

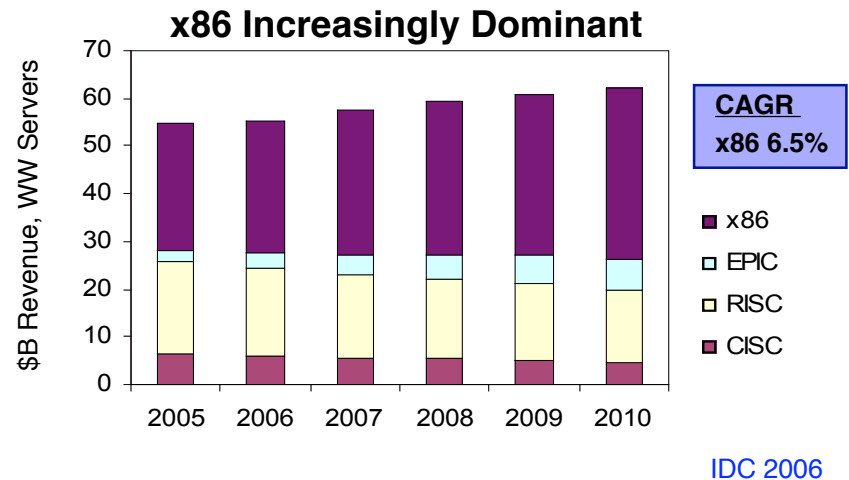
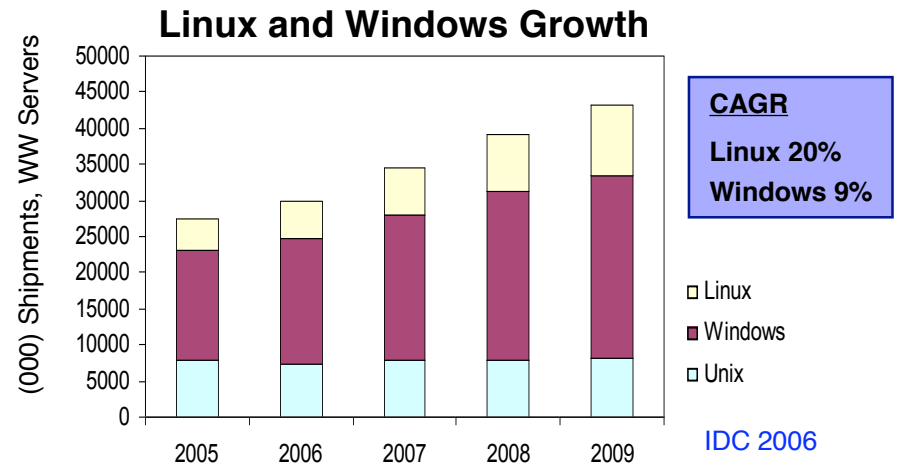
The Need For Speed

Rick Reid
Principal Engineer
SGI



Commodity Systems

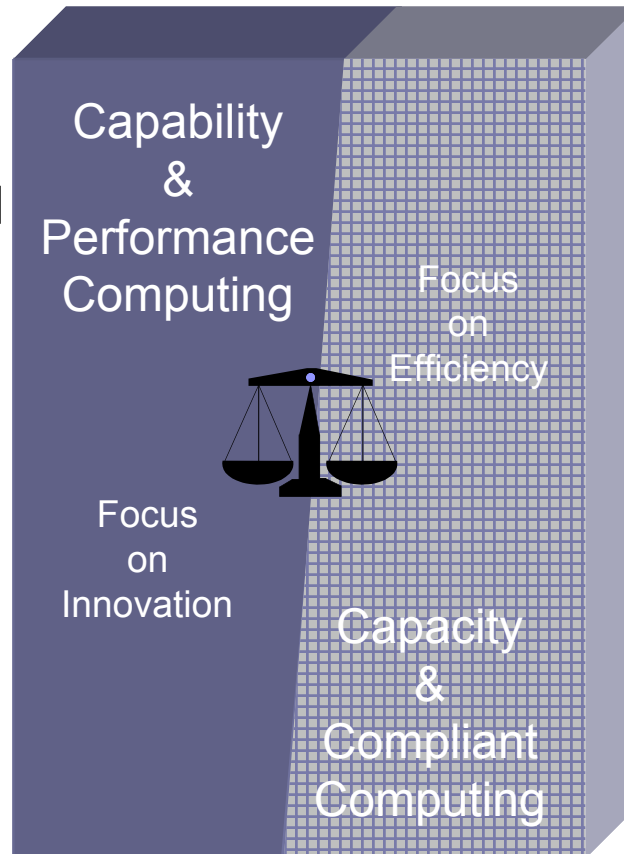
- Linux
 - Red Hat
 - SUSE
 - SE-Linux
- X86-64
 - Intel Xeon
 - AMD
- Scalable Programming Model
 - MPI
- Global Data Access
 - NFS
 - NFS/RDMA
 - pNFS



Servers: One Size Doesn't Fit All!

Workflow Characteristics

- Performance Oriented
- Data-Intensive
- Product Timelines
- Interactivity
- Rapid development cycles

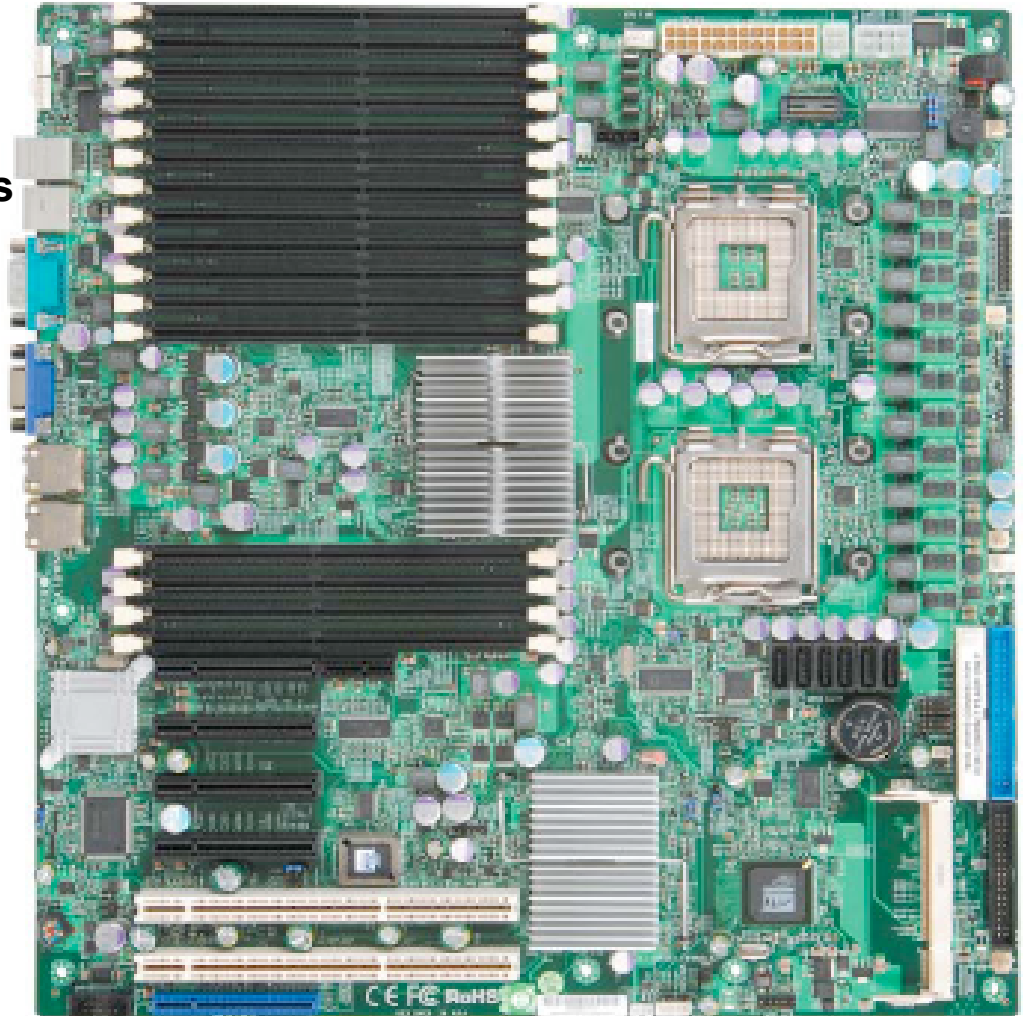


Workflow Characteristics

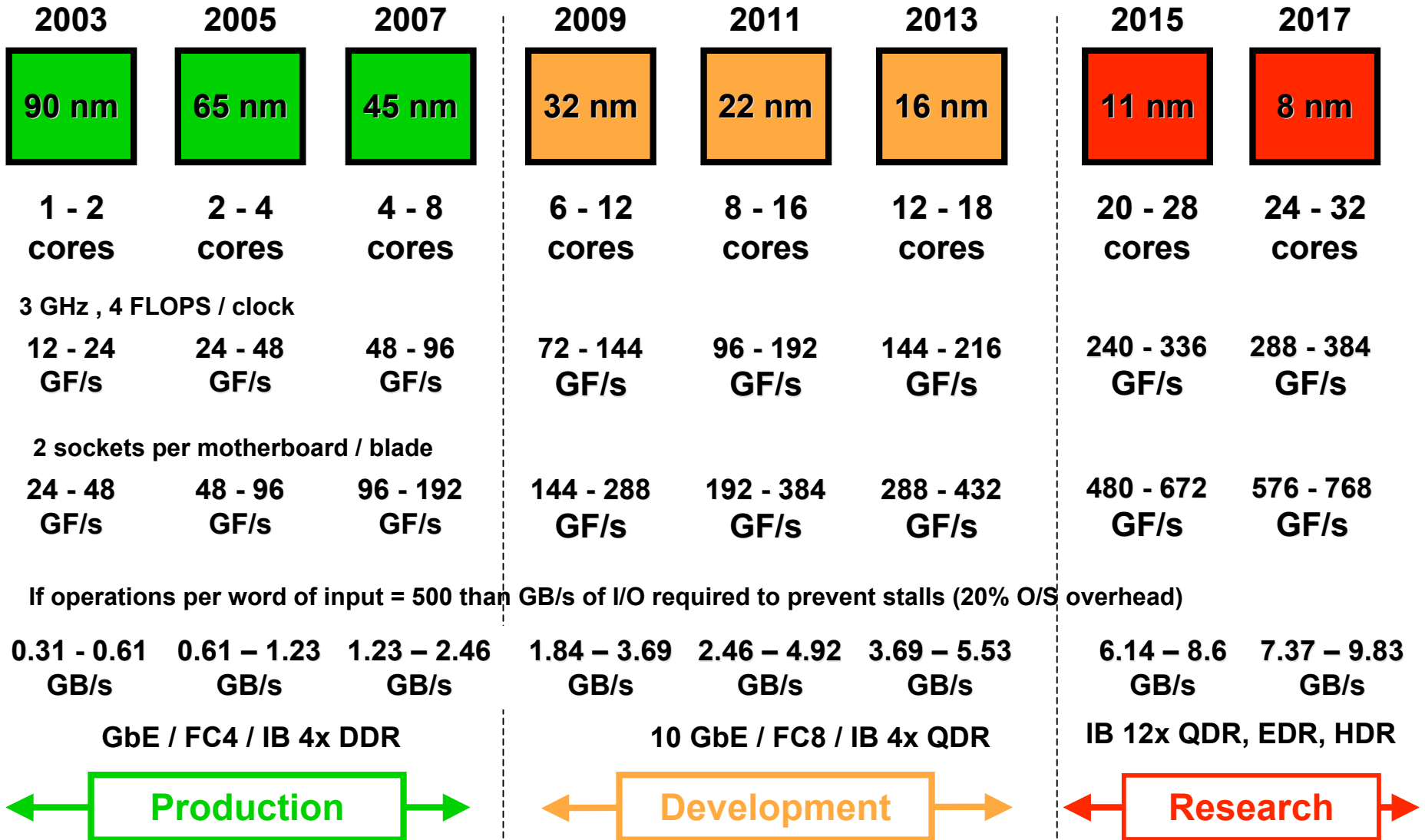
- Price-performance key
- Little data sharing
- Predictable Workloads
- Non-interactive
- Mature Apps

Commodity Xeon Motherboard – Today

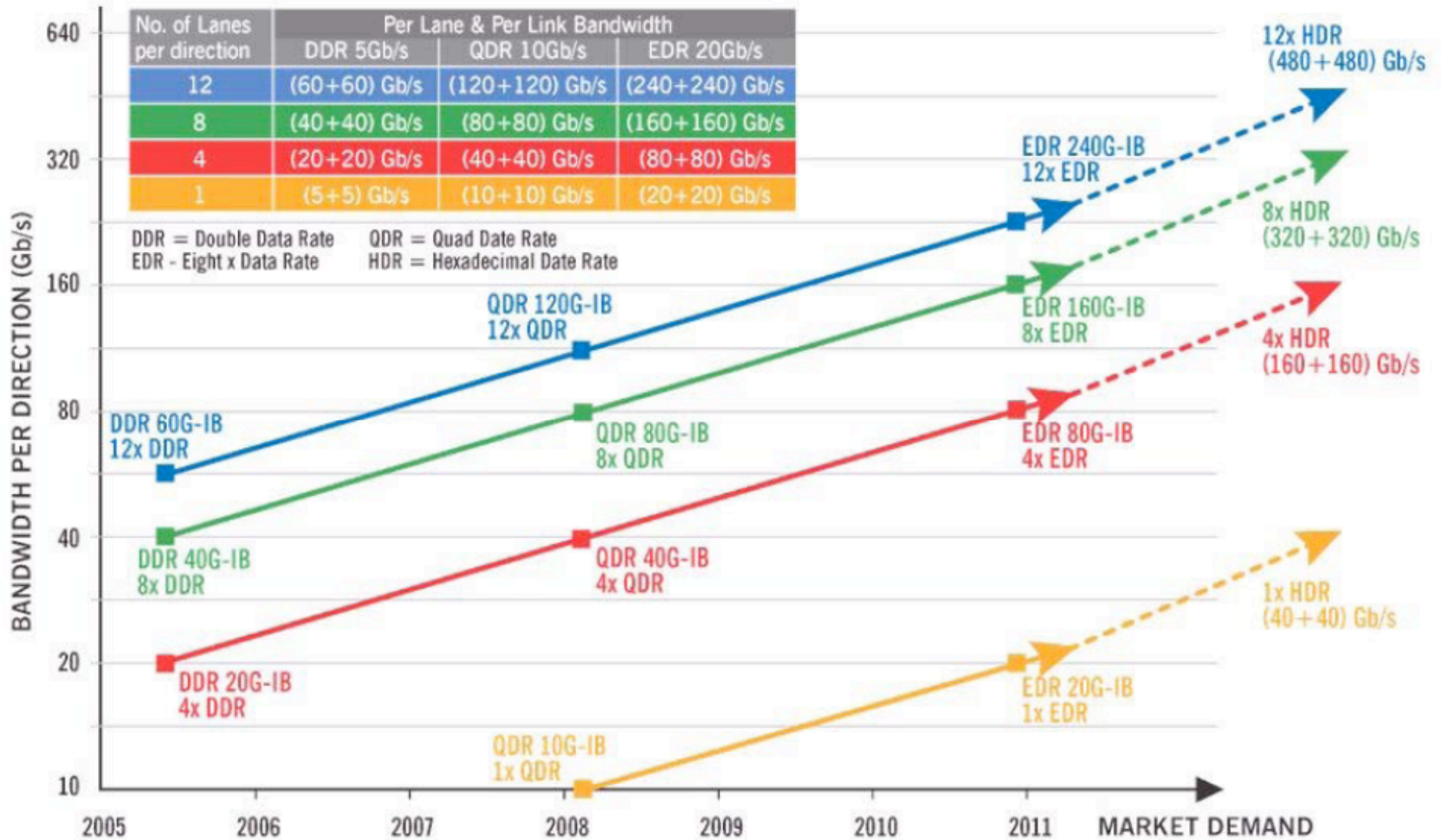
- Seaburg chipset
 - Wolfdale and Harpertown
 - 1600 MHz or 1333 MHz front side bus
- 16 FB DIMM slots
 - 64 GB ECC DDR2 FBD
- 8 SATA2 or SAS drives
 - H/W RAID 0, 1, 5, 6, 10, 50, 60
- ATI 32 MB DDR2 Graphics
- Dual GbE interfaces
- Redundant power supplies
- 5 external I/O slots
 - 2 x PCIe x8 gen 2
 - 1 x PCIe x8 gen 1
 - 1 or 2 x PCIe x4 gen 1
 - 1 or 2 x PCI-x 133/100



Fabrication Technology Drives the Need



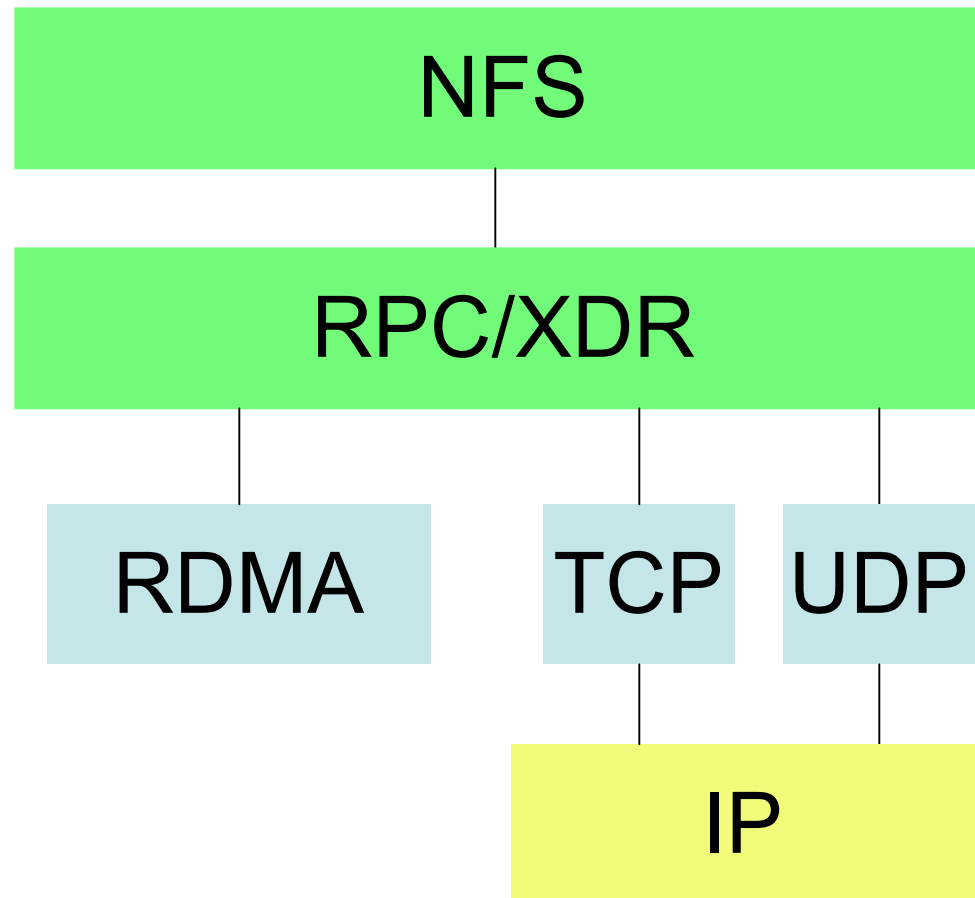
Infiniband Roadmap



NFS/RDMA - A new NFS standard

- NFS/RDMA is provided by transport switch in the RPC layer which allows data to be transported via TCP/IP, UDP/IP or RDMA over Infiniband.
- The transport switch allows NFS to use the efficient, high bandwidth native Infiniband protocol.
- NFS/RDMA can be simultaneously intermixed with NFS over IP-based protocols on the same server and even the same files
- NFS/RDMA transfers the data from the Servers memory to the client memory without host CPU intervention
 - This allows the NFS server to serve more clients
 - This improves NFS bandwidth
 - This allows the NFS client to dedicate more CPU resources to the application
- Latest OpenFabric drivers provide most of the capabilities required
 - Current version: OpenFabrics Enterprise Distribution (OFED) 1.3

What is NFS over RDMA?



NFS/RDMA

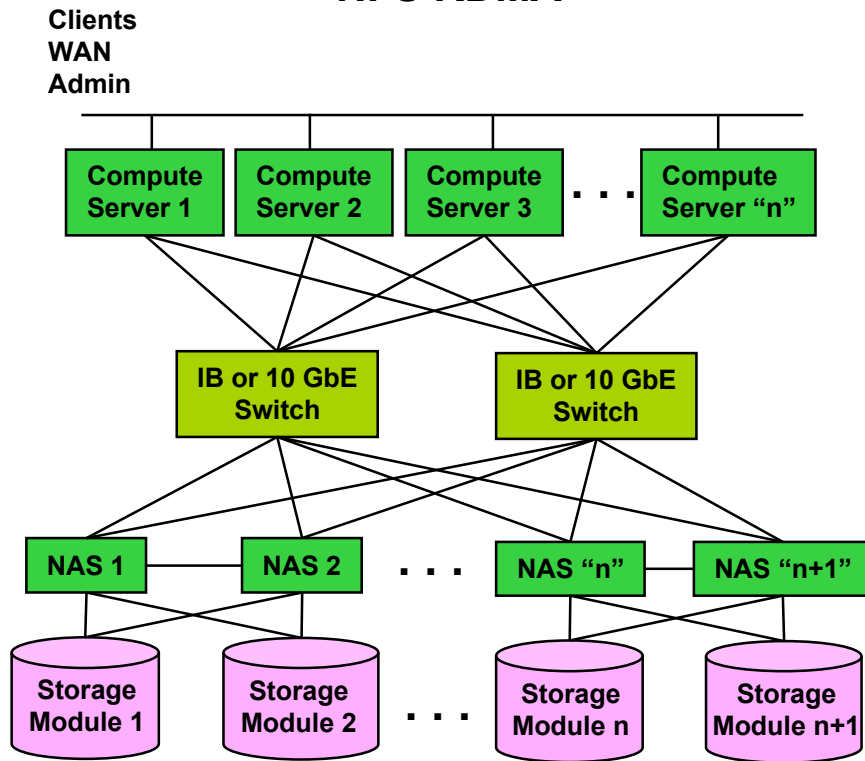
- **IB or 10 GbE connectivity between NFS server and NFS clients**
- **IB, SAS or FC backend to storage**
- **Performance**
 - 10 RAIDs, 480 15K SAS drives, 96 GB RAM, 2 IB interfaces
 - NFS/RDMA IB Read: 3.5 GB/s
 - NFS/RDMA IB Write: 2.25 GB/s
 - NFS 10 GbE Read: 800 MB/s
 - NFS 10 GbE Write: 400 MB/s
- **Higher aggregate bandwidth**
- **Lower CPU overhead - ~10% CPU load/IB interface**

pNFS - Future Technology

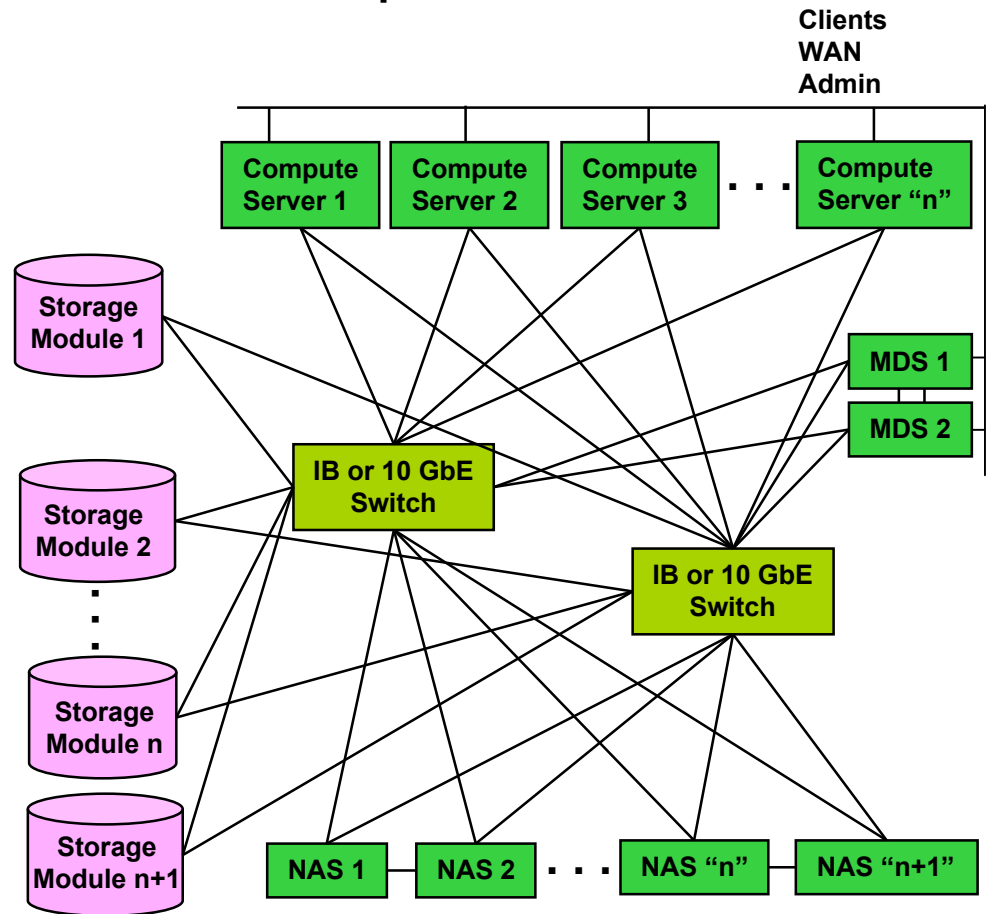
- Parallel NFS: Extension of NFS 4.x Standard
 - Expected draft release in NFS 4.1
- Dr. Garth Gibson, professor at CMU, wrote the Internet Draft:
 - “pnfs-problem-statement”
 - www.pdl.cmu.edu/pNFS/archive/gibson-pnfs-problem-statement.html
- File base sharing
- Uses metadata servers
 - NFS semantics for metadata access
 - Metadata server points client to the storage node
 - Supports Block, File System, or Object Storage backend
 - Expected to be able to serve 100K+ of clients
- Multiple storage and system suppliers have stated their plans to support pNFS

NAS Storage

NFS RDMA



pNFS



Questions

Thank You

