



High Performance Computer Hardware

G Burton – ICG – May18 – v1.1



Overview

- What is an HPC
- Fundamentals of executing a program
- Commodity clusters
- Shared Storage
- Interconnects and Networks
- Login layer
- GPU's and Intel Phi's

What is a High Performance Computer ?

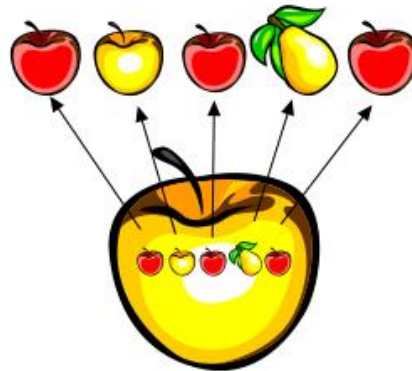
- This depends on who you are and what you want to do ...
 - Are you a graphics designer that needs high end graphics on a work station.
 - Are you a financial trader that requires super fast networking.
 - Are you my Mum who thinks her iPhone is the most powerful computer in the world.
 - Are you weather forecaster that needs to run a weather simulation.
 - Are you a Scientist wanting to experiment with a quantum computer.
 - Are you a banker requiring mission critical transaction guarantee.
 - Are you mining for bit coins.
- We need to narrow things down

Specifically

High Performance Computing
(HPC) Hardware for
Scientific Numerical Analyses

Demystifying the Techno Babble

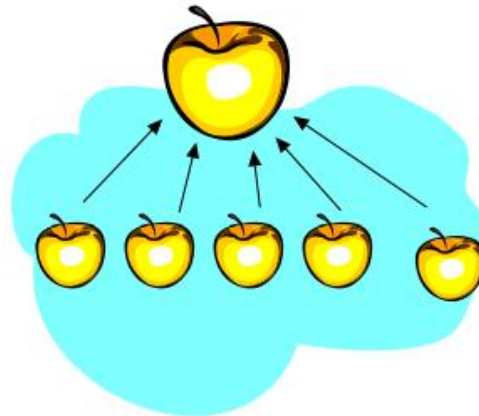
Virtualisation - One machine appears as many



Virtualisation maximises hardware resources through virtual machines (VM's)

Host machine may support VM's with different OS's

Clouds - Many machines appear as one



Clouds are all about scalability

Being able to switch (automatically) resources in or out.

Sometimes called SOA
Service Orientated Architecture.
IaaS - Infrastructure as a Service
PaaS - Platform as a Service
SaaS - Software as a Service

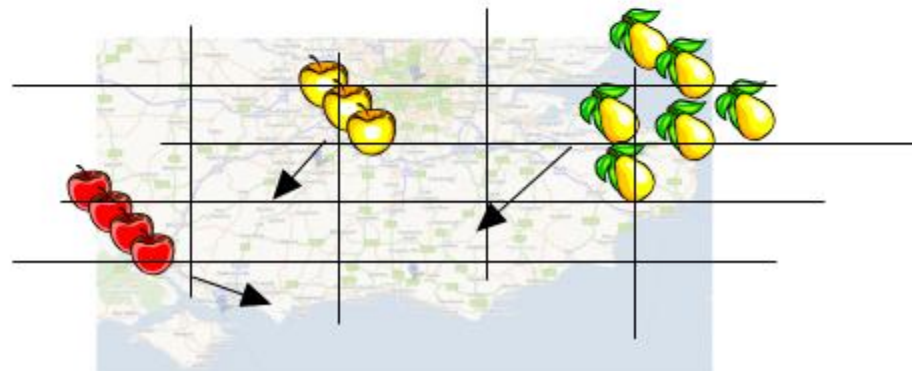
Clouds use virtualisation.

Demystifying the Techno Babble (2)

Clusters - Many homogeneous machines at one location

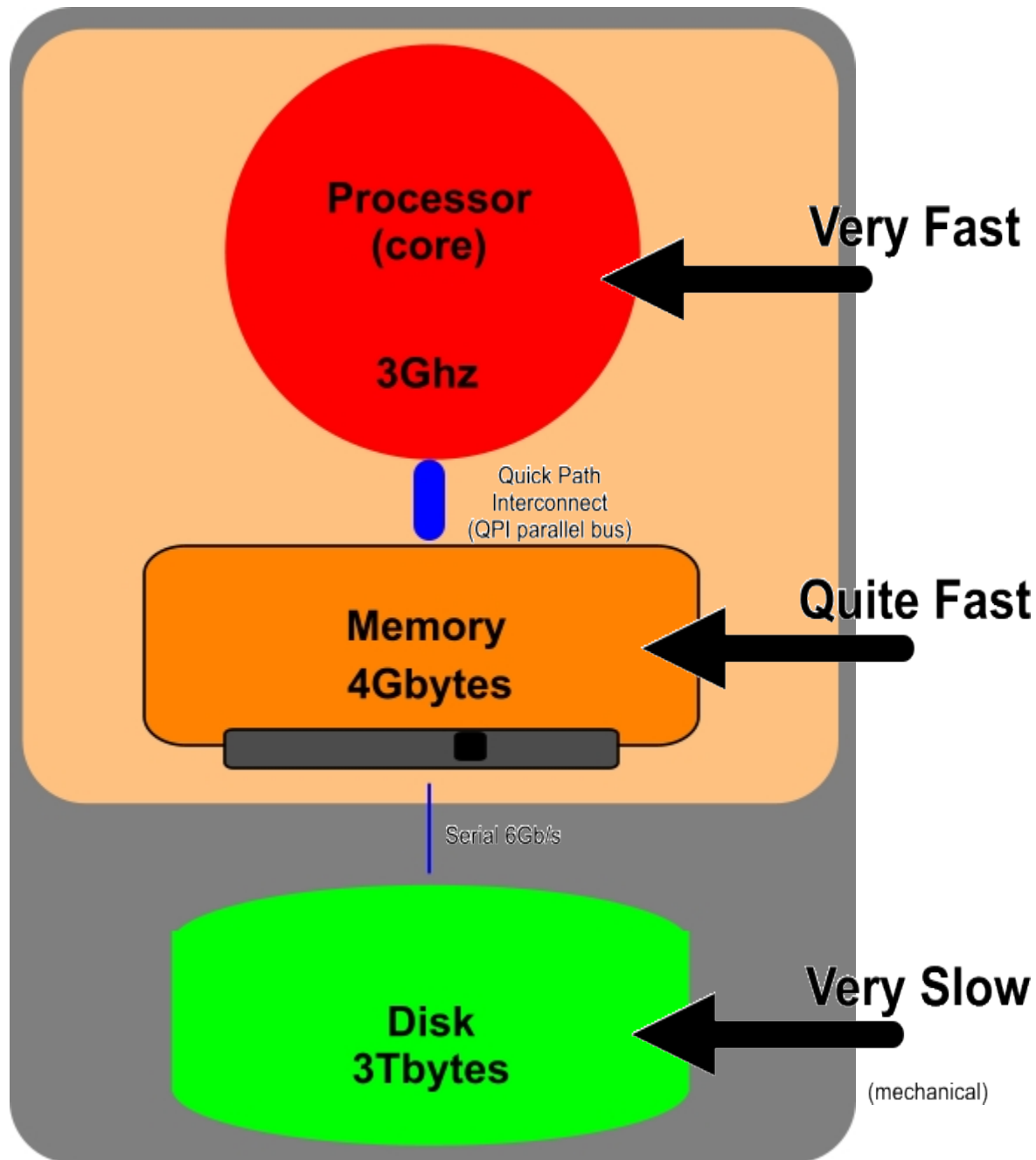


Grid - Many heterogeneous machines at disparate locations

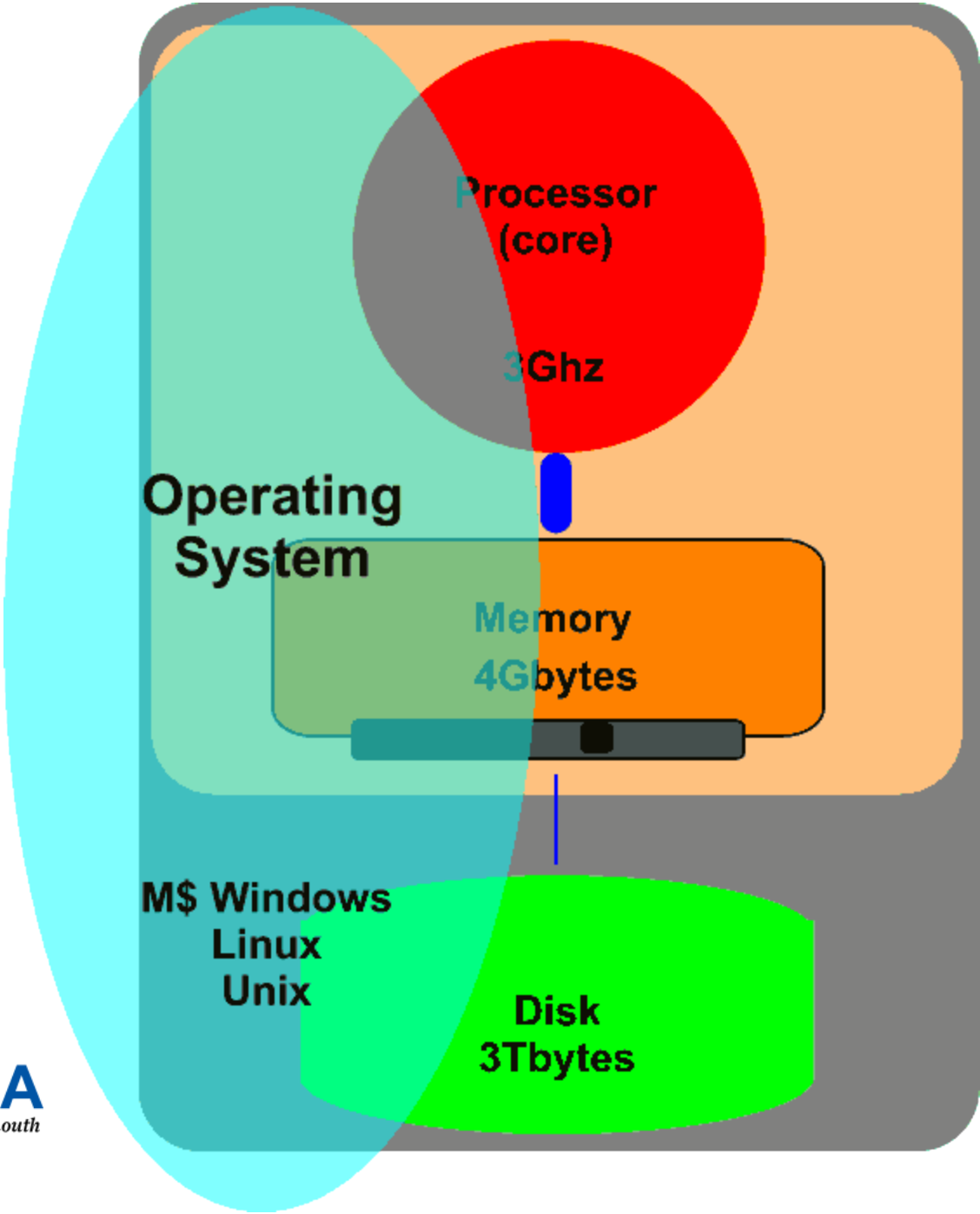


Fundamentals of executing a program

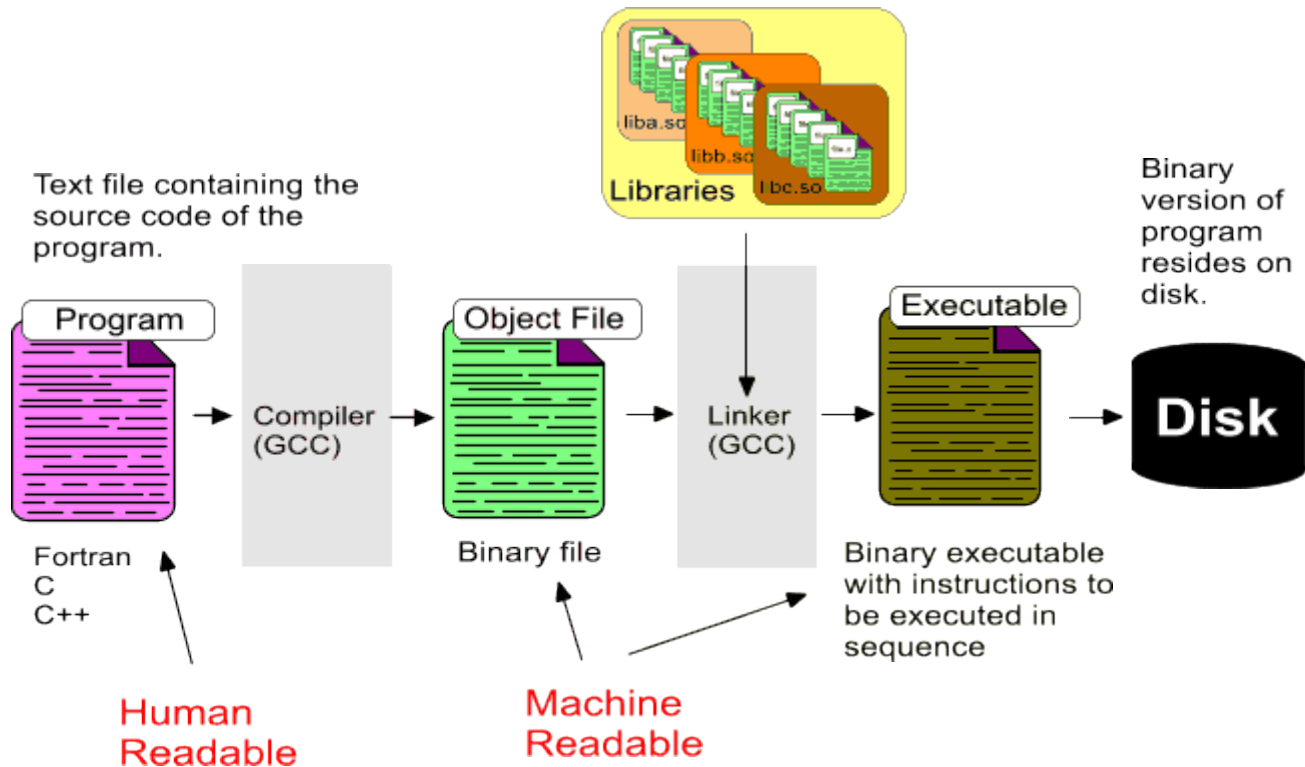
Back to Basics

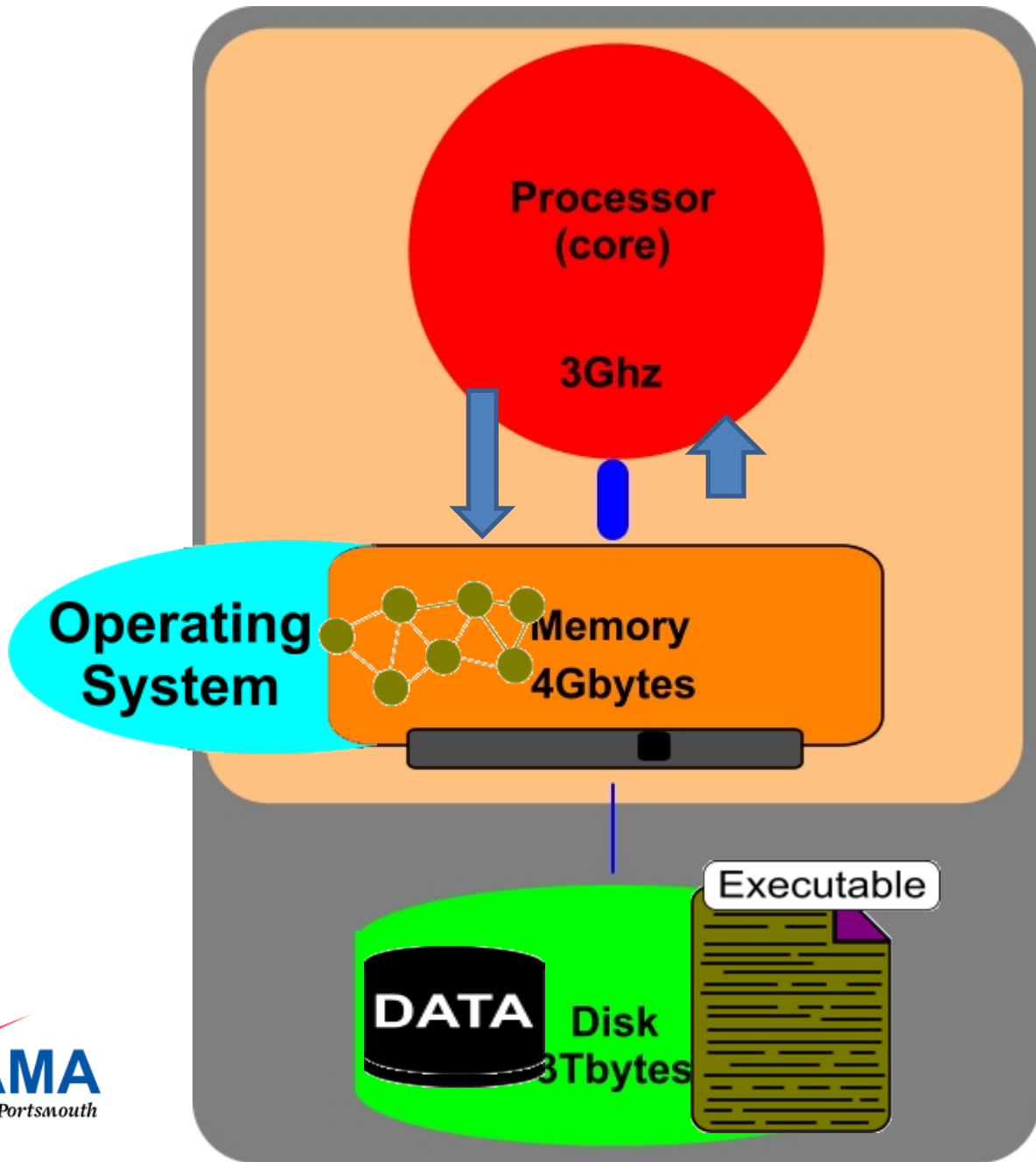


Back to Basics



In order to solve our problem we need a “Program” to run.

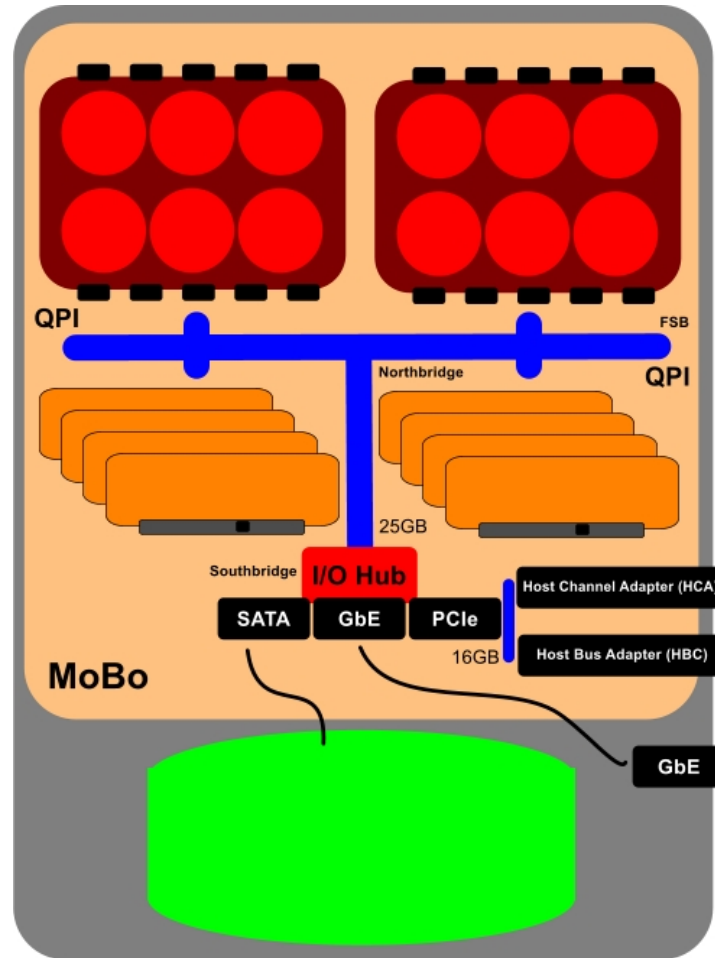




Paging
Swapping
OOM killer

... these days much more packed into the same space ... but basically the same!

Its doesn't matter how many cores the standard executable will execute sequentially one instruction at a time.



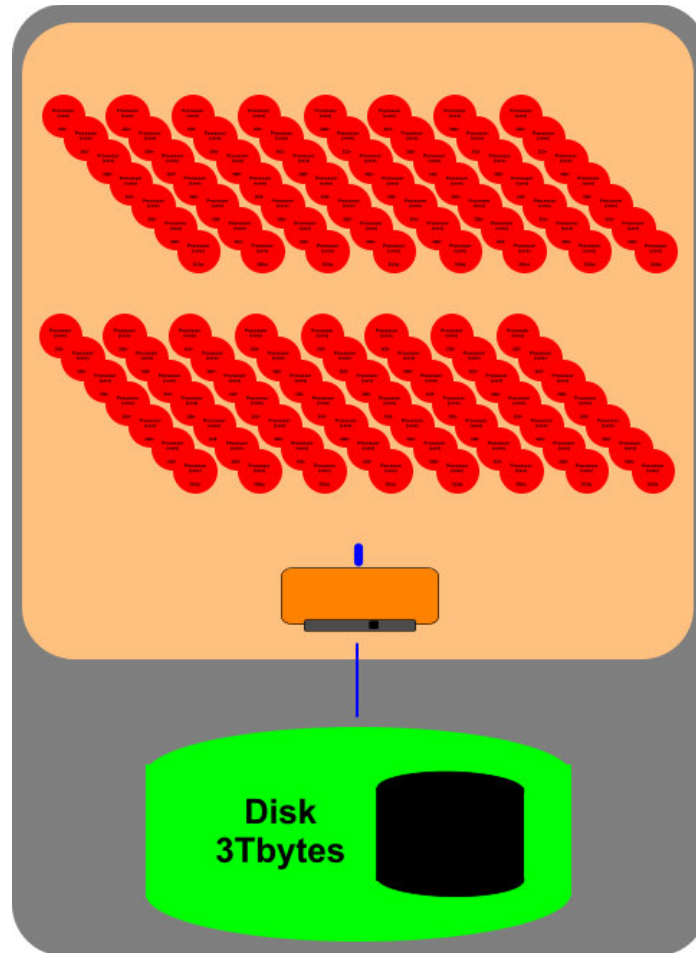
Programming Frameworks

- Shared memory (cooperating threads – or **multi-threading**)
 - OpenMP
 - POSIX Threads
 - Cilk
 - Threaded Building Blocks
 - etc, etc
- Distributed memory (cooperating processes)
 - MPI (PVM, etc)
 - Co-array Fortran, UPC, etc
 - Global Array Toolkit (etc)
 - Adlib and HPspmd?!
 - etc, etc

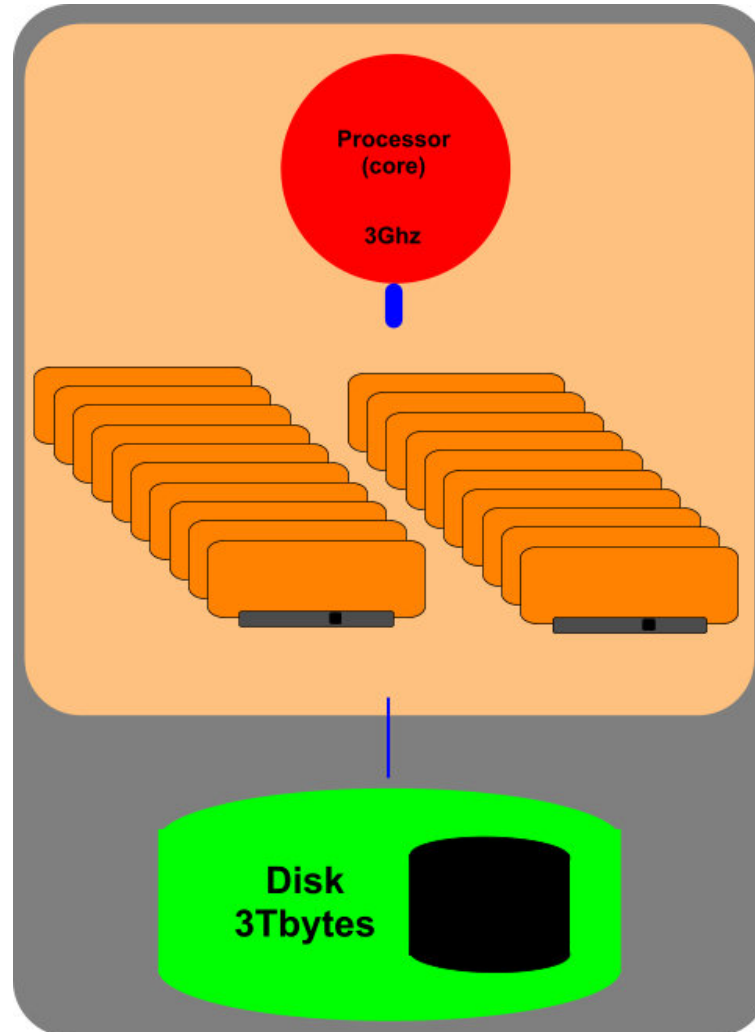
How effective.....
Depends on the nature
of what you are trying
to solve.

Not always
Possible !!

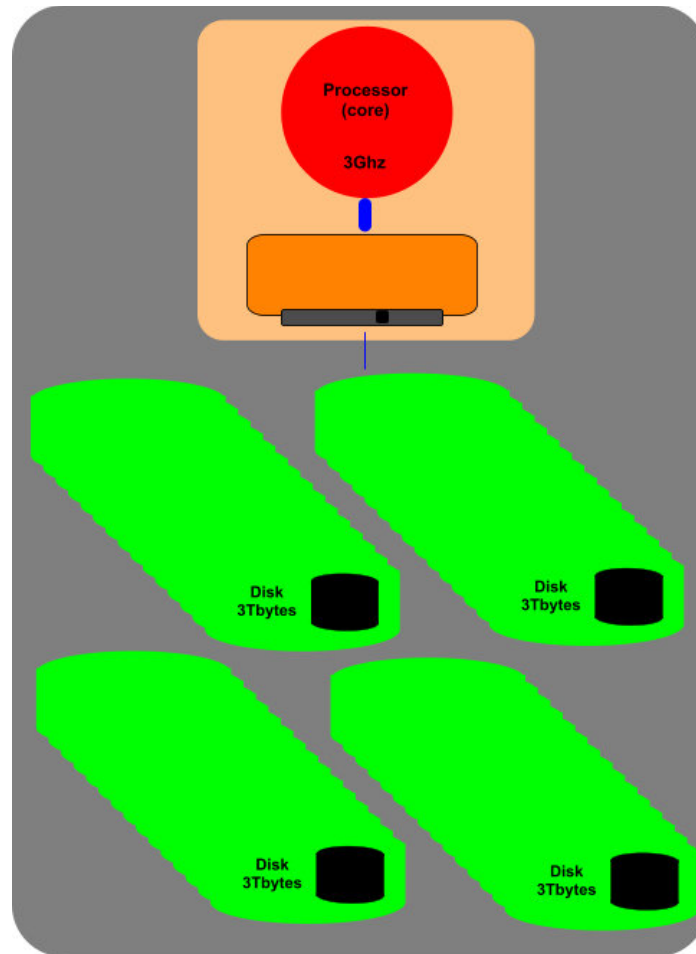
In Reality – A program may require many 100's of cores



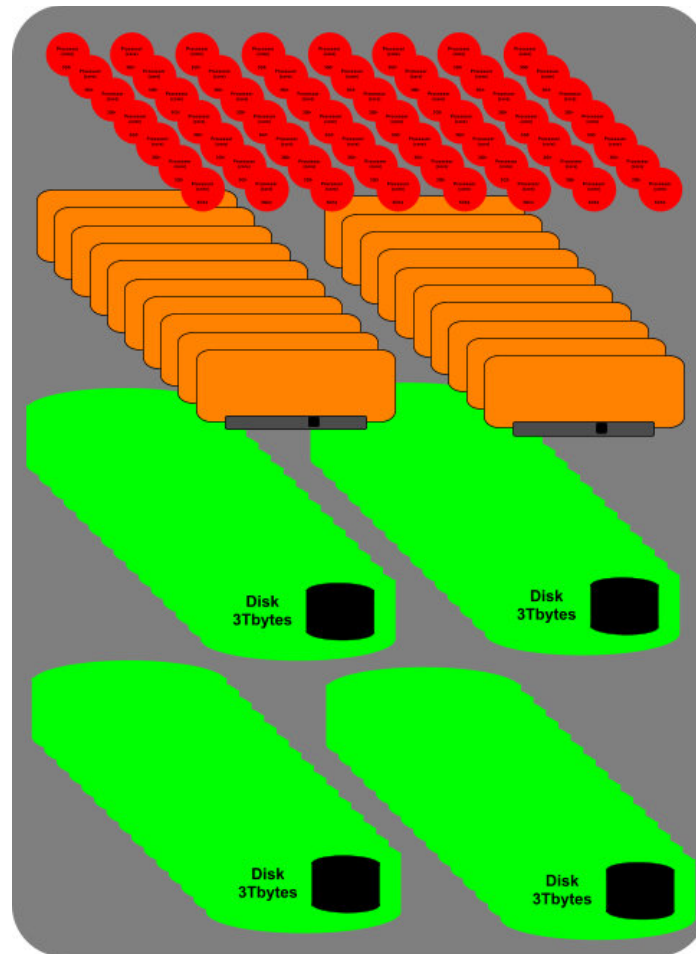
... or many Gbytes of memory



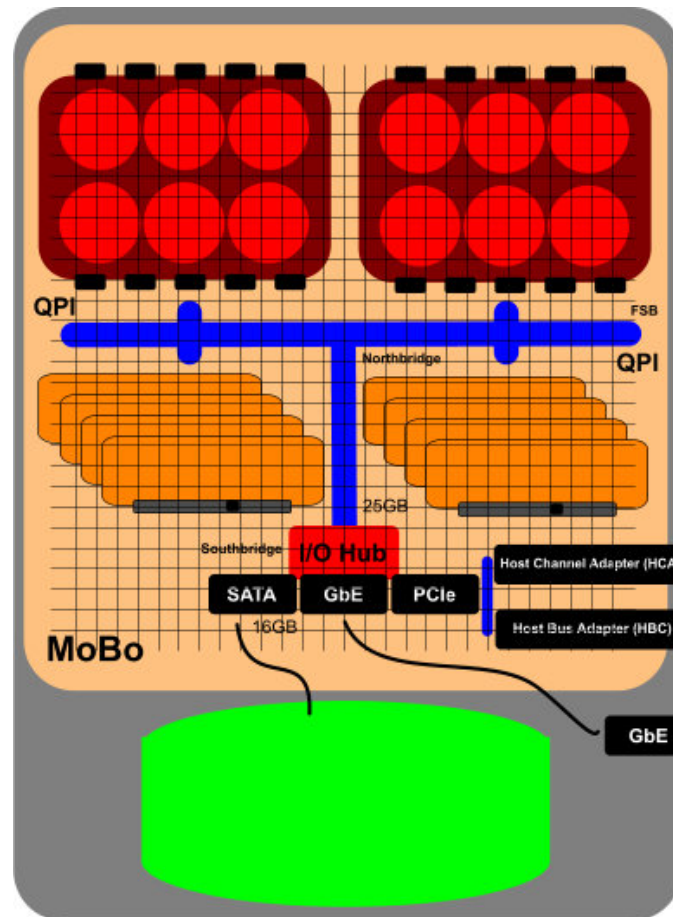
....or a lot of disk space



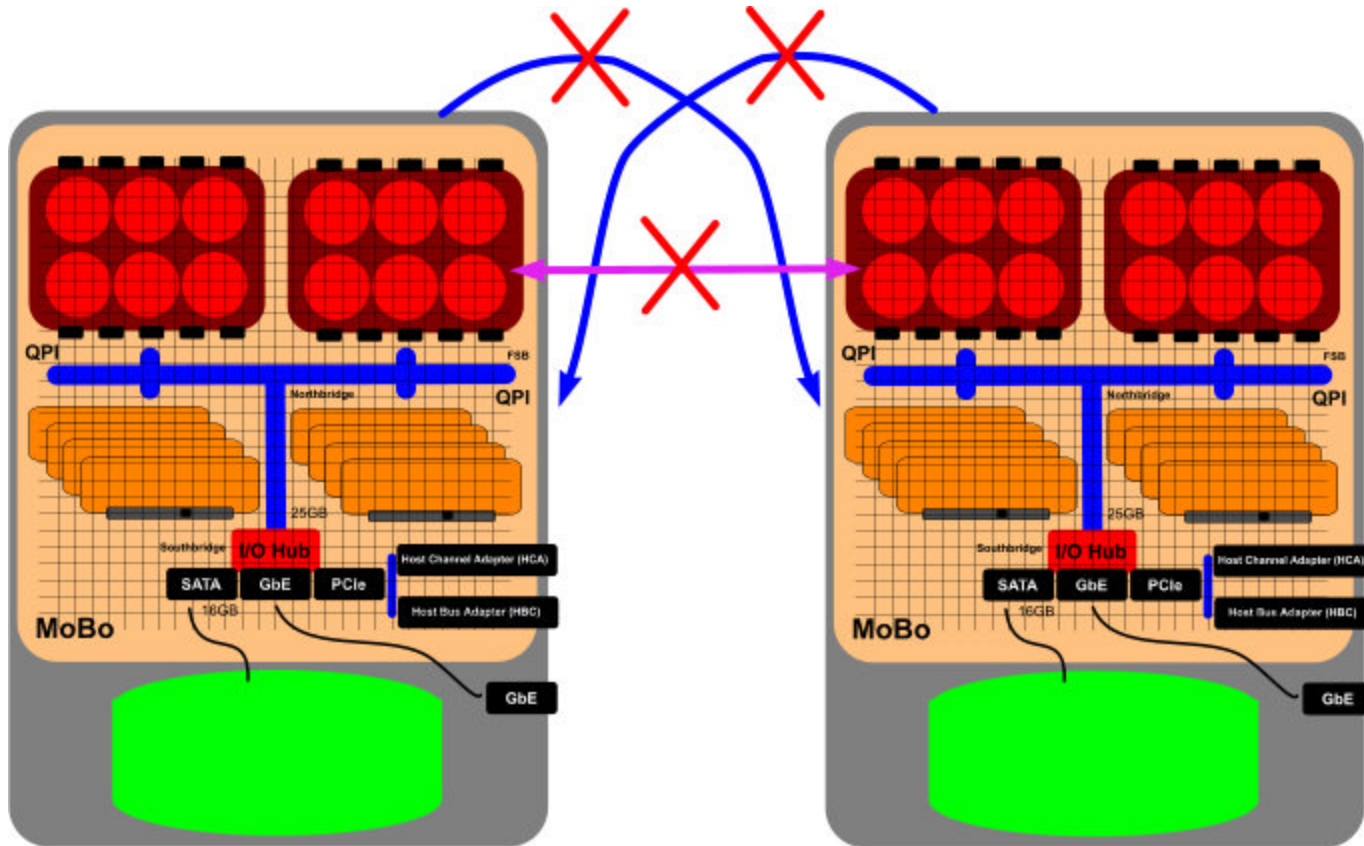
...or a combination of all.



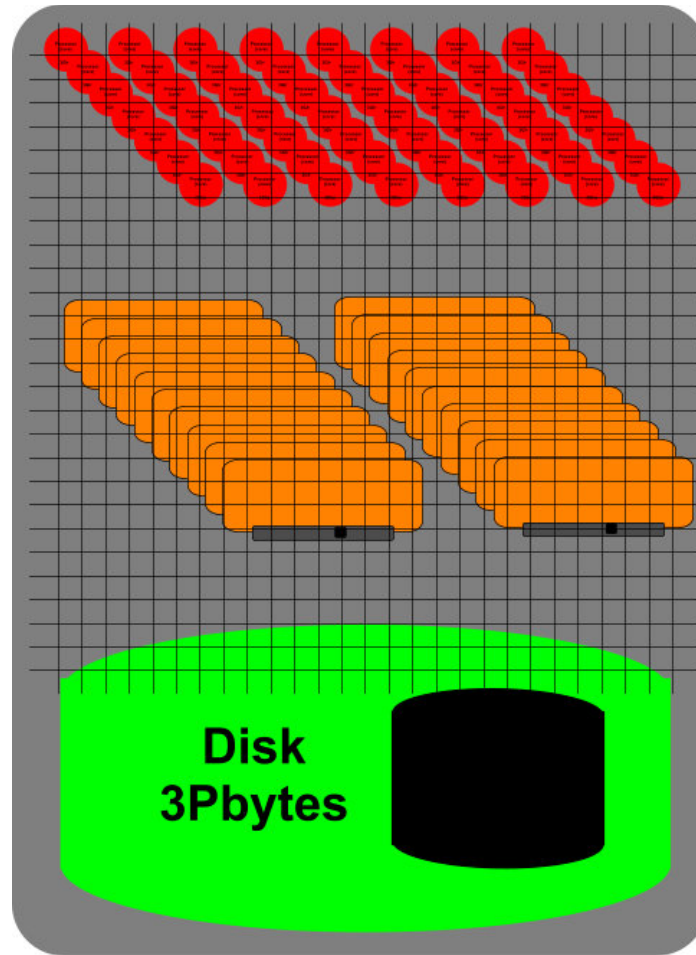
Internal Communications Possible



No direct comms between Pc's

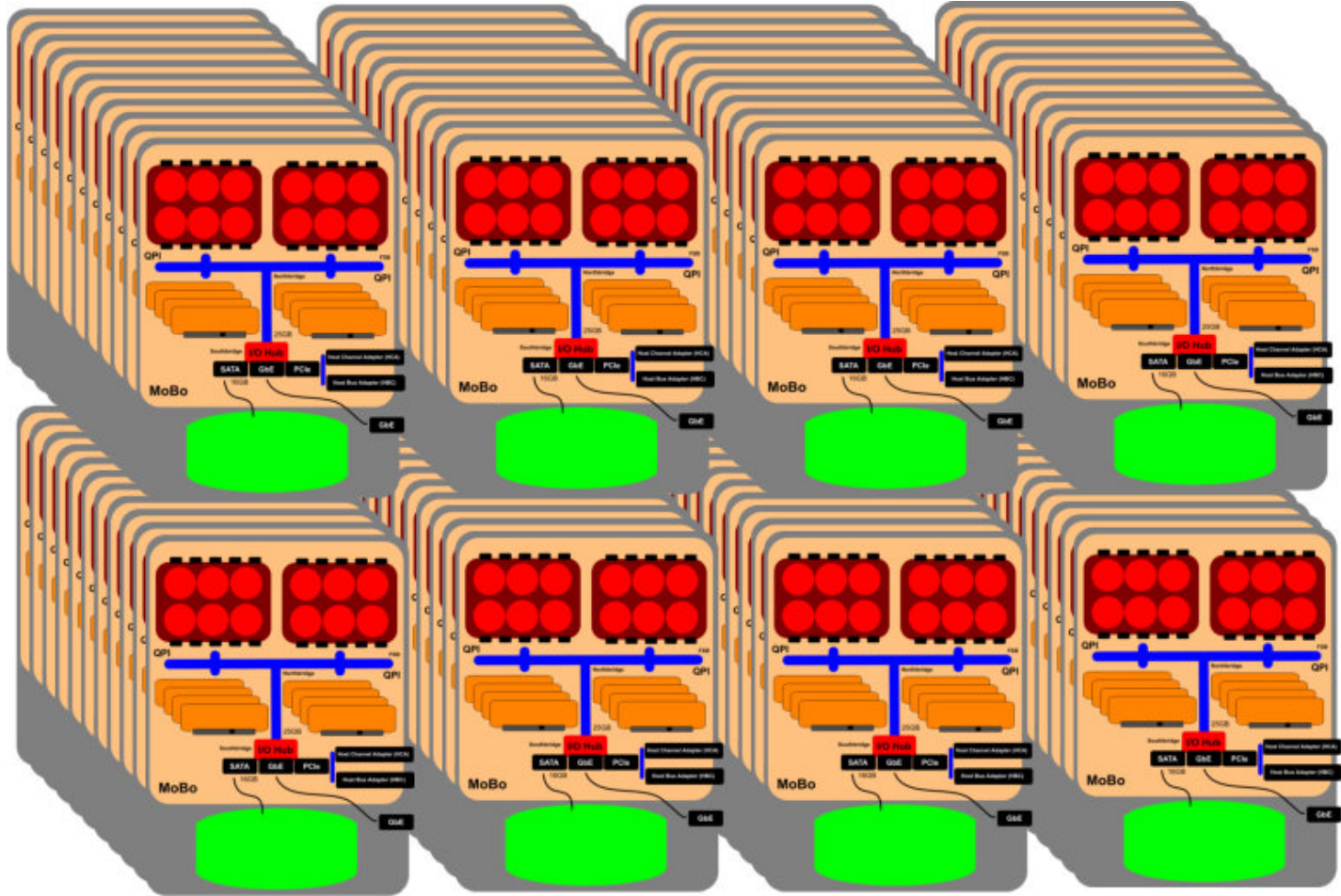


We could build a very expensive machine

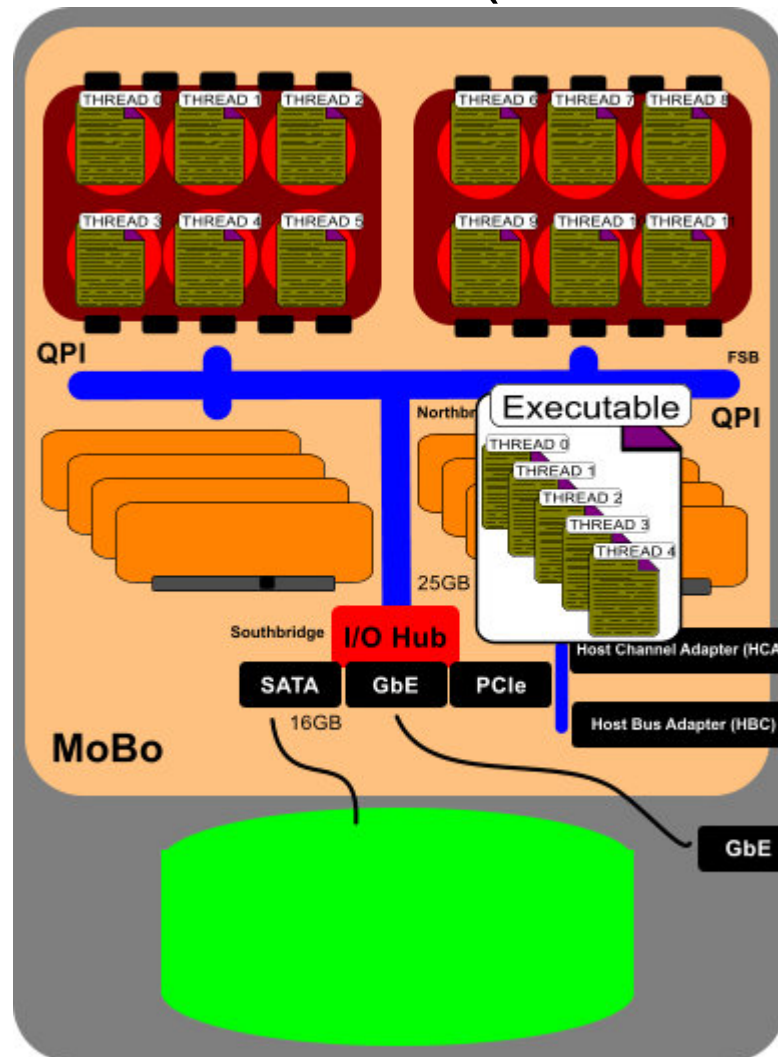


Referred to as a Shared Memory Machine.

Commodity cluster using high end PCs

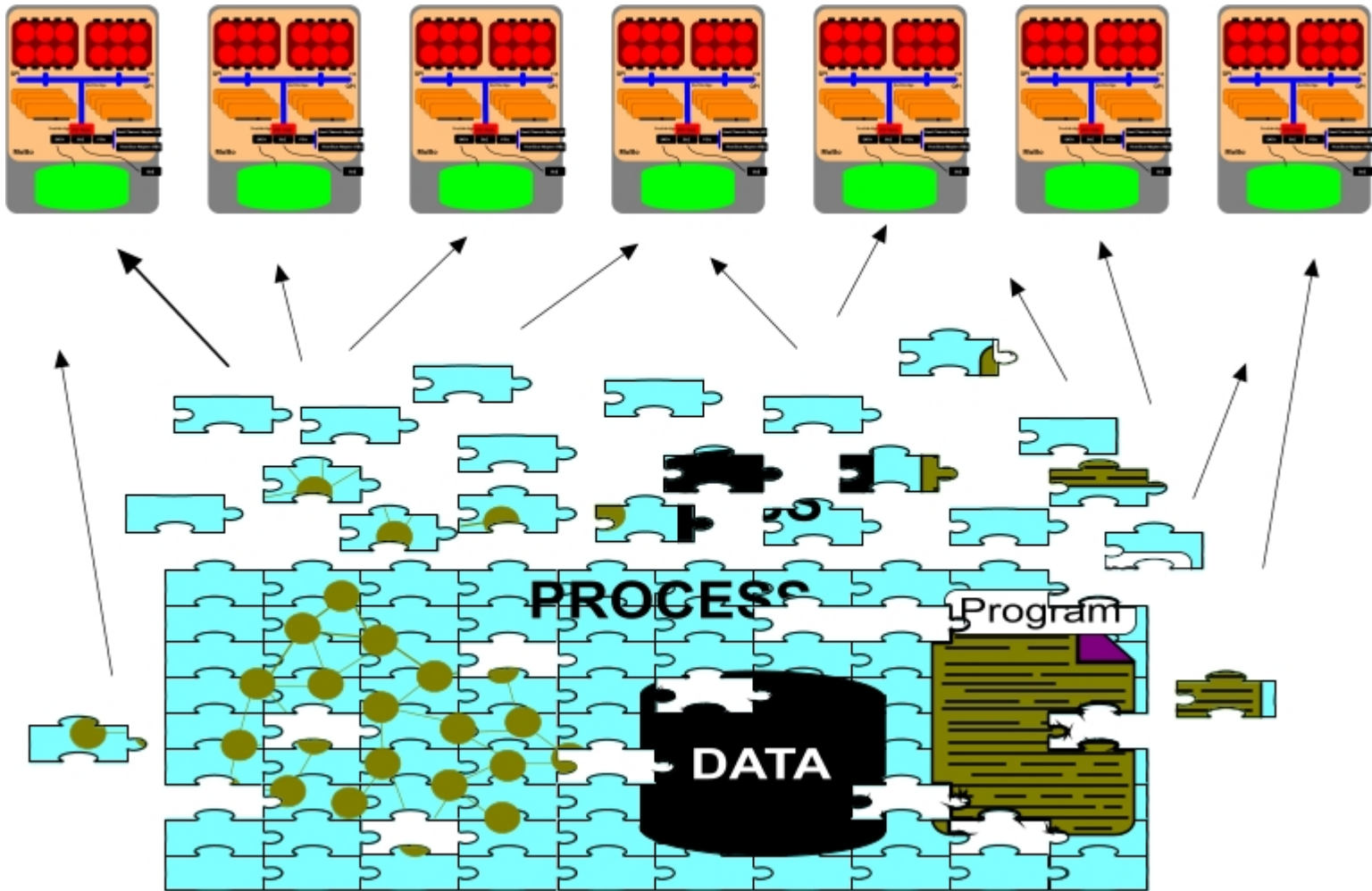


Multithreading takes advantage of onboard cores...(but not across PC's)



Limited to number
Of cores on Mobo.

Distributed Memory works between PC's

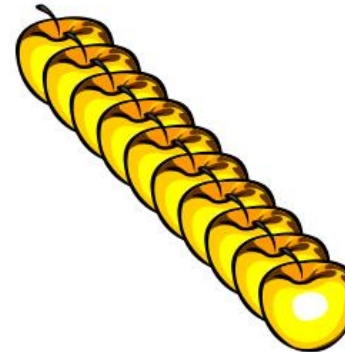


Summary - HPC alternatives for Scientific Computing

HPC through
Specialised Hardware



HPC using
Commodity Hardware



Examples:-

Highend Graphics Workstation
for 3D modelling.

Shared Memory or Symmetric
Multiprocessor (SMP) machines
like COSMOS & Universe

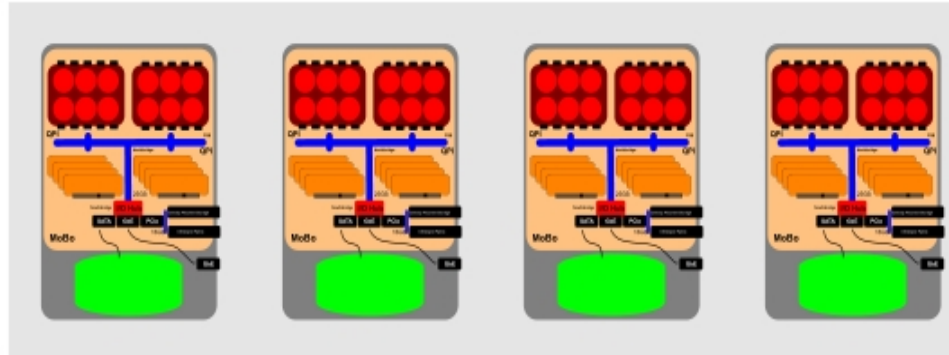
Vector Processors using, for
example, expensive Cray
or Convex Mainframes.

High Performance Computer
=
Commodity Cluster

Commodity Clusters

- Made from commodity (off-the-shelf) components.
- Consequently (relatively) cheap.
- Usually Linux based
- High availability storage (no single point of failure)
- Generic compute pool (cloned servers that can easily be replaced).

Dell 6100 Front View



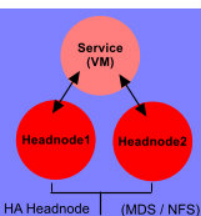
Google
Amazon
Ebay



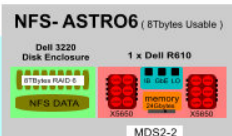
Dell 6100 Rear View

High density - Google / Amazon

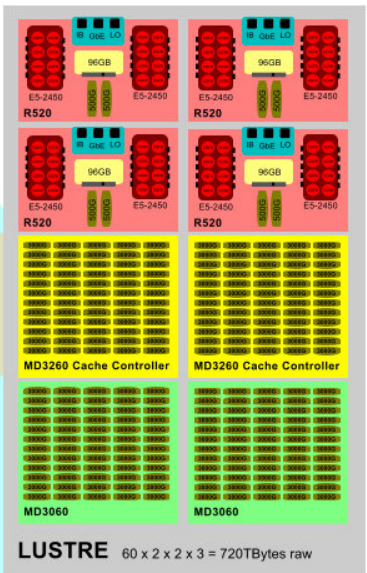
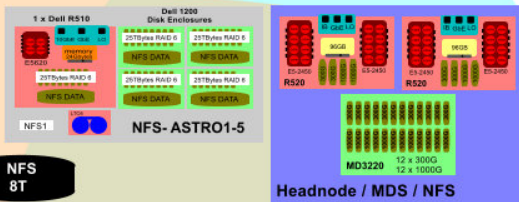




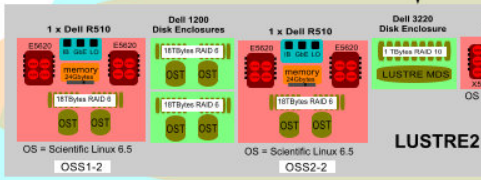
James Watson DC



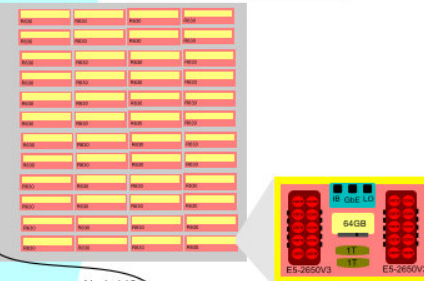
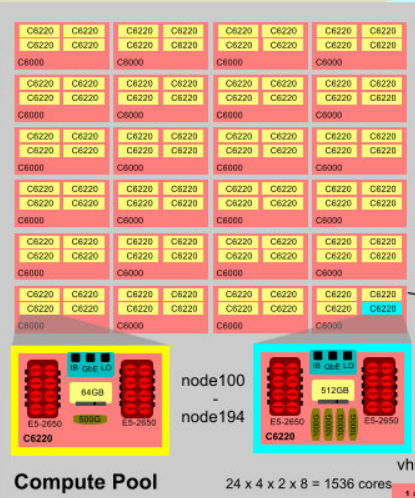
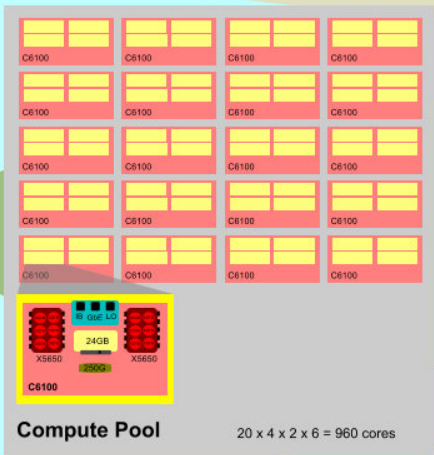
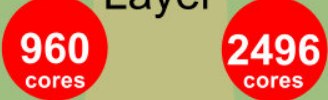
Anglesea DC



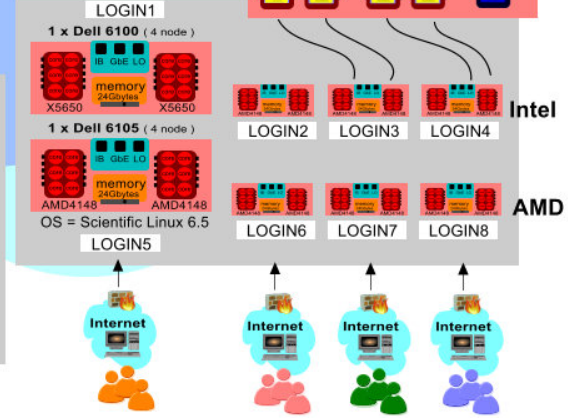
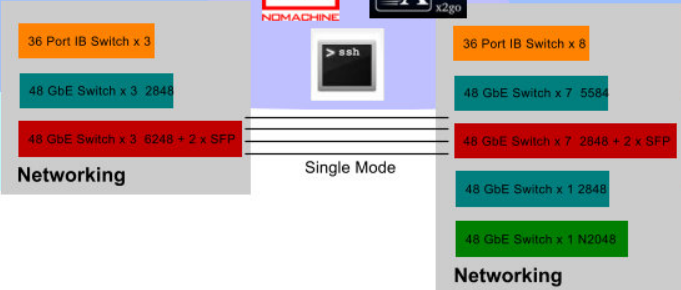
Shared Storage Layer



Compute Layer



Login Layer



Processors

- X5650 = 2.66GHz
- E5620 = 2.4GHz
- AMD4148 = 2.8GHz
- E5-2650 = 2.6GHz
- E5-2650V3 = 2.3GHz

Network

- IB = Infiniband
- GbE = Gigabit Ethernet
- 10GbE = 10G Ethernet
- LO = Lights Out (BMC)

OS

Scientific Linux 6.5



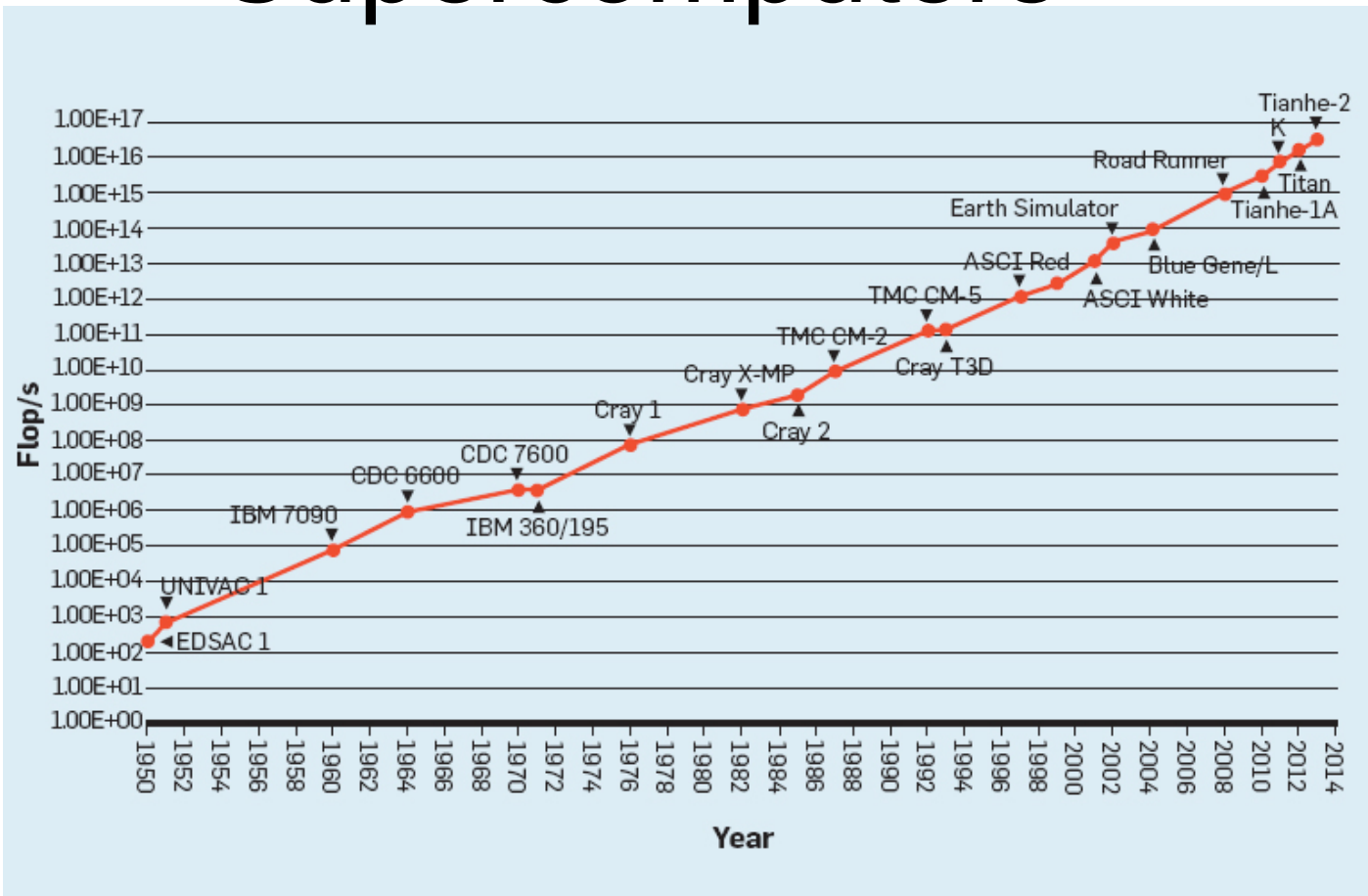
HPC Facts

- China's Sunway TaihuLight now number 1 in top 500 with 10.6 million Cores.
- 125 petaFLOPS (Floating Point Operations / Second). Sciama 10 teraFLOPS.
- China now has 202 in top 500 with US second with 143.
- Japan=35, Germany=20, France=18, UK=15
- Exascale expected 2020. Race between Japan, France, China, USA
- Uk – top 5 to do with weather, followed by Cambridge and Edinburgh m/cs

Computer performance

Name	FLOPS
yottaFLOPS	10^{24}
zettaFLOPS	10^{21}
exaFLOPS	10^{18}
petaFLOPS	10^{15}
teraFLOPS	10^{12}
gigaFLOPS	10^9
megaFLOPS	10^6
kiloFLOPS	10^3

Moore's Law for Supercomputers

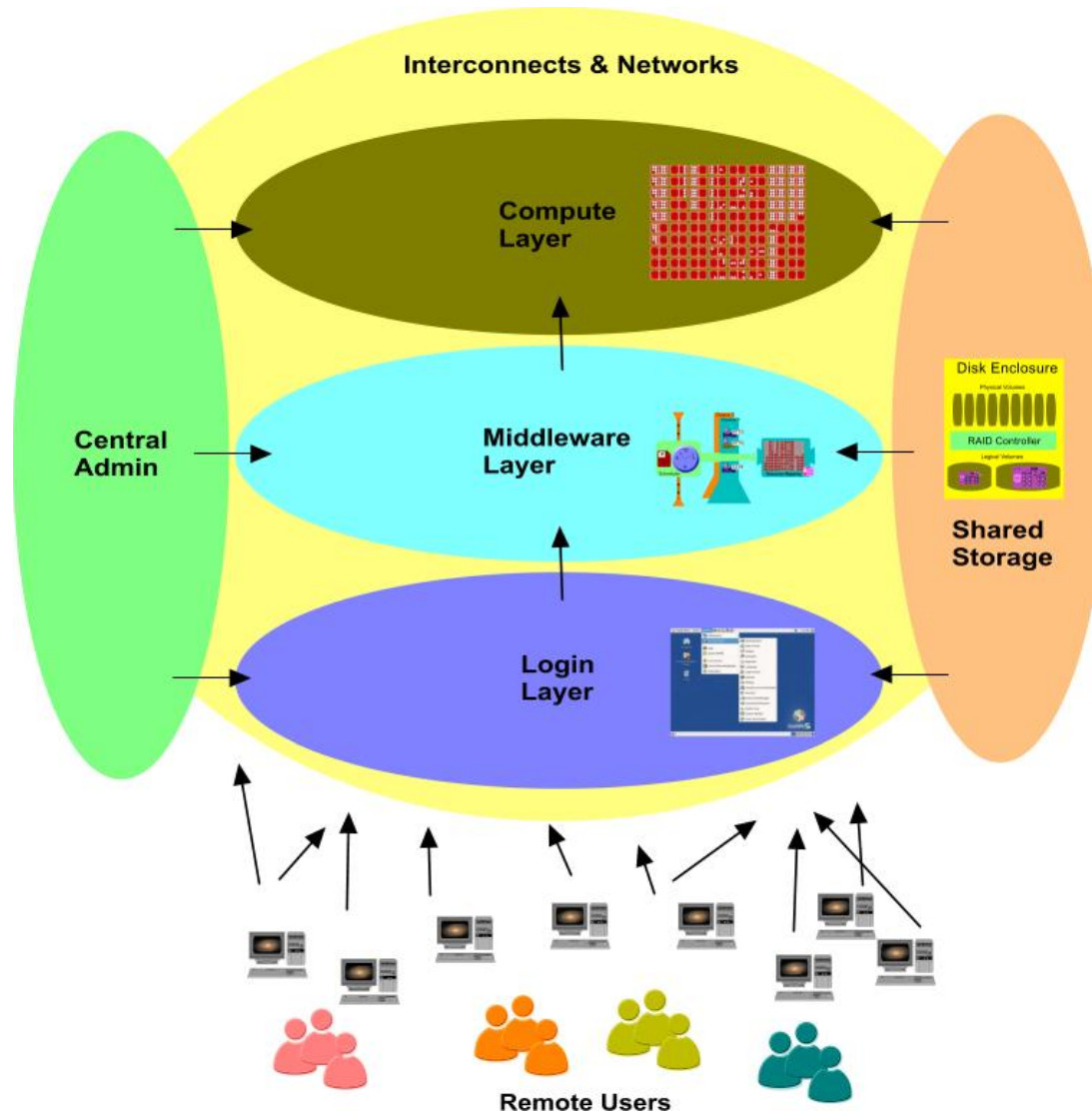




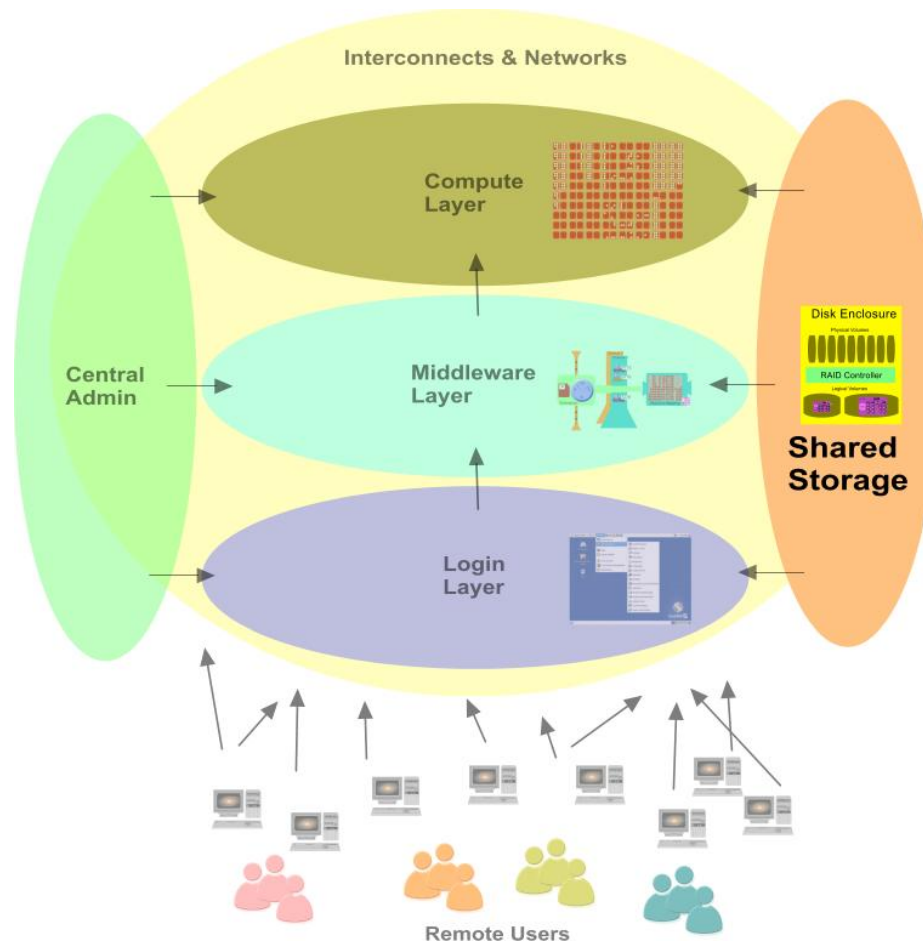
Common Misconception

- A super computer does not necessarily mean a program will run faster.
- Commodity processors typically slower than average desktop / laptop (2.6GHz vs 3.0GHz)
- Unless a program can be parallelised it may run slower.
- Disk access may also be slower (directly attached).
- However super computer much more stable for long runs.

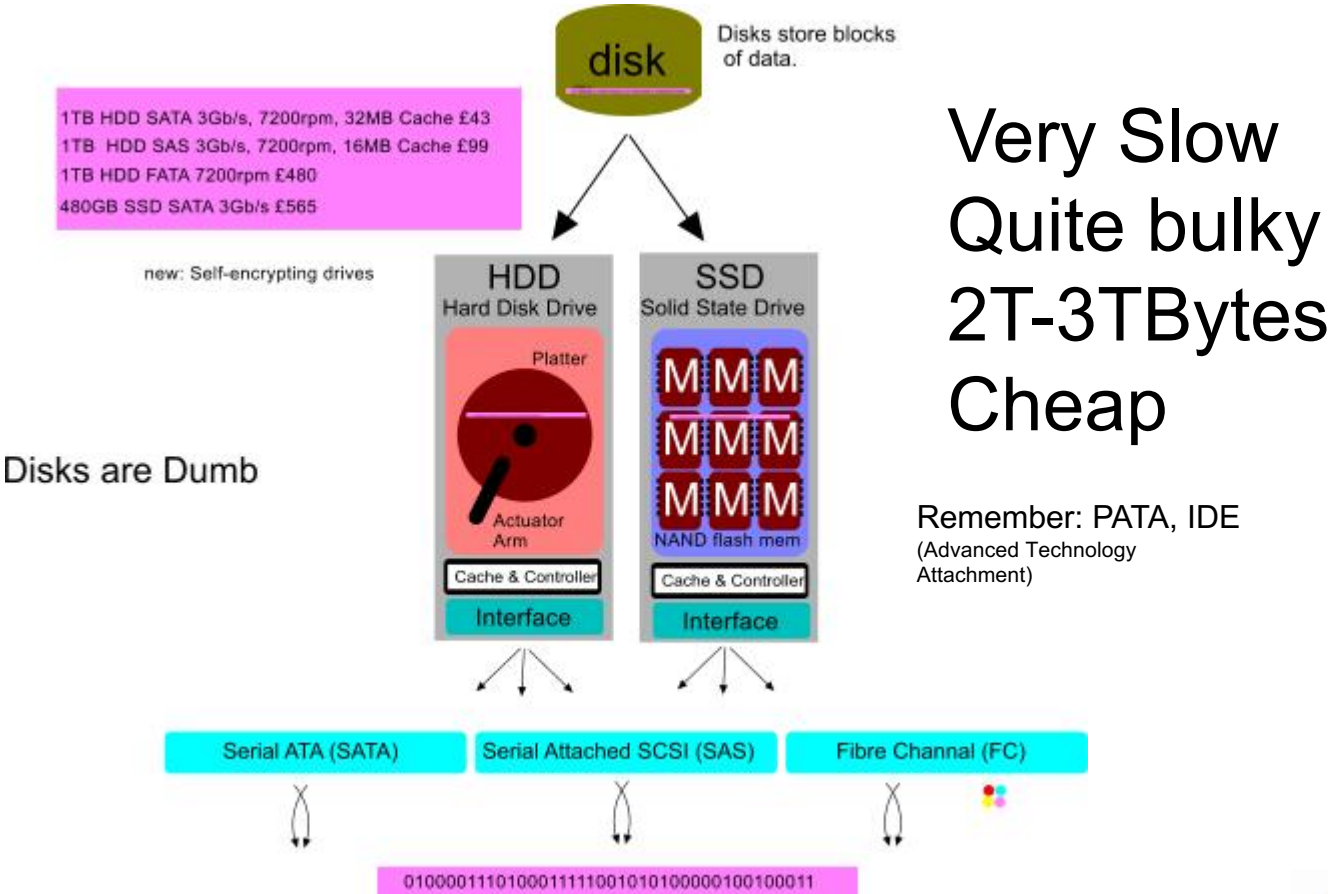
Cluster Concept

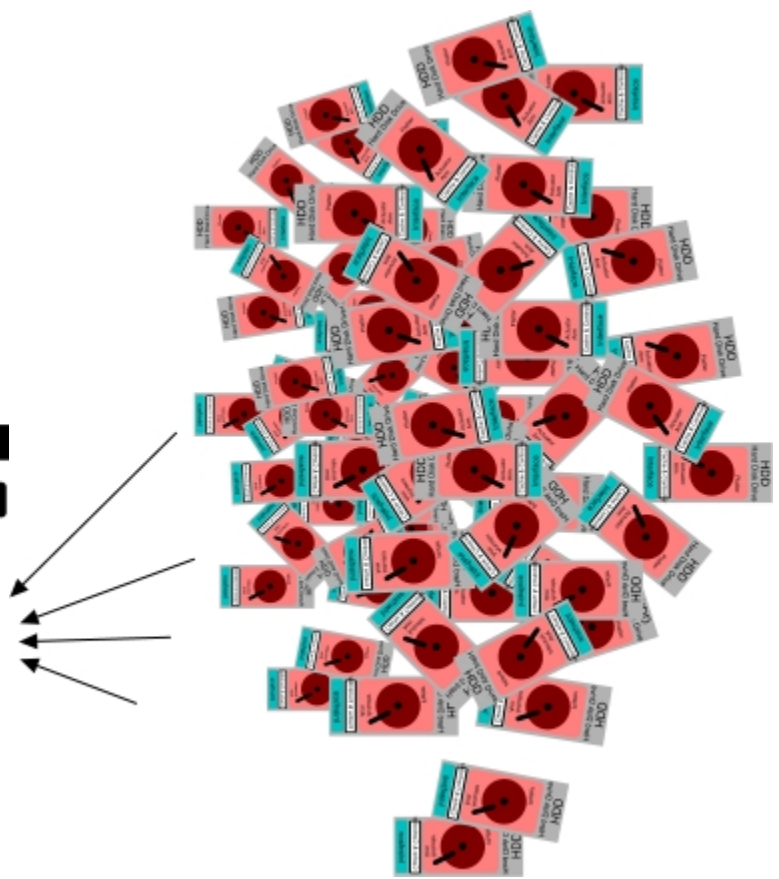
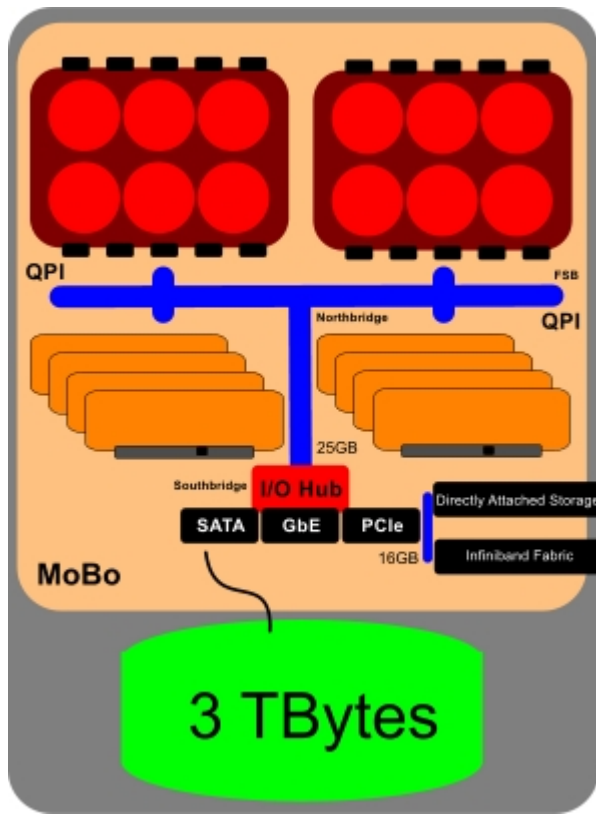


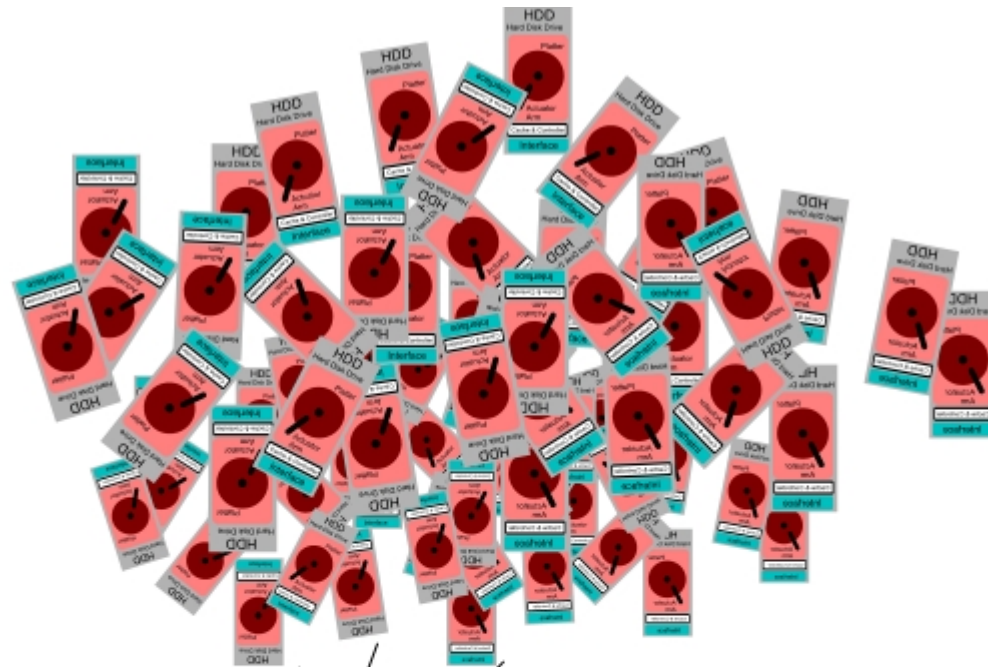
Shared Storage



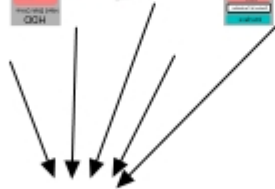
Raw Disks are Dumb





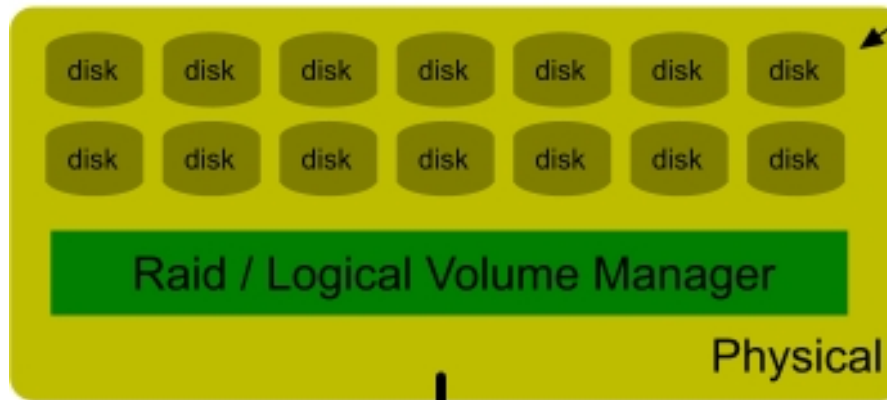


Commodity Disks

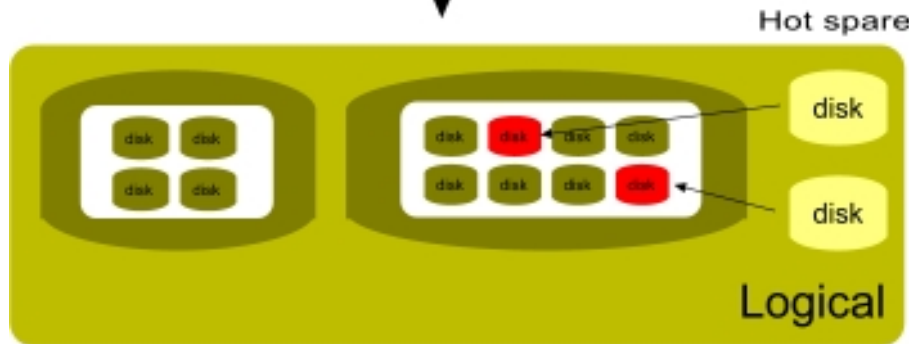


In HPC internal disks usually only contain the Operating System.
Some times two disks are "mirrored" for security.

Single Disk of little use for data :-



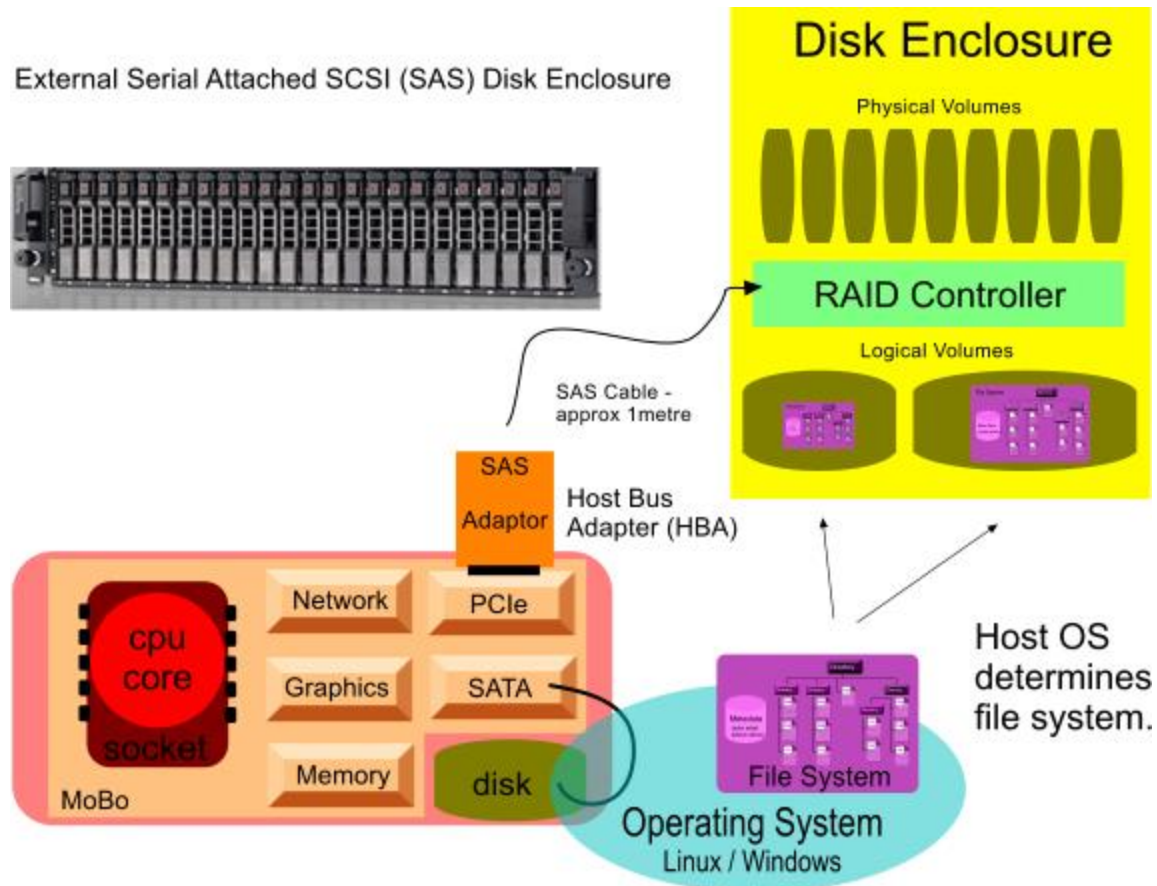
- Limited size.
- Limited performance
- Limited fault tolerance



Logical Volumes

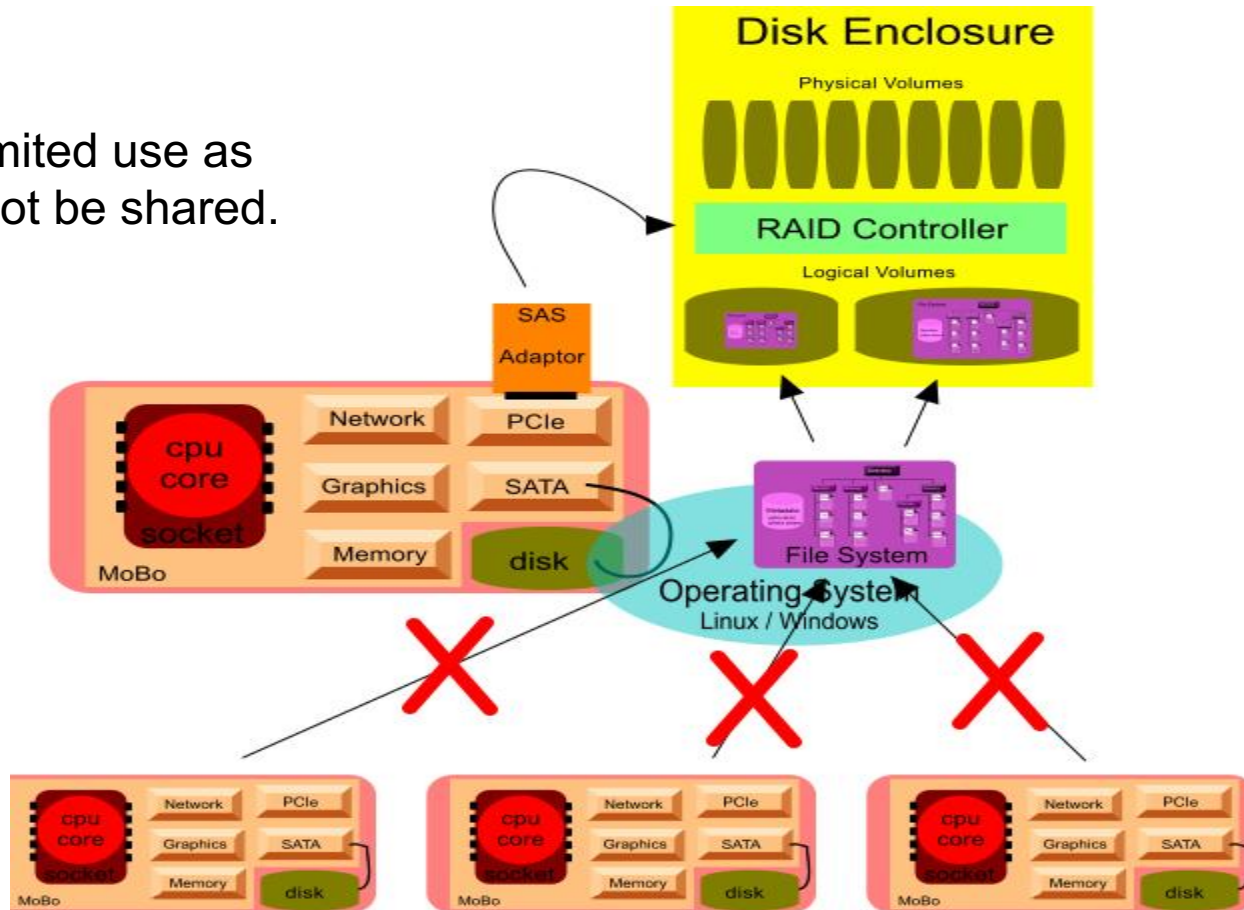
- Increased size
- increased performance through striping.
- Increased resilience through parity and mirroring.

Directly Attached Storage (DAS)

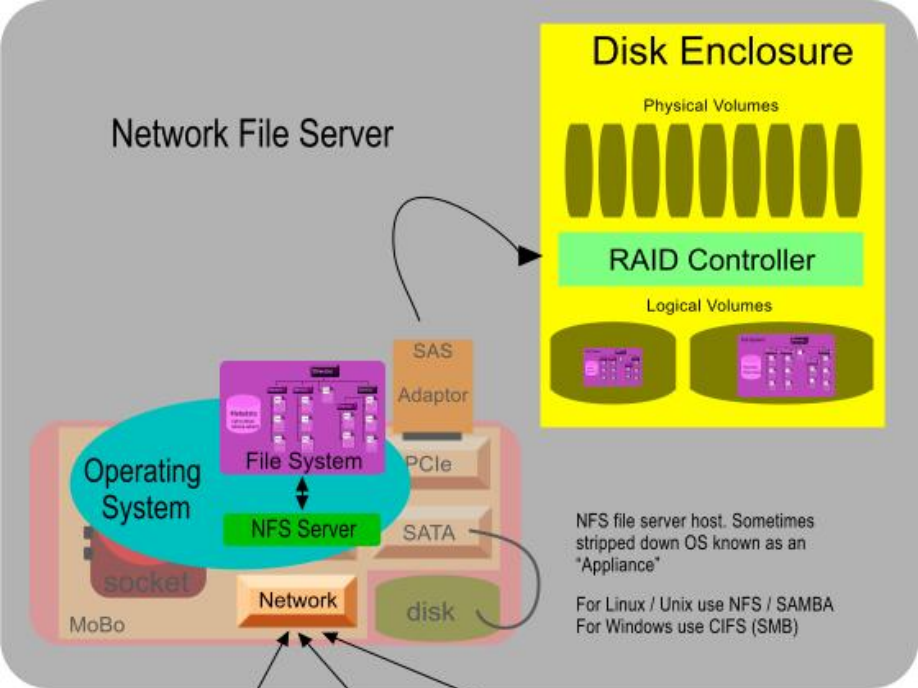


Directly Attached Storage (DAS)

Of limited use as cannot be shared.



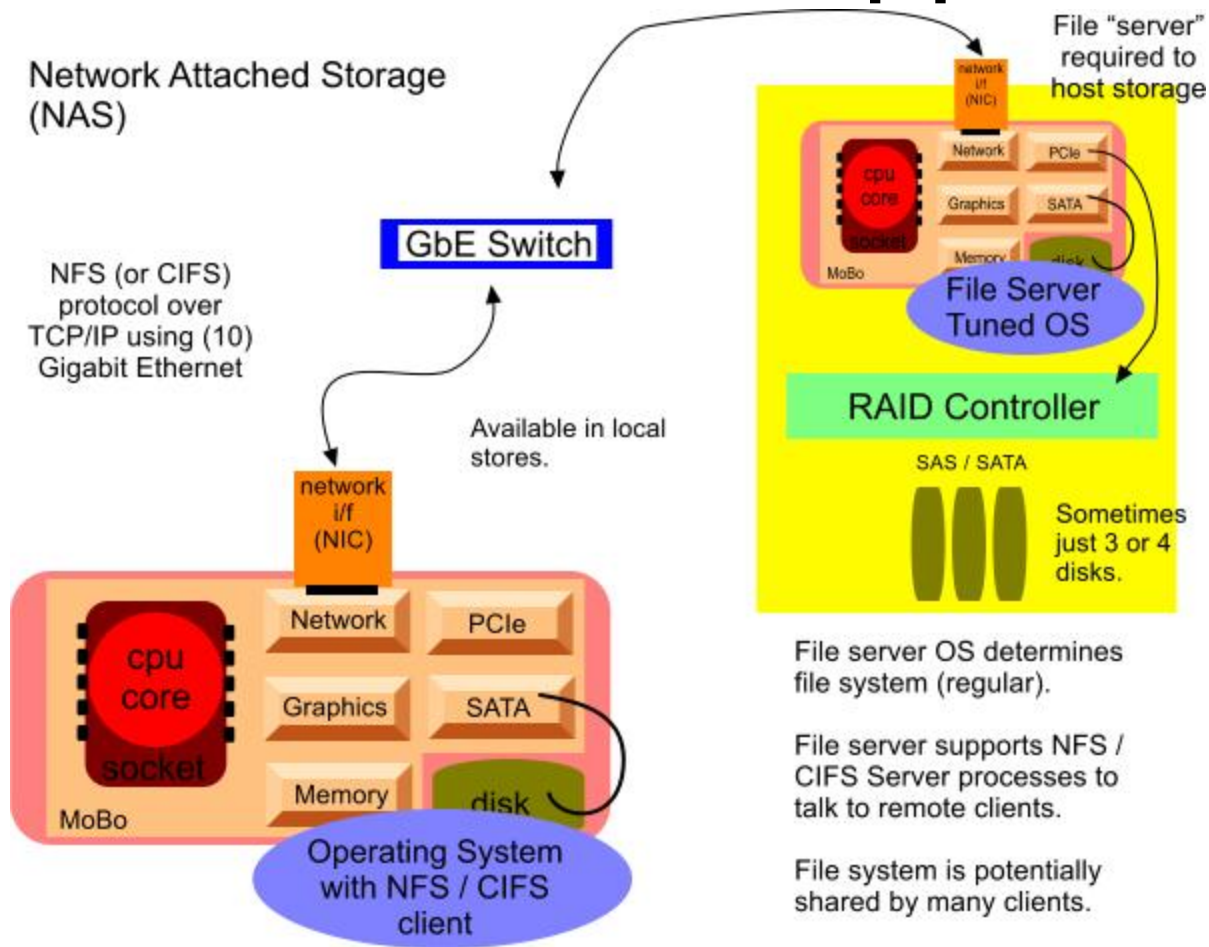
Network Attached Storage (NAS)



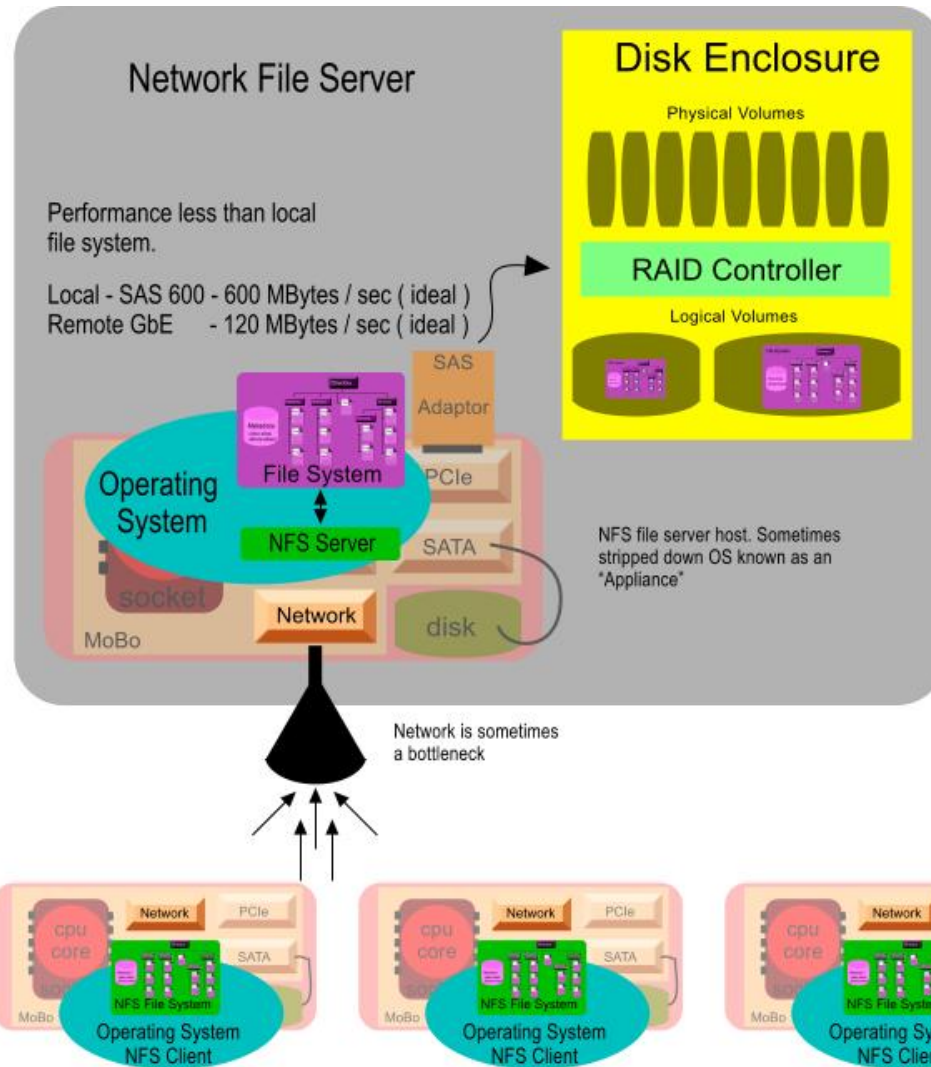
Remote File System is mounted over normal TCP/IP network.



NAS or Network Appliance



Network BW is often the bottleneck



Distributed File Systems – quick glance from Wikipedia

Examples [\[edit \]](#)

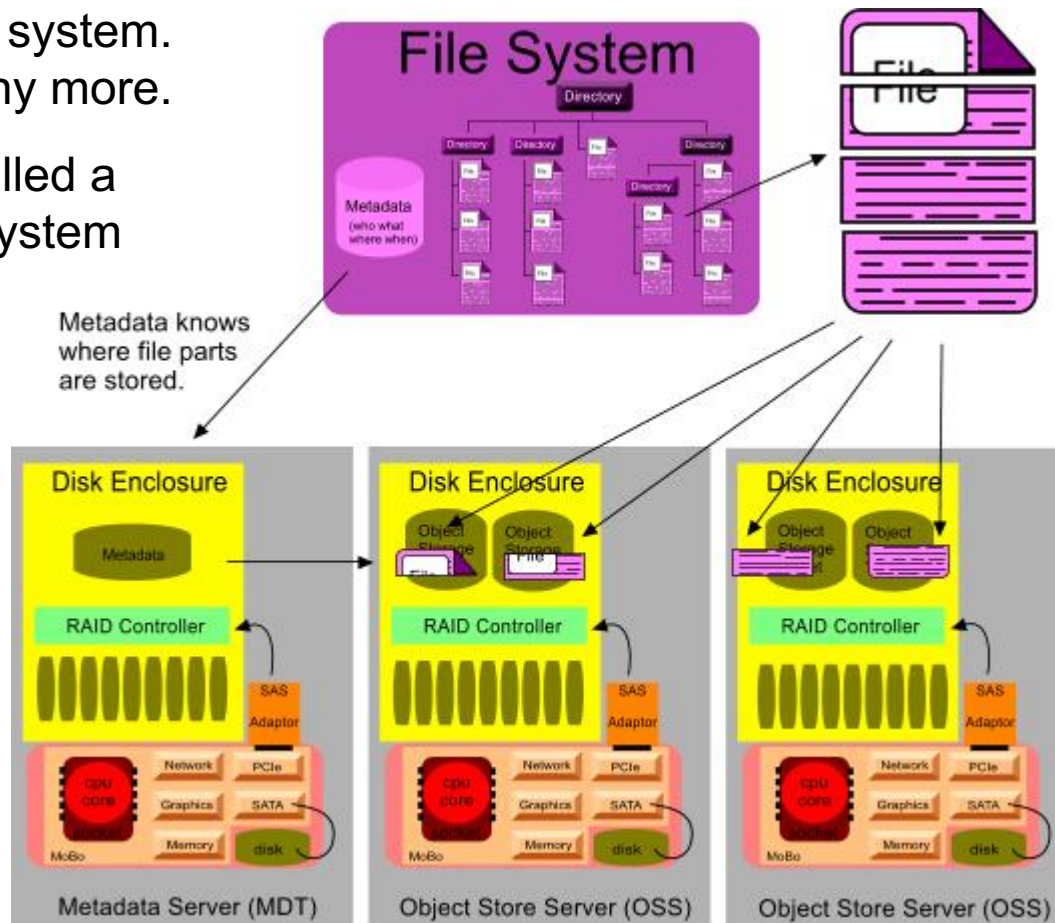
Main article: [List of distributed file systems](#)

- [Alluxio](#)
- [BeeGFS](#) (Fraunhofer)
- [Ceph](#) (Inktank, Red Hat, SUSE)
- [Windows Distributed File System \(DFS\)](#) (Microsoft)
- [Infini](#)
- [GfarmFS](#)
- [GlusterFS](#) (Red Hat)
- [GFS](#) (Google Inc.)
- [HDFS](#) (Apache Software Foundation)
- [IPFS](#)
- [iRODS](#)
- [LizardFS](#) (Skytechnology)
- [MapR FS](#)
- [MooseFS](#) (Core Technology / Gemius)
- [ObjectiveFS](#)
- [OneFS](#) (EMC Isilon)
- [OpenIO](#)
- [OrangeFS](#) (Clemson University, Omnibond Systems), formerly [Parallel Virtual File System](#)
- [Panfs](#) (Panasas)
- [Parallel Virtual File System](#) (Clemson University, Argonne National Laboratory, Ohio Supercomputer Center)
- [RozoFS](#) (Rozo Systems)
- [Torus](#) (CoreOS)
- [XtreemFS](#)

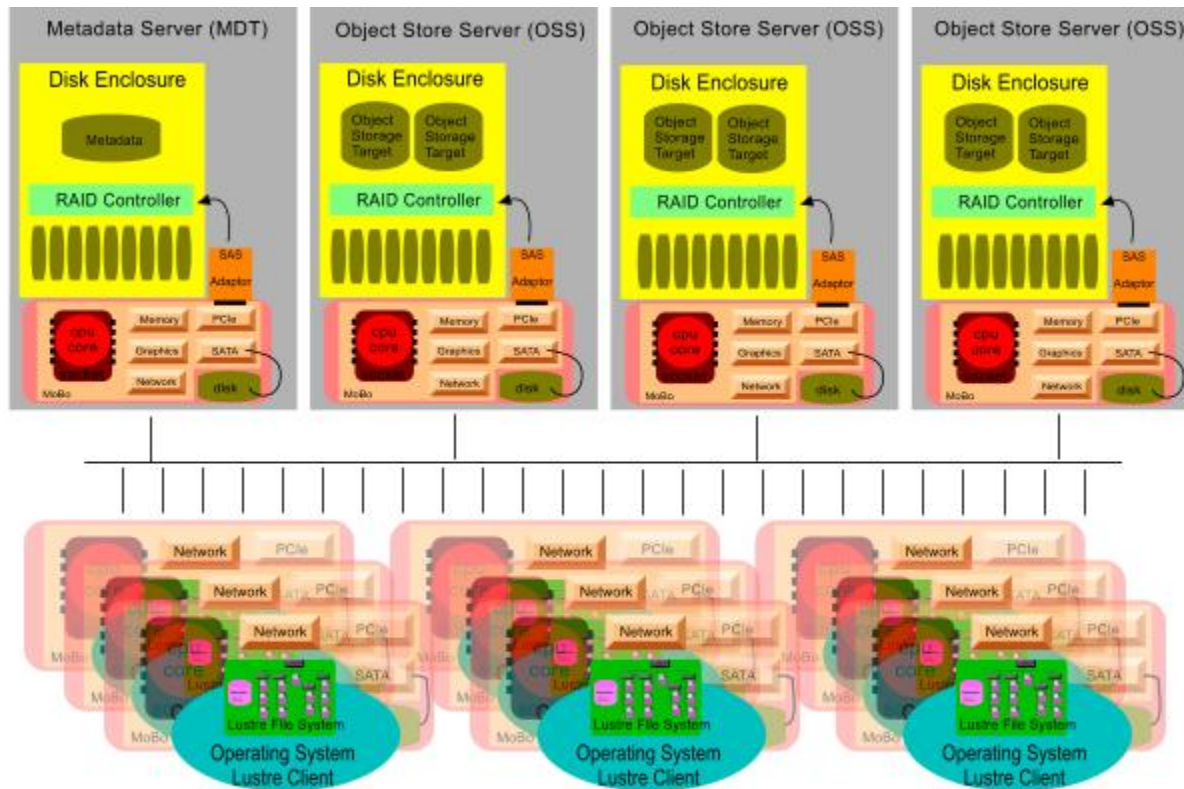
NAS - Lustre File System

Lustre is an example of a distributed file system. There are many more.

Sometimes called a “Cluster” file system



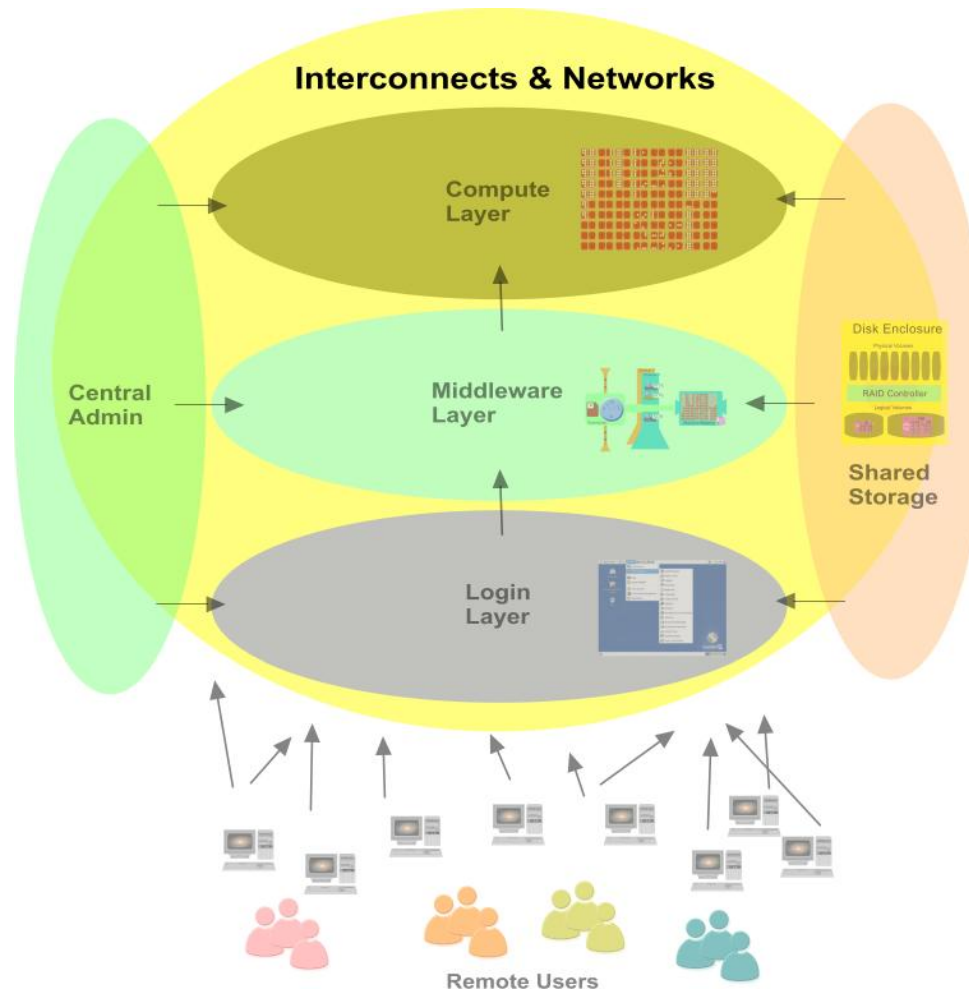
NAS – Lustre



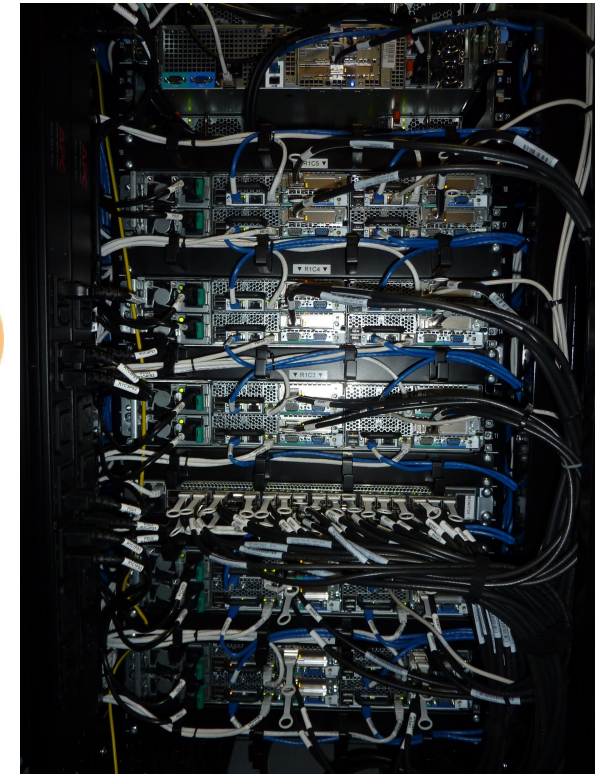
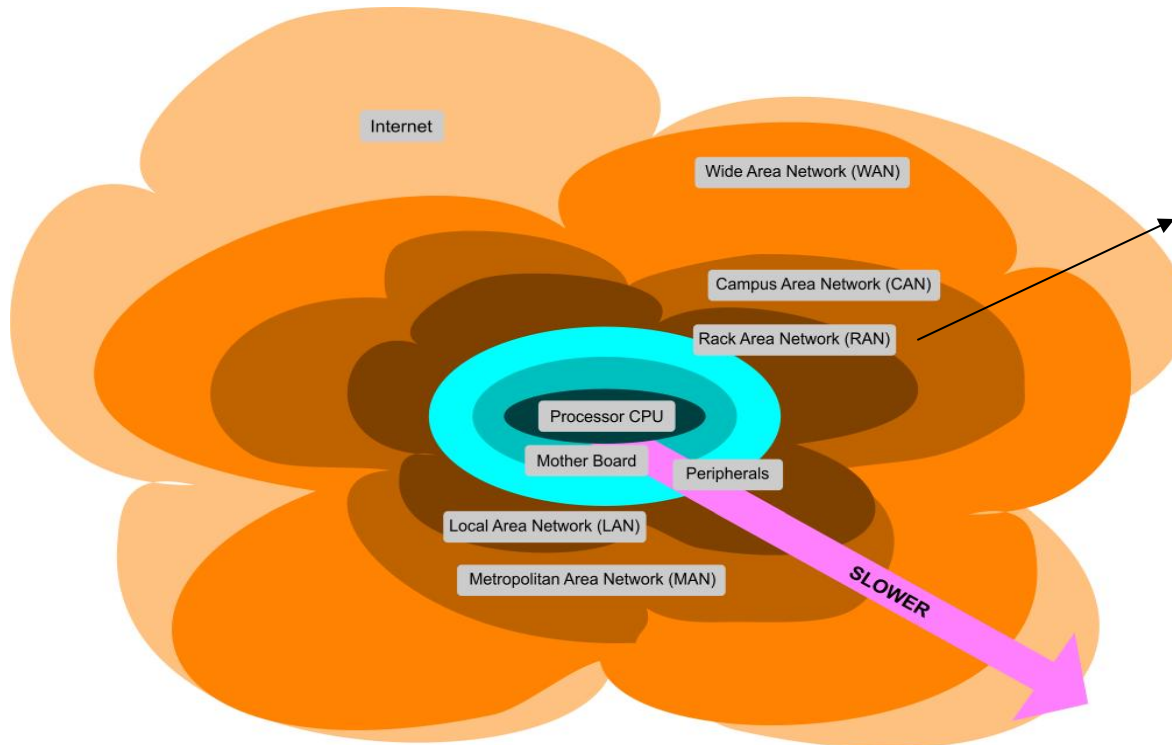
Often used with an Infiniband fabric.

LUSTRE - Massively Parallel Distributed Shared File System

Interconnects and Networks

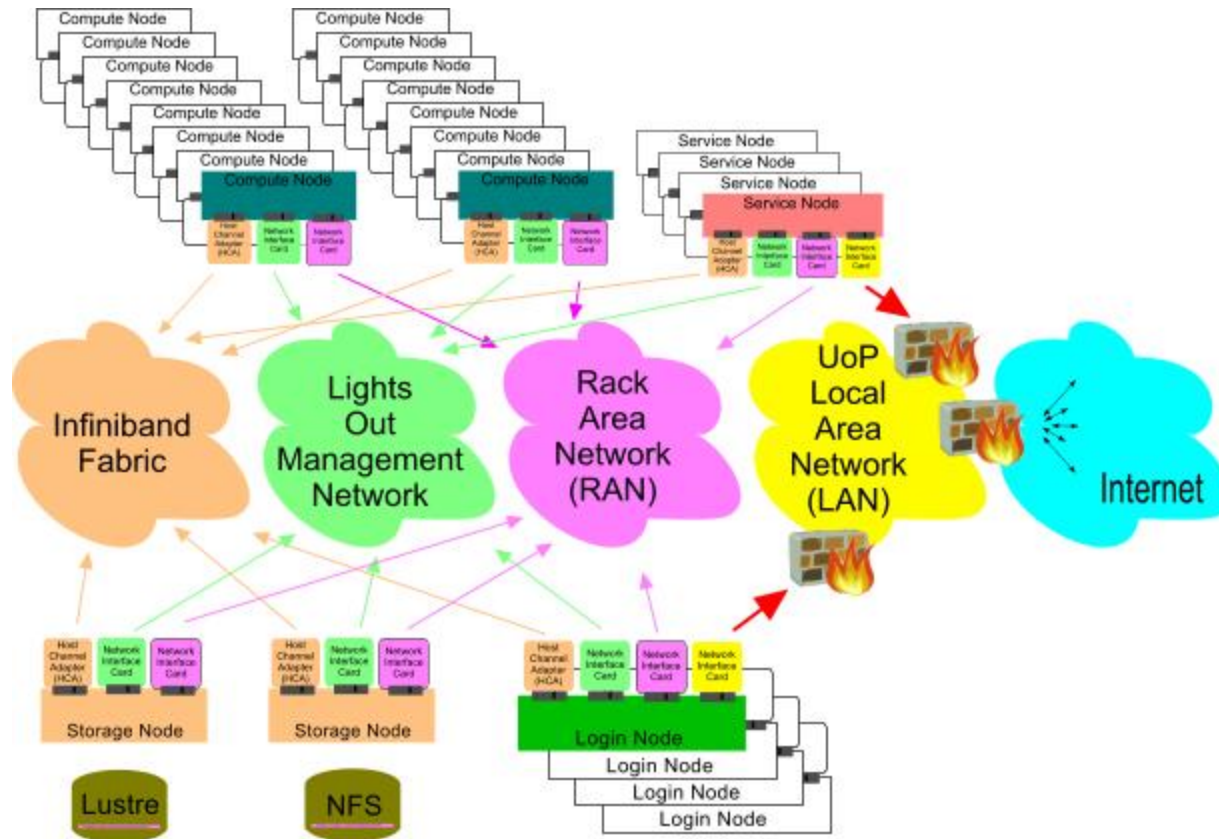


Interconnects and Networks



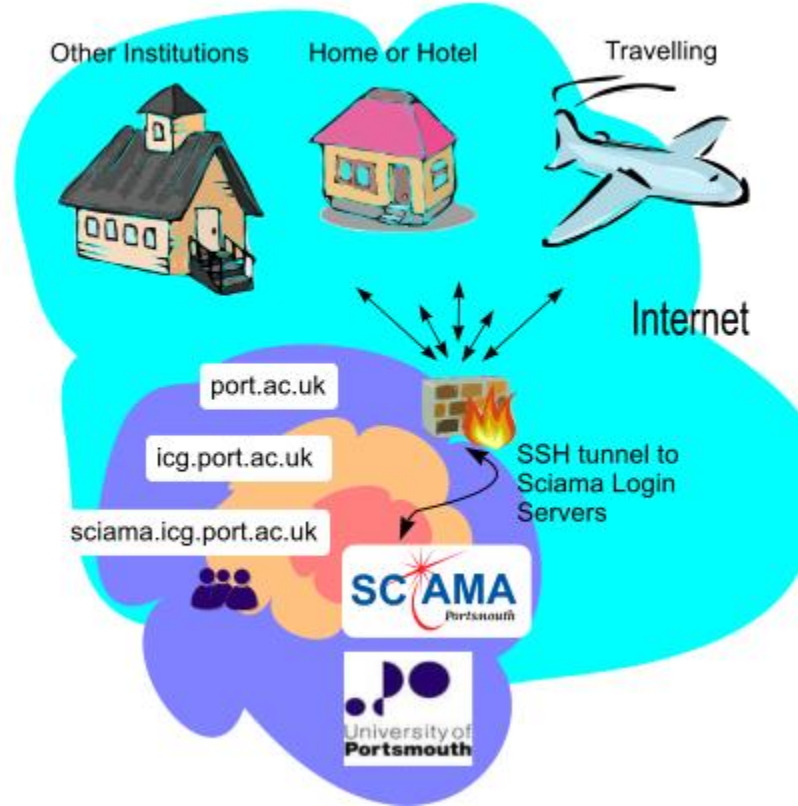
Moving away from the Processor towards the Internet you get slower and slower due to Increased Latency and Reduced Bandwidth

Commodity Cluster Networks

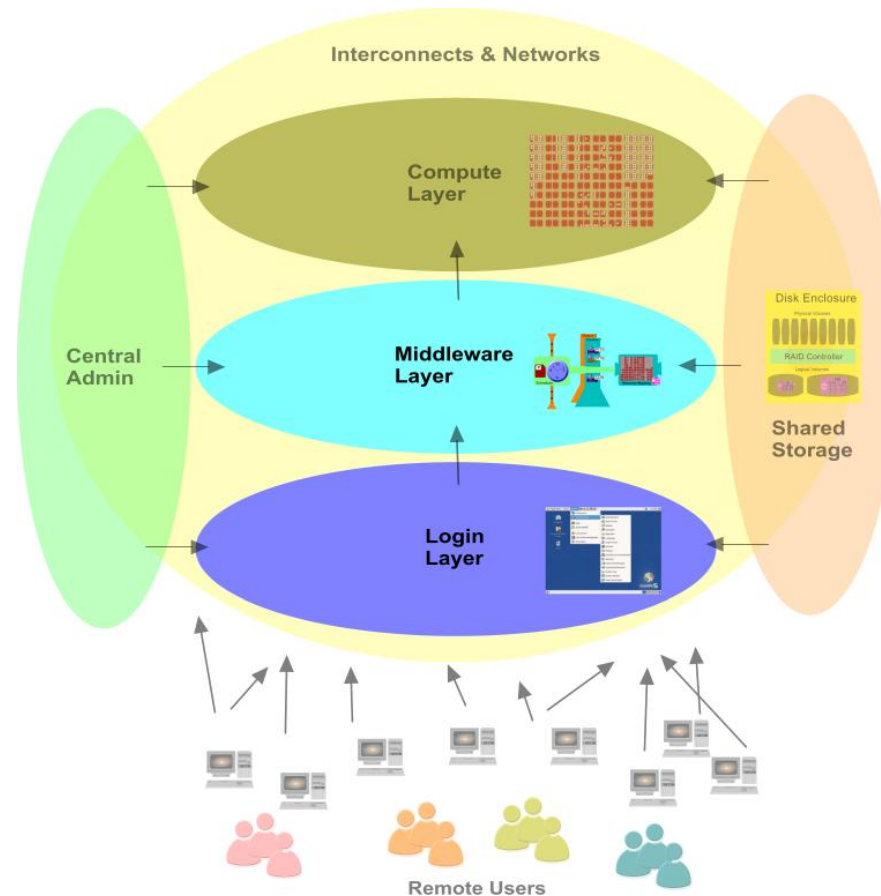


Connecting to an HPC

Remote Access to Sciama



Login Layer

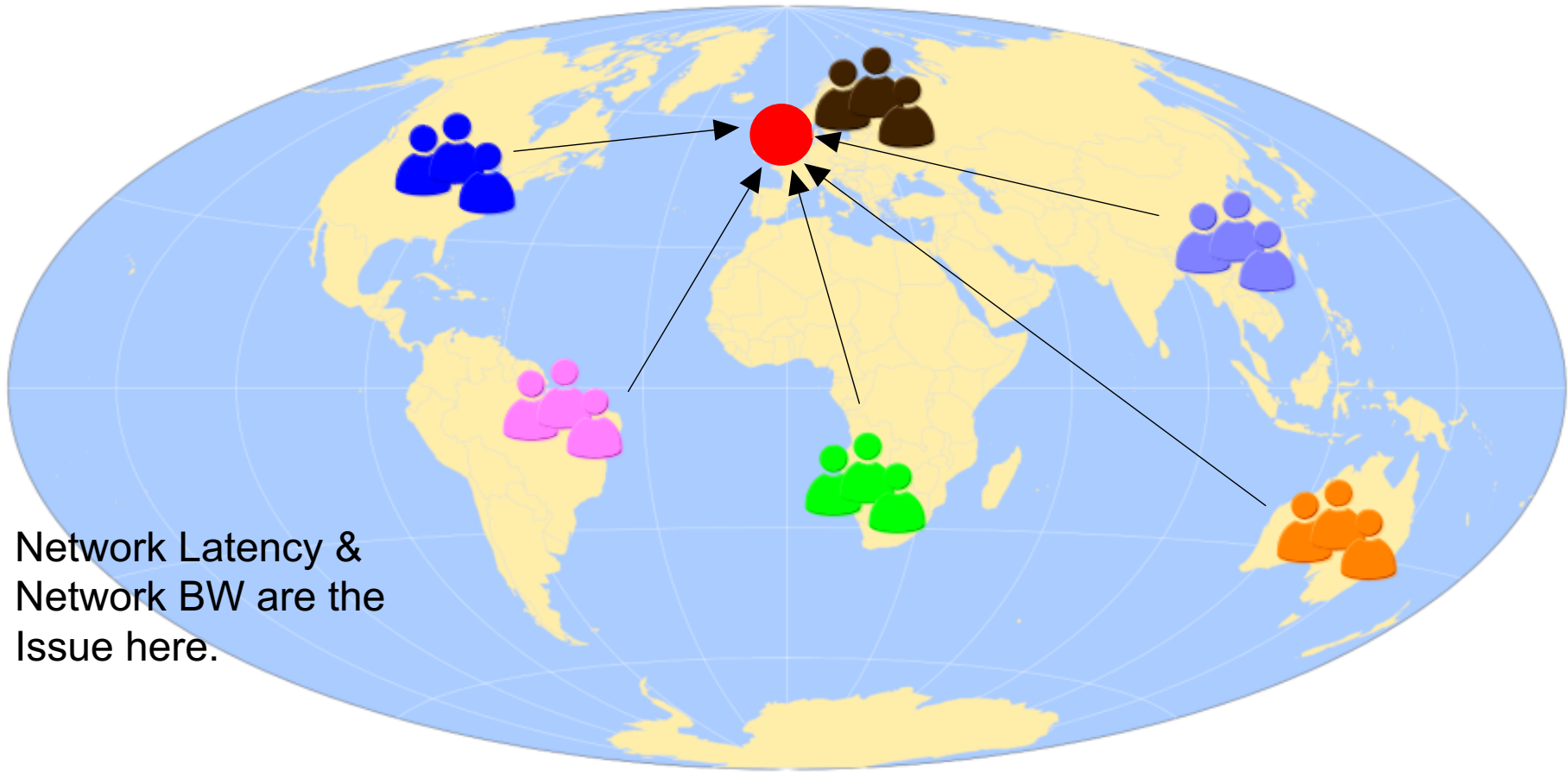


Why Login Servers

- Login servers will provide the gateway to the cluster.
- Users can remotely login into the servers using “ssh” or a Remote Desktop Client.
- A desktop client gives a full working desktop in the environment (can full screen)



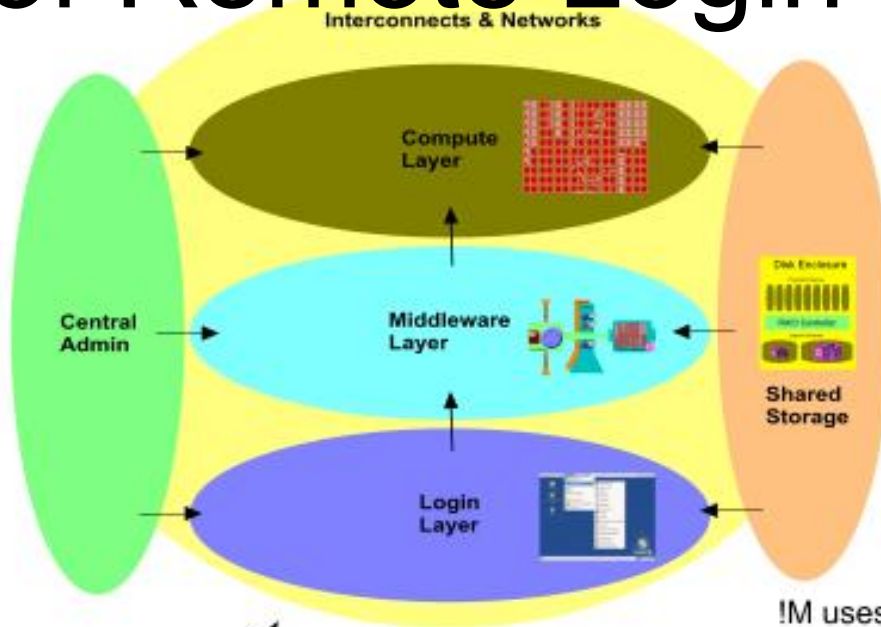
Global User Base



Network Latency &
Network BW are the
Issue here.

Use of Remote Login Client

1.) User downloads IM client from Internet. Client is available for Linux, Windows and MAC. The client is free.



IM uses a very lean protocol for use over high latency low bandwidth links.

2.) Once installed and configured the IM client can connect to the IM server running in the cluster.

3.) the user is returned a login prompt or a suspended session.



Remote Users





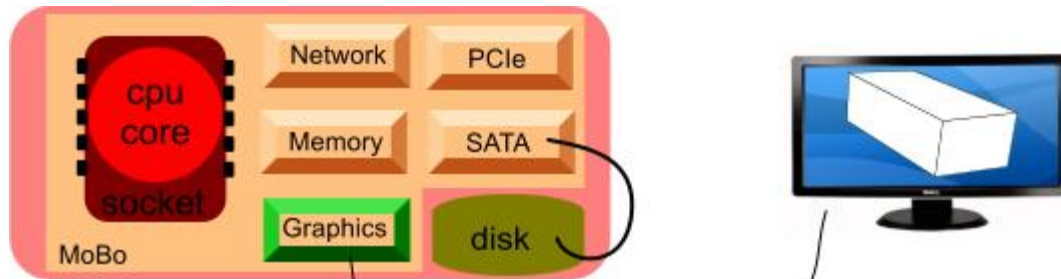
GPGPU's

General Purpose Graphics Processing Units

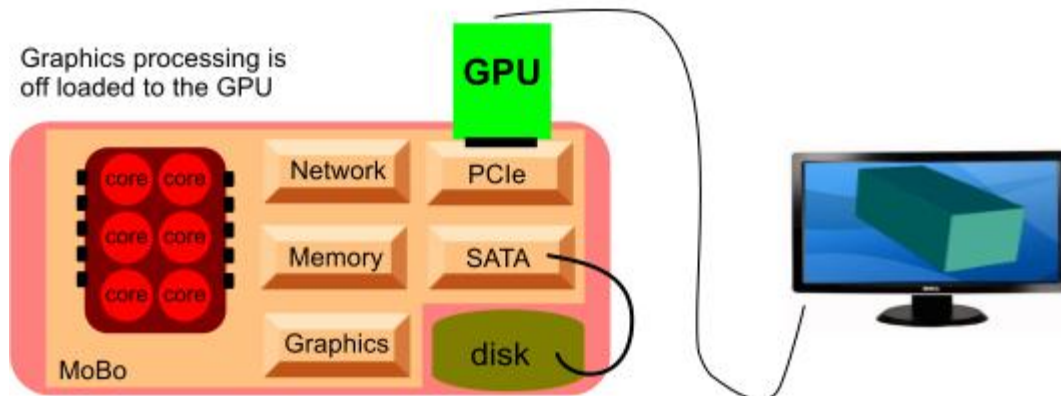
G Burton – ICG – Oct12 – v1.0



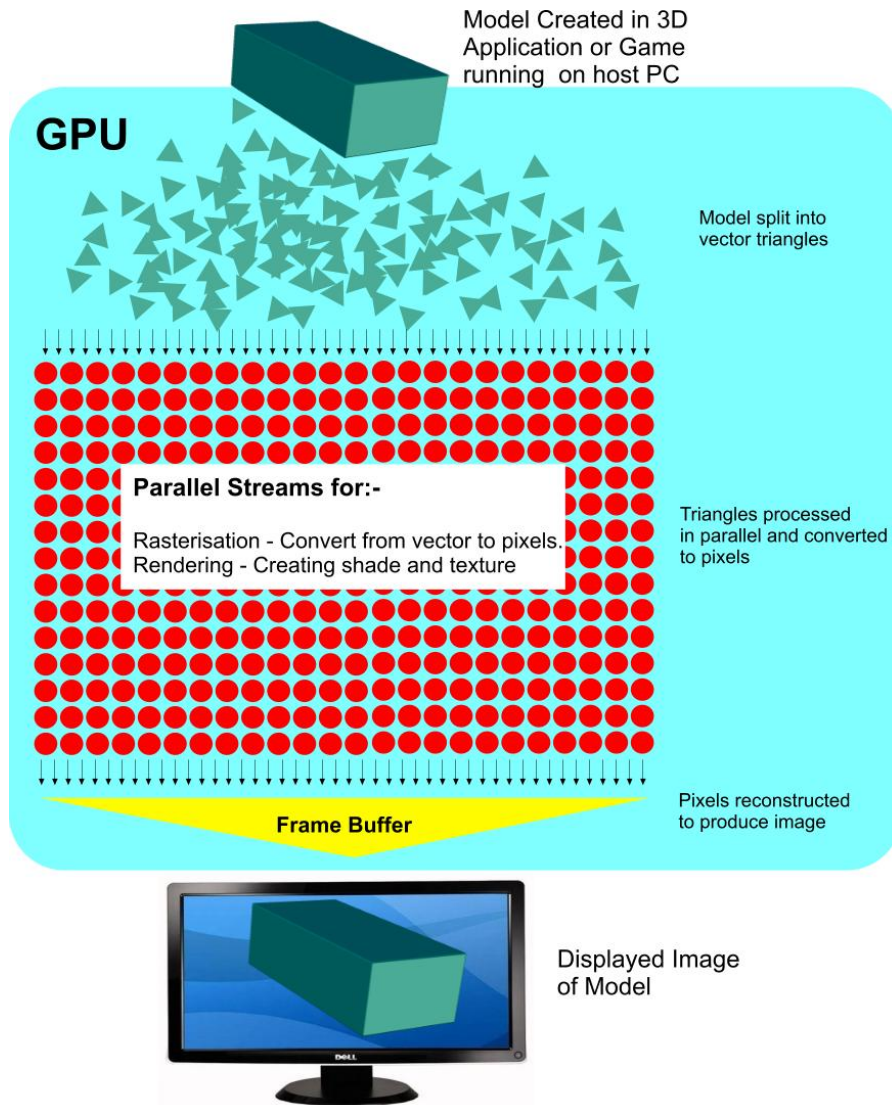
Graphics Processing



Graphics processing is handled by the host CPU.



Graphics processing is off loaded to the GPU



Cheap commodity hardware mainly from Nvidia, AMD and Intel for home PC's.

Driven by gaming.

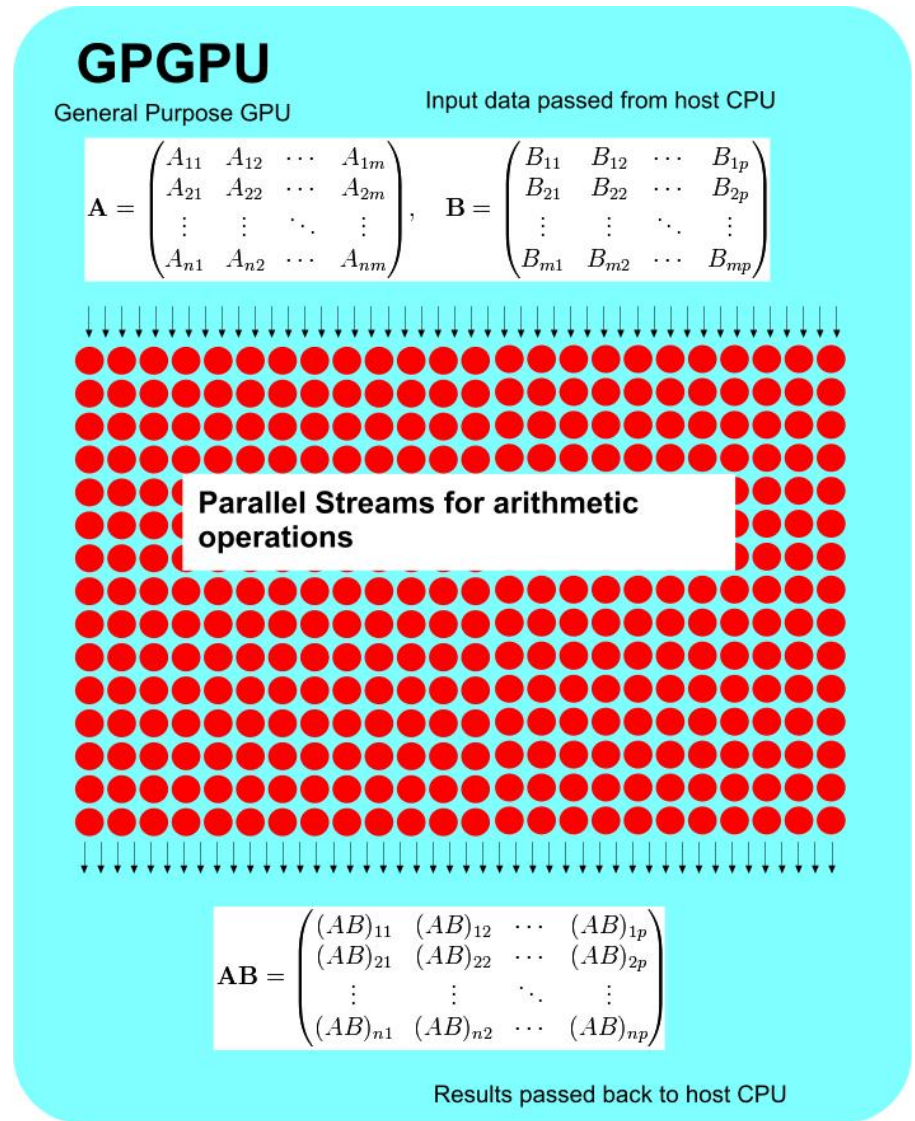
Bespoke architecture replaced by generic programmable architecture.

The Birth of Cuda

(compute Unified Device Architecture) **and**
OpenCL.

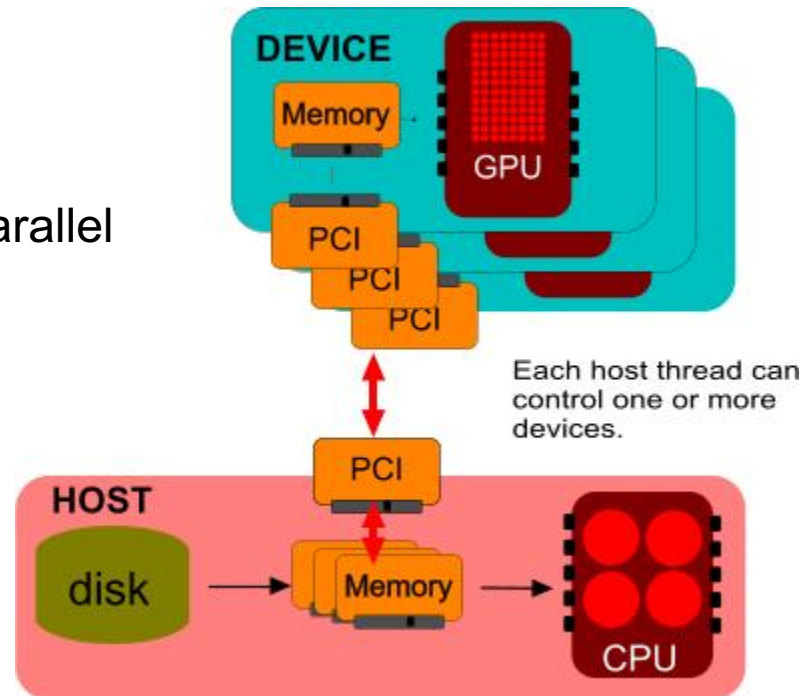
Cuda (Nvida) is
cutting edge,
OpenCL follows.

If you know Cuda
then OpenCL
Is easy.

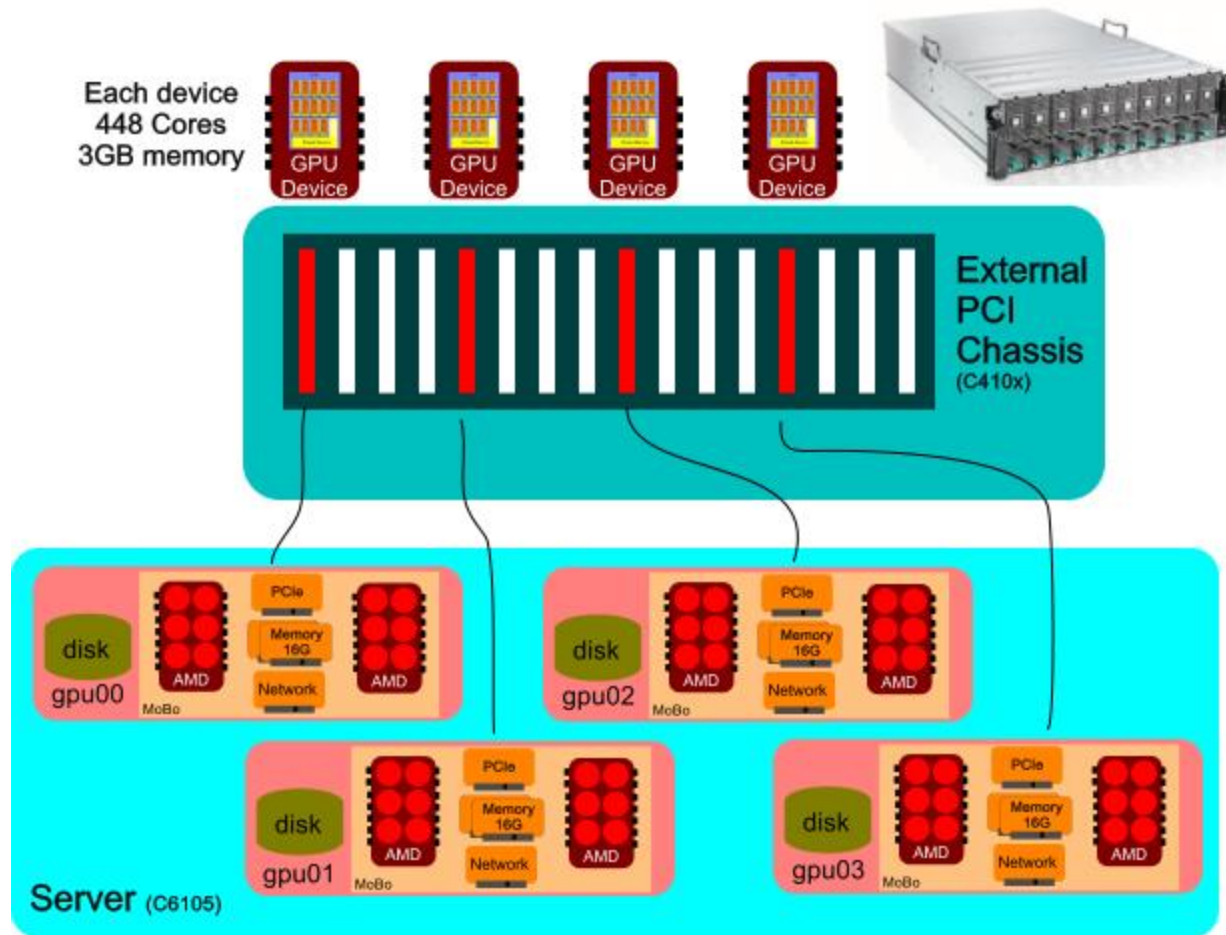


Several GPU cards can be connected in Parallel

Allows 1000's of GPU cores for massively parallel computation.



Sciama Specifics



Sciama GPU Specification

CUDA Device Query (Runtime API) version (CUDA static linking)

Found 1 CUDA Capable device(s)

Device 0: "Tesla M2050"

CUDA Driver Version / Runtime Version 4.20 / 4.0

CUDA Capability Major/Minor version number: 2.0

Total amount of global memory: 2687 MBytes (2817982464 bytes)

(14) Multiprocessors x (32) CUDA Cores/MP: 448 CUDA Cores

GPU Clock Speed: 1.15 GHz

Memory Clock rate: 1546.00 Mhz

Memory Bus Width: 384-bit

L2 Cache Size: 786432 bytes

Max Texture Dimension Size (x,y,z) 1D=(65536), 2D=(65536,65535), 3D=(2048,2048,2048)

Max Layered Texture Size (dim) x layers 1D=(16384) x 2048, 2D=(16384,16384) x 2048

Total amount of constant memory: 65536 bytes

Total amount of shared memory per block: 49152 bytes

Total number of registers available per block: 32768

Warp size: 32

Maximum number of threads per block: 1024

Maximum sizes of each dimension of a block: 1024 x 1024 x 64

Maximum sizes of each dimension of a grid: 65535 x 65535 x 65535

Maximum memory pitch: 2147483647 bytes

Texture alignment: 512 bytes

Concurrent copy and execution: Yes with 2 copy engine(s)

Run time limit on kernels: No

Integrated GPU sharing Host Memory: No

Support host page-locked memory mapping: Yes

Concurrent kernel execution: Yes

Alignment requirement for Surfaces: Yes

Device has ECC support enabled: Yes

Device is using TCC driver mode: No

Device supports Unified Addressing (UVA): Yes

Device PCI Bus ID / PCI location ID: 15 / 0

Compute Mode:

< Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 4.20, CUDA Runtime Version = 4.0, NumDevs = 1,

Device = Tesla M2050

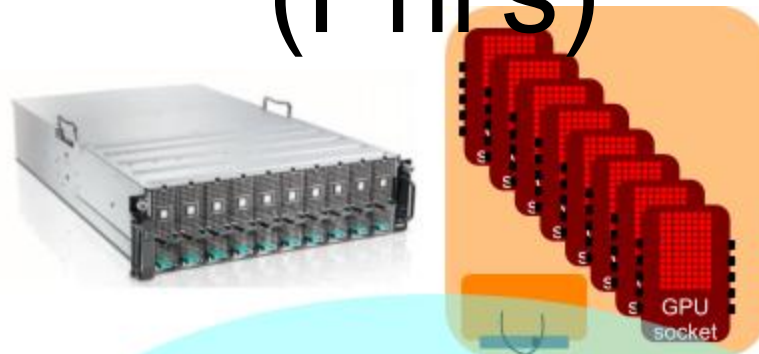
[gputest] test results...

PASSED

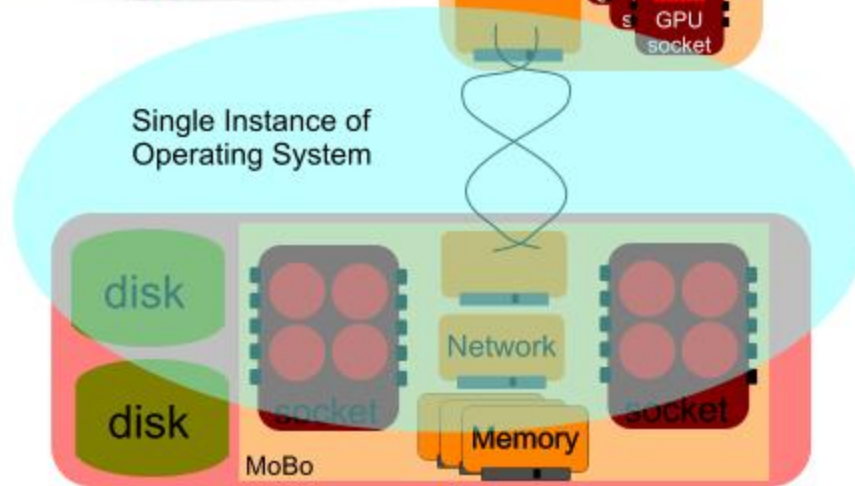
Graphical Processing Units (GPU's) and Intel CoProcessors (Phi's)

Three players:-

Intel
AMD
Nvidia



Special programming language. CUDA and OpenCL



Cpu – multiple cores
Gpu – 100's of cores

CPU's still in charge