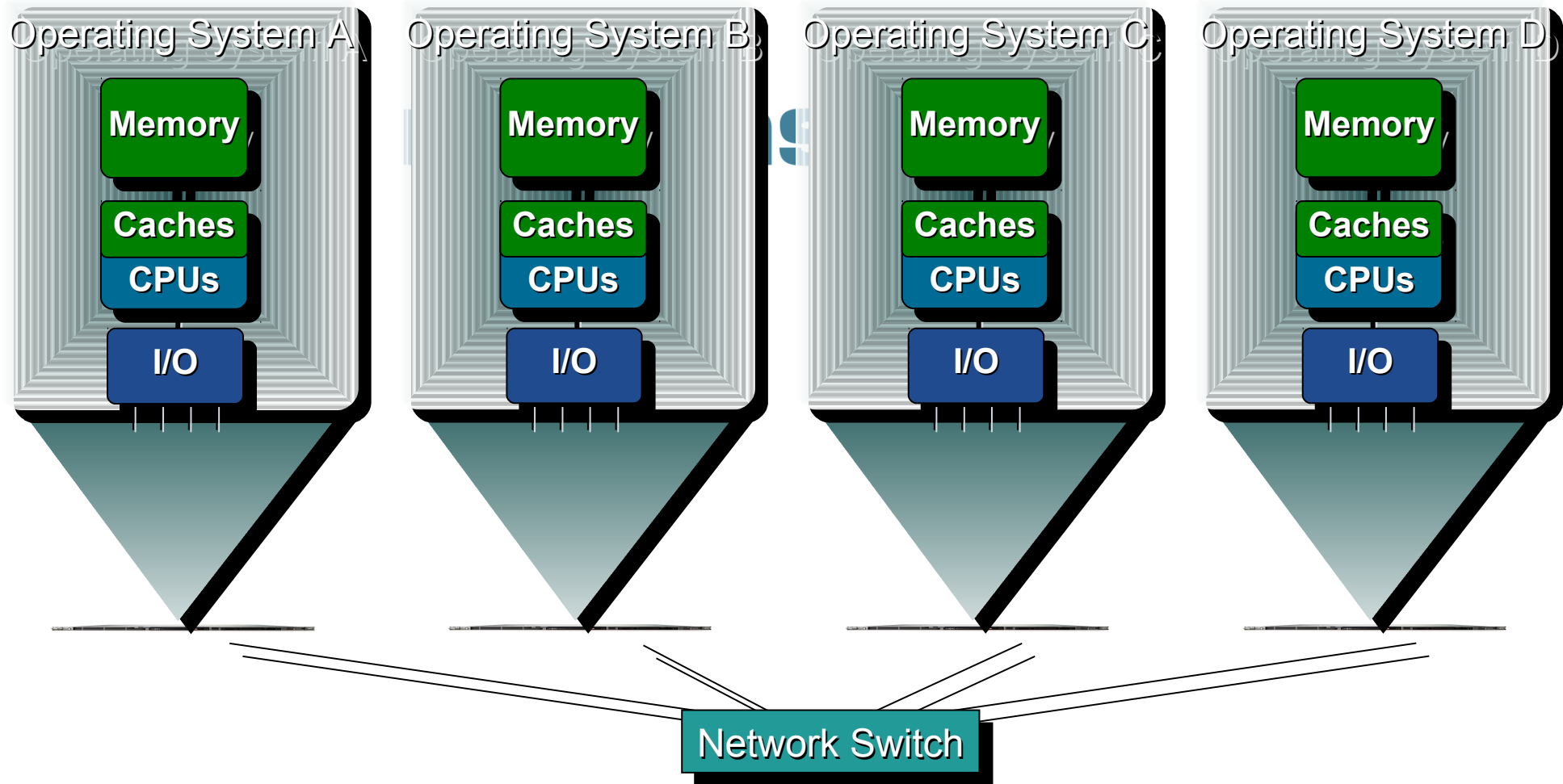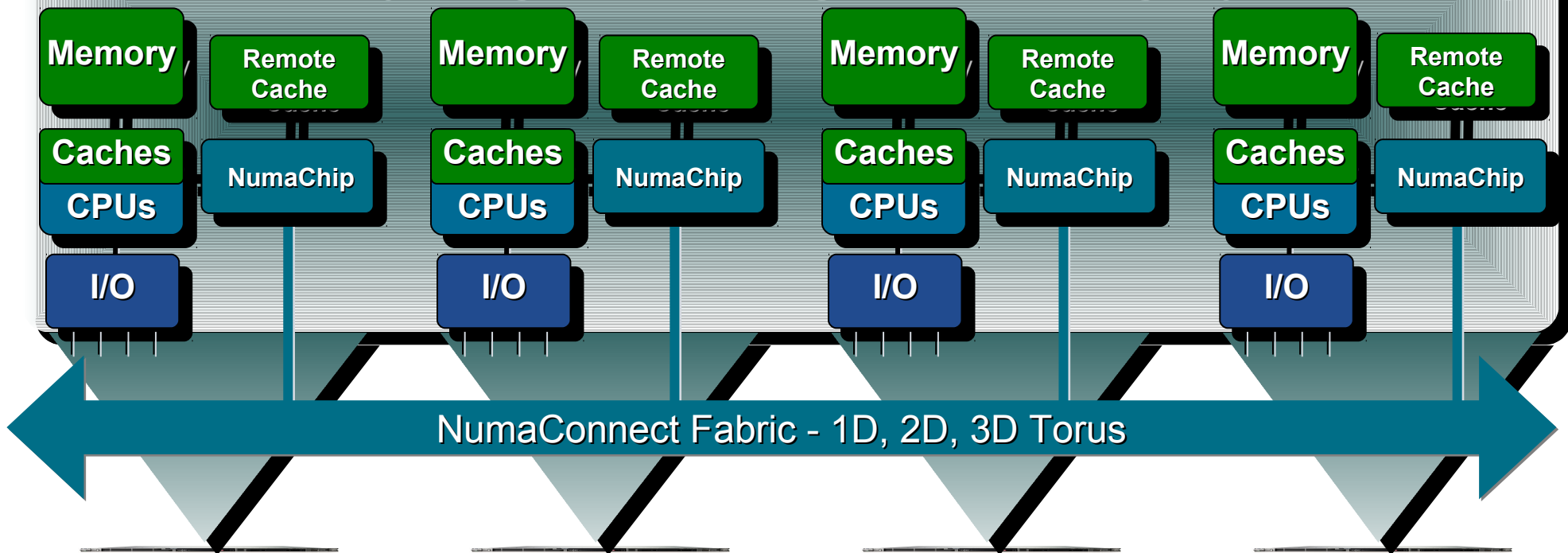# Clusters - NO Shared Resources

## Individual Instances of the Operating System
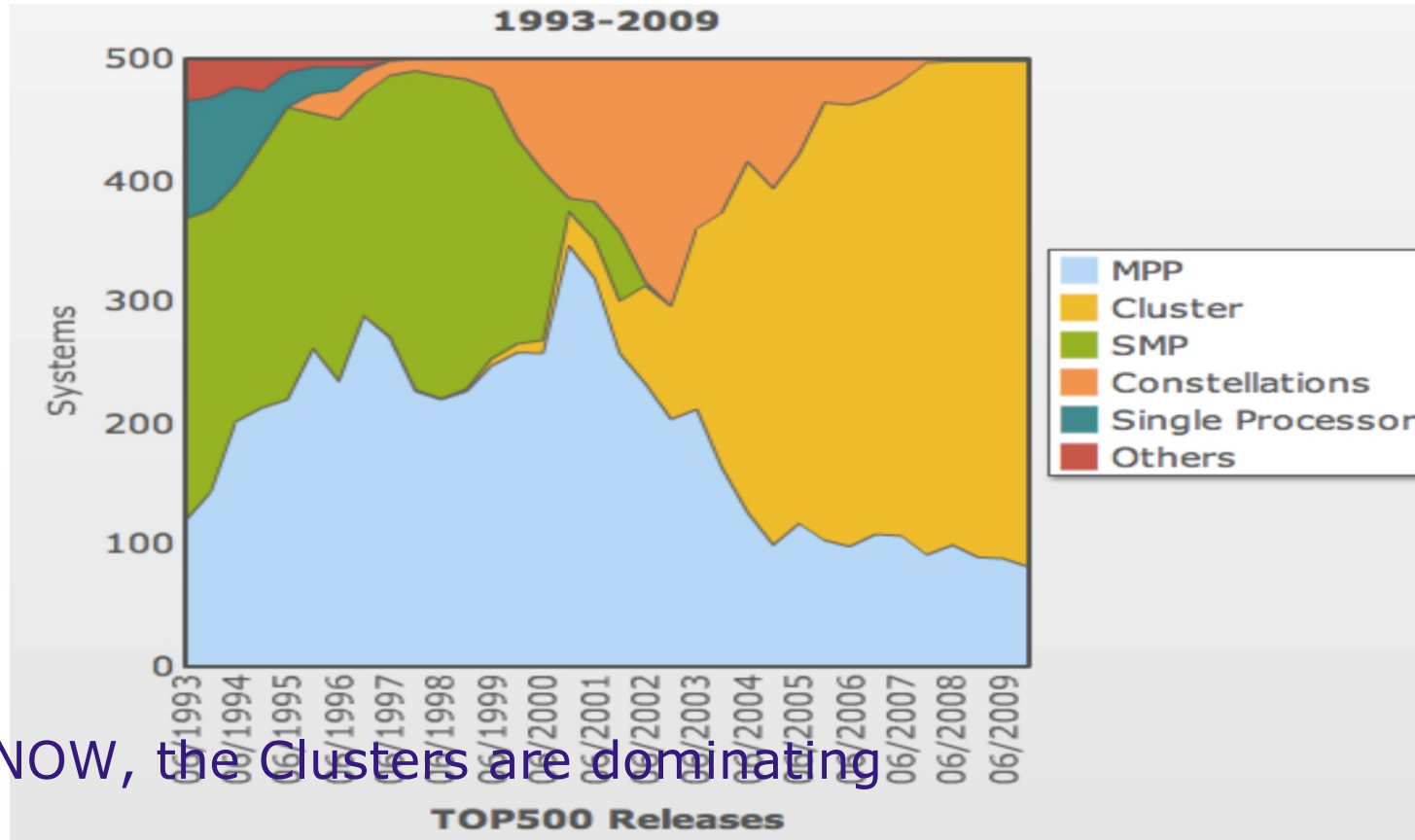
# Cache Coherent Shared Memory



Shared Everything - One Single Operating System Image

| Memory | Remote Cache | Memory | Remote Cache | Memory | Remote Cache | Memory | Remote Cache |

| Caches | NumaChip | Caches | NumaChip | Caches | NumaChip | Caches | NumaChip |

CPUs / CPUs / CPUs / CPUs

I/O / I/O / I/O / I/O

NumaConnect Fabric - 1D, 2D, 3D Torus

## Capabilities like Mainframe - Price like Cluster

**numascale**

- The expensive SMPs used to rule:
  - Cray XMP, Convex Exemplar, Sun ES



1993-2009

Legend: MPP, Cluster, SMP, Constellations, Single Processor, Others
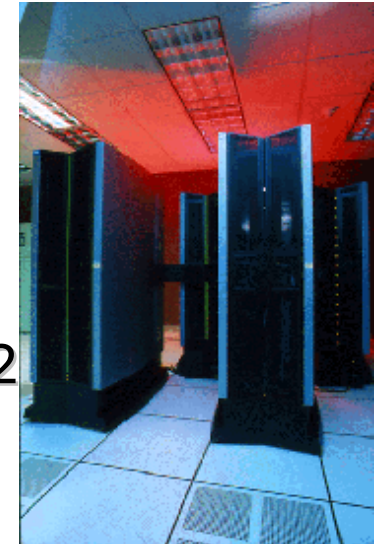
TOP500 Releases

- NOW, the Clusters are dominating

# numascale

**Convex Exemplar (Acquired by HP)**

- First implementation of the ccNUMA architecture from Dolphin in 1994

**Data General Aviion (Acquired by EMC)**

- Designed in 1996 with deliveries from 1997 - 2002
- Used Dolphin's chips with 3 generations of processor/memory buses

**I/O Attached Products for Clustering OEMs**

- Sun Microsystems (SunCluster)
- Siemens RM600 Server (IO Expansion)
- Siemens Medical (3D CT)
- Philips Medical (3D Ultra Sound)
- Dassault/Thales Rafale

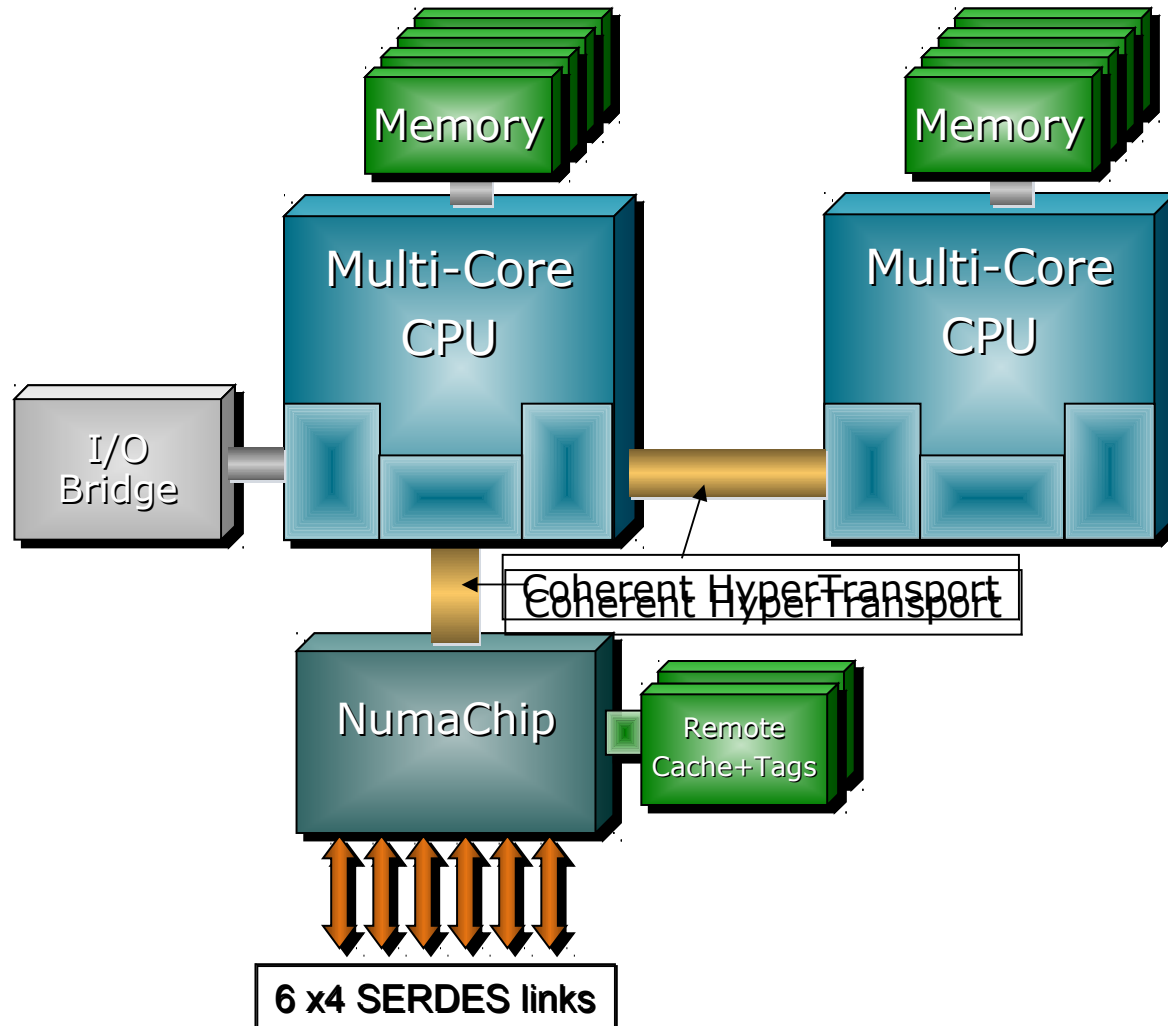Dolphin's Cache Chip

Dolphin's Low Latency Clustering HW

**HPC Clusters (WulfKit w. Scali)**

- First Low Latency Cluster Interconnect
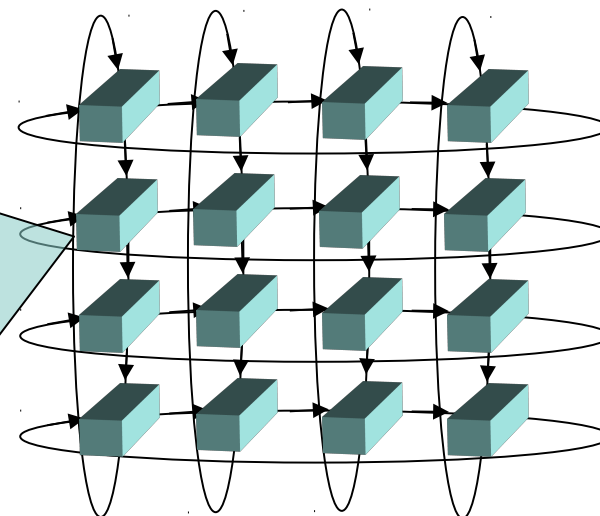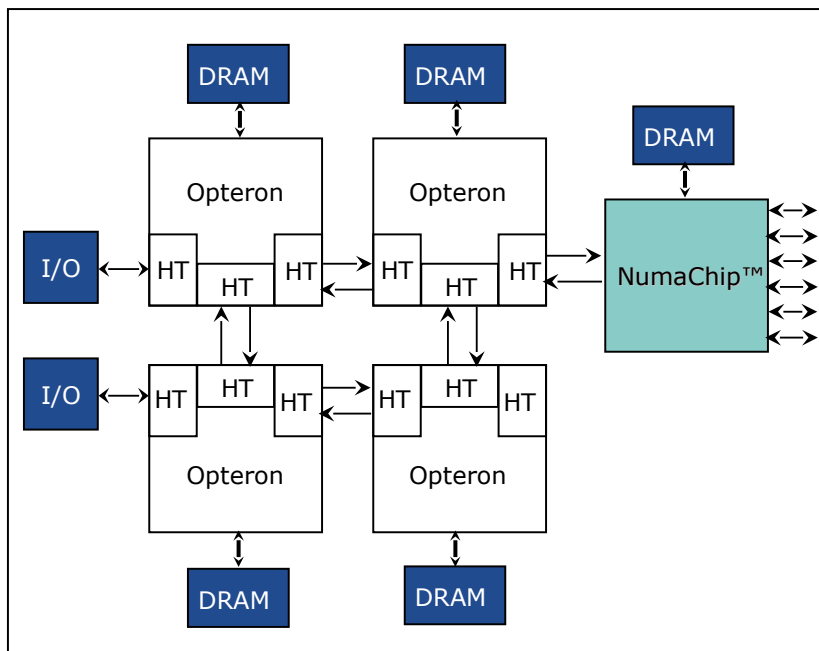
# Company Background

numascale

- Founded in Oslo in 2008 as a spin-out from Dolphin Interconnect Solutions

- Technology from Norsk Data 1987→
  - Dolphin Interconnect 1992

- 24 Experienced Staff Members
  - Interconnects
  - Processor Architecture
  - Supercomputing
  - Data Acquisition

- Main Owners:
  - Investinor (32.1%)
  - ProVenture Seed (20.8%)
  - Statoil ASA (17.5%)
  - Helge B. Risnes (Ex. InfoCare) (7.3%)
  - Svein A. Tunheim (Ex. ChipCon) (7.3%)
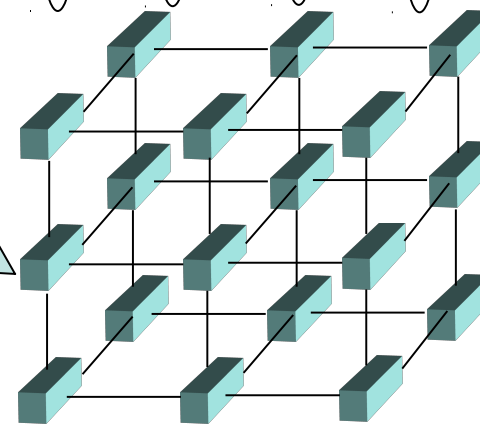
# NumaConnect Main Features

numascale

- 12 bits Node ID = 4,096 nodes, (lots and lost of cores), >19,000 cores

- 48 bits node physical address space = 256TBytes

- Scalable, directory based cache coherency protocol

- Scalable On-Chip switch fabric (2-D, 3-D Torus)

- Configurable Cache for remote data (1 - 16GB/node)

- System-wide cache coherency in hardware

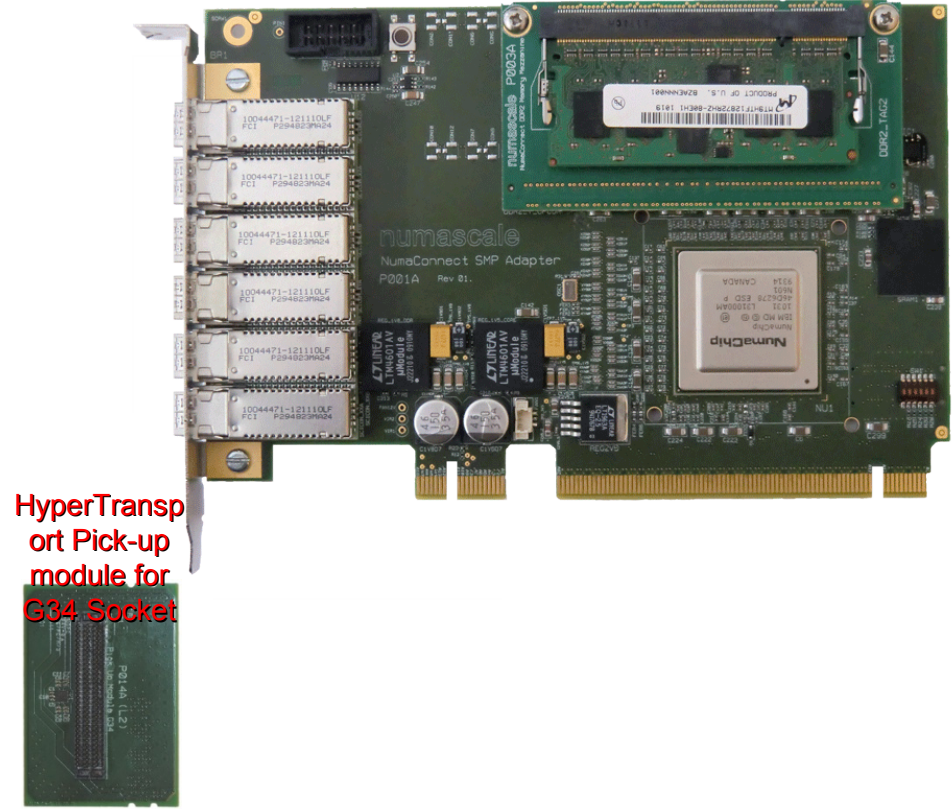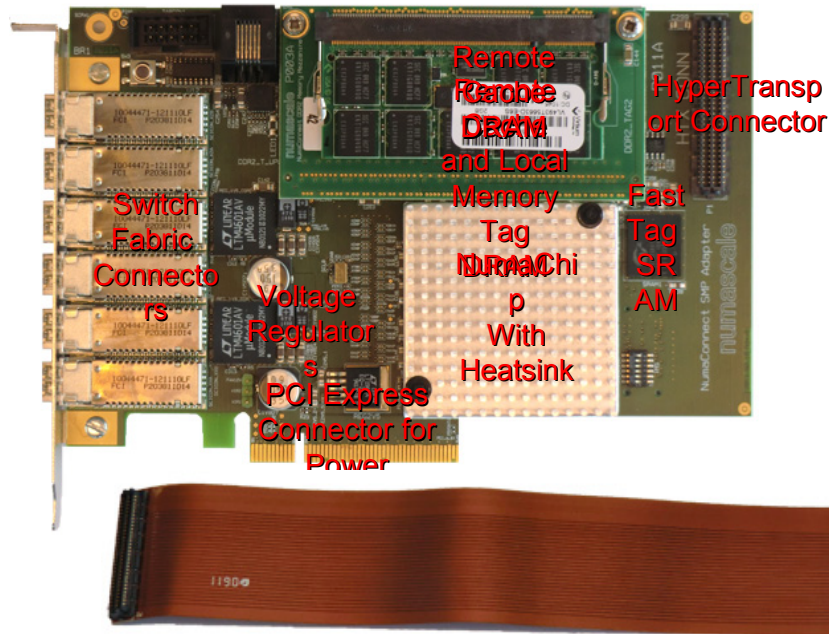- 64Byte cache line granularity same as x86 CPUs

# NumaChip™ Server Configuration

# NumaChip™ System Architecture

**numascale**

**Multi-socket Node**



**2-D Torus**

**3-D Torus**

**6 links allow flexible system configurations in multi-dimensional topologies**

# NumaConnect Cards HTX & PCIe

Remote
Remote
DRAM
and Local
Memory
Tag
NumaChi
p
With
Heatsink

HyperTransp
ort Connector

Fast
Tag
SR
AM

HyperTransp
ort Pick-up
module for
G34 Socket

Switch
Fabric
Connecto
rs

Voltage
Regulator
s
PCI Express
Connector for
Power

# Shorter Time to Performance

**numascale**

A NumaConnect system can be programmed just as an "ordinary" computer!

The full memory range is available to all applications

You can run "top" on a 1.5TB NumaConnect system

```
top - 09:29:57 up 20 days, 20:07,  6 users,  load average: 8.32, 8.31, 8.32
Tasks: 3068 total,   3 running, 3065 sleeping,   0 stopped,   0 zombie
Cpu(s):  1.4%us,  0.1%sy,  0.0%ni, 98.5%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Mem:  1583148160k total, 1185082348k used, 398065812k free,      4k buffers
Swap: 33046524k total,         0k used, 33046524k free,  4055576k cached

  PID USER      PR  NI  VIRT  RES  SHR S %CPU %MEM    TIME+  COMMAND
202867 root     20   0 1119g 1.1t  524 R  815 74.1  3758:34 x.mod2as
197024 root     20   0 19160 6560 1204 R   37  0.0 403:40.55 htop
206397 root     20   0 13512 3704  924 R   14  0.0   0:04.48 top
197023 root     20   0 13512 3700  924 S   13  0.0 177:20.24 top
 1399 root      20   0     0    0    0 S    3  0.0   1:03.79 ksoftirqd/348
 2187 root      20   0     0    0    0 S    3  0.0   0:50.98 ksoftirqd/545
   10 root      20   0     0    0    0 S    2  0.0 290:26.05 rcu_sched
11190 root      20   0     0    0    0 S    1  0.0  68:02.65 kworker/40:1
11177 root      20   0     0    0    0 S    1  0.0 173:47.85 kworker/24:1
11241 root      20   0     0    0    0 S    1  0.0  86:23.99 kworker/44:1
11668 root      20   0     0    0    0 S    1  0.0  53:37.71 kworker/348:1
11688 root      20   0     0    0    0 S    1  0.0  59:31.14 kworker/336:1
12035 root      20   0     0    0    0 S    1  0.0  20:37.15 kworker/545:1
11197 root      20   0     0    0    0 S    0  0.0  96:09.80 kworker/36:1
11203 root      20   0     0    0    0 S    0  0.0 148:15.71 kworker/32:1
11209 root      20   0     0    0    0 S    0  0.0 166:31.05 kworker/28:1
11233 root      20   0     0    0    0 S    0  0.0  82:54.52 kworker/48:1
11626 root      20   0     0    0    0 S    0  0.0  32:13.59 kworker/308:1
11710 root      20   0     0    0    0 S    0  0.0  26:04.65 kworker/357:1
11714 root      20   0     0    0    0 S    0  0.0  31:09.39 kworker/354:1
47717 root      20   0  7764  576  484 S    0  0.0  22:06.24 tail
47736 root      20   0  7764  576  484 S    0  0.0  22:13.57 tail
```

# Large Shared Memory, all threads available

# Linux

**numascale**

NumaConnect Architecture Supported in Linux kernel

  Interprocessor Interrupt (APIC extension HW)

Runs with standard kernel

Tuned kernel recommended

  Especially for large systems >8 servers

  "Custom Kernel" with recommended options

Patches

  Queue-Based Spin Locks (Scalability)

  Optimized Timing Framework for NumaConnect Fabric

# Memory Bandwidth-Stream

**numascale**

This system uses 8 bytes per DOUBLE PRECISION word.

----------------------------------------------------------------

Array size = 180000000000, Offset = 0

Total memory required = 4119873.0 MB.

Each test is run 10 times, but only

the *best* time for each is used.

----------------------------------------------------------------

Number of Threads requested = 828

----------------------------------------------------------------

| Function | Rate (MB/s) | Avg time | Min time | Max time |
|----------|-------------|----------|----------|----------|
| Copy: | 1599317.5028 | 1.9224 | 1.8008 | 2.1393 |
| Scale: | 1468219.1643 | 2.0954 | 1.9616 | 2.2290 |
| Add: | 1664455.1221 | 2.8375 | 2.5954 | 3.0947 |
| Triad: | 1492414.0721 | 3.0478 | 2.8946 | 3.3267 |

University of Oslo:
- 72 nodes - IBM x3755
- 1 728 cores
- 4.6 TBytes Shared Memory

Human readable numbers:
Array size = 180 GB
Copy = 1.6 TB/s
Scale =1.5 TB/s
Add = 1.7 TB/s
Triad = 1.5 TB/s

# LMbench - Chart

Latency - LMbench

RemoteMemory

NumaCache(L4)

Cpu Caches

L3

L2

L1

1087,576
628,583
308,61
304,582
287,01
147,125
16,732
6,421
1,251

Nanoseconds

Array Size (MBytes)

# RWTH TrajSearch

**numascale**

Code performance and scaling results

Dipl.-Inform. Dirk Schmidl

    High Performance Computing

    RWTH Aachen

    Center for Computing and Communication

Numascale Demo System

    8 Supermicro 1042G-LTF+ Servers with NumaConnect

    Each Server

        128GB RAM

        Two 16 cores AMD 6380 processors

        512GB SSD per node

    Total 256 cores sharing 1TB

# RWTH TrajSearch code

## numascale

The graph shows that the application has the **most speedup on NumaConnect**, even if it was originally adapted to ScaleMP.



Legend:
- SGI Altix UV: Nehalem EX
- SCALEMP: Nehalem EX
- Bull BCS: Nehalem EX
- SCALEMP: SandyBridge EP
- NUMASCALE
- SGI Altix UV: Nehalem EX-Speedup
- SCALEMP: Nehalem EX-Speedup
- Bull BCS: Nehalem EX-Speedup
- SCALEMP: SandyBridge EP-Speedup
- Numascale-Speedup

# RWTH TrajSearch code

# RWTH TrajSearch code

# NAS Parallel benchmarks MPI

**numascale**

**If you can get scalable OpenMP and MPI performance, ease of programming and ease of administration at commodity cluster price points, why limit yourself to an MPI cluster?**

## NPB-SP MPI CLASS D
### Time in seconds

- ■─ Numascale (2,5 GHz AMD Opteron 6380)
- ◆─ FDR Infiniband System (Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz)

Runtime [sec] vs Number of Processes

- 2818.76
- 1351.28
- 814.49
- 476.2
- 433.06

X-axis: 16, 36, 64, 121

## NPB-LU MPI CLASS D
### Time in seconds

- ■─ Numascale (2,5 GHz AMD Opteron 6380)
- ◆─ FDR Infiniband System (Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz)

Runtime [sec] vs Number of Processes

- 3049.51
- 2325.17
- 1607.2
- 1036.75
- 425.7
- 381.25
- 241.9
- 215.14

X-axis: 16, 32, 64, 128

## NPB-BT MPI CLASS=D
### Time in seconds

- ■─ Numascale (2,5 GHz AMD Opteron 6380)
- ✳─ FDR Infiniband System (Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz)

Runtime [sec] vs Number of processes

- 2289.97
- 1922.7
- 1049.93
- 836.79
- 627.6
- 561.46
- 382.63
- 350.87

X-axis: 16, 36, 64, 121

The **overhead** introduced by MPI is not needed when we are running on a **Shared Memory System**

➤ The NAS Parallel Benchmarks (NPB)

- ➢ evaluate the performance of parallel supercomputers
- ➢ derived from computational fluid dynamics (CFD) applications
- ➢ LU is a simulated uses symmetric successive over-relaxation (SSOR) method to solve a seven-block-diagonal system resulting from finite-difference discretization of the Navier-Stokes

## NPB-NC-OMP LU E: Time in Seconds
## AMD Opteron(tm) Processor 6174
## 72 NumaConnect Nodes

Legend: NPB-MPI LU CLASS=E Time in Seconds | NPB-NC-OMP LU E: Time in Seconds

Y-axis: Runtime [sec] — 16000, 8000, 4000, 2000, 1000

X-axis: Number of processes — 128, 256, 512, 1024

Data points (NPB-MPI LU CLASS=E): 14509.94, 7742.66, 3165.47, 2381.32

Data points (NPB-NC-OMP LU E): 11169.93, 5318.42, 2691.01, 1628.4

# CD-adapco STAR-CCM+ (MPI)

**numascale**

- STAR-CCM+ is a technology leading Computational Fluid Dynamics (CFD) package unrivalled in its ability to tackle problems involving multi-physics and complex geometries,
  http://www.cd-adapco.com/products/star-ccm®

- The NumaConnect Shared Memory test system used to conduct the tests has:
- 1TB of memory
- 256 cores
- It utilizes 8 servers each equipped with:
  - 2 x AMD Opteron 2,5 GHz 6380 CPUs
  - 16 cores in each CPU
    -

## STAR-CCM+
## "Time per Iteration [seconds]"
## Lower is better



Time per Iteration [seconds] vs Number of processes:
- 32: 141.77
- 64: 75.15
- 128: 43.56

# GROMACS

> **GROMACS** is a versatile package to perform molecular dynamics, i.e. simulate the Newtonian equations of motion for systems with hundreds to millions of particles. It is primarily designed for **biochemical molecules** like proteins, lipids and nucleic acids that have a lot of complicated bonded interactions, but since GROMACS is extremely fast at calculating the nonbonded interactions (that usually dominate simulations) many groups are also using it for research on non-biological systems, e.g. polymers.

> The NumaConnect Shared Memory test system used to conduct the tests has:
> 1TB of memory
> 256 cores
> It utilizes 8 servers each equipped with:
> > 2 x AMD Opteron 2,5 GHz 6380 CPUs
> > 16 cores in each CPU
> > 128GB

## GROMACS with NC-OpenMPI [ns/day] (higher is better) case: Test-performance_protein-water-membrane.tpr

# Shared GPU, GPGPU

3 AIC Octans

3 NVIDIA GeForce GT 640 2GB

Numascale Shared Memory System

- Cache Coherent Global Shared Memory and Shared IO

- All GPUs providing aggregated TFLOPS

- Running N-body CUDA application

# Numa-Q

## In-memory analytics appliance



4 x Dell R815s

Running 1Billion row, Spark regression benchmark, 4X gain over cluster

**Advanced Analytics & Visualization**
**Simplified Management With NumaManager**
**Terabytes of Memory**

**Thousands of Cores**

**Single Linux Instance**

# Spark Benchmark

Apache Spark™ Benchmark

1B rows, 10 variables, Logistic Regression

4 node distributed cluster vs 4 node NumaQ

|  | 4 nodes Cluster 256GB RAM 32 cores each | NumaQ 1TB RAM 128 cores |
|---|---|---|
| Logistic Regression 1B rows 10 variables | 108 sec | 27 sec |

# Competition - differentiation

**numascale**

| | Perfor-mance | Shared Memory | Price | Comments |
|---|---|---|---|---|
| Software solutions with InfiniBand or 10Gbe (ScaleMP)<br>**ScaleMP** | ? | ✓ | 2X | Software emulation Non-standard Operating System - Virtualization Layer |
| Mainframes (SGI, HP, IBM, Oracle (Sun))<br>**sgi** | High | ✓ | 10-30X | • "Max" performance – shared memory<br>• 50TB limit for SGI<br>• Limited Scalability |
| High-end interconnect for clusters (InfiniBand)<br>**Mellanox TECHNOLOGIES** | High | | 1X | • Pure message passing only |
| YarcData - uRiKA<br>**CRAY** THE SUPERCOMPUTER COMPANY | High | ✓ | 10X? | • Complete system solution – Big Data Appliance, Proprietary architecture |
| Numascale<br>**numascale** | High | ✓ | 1X+ | • Independent hardware vendor<br>• Commodity server hardware |

# Shared Memory returning - Why?

**numascale**

Compelling programming model
- Less code
- Large memories - less effort, no data domain decomposition
- Flexibility

System Utilization
- More efficient utilization of Resources, up to 90% Hitachi mainframe Cambridge University, versus 50% cluster University of Oslo.
- Reduced sysadmin', single OS
- Data Center Fabric

NumaConnect
- Turns COTs servers into SSI - ccNUMA , CHEAPLY!

With ASIC transistor density, and HT, it can now be done cheaply!

# What's best for scientific research?

**numascale**

 Cluster
- Look at the size of my Linpack!
- # nodes/cores
- Grand Challenge benchmarks
- Interesting Computer Science
- Poor ROI

 Shared Memory
- High system utilization
- Easier sysadmin'
- Easier programming
- More research papers produced
- Better science

 Intel

– XEON based solution

 AMD

– Higher Link Speeds
– Larger node memory
– Higher core count
– Hybrid CPU/GPU