

CS 380 - GPU and GPGPU Programming

Lecture 2: Introduction, Pt. 2

Markus Hadwiger, KAUST

Reading Assignment #1 (until Sep 4)



Read (required):

- Orange book, chapter 1 (*Review of OpenGL Basics*)
- Orange book, chapter 2 (*Basics*)

What are GPUs?



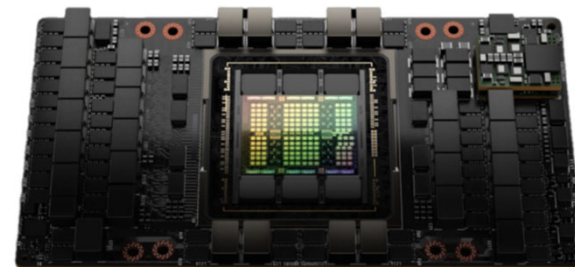
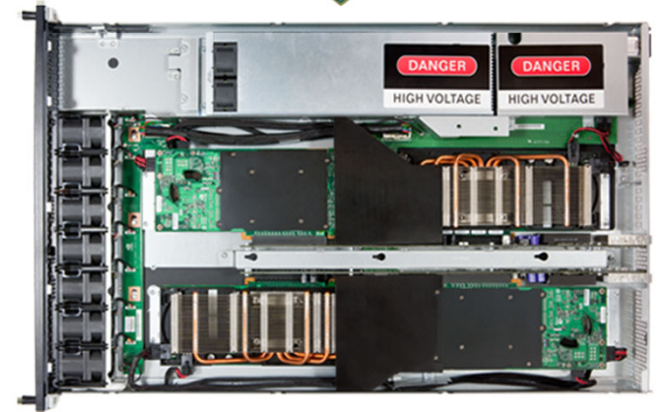
Graphics Processing Units

But evolved toward

- Very flexible, massively parallel floating point co-processors
- But not entirely programmable!
- Fixed-function parts have definite advantages (e.g., texture filtering, z-buffering)

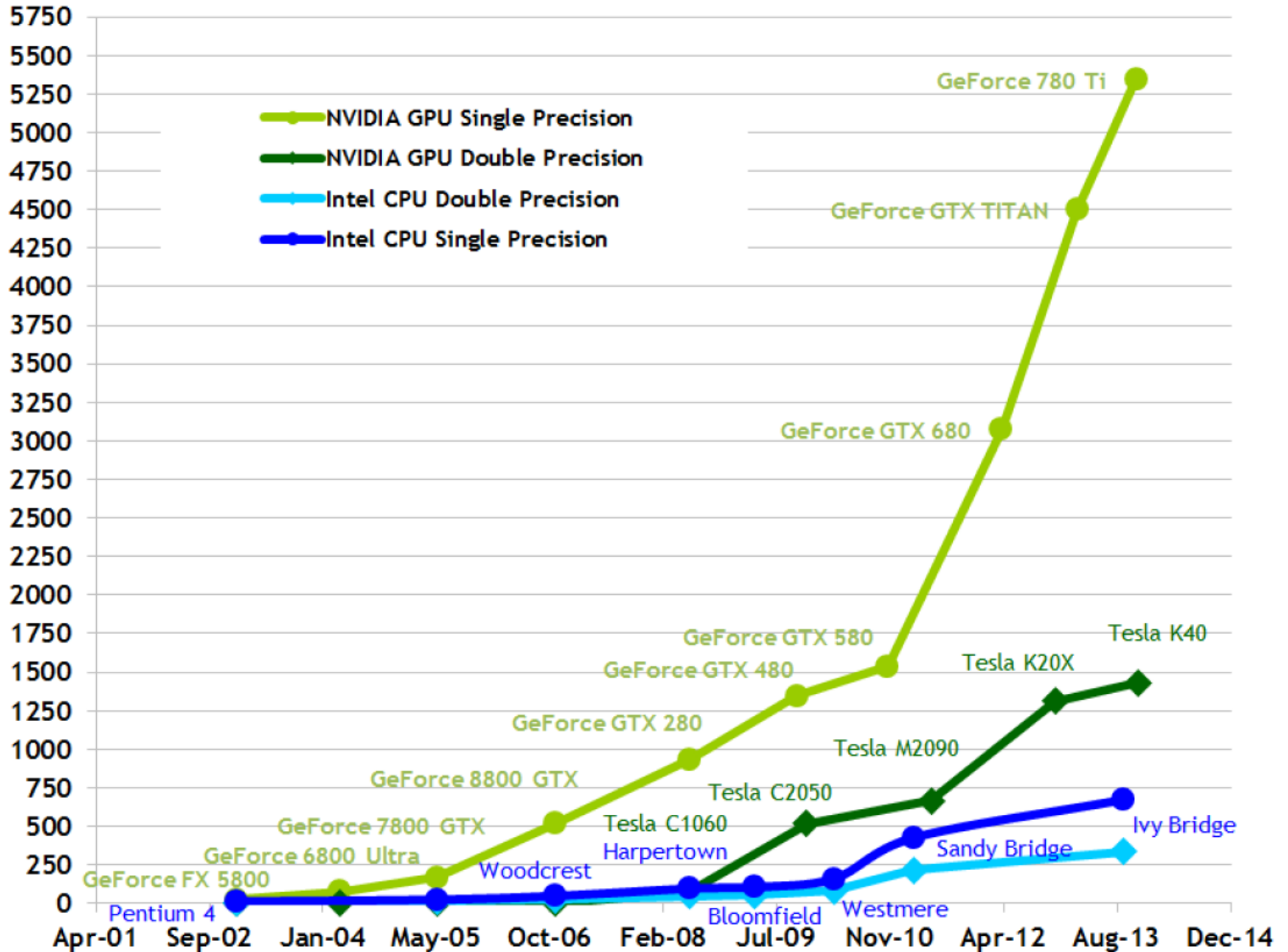
We will cover both perspectives

- GPUs for graphics
- GPU computing (GPGPU – general purpose computation on GPU)



Theoretical GFLOP/s

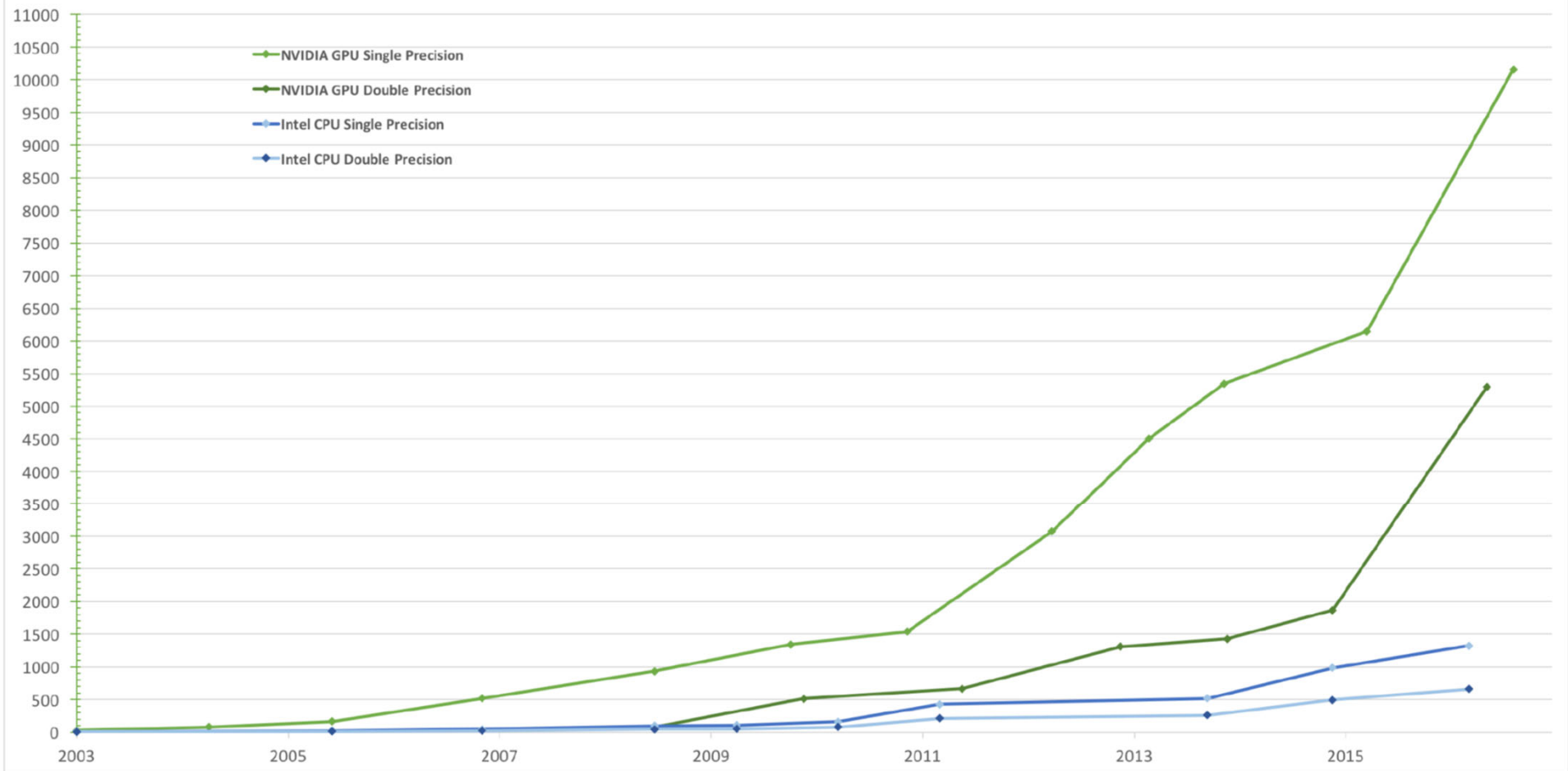
Peak Performance



Peak Performance

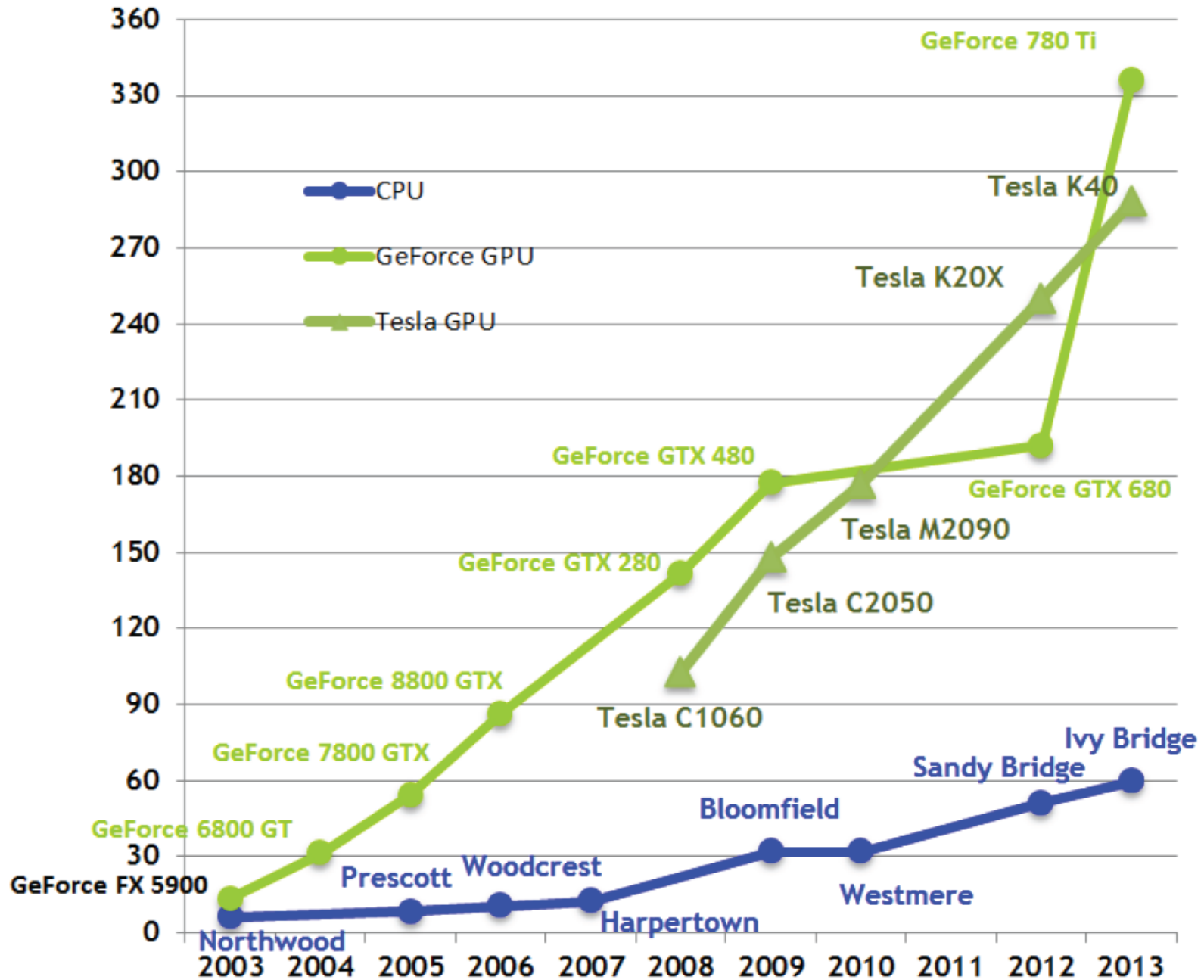


Theoretical GFLOP/s at base clock

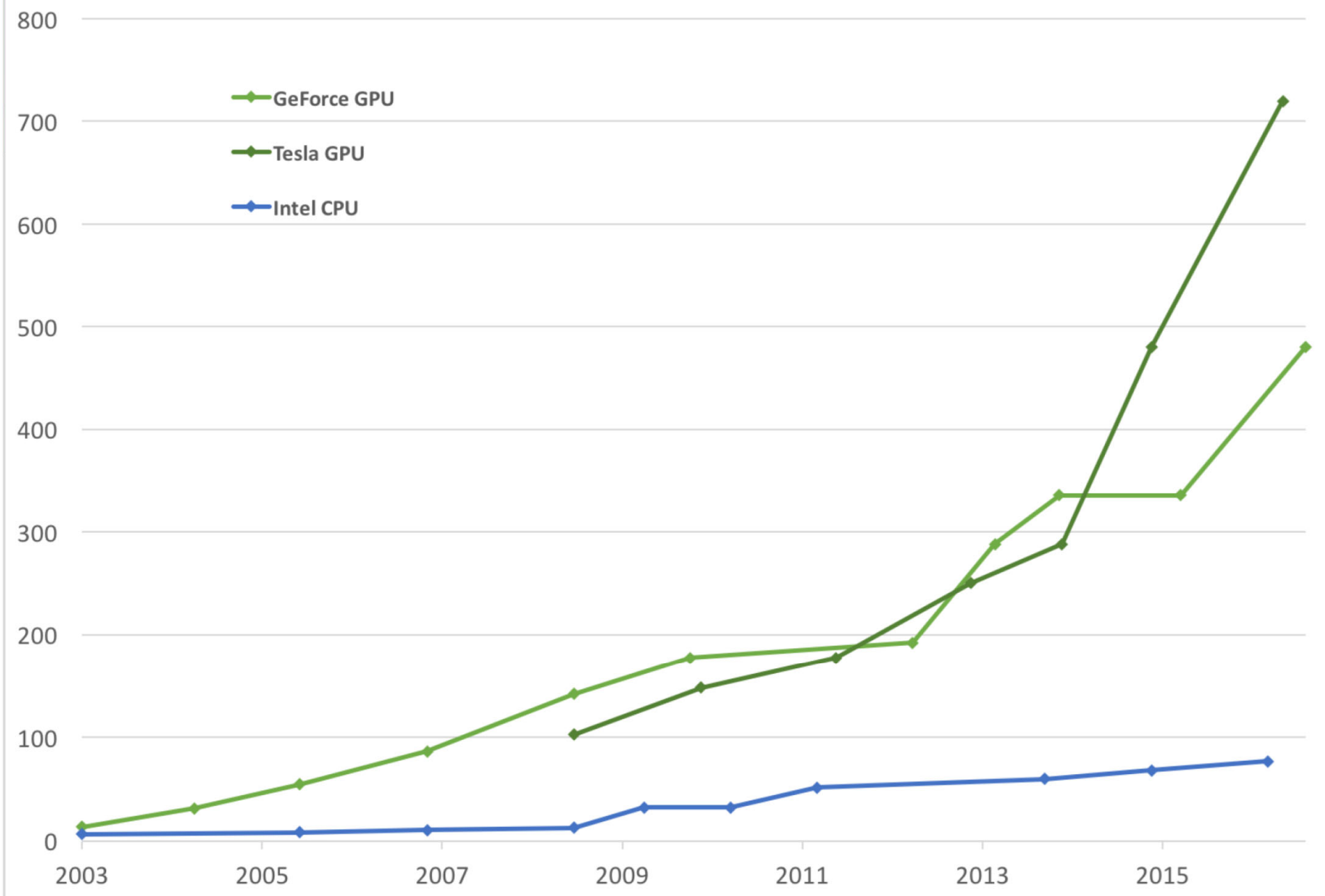


Theoretical GB/s

Peak Bandwidth

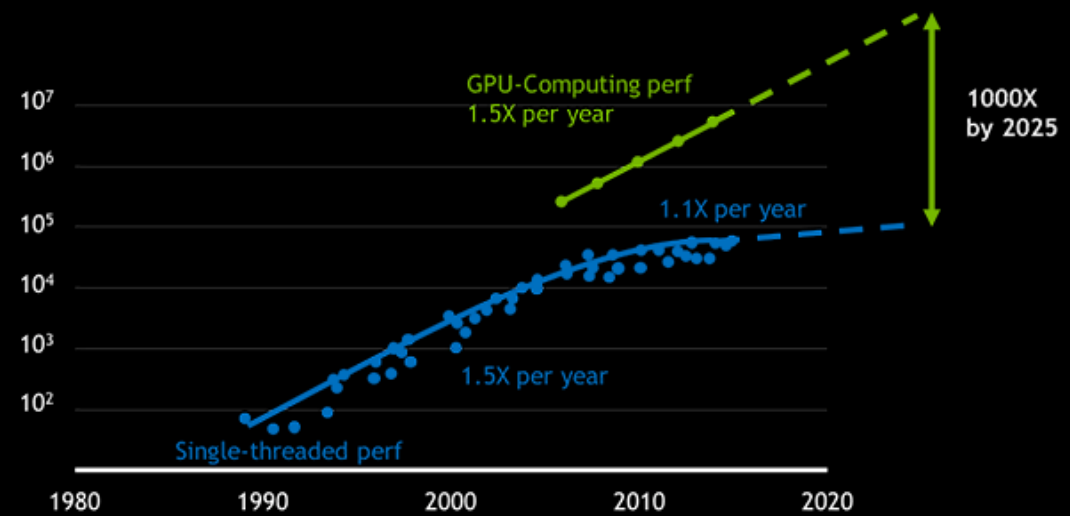
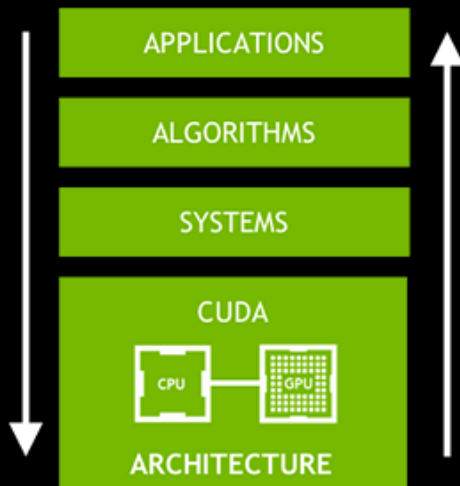


Theoretical Peak GB/s





RISE OF GPU COMPUTING

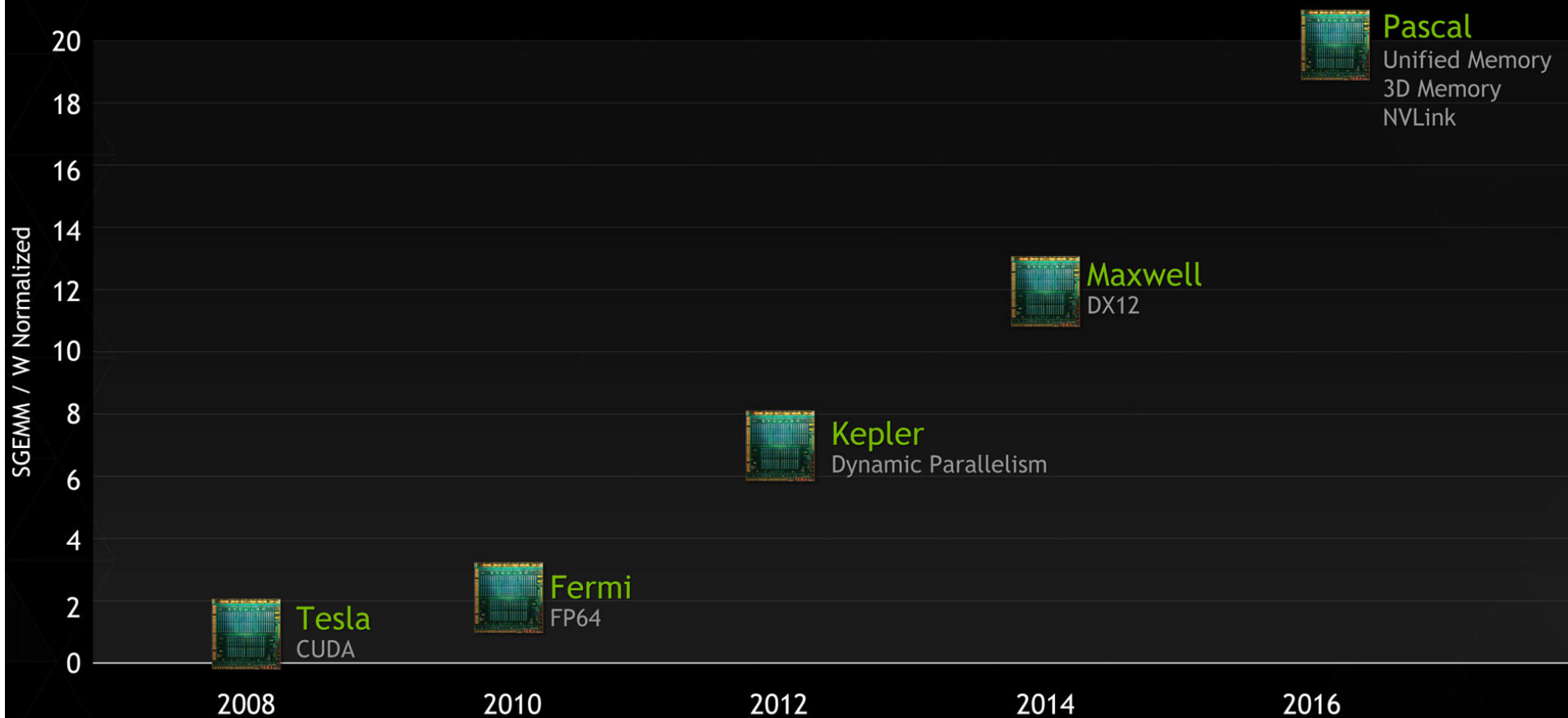


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten. New plot and data collected for 2010-2015 by K. Rupp.

GPU Architectures Over the Years



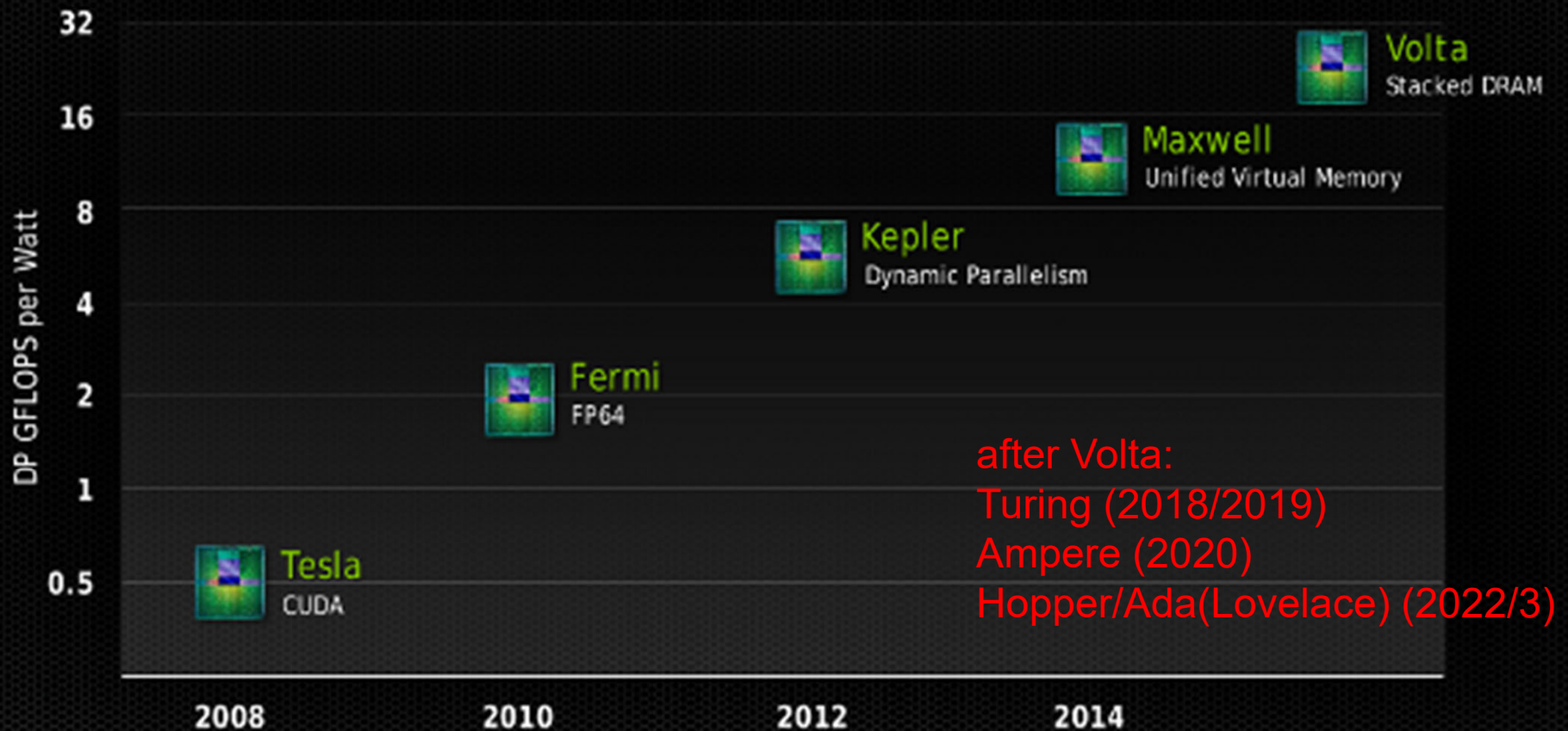
GPU Roadmap



GPU Architectures Over the Years



GPU Roadmap



Recent Updates (1)



NVIDIA Ampere architecture (2020)

[https://en.wikipedia.org/wiki/Ampere_\(microarchitecture\)](https://en.wikipedia.org/wiki/Ampere_(microarchitecture))

Initial presentation 2020:

<https://blogs.nvidia.com/blog/2020/05/14/gtc-2020-keynote/>

Geforce 30-series (Ampere):

<https://nvidia.com/en-us/geforce/graphics-cards/30-series/>

RTX 3090 has 10,496 FP32 CUDA cores (max arch would be 10,752)

A100 (Ampere):

<https://www.nvidia.com/en-us/data-center/a100/>

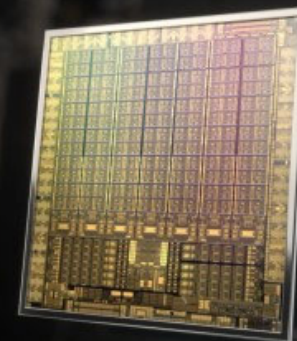
A100 has 6,912 FP32 CUDA cores (max arch would be 8,192)

Recent Updates (1)



NEW AMPERE ARCHITECTURE 2nd Generation RTX

28 Billion Transistors
30 Shader-TFLOPS | 58 RT-TFLOPS | 238 Tensor-TFLOPS
Micron G6X – World's Fastest Memory
Samsung 8N NVIDIA Custom Process



Recent Updates (2)



NVIDIA Hopper architecture (2022)

[https://en.wikipedia.org/wiki/Hopper_\(microarchitecture\)](https://en.wikipedia.org/wiki/Hopper_(microarchitecture))

Presentation [27:02] from NVIDIA GTC keynote, March 2022:

<https://www.youtube.com/watch?v=39ubNuxnrK8>

Hopper Whitepaper:

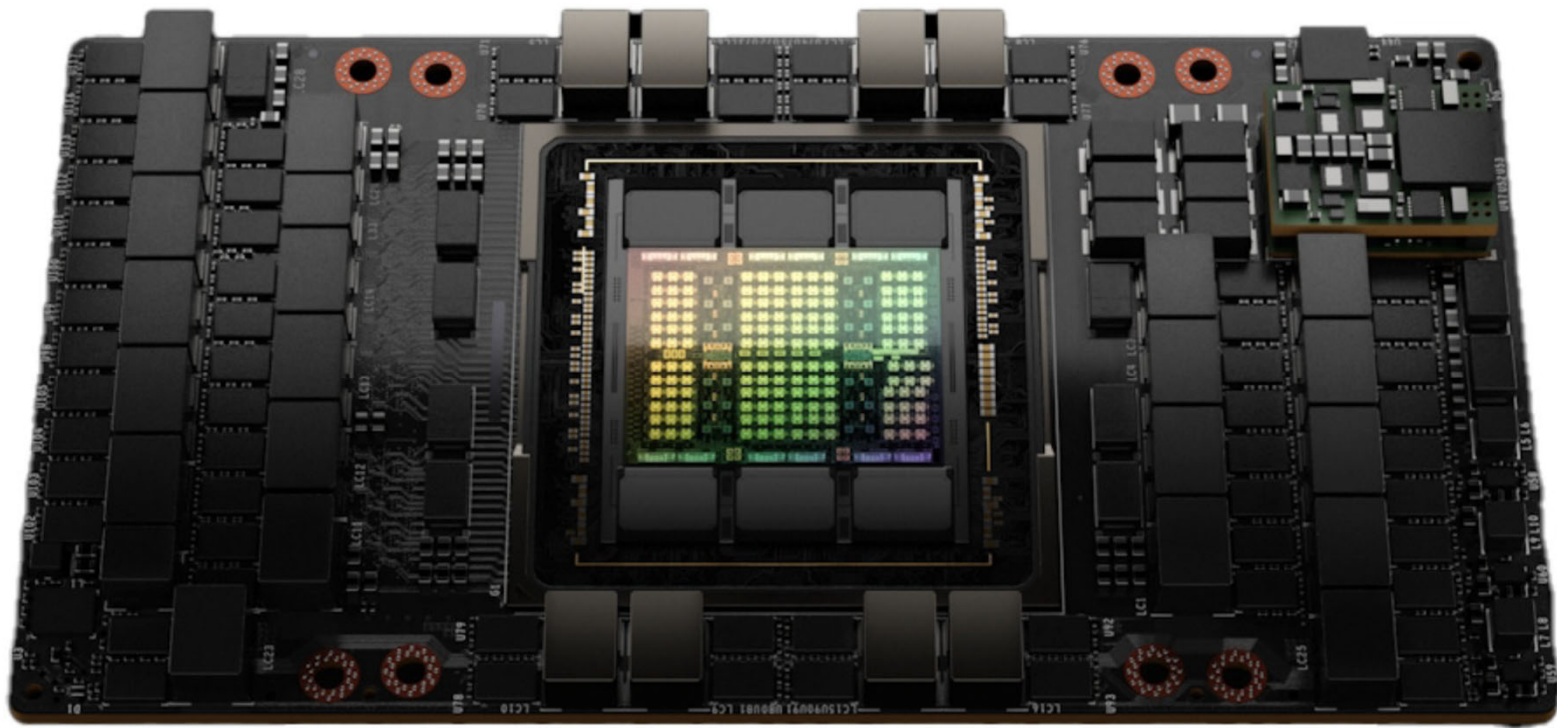
<https://resources.nvidia.com/en-us-tensor-core/gtc22-whitepaper-hopper/>

H100 (Hopper):

<https://www.nvidia.com/en-us/data-center/h100/>

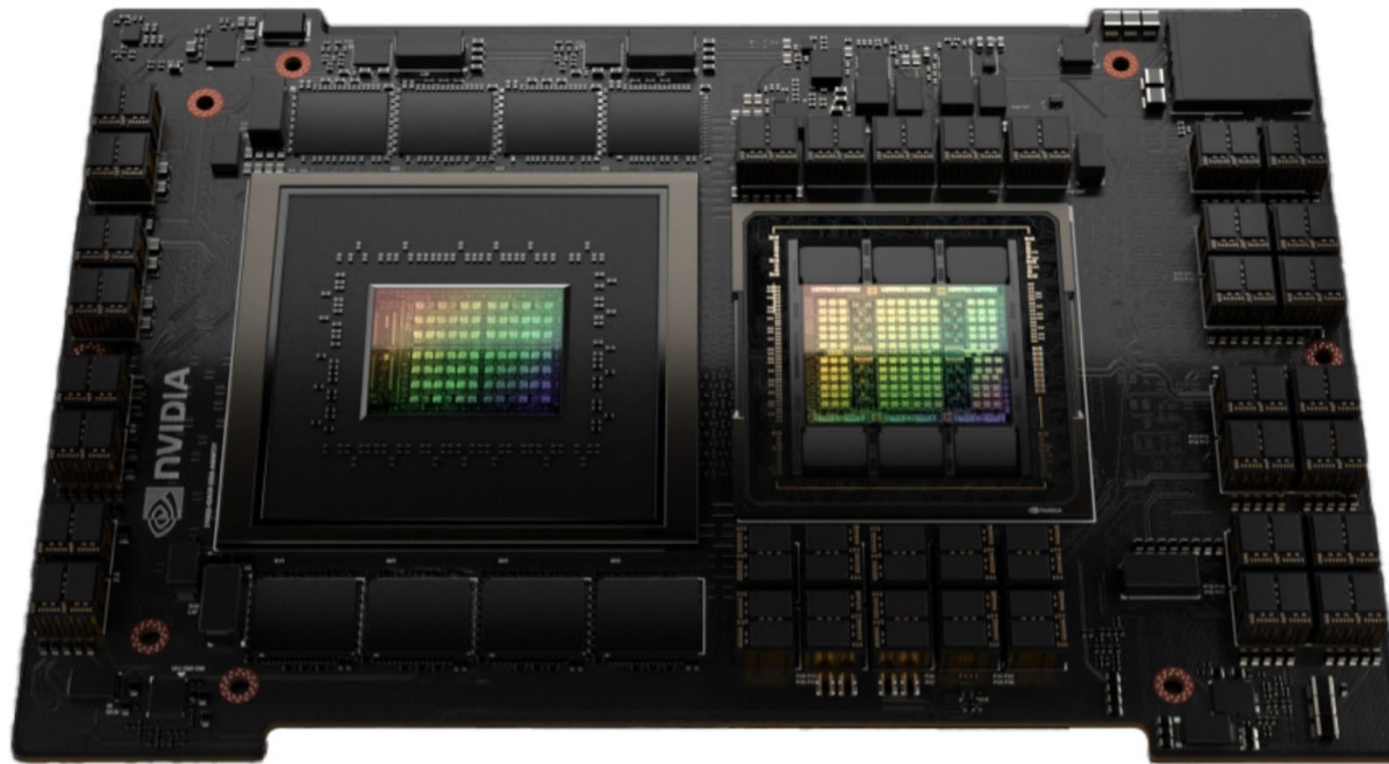
*H100 has up to 18,432 FP32 CUDA cores (max arch)
(H100 SXM5: 16,896; H100 PCIe: 14,592)*

Recent Updates (2)



NVIDIA H100 SXM5

Recent Updates (2)



NVIDIA Grace Hopper Superchip (Grace CPU + Hopper GPU)

NVIDIA Hopper GH100 Architecture (2022)



GH 100 (H100 Tensor Core GPU)

Full GPU: 144 SMs (in 8 GPCs/72 TPCs)



NVIDIA GH100 SM

CC 9.0 Multiprocessor

- 128 FP32 + 64 INT32 cores
- 64 FP64 cores
- 4x 4th gen tensor cores
- ++ thread block clusters, DPX insts., FP8, TMA

4 partitions inside SM

- 32 FP32 + 16 INT32 cores
- 16 FP64 cores
- 8x LD/ST units each
- 1x 4th gen tensor core each
- Each has: warp scheduler, dispatch unit, 16K register file



NVIDIA Hopper GH100 Architecture (2022)



GH 100 (H100)

Full GPU: 144 SMs (in 8 GPCs/72 TPCs)

- 64K 32-bit registers / SM = 256 KB register storage per SM
- 256 KB shared memory / L1 per SM

For 144 SMs on full GPU [SXM5: 132; PCIe: 114]

- 36 MB register storage, 36 MB shared mem / L1 storage = **72 MB context+”shared context” storage !**
- L2 cache size on H100: 50 MB
- 18,432 FP32 cores (128 FP32 cores per SM)
- 294,912 max threads in flight (max warps / SM = 64)

Recent Updates (3)



NVIDIA Ada (Lovelace) (2022/2023)

Upcoming GTC presentation September 19-22, 2022:

<https://www.nvidia.com/gtc/>

Geforce 40-series (speculative for now):

https://en.wikipedia.org/wiki/Draft:GeForce_40_series/

RTX 40xx has (not announced yet) FP32 CUDA cores

Overviews and Specs



Wikipedia has many comprehensive lists of architectures and specs

`https://en.wikipedia.org/wiki/
List_of_Nvidia_graphics_processing_units`

`https://en.wikipedia.org/wiki/
List_of_AMD_graphics_processing_units`

Thank you.