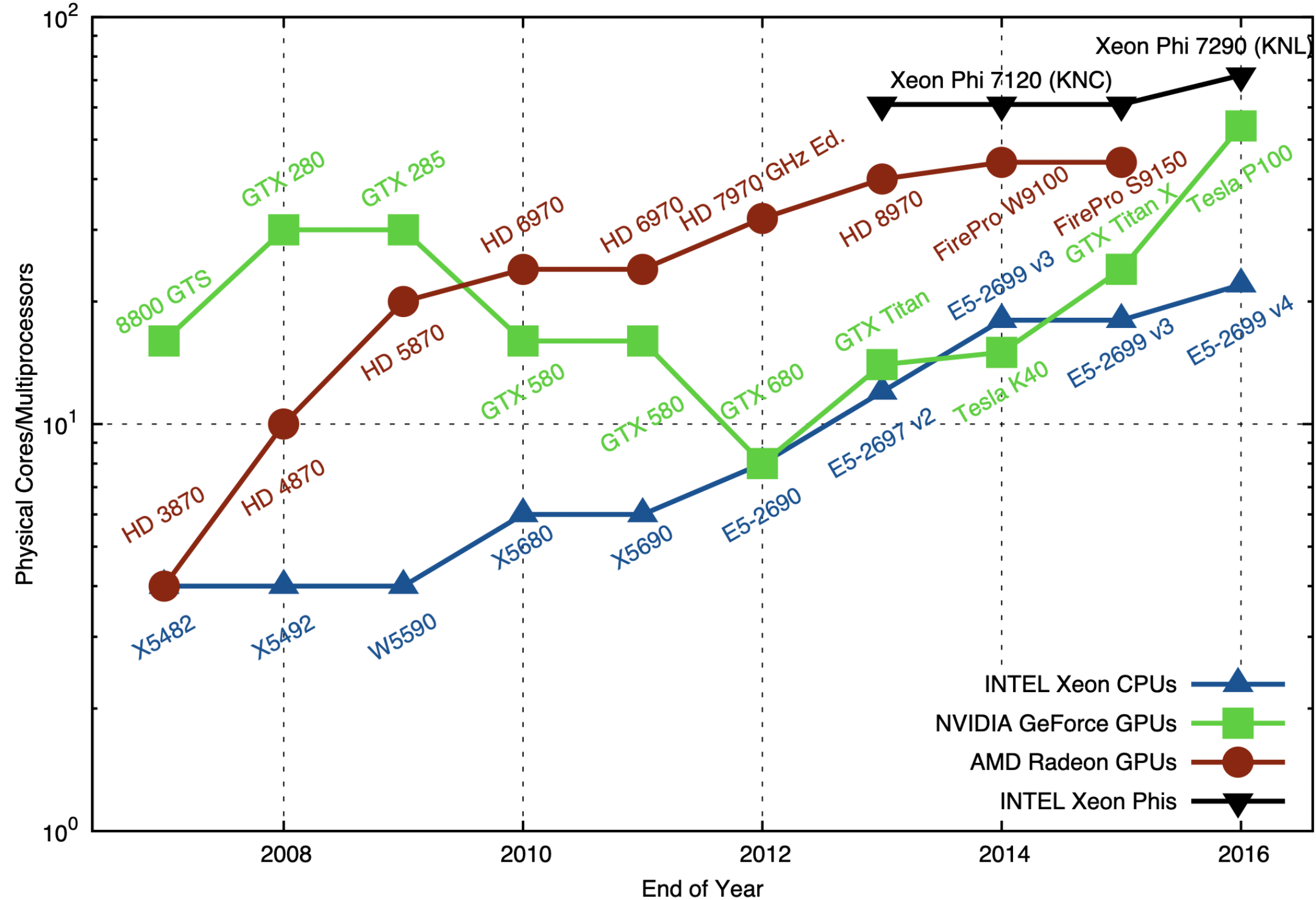
## Reference

<https://github.com/karlrupp/cpu-gpu-mic-comparison>

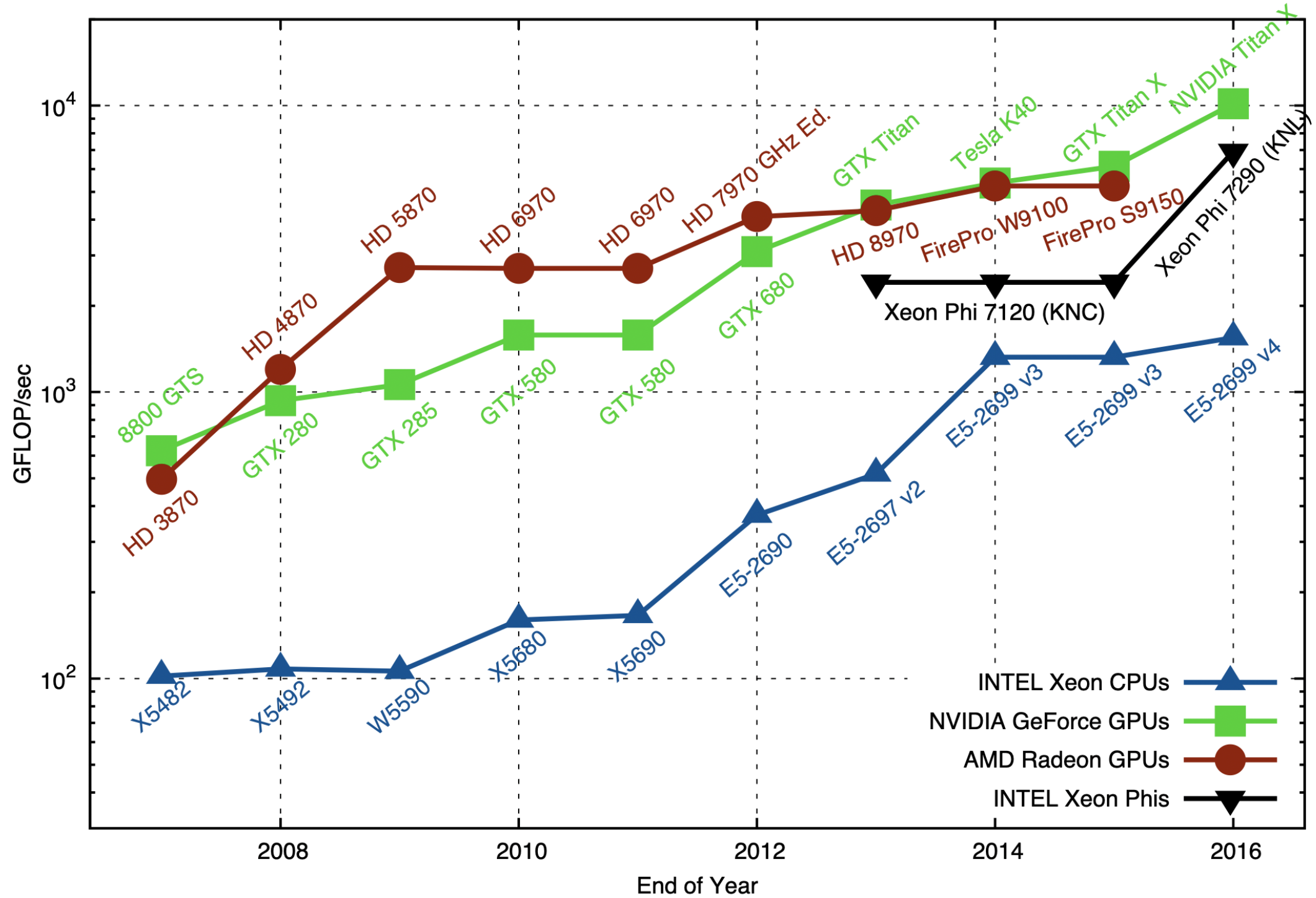
<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

<https://www.karlrupp.net/2015/06/40-years-of-microprocessor-trend-data/>

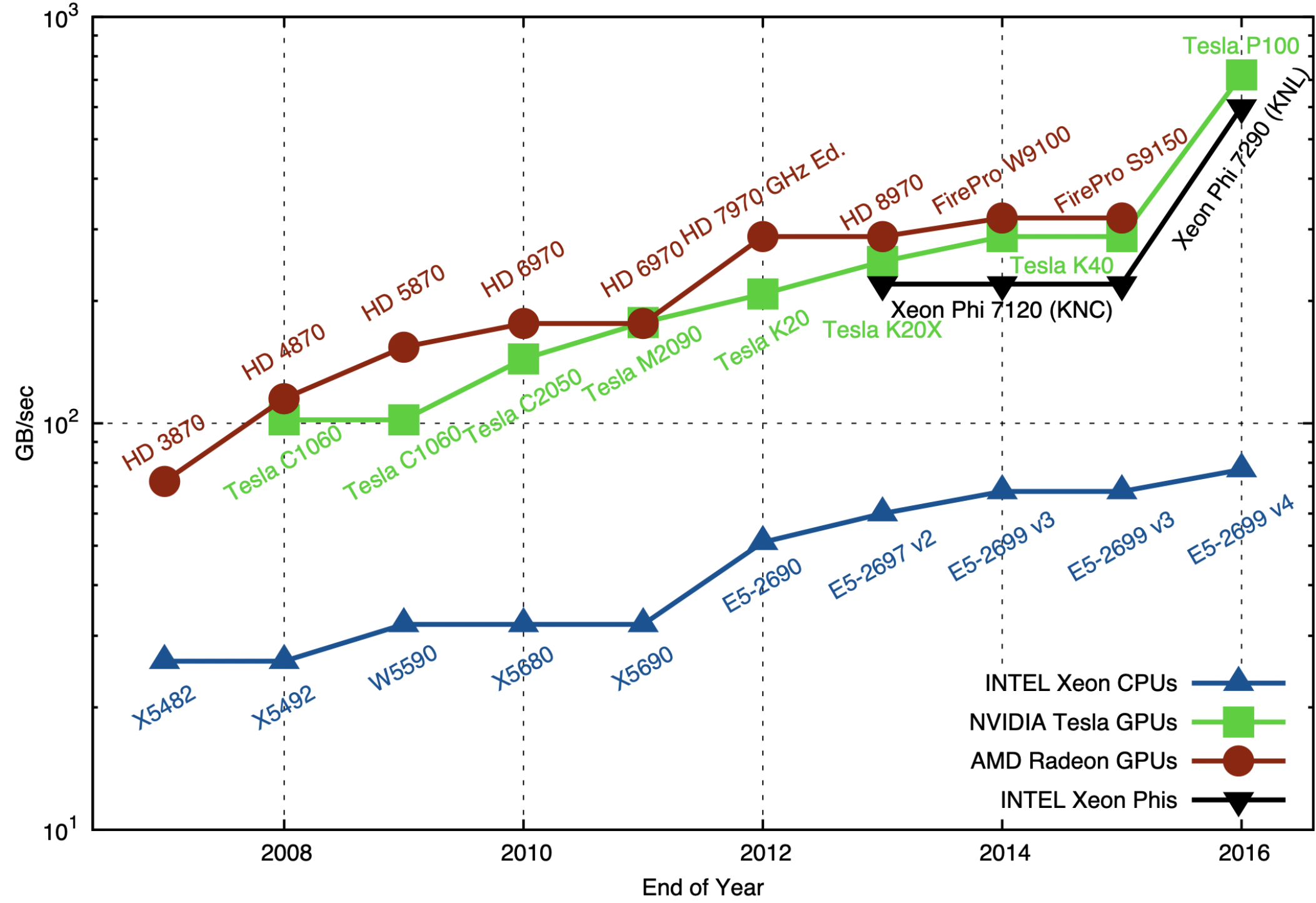
Number of Physical Cores/Multiprocessors, High-End Hardware



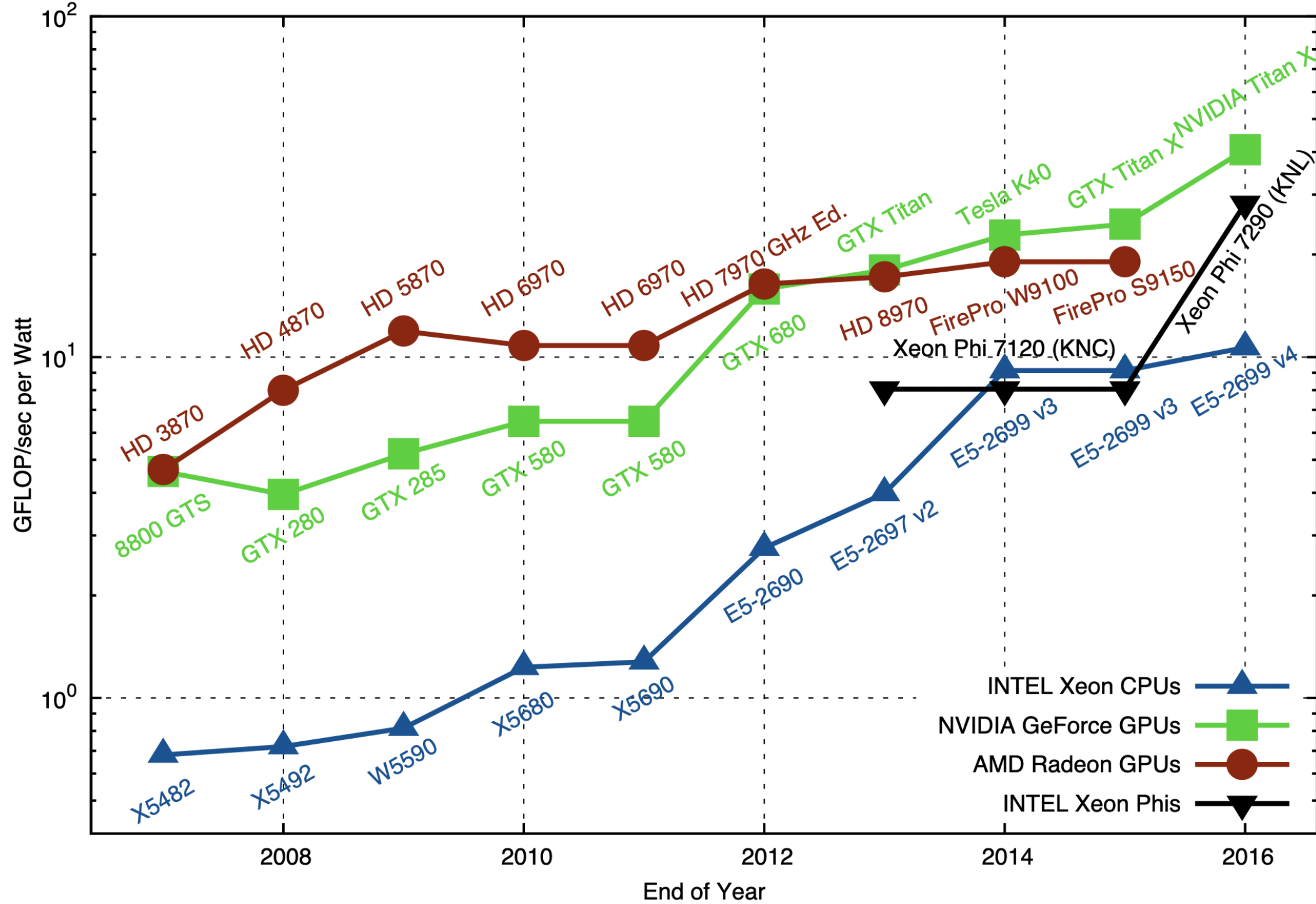
Theoretical Peak Performance, Single Precision



Theoretical Peak Memory Bandwidth Comparison

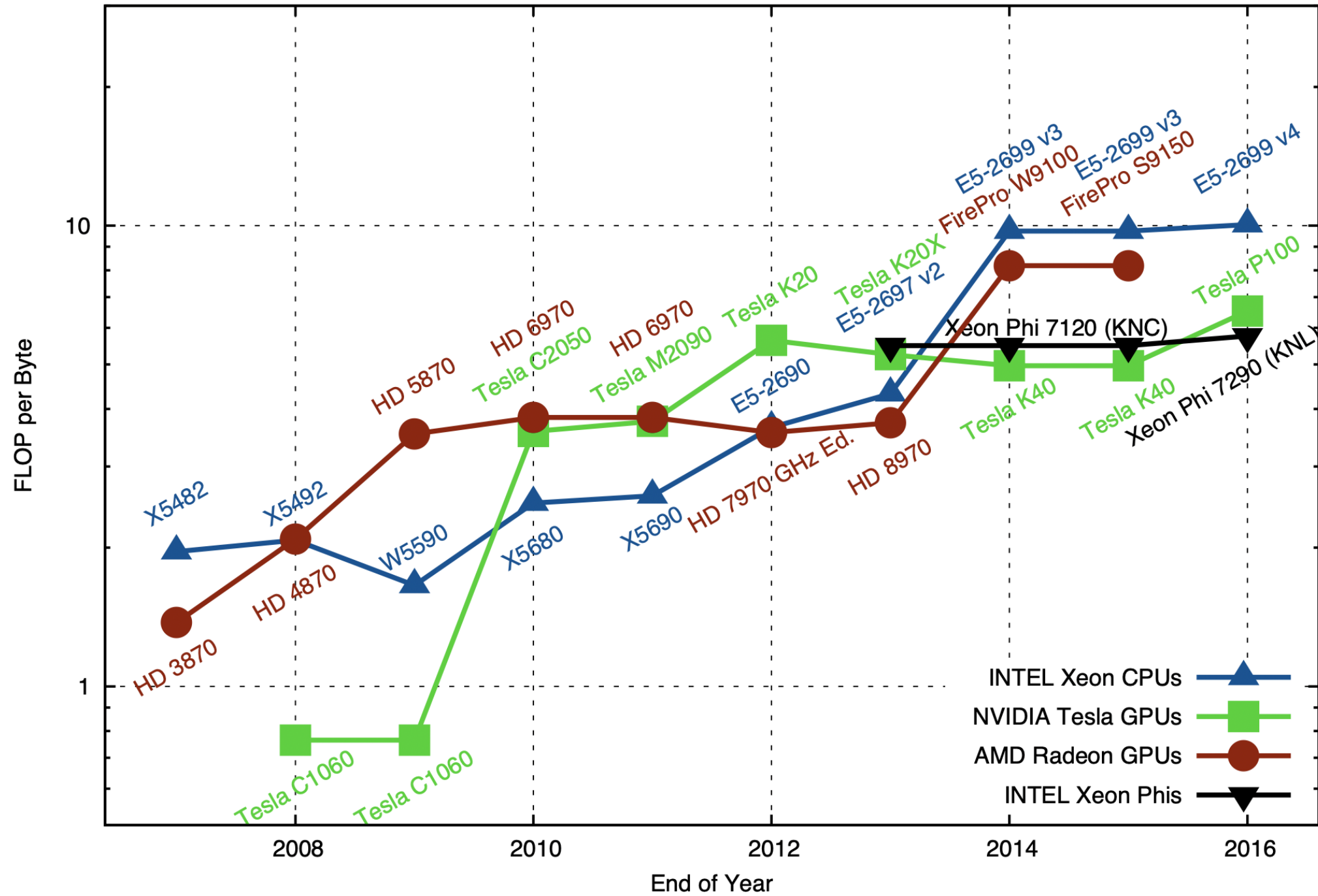


Theoretical Peak Floating Point Operations per Watt, Single Precision





Theoretical Peak Floating Point Operations per Byte, Double Precision



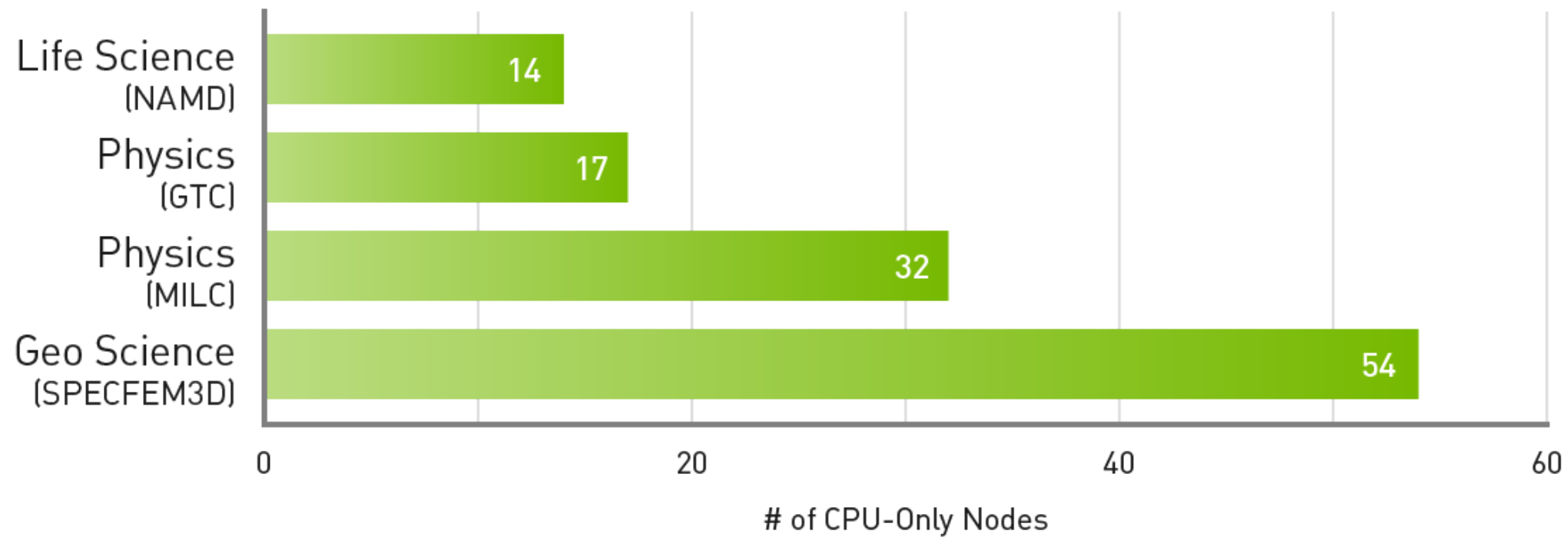
## Example: Volta V100

- 8.2 teraflops double-precision performance
- 16.4 teraflops single-precision performance
- 130 teraflops for tensor (deep learning)
- 1134 GB/sec memory bandwidth
- 250 Watts power



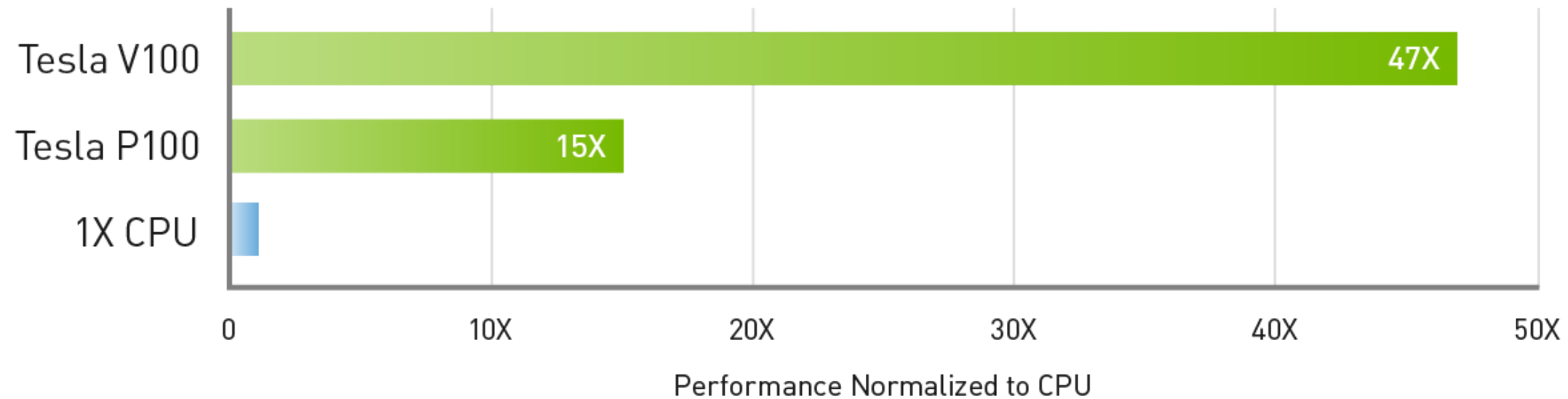
# 1 GPU Node Replaces Up To 54 CPU Nodes

Node Replacement: HPC Mixed Workload



CPU Server: Dual Xeon Gold 6140@2.30GHz, GPU Servers: same CPU server w/ 4x V100 PCIe | CUDA Version: CUDA 9.x | Dataset: NAMD (STMV), GTC (mpi#proc.in), MILC (APEX Medium), SPECFEM3D (four\_material\_simple\_model) | To arrive at CPU node equivalence, we use measured benchmark with up to 8 CPU nodes. Then we use linear scaling to scale beyond 8 nodes.

## 47X Higher Throughput Than CPU Server on Deep Learning Inference



Workload: ResNet-50 | CPU: 1X Xeon E5-2690v4 @ 2.6 GHz | GPU: Add 1X Tesla P100 or V100

What is the technology behind GPU processors?

It's started with 3D graphics



GPUs were designed to perform the same instruction on a large amount of data

Graphics processing == scientific computing

## The multicore processor

Fast processor because:

- Pipelining of execution
- Out of order execution
- Branch prediction
- Pre-fetching
- Large amount of multilevel cache
- ...





The GPU processor

Many computing units operating in parallel

Ideal for simple but repetitive tasks



Great for:

- Dense linear algebra
- Finite-difference
- Neural network

Bad when calculation involves branching

## The secret behind the magic

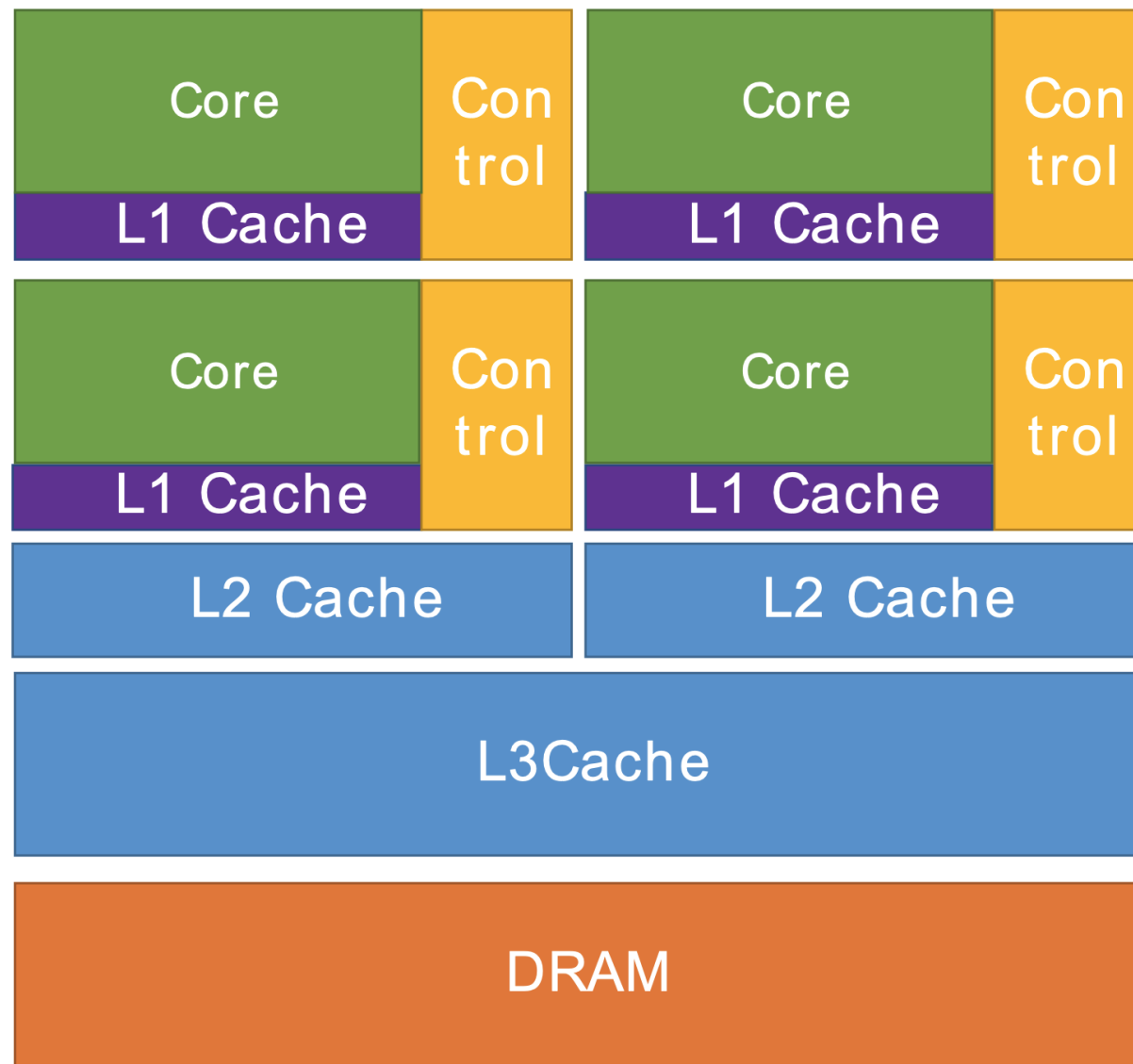
Striking the right trade-off

More computing units means you need to give up something

- No processor space dedicated to complex optimizations, e.g., out-of-order execution
- Cache/memory is limited because it has to be shared amongst the threads

- Light threads: hardware supports the ability to switch threads every cycle
- Limited logic for program control: 32 threads are grouped into warps; all threads in warp execute the same instruction at the same time.

Each computing unit is less powerful but there are more of them.



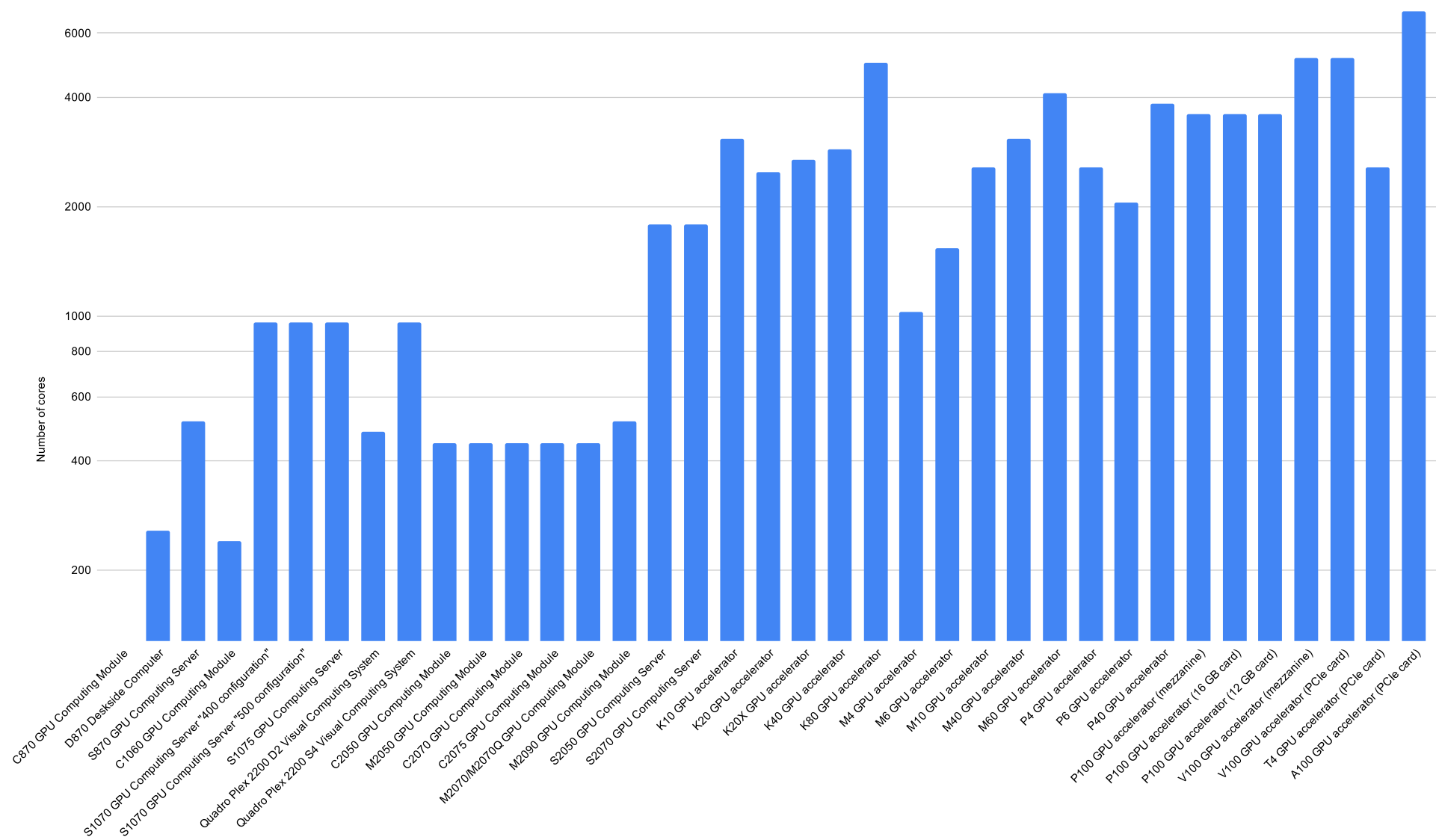
CPU



GPU

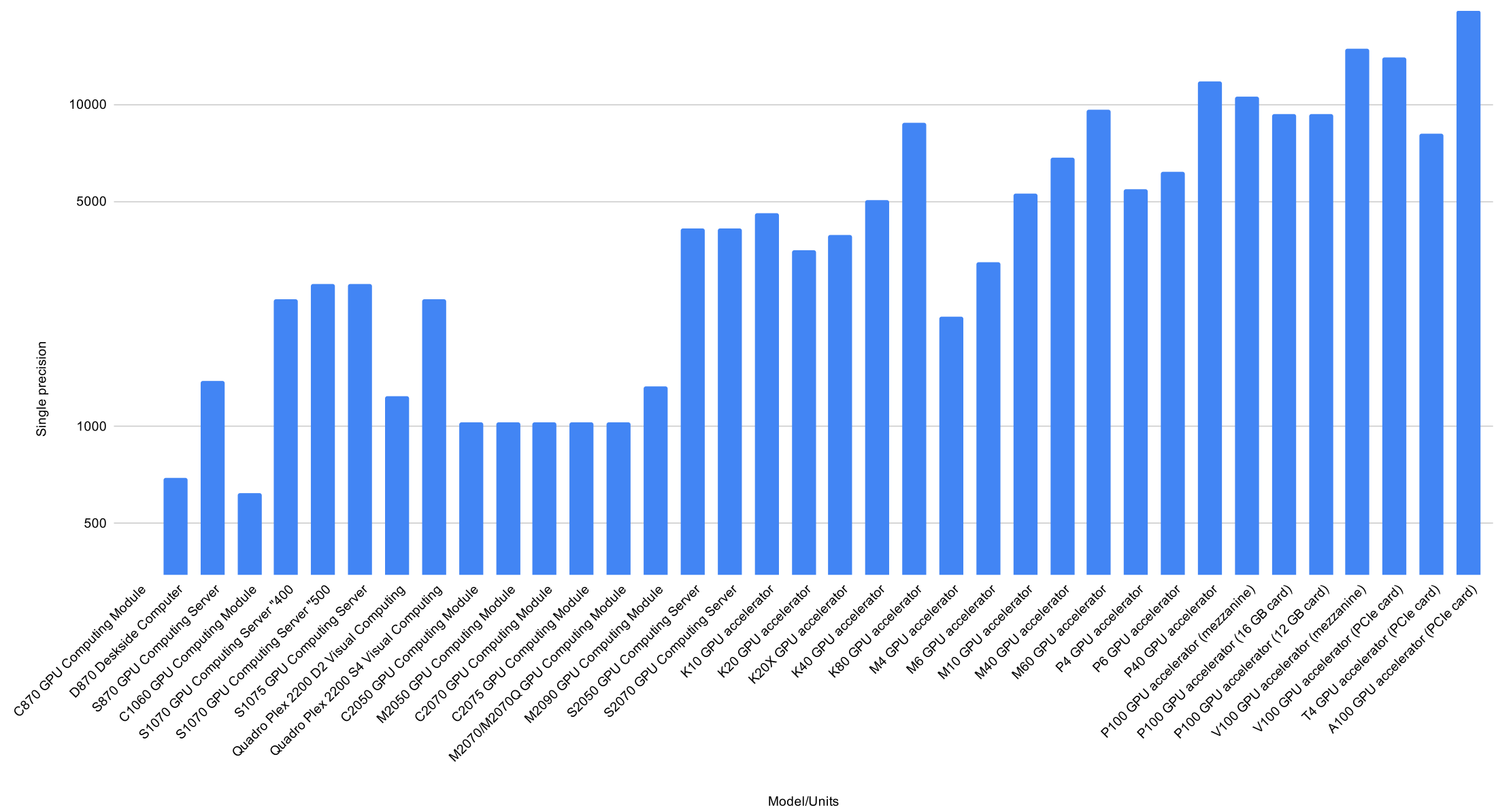
Micro-architecture	Release	Compute Capability	GPU code name
G70	2005		
Tesla	2006	1.0-1.3	GXX, GT2XX
Fermi	2010	2.0-2.1	GFXXX
Kepler	2012	3.0-3.7	GKXXX
Maxwell	2014	5.0-5.3	GMXXX
Pascal	2016	6.0-6.2	GPXXX
Volta	2017	7.0-7.2	GVXXX
Turing	2018	7.5	TUXXX
Ampere	2020	8.0-8.6	GAXXX
Lovelace	?	?	?
Hopper	?	?	?





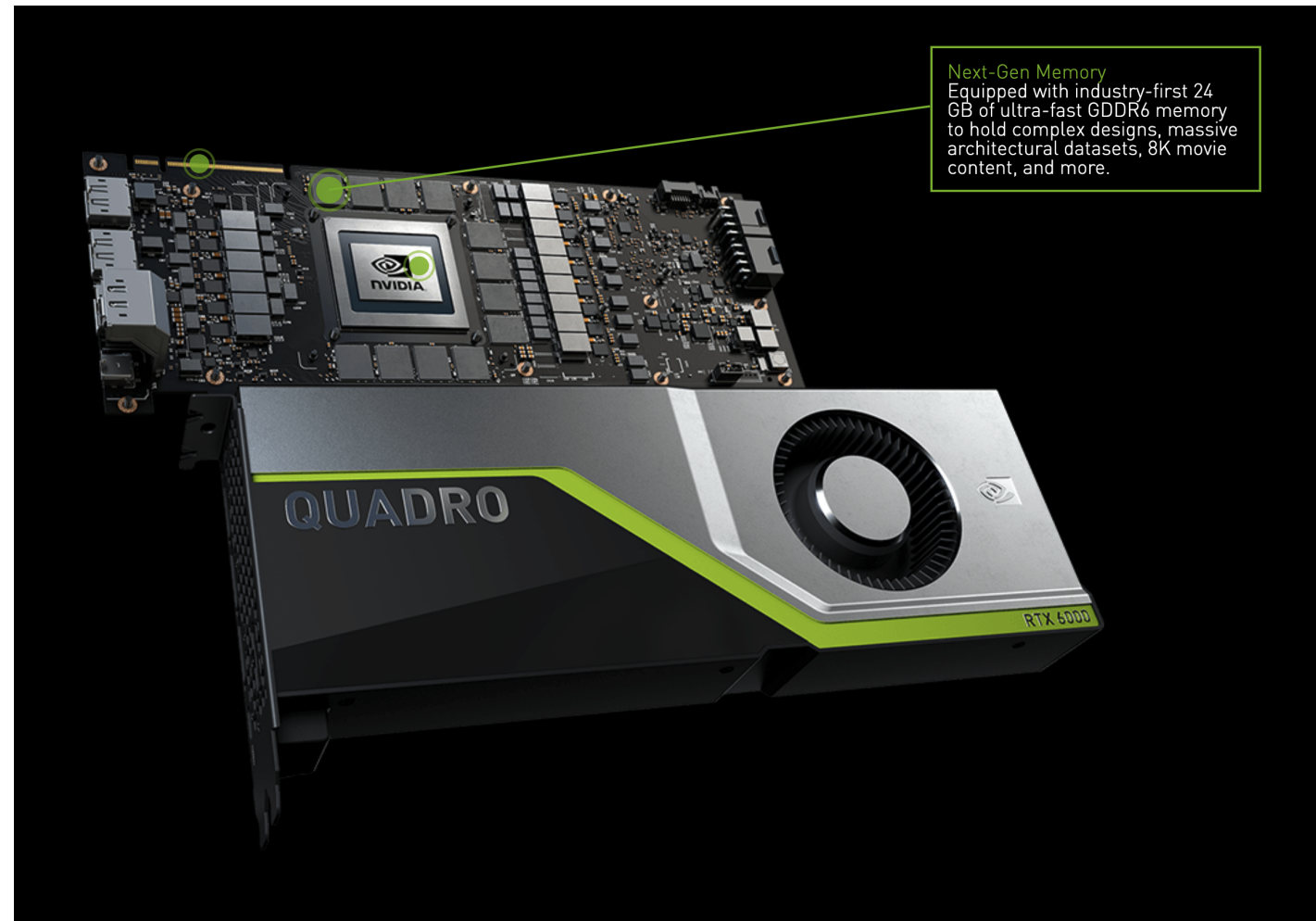
Tesla GPUs

Single precision GFLOPS vs. Model/Units





## Quadro RTX 6000 in icme-gpu



4,608 parallel processing cores; Turing architecture

- 72 Streaming multiprocessors (6 graphics processing clusters)
- 4,608 cores (64 per SM)
- Total global memory: 24 GB GDDR6
- Bandwidth: 672 GB/sec
- Peak performance: single 16.3 TFLOPS; double 510 GFLOPS (1/32)
- Compute capability: 7.5
- Power: 295 W
- Architecture: Turing
- L1 Cache: 96 KB (per SMX); L2 Cache: 6,144 KB

Quadro RTX 6000

Turing architecture

Whitepaper

Specs

NVIDIA GPUs are different from conventional processors.

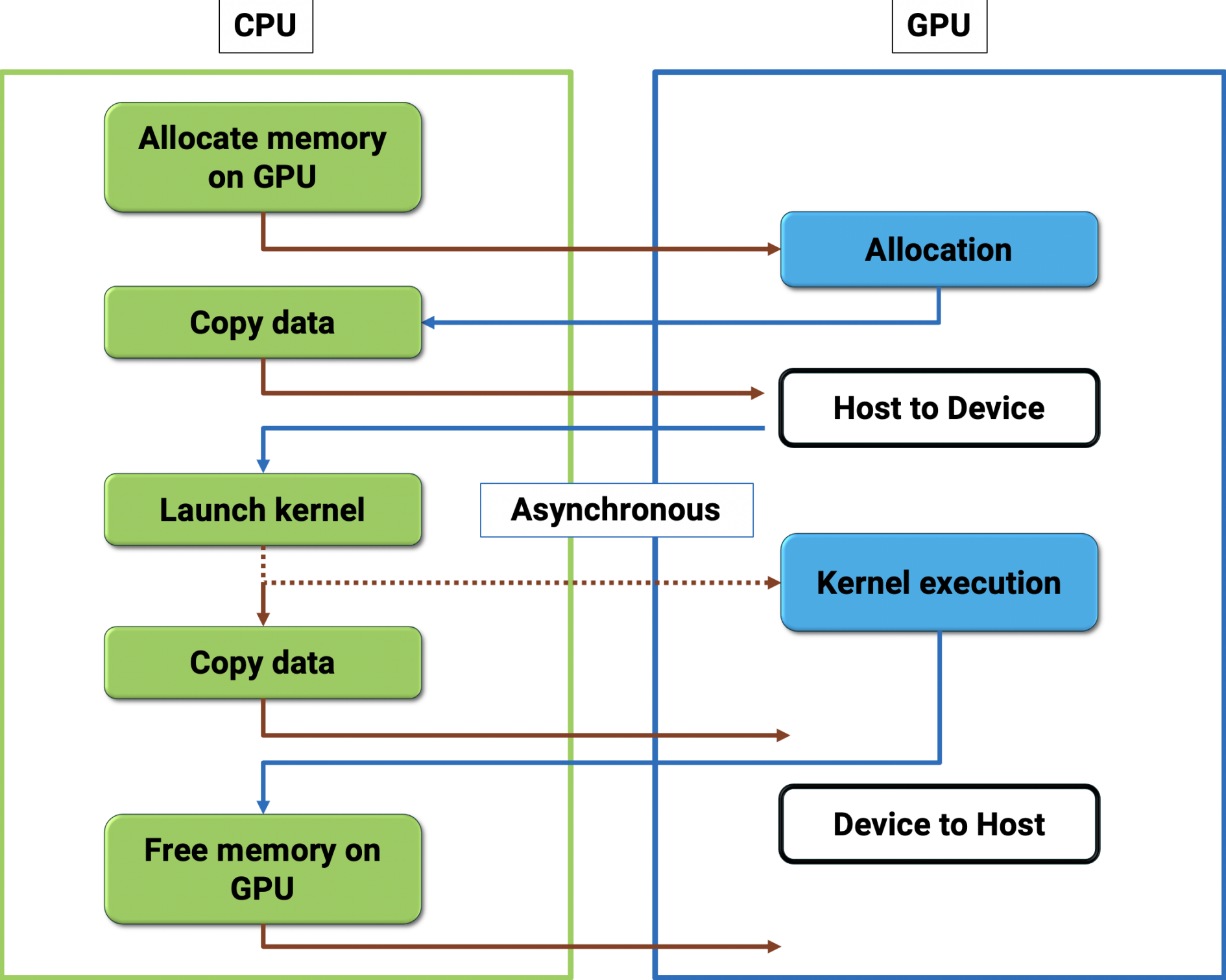
They only work as co-processors.

This means you need a host processor (e.g., Intel Xeon).

Your program runs on the host and uses the CUDA API to move data back and forth to the GPU and run programs on the GPU.

You cannot log on the GPU directly or run an OS on the GPU.





You cannot program a GPU without understanding the architecture of the processor.

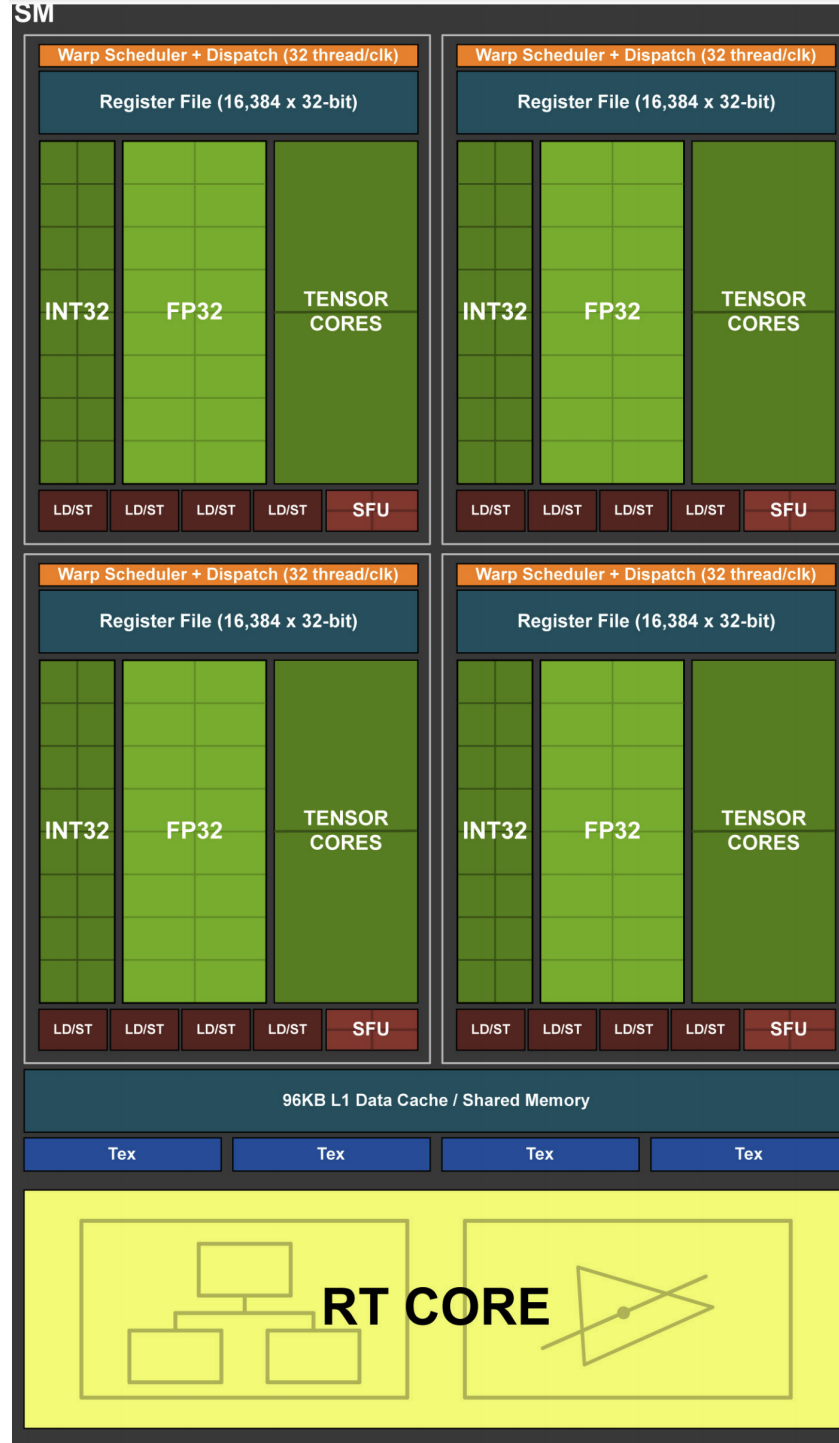
NVIDIA Turing

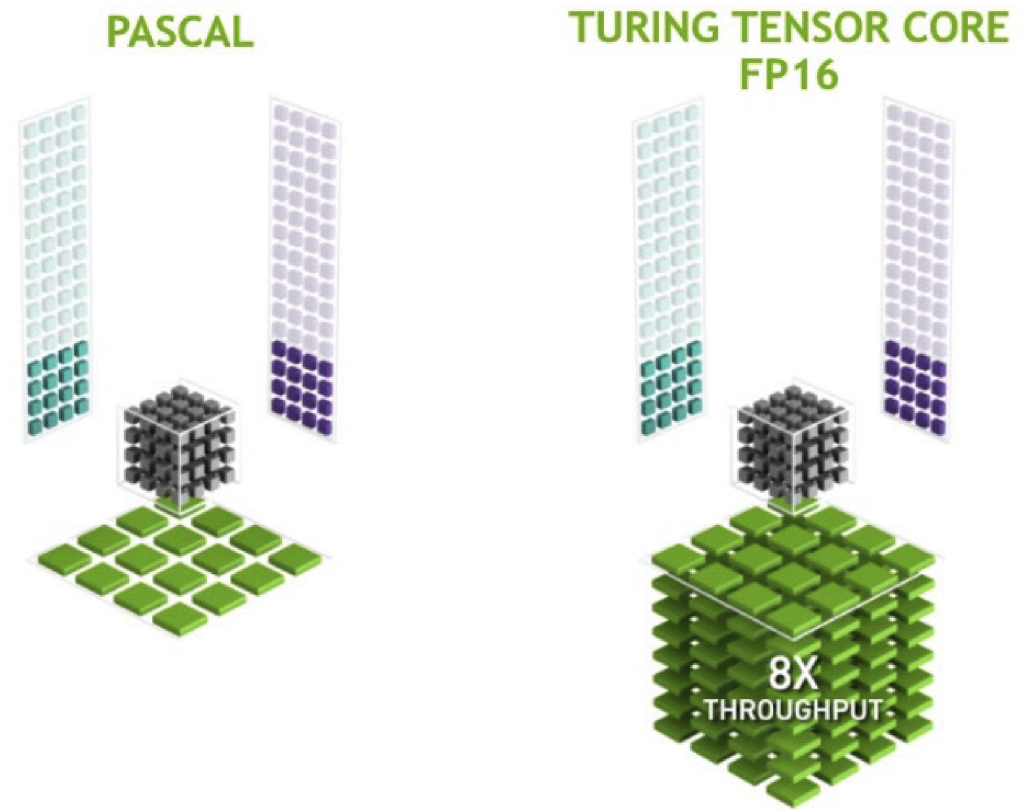
TU102 GPU architecture



- GPC: graphics processing cluster
- TPC: texture processing cluster
- SM: Streaming Multiprocessor
- RT core: ray-tracing core
- Tensor cores

SM





Tensor Core 4x4 Matrix Multiply and Accumulate

## Hierarchical breakdown of threads



## Warp

Bottom-most level: 32 cores = 32 threads, grouped in a warp.

The hardware is optimized to have all threads in a warp execute the same instruction at the same time.

SIMT (single instruction multiple threads) model.

## Block

Groups of warps form a block. A block executes on 1 SM (streaming multiprocessor).

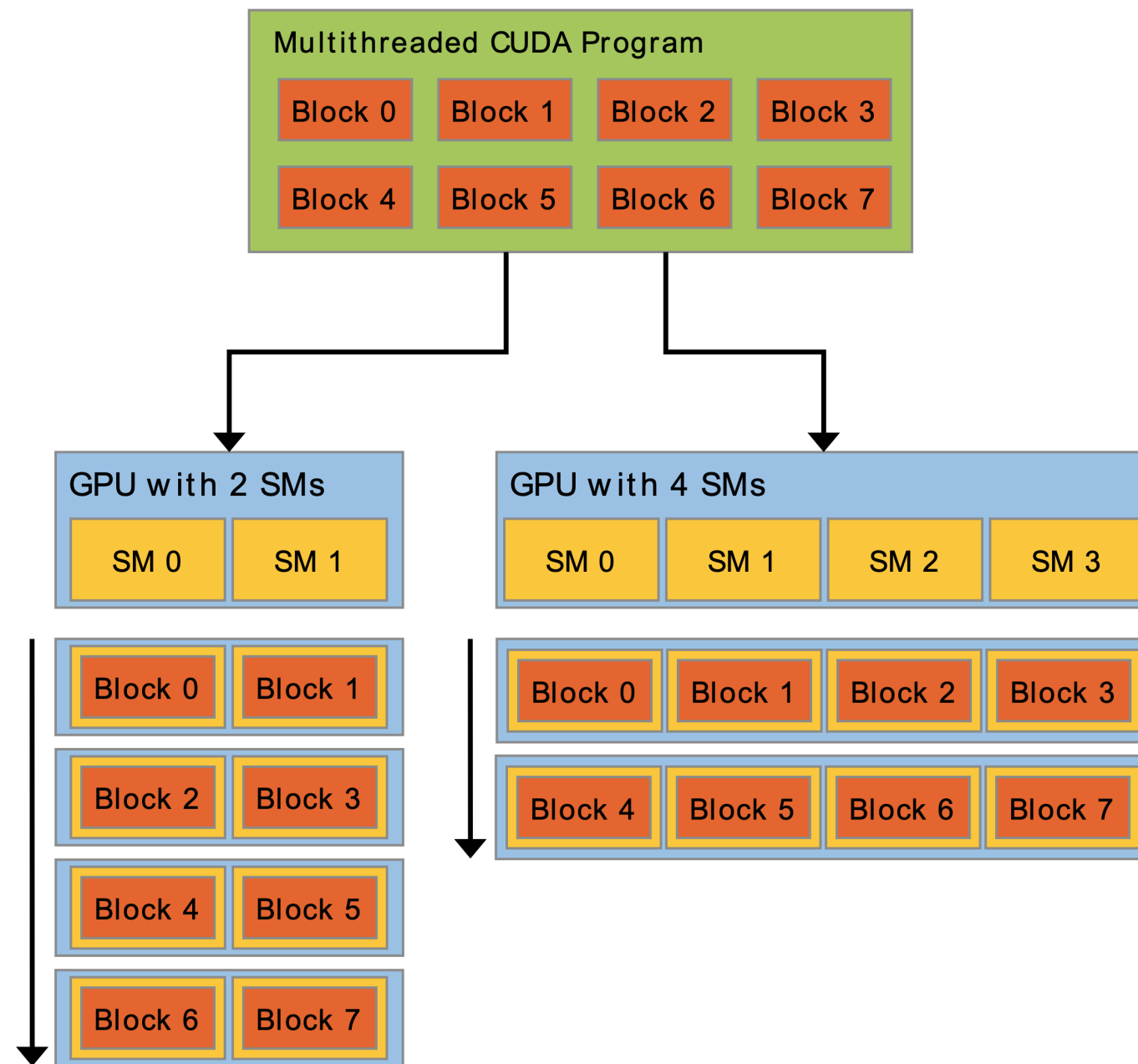
Threads within a block have some ability to exchange data and synchronize.



## Grid

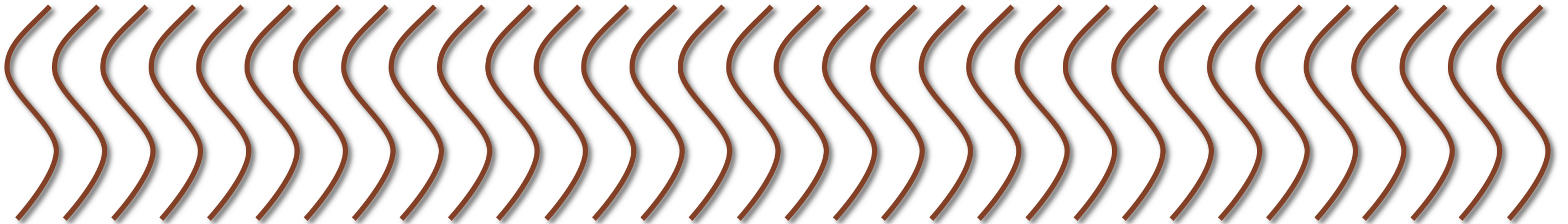
Collection of many blocks constitute the entire "dataset" that will be operated on by a kernel.





# Execution model

Input



Output