

## Lecture 8: looking to the future

Prof. Mike Giles

mike.giles@maths.ox.ac.uk

Oxford University Mathematical Institute

Important in scientific computing to keep an eye on what is happening with both hardware and software

(I am self-taught through reading lots of blogs and websites, as well as academic papers on scientific computing)

Remember: at times the business aspects are as important as the technical in thinking about how things are developing

Feb 2025 market capitalization (i.e. company value)

- **NVIDIA:** \$ 3.18 trn
- **AMD:** \$ 175 bn
- **Intel:** \$ 82 bn

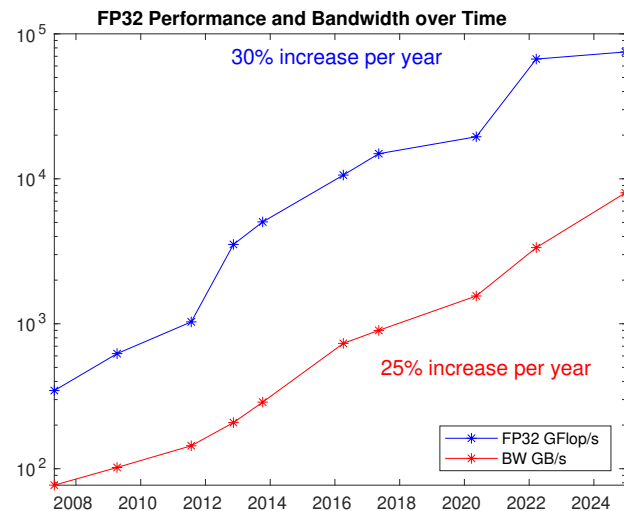
10 years ago the order would have been reversed!

Lecture 8 – p. 1/28

Lecture 8 – p. 2/28

## Hardware trends

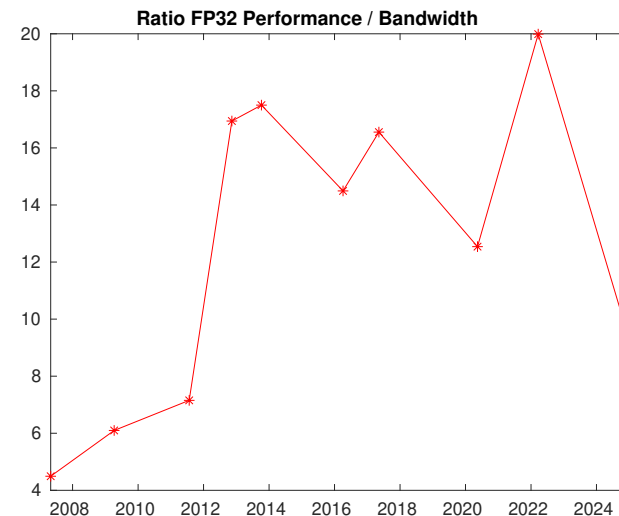
NVIDIA high-end GPU performance and bandwidth



Lecture 8 – p. 3/28

## Hardware trends

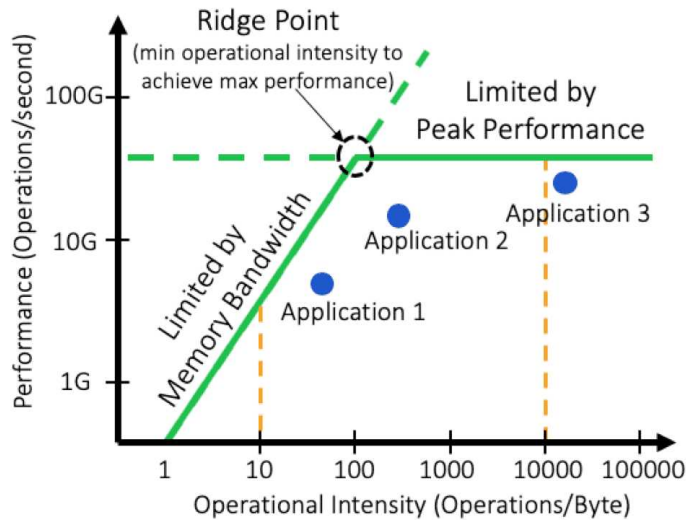
Compute / bandwidth ratio



Lecture 8 – p. 4/28

# Hardware trends

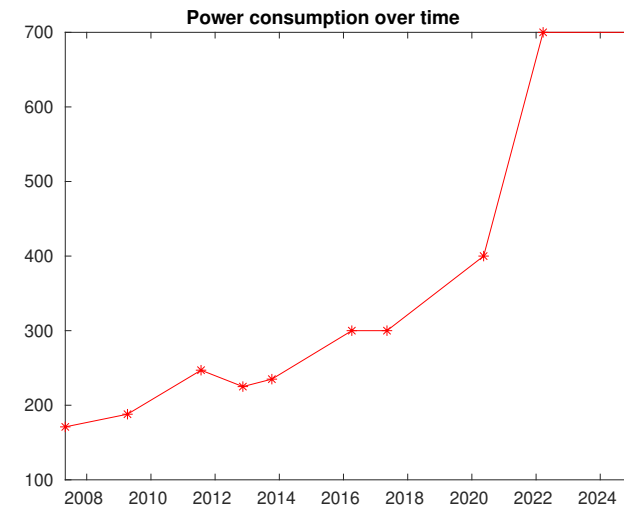
Roofline model (image copyright Rambus Inc.)



Lecture 8 – p. 5/28

# Hardware trends

Increasing energy consumption by NVIDIA GPUs – moving to chilled-water cooling blocks



Lecture 8 – p. 6/28

# NVIDIA

Market Summary > NVIDIA Corp

3,18 trillion USD

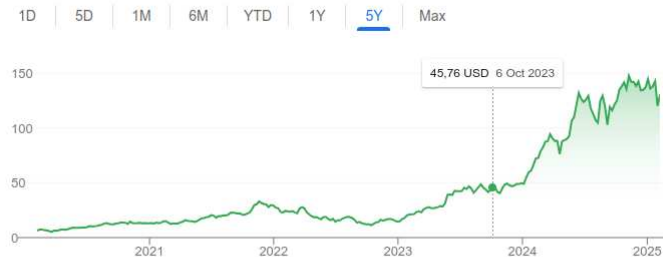
Market capitalisation

129,84 USD

+123.55 (1,964.23%) ↑ past 5 years

Closed: 07 Feb, 20:00 GMT-5 • Disclaimer

After hours 128,73 -1,11 (0,86%)



Open	129,22	Mkt cap	3,18T	52-wk high	153,13
High	130,37	P/E ratio	51,16	52-wk low	66,25
Low	125,00	Div yield	0,031%		

Lecture 8 – p. 7/28

# NVIDIA

Ampere came out in 2020:

- A100 with 108 SMs, 40-80 GB HBM2 memory
- wide range of “tensor core” capabilities
- NVIDIA DGX A100 Deep Learning server
  - <https://www.nvidia.com/en-us/data-center/dgx-a100/>
  - 8 NVIDIA A100 GPUs, each with 80GB HBM2
  - 2 × 64-core AMD “Rome” CPUs
  - 2 TB DDR4 memory, 30 TB SSD
  - 600GB/s NVlink interconnect between the GPUs

Lecture 8 – p. 8/28

# NVIDIA

- Hopper came out in 2023:
  - H100 for HPC
  - 228-264 SMs
  - 80GB HBM3 memory
  - 40MB L2 cache
  - NVlink improvements – up to 50% faster, 900GB/s
  - PCIe v5.0 – 2× improvement
- Grace CPU also arrived in 2023:
  - Arm-based
  - up to 72 cores
  - 550GB/s bandwidth to LPDDR5X memory
  - 900GB/s NVlink connection to Hopper GPU in GB100 “superchip”

Lecture 8 – p. 9/28

# NVIDIA

- A Hopper refresh came out in 2024:
  - H100NVL – two GPUs on one card
  - 2 × 96GB HBM3 memory
  - 2 × 3.9TB/sec bandwidth
  - aimed particularly at LLMs needing lots of memory and bandwidth
- Late 2024, the Blackwell B100 came out:
  - 192 GB memory, 8192-bit bus, 8TB/s bandwidth
  - modest improvement in compute except for low-precision tensor cores
  - persistent rumours of over-heating

Lecture 8 – p. 10/28

# NVIDIA

Current status:

- big AI companies are competing to buy huge numbers (10,000+) of Hopper H100 and Blackwell B100 GPUs – many orders are worth over \$1bn
- supply is limited, prices have become inflated, and it's very difficult for academics to get any
- emergence of Grace CPU is significant – gives NVIDIA freedom to design their own combined CPU/GPU offerings with high bandwidth interconnect

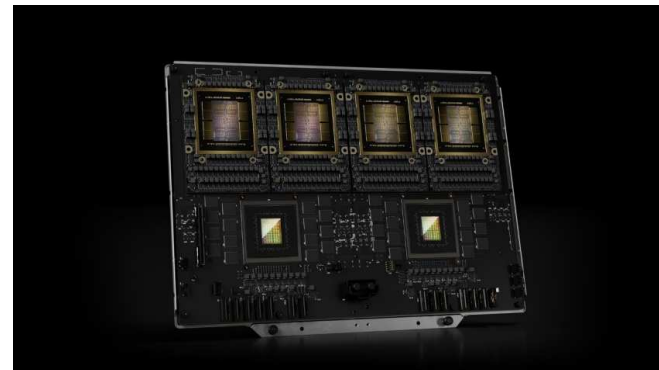
(also part of ARM breakthrough into the server market? hyperscalers are responsible for 50% of the server market, and 50% of their CPUs are Arm-based)

Lecture 8 – p. 11/28

# NVIDIA

Coming later this year:

- GB200 superchip combining 1 Grace CPU / 2 B200s
- [GB200 NVL72](#) rack with 36 GB200s (72 B200 GPUs) connected by an NVLINK switch
- also [GB200 NVL4](#) card with 2 Grace CPUs / 4 B200s



Lecture 8 – p. 12/28

# NVIDIA

Coming later this year – NVIDIA DIGITS system:



Lecture 8 – p. 13/28

# NVIDIA



“A Grace-Blackwell AI Supercomputer on your desk”:

- complete small desktop system – \$3000 in US
- GB10 combined CPU/GPU superchip
- 20 Arm cores in CPU with 128GB of DDR5X memory
- 1 PFlop FP4 (!) Blackwell GPU – 7% of B100
- 4 TB SSD and Linux operating system

I want one for Christmas!

Lecture 8 – p. 14/28

# AMD

Market Summary > Advanced Micro Devices Inc

174,55 billion USD

Market capitalisation

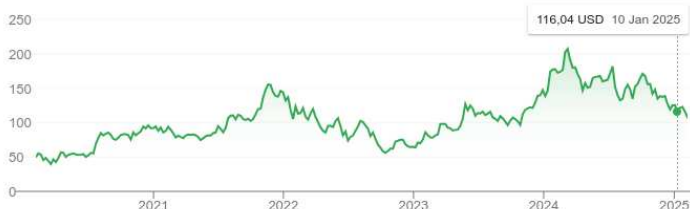
107,56 USD

+57.83 (116.29%) ↑ past 5 years

Closed: 07 Feb, 19:59 GMT-5 • Disclaimer

After hours 106,98 -0,58 (0,54%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max



Open	109,13	Mkt cap	174,55B	52-wk high	227,30
High	109,92	P/E ratio	107,30	52-wk low	106,50
Low	106,79	Div yield	-		

Lecture 8 – p. 15/28

# Top500

Top 5 on Top500 list, November 2024:

- #1 El Capitan (DoE/LLNL, USA)
  - HPE: 44,000 AMD MI300A GPUs
- #2 Frontier (DoE/ORNL, USA)
  - HPE: 40,000 AMD MI250X GPUs
- #3 Aurora (DoE/ALC, USA)
  - Intel/HPE: 54,000 Intel Max GPUs
- #4 Eagle (Microsoft Azure)
  - Microsoft: NVIDIA H100 GPUs
- #5 HPC6 (Eni, Italy)
  - HPE: 14,000 AMD MI250X GPUs

Lecture 8 – p. 16/28



**El Capitan:** #1 supercomputer based on Linpack performance

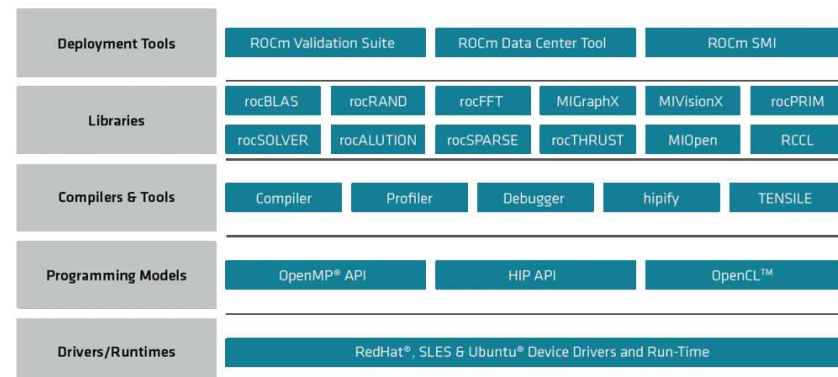
- sited at Lawrence Livermore National Laboratory (DoE)
- 1.7 Exaflops, 30 MW
- system from HPE; CPUs and GPUs from AMD
- 11,136 compute nodes, each with 4 MI300A GPUs

- over past decade AMD has had excellent CPUs and GPUs (and pioneered chiplet packaging) but has not invested enough in software – that is changing
- hired lots of software specialists in the past 2 years, including many of the NAG team responsible for ACML (AMD’s version of Intel’s MKL libraries)
- “Genoa” Zen4 EPYC CPUs:
  - up to 64 cores with vector units and 384MB L3
  - now getting about 20% share of server market
- Frontier has previous generation “Trento” Zen3 EPYC CPUs

• **Instinct GPUs:**

- MI250X has 220 Compute Units, 128 GB HBM2e, 3.2 TB/s: comparable to A100 GPU [for PyTorch](#)
- MI300A has 228 Compute Units, 912 Matrix Cores, 128 GB HBM2e, 5.3 TB/s: comparable to H100?
- MI325X has 304 Compute Units, 1216 Matrix Cores, 256 GB HBM3e, 6 TB/s: comparable to B100?
- programmed using AMD’s ROCm (similar to CUDA) with extensive library support
- portability provided through [HIP](#) (Heterogeneous computing Interface for Portability) with compilation to either CUDA or AMD’s ROCm

AMD’s ROCm eco-system:



## AMD's HIP – some example code:

```
char* inputBuffer;
char* outputBuffer;

hipMalloc((void**)&inputBuffer, (strlen+1)*sizeof(char));
hipMalloc((void**)&outputBuffer, (strlen+1)*sizeof(char));

hipMemcpy(inputBuffer, input, (strlen+1)*sizeof(char),
          hipMemcpyHostToDevice);

hipLaunchKernelGGL(helloworld, dim3(1),dim3(strlen), 0, 0,
                  inputBuffer, outputBuffer );

hipMemcpy(output, outputBuffer, (strlen+1)*sizeof(char),
          hipMemcpyDeviceToHost);

hipFree(inputBuffer);
hipFree(outputBuffer);
```

Lecture 8 – p. 21/28

- ROCm and HIP look very similar to CUDA – probably required to win the major DoE and EU contracts
- pricing and availability of GPUs are both much better than NVIDIA currently, especially for academics  
(major AI companies are placing \$1bn orders with NVIDIA so no GPUs left for us!)
- AMD's software eco-system is still maturing – will take at least another 5 years to get close to CUDA
- still, very good to see competition in the marketplace

Lecture 8 – p. 23/28

## Now for some kernel code:

```
__global__ void helloworld(char* in, char* out)
{
    int num = hipThreadIdx_x + hipBlockDim_x * hipBlockIdx_x;
    out[num] = in[num] + 1;
}
```

Can see why it is fairly easy for AMD's HIPIFY tool to convert most simple CUDA code to HIP – this is another reason to avoid “exotic” CUDA features as much as possible.

Warning: AMD GPUs have a warp size of 64, not 32, so use `warpSize` variable in your code rather than hard-coding a warp size of 32.

Lecture 8 – p. 22/28

Market Summary > Intel Corp

**82,38 billion** USD

Market capitalisation

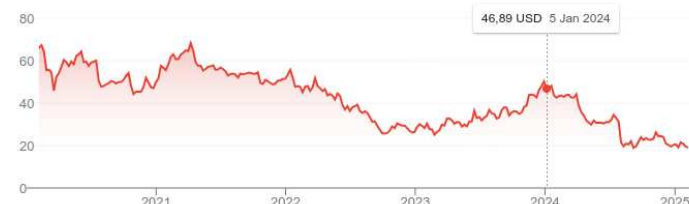
**19,10** USD

-46.92 (-71.07%) ↓ past 5 years

Closed: 07 Feb, 20:00 GMT-5 • Disclaimer

After hours 19,08 -0,020 (0,10%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max



Open	19,35	Mkt cap	82,38B	52-wk high	46,63
High	19,36	P/E ratio	-	52-wk low	18,51
Low	19,03	Div yield	2,62%		

Lecture 8 – p. 24/28

## Intel

- current “Granite Rapids” Xeon-SP CPUs:
  - up to 128 cores, each with one or two 512-bit AVX-512 vector units per core (512 bits = 16 floats)
  - up to 500MB L3 (shared), 2MB L2 per core
  - up to 600 GB/s memory bandwidth with new DDR5 MRDIMM memory
- “Ponte Vecchio” a.k.a. Data Center GPU Max:
  - 128 Xe cores, each with  $8 \times 256$ -bit vector units and 8 tensor cores
  - 408MB L2 cache, 128GB HBM2e with 8192-bit bus
  - development of “Rialto Bridge” successor ended; not clear what Intel will do now

Lecture 8 – p. 25/28

## Outlook

My current software assessment:

- CUDA is dominant in HPC because of
  - ease-of-use
  - NVIDIA dominance of hardware, with huge sales in machine learning in particular
  - extensive library support
  - support for many different languages (Fortran, Python, R, MATLAB, etc.)
  - extensive eco-system of tools
- HIP is a real threat to that dominance by offering platform independence with compilation to both CUDA and AMD’s ROCm

Lecture 8 – p. 27/28

## Others

Special designs, solely for the needs of Machine Learning:

- Google: Tensor Processing Unit (TPU)
- Graphcore: Colossus Intelligent Processing Unit
- Cerebras: in-memory computing (lots of computing elements interspersed within a huge amount of memory in wafer-scale chips)

It seems unlikely that Google will get into the hardware business in a big way, and if any startup makes real progress they’ll be bought out by NVIDIA, AMD or Intel.

Lecture 8 – p. 26/28

## Final thoughts

- NVIDIA holds a dominant market position, maybe hard to justify their huge market valuation but they’re the leader for a good reason – they have excellent hardware and software, and focussed early on the needs of AI/ML  
Original gaming market no longer significant, the auto market is the next big one they’re working on
- By addressing their software weakness, AMD is back in the game for both HPC and AI/ML – great to have competition again
- I remain unconvinced by Intel’s new hardware and software products, though traditional Xeon CPUs remain powerful and sell well
- Other vendors are unlikely to break through significantly

Lecture 8 – p. 28/28