

Introductory meeting

Sep 11th 2023

Bible verse of the day

- You, O God, do see trouble and grief; you consider it to take it in hand. The victim commits himself to you; you are the helper of the fatherless. —[Psalm 10:14](#)



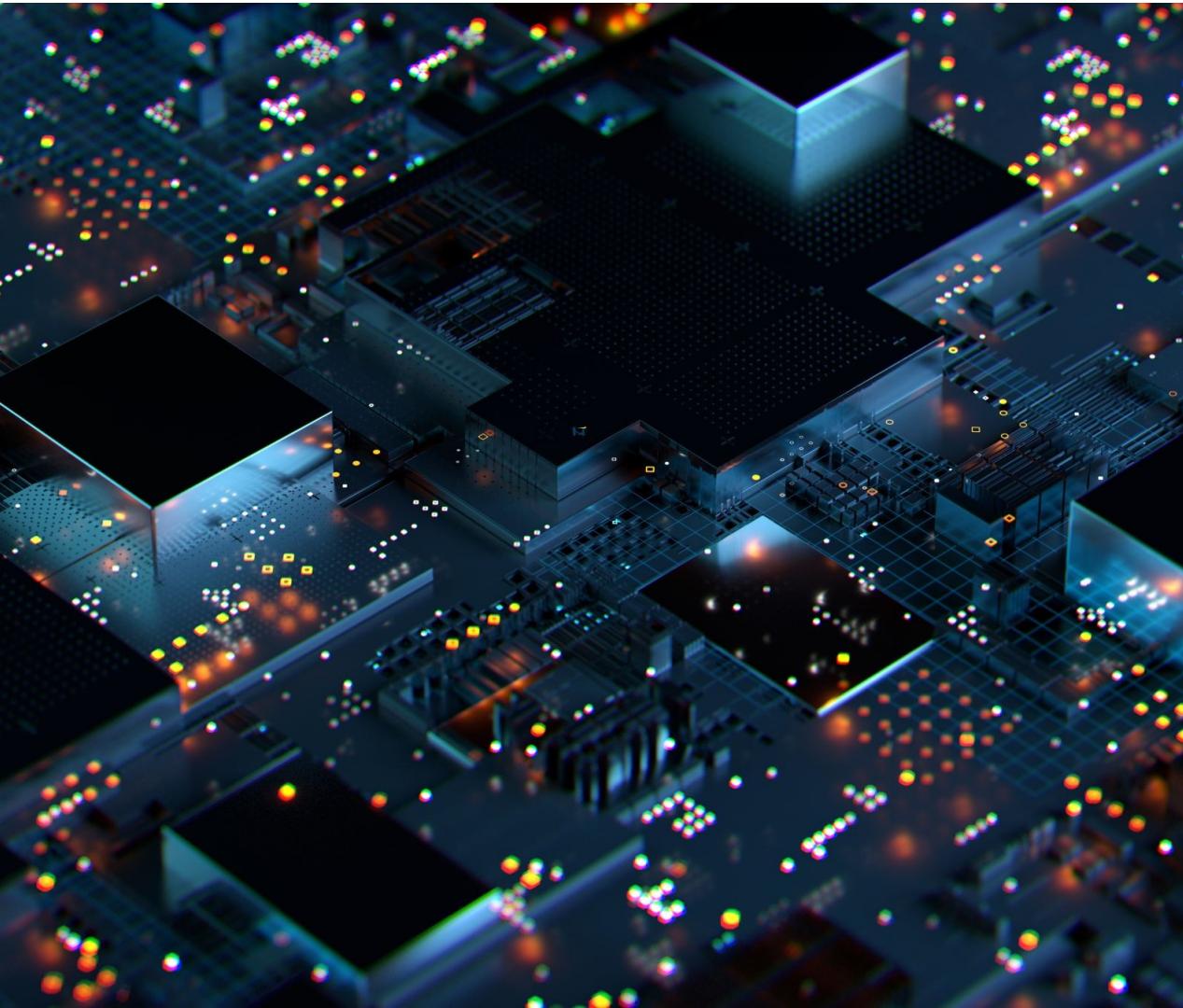
Introduction for HPC concepts

- Hardware
- OS
- Software
- Networking

What is an HPC?

- High Performance Computing
- Primarily associated with fast computers
- Differences occur in multi threaded vs single threaded calculations

History of HPC ~A rough overview 80s-2ks



- Cray Super Computers
 - Primarily focused on really fast single threaded computers
- 90's led to Beowulf clusters
 - Cheap consumer hardware networked together
 - Tasks at this point were becoming very multi-threaded
- 2000's-present
 - The majority of super computers have both a gpu and cpu presence
 - Multicore processors led to a huge increase in performance



What is a computer cluster?

- Two or more computers networked together to achieve some task
 - At its simplest form a LAN with two computers can be considered a cluster



Multi-threaded vs Single-threaded

- If the process can use more than one core it is multi-threaded
- As we have encountered a blockage in recent history preventing cpus from massively gaining single threaded speed to gain more performance software relies on being programmed to utilize more than one core

How to get more cores?

- Beowulf clusters
 - Take a bunch of cheap consumer computers and network them together
 - A lot of research was done on this in the 90's



Different methods of achieving multi core over a network

MPI

- Openmpi
- Mpicc

Mosix

- Not free

kerrighed

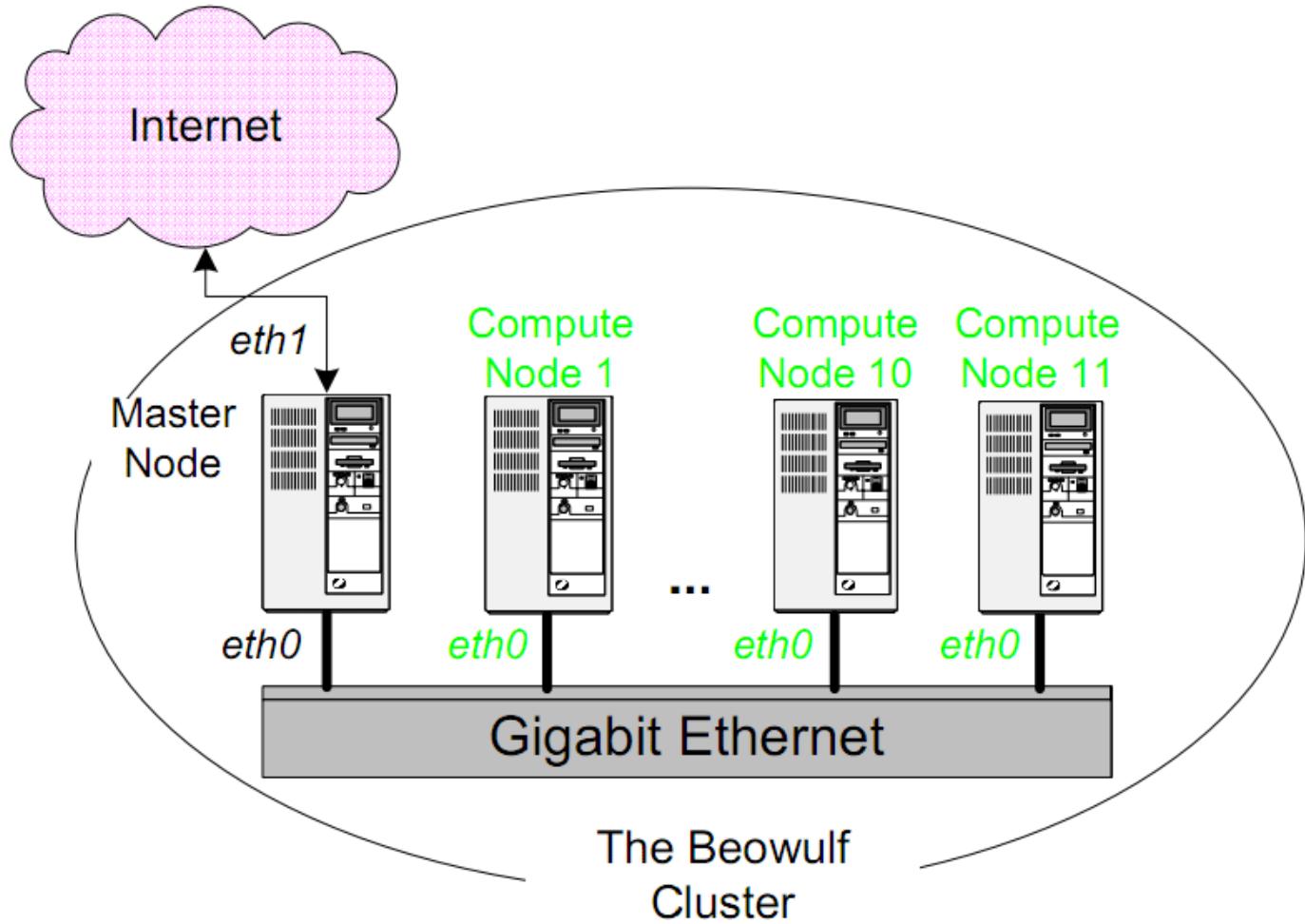
- Not supported



We primarily use MPI

- MPI is easy to setup
- There is much support for it
- You build your program with MPI support

Beowulf cluster network diagram



Back to the building blocks part 1

- Hardware
- OS
- Software
- Networking

Back to the building blocks part 2

- Hardware (cheap computers with cpu cores)
- OS
- Software
- Networking

Back to the building blocks part 3

- Hardware (cheap computers with cpu cores)
- OS (rocky linux, although can be any operating system)
- Software
- Networking

Back to the building blocks part 4

- Hardware (cheap computers with cpu cores)
- OS (rocky linux, although can be any operating system)
- Software (MPI, openMPI)
- Networking

Back to the building blocks part 5

- Hardware (cheap computers with cpu cores)
- OS (rocky linux, although can be any operating system)
- Software (MPI, openMPI)
- Networking (Ethernet switches, high end uses infiniband)



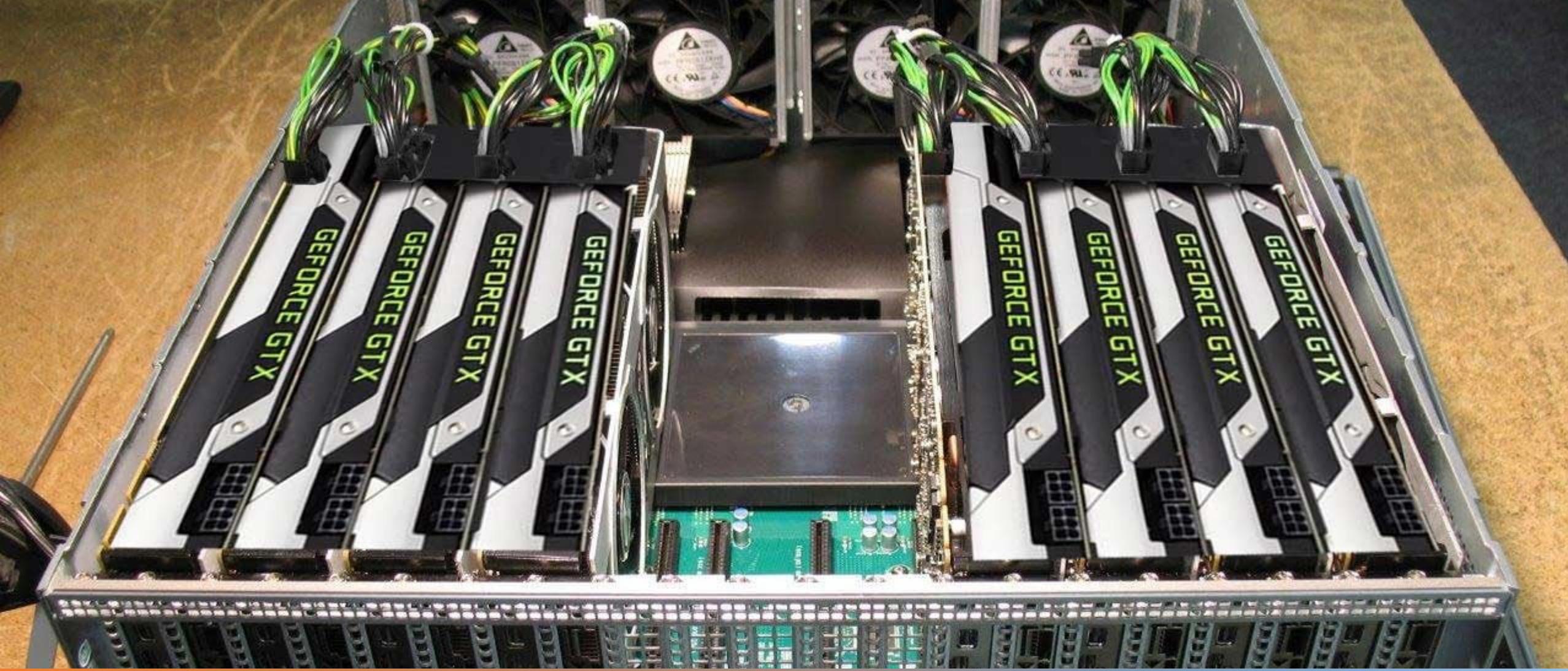
Now that you know what we deal with,
let's go over what is outside of this club's
budget

- GPU based clusters
- RDMA supported clusters, cause infiniband

What are GPU Clusters

What we have covered so far has been CPU based, but GPU has same concept

Instead of maxing out number of cores you maximize number of GPUs

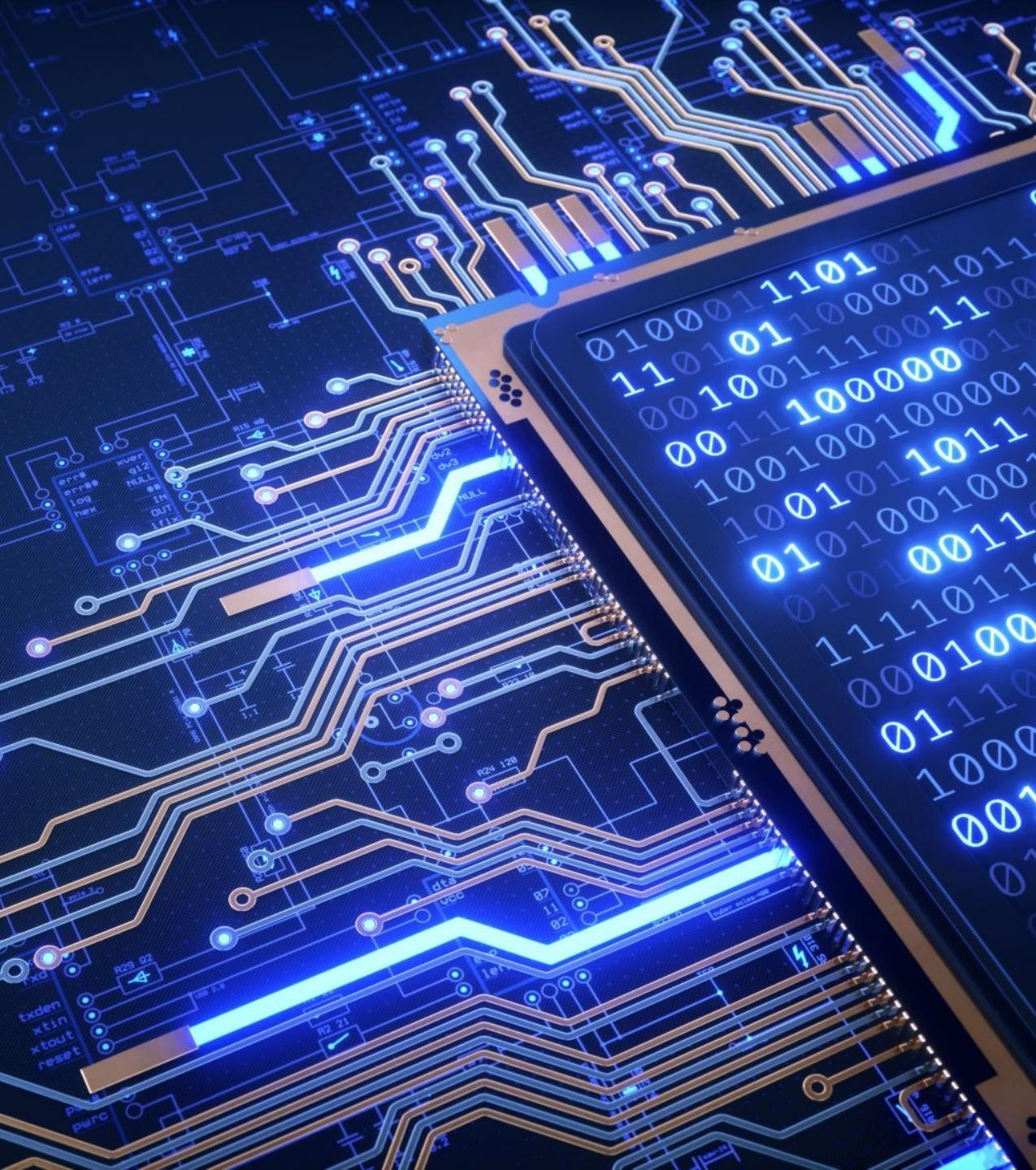


Example photo of a node in a GPU cluster

GPU more \$\$ than CPU clusters

- Researchers discovered in the early 2010's that if the problem they are trying to solve can be reduced to matrix math a GPU can perform it faster than a CPU for the same \$\$
- This led to researchers building heavy GPU based clusters, and the rise of NVIDIA in that sector
- AI is heavily matrix math based, QED





What about problems that are not matrix math based?

- These, assuming they can be made to support multiple cores, are the reason CPU clusters are still used
- Not all problems can run on GPUs and not all can benefit from GPUs and this all has to do with instruction sets and the problem itself.
- A lot of High end computer clusters that deal with CPU's are ARM based rather than x86 based because of energy efficiency

ARM vs X86 rough overview



ARM think apple silicon



X86 think intel



That's all you really need
from an overview stand
point

Other things that are also ARM



Raspberry pi



Majority of all smartphones android and ios



Most chromebooks

Networking, What is infiniband?

- Instead of using copper wires, InfiniBand uses lasers through a tube
- This has much lower latency and a larger bandwidth than ethernet
- New IB is expensive tho

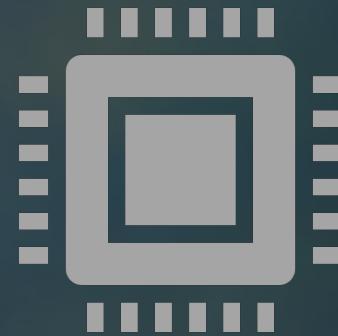
What can you do with IB?

- RDMA (Remote Direct Memory Access)
- Allows for a computer on the network to directly access the RAM or VRAM of a secondary computer without going through that computers CPU
- This speeds up multi core and multi GPU operations

In conclusion



High performance computing is basically taking a bunch of smaller things and making a bigger thing



Many different ways to achieve it and the computer should be built around the task that is needed to be completed



Questions?

