BLUE WATERS SUSTAINED PETASCALE COMPUTING

Power 7 for Scientific Computing

Public Information











ONE CORE





Power 7 Core

- Execution Units
 - 2 Fixed point units
 - 2 Load store units
 - 4 Double precision floating point
 - 1 Branch
 - 1 Condition register
 - 1 Vector unit
 - 1 Decimal floating point unit
 - 6 wide dispatch
- Recovery Function Distributed
- 1,2,4 Way SMT Support
- Out of Order Execution
- 32KB I-Cache
- 32KB D-Cache
- 256KB L2
 - Tightly coupled to core



I

TEM





8 cores

ONE CHIP





Power 7 Chip



IEM

GREAT LAKES CONSORTIUM

I

NESA





GREAT LAKES CONSORTIUN

TRM



RAM Technologies

Conventional	IBM ASIC	IBM Custom	Custom	Custom
Memory DRAM	eDRAM	eDRAM	Dense SRAM	Fast SRAM
Dense, low por Low speed/band	wer width	Off On uP uP Chip Chip	Hi	High Area/power gh speed/bandwidth
				and an interview of the states
Conventional	Large, Off-chip	On-processor	On-processor	Private core
Memory DIMMs	30+ MB Cache	30+ MB Cache	Multi-MB Cache	Sub-MB Cache

- DIMMs, Dense and Fast SRAM are industry standard; eDRAM is IBM technology
 - Used in Power 4,5,6 for off-chip L3 cache
 - Used in Power 7 for on-chip cache (to avoid pin limitations and support BW of 8 cores)



L3 Cache



NESA IBM

GREAT LAKES CONSORTIUM



L3 Cache Requirements

Cache Hierarchy Rqmt for POWER Servers

Challenge for Multi-core POWER7

I

IRM

GREAT LAKES CONSORTIUM





Solution: "Hybrid" L3 – Fluid Cache Structure

TRM

GREAT LAKES CONSORTIUN





Solution: "Hybrid" L3 – Fluid Cache Structure

TRM

GREAT LAKES CONSORTIU



- Keeps multiple footprints at ~3X lower latency than local memory.
- Automatically migrates private footprints (up to 4M) to fast local region (per core) at ~5X lower latency than full L3 cache.
- Automatically clones shared data to multiple private regions.



Solution: "Hybrid" L3 – Fluid Cache Structure

I

TRM

GREAT LAKES CONSORTIUM



- Enables a subset of the cores to utilize the entire large shared L3 cache when the remaining cores are not using it.





Local L3 Region



IBM

GREAT LAKES CONSORTIUN

NESA



L2 Still Needed



I

TEM

GREAT LAKES CONSORTIUM



L2 Turbo Cache



I

TEM

GREAT LAKES CONSORTIUM

- L2 "Turbo" cache keeps a tight 256K working set with extremely low latency (~3X lower than local L3 region) and high bandwidth, reducing L3 power and boosting performance.



L2 Cache



IBM

GREAT LAKES CONSORTIUM

I

NESA



Cache Hierarchy Summary



I

NESA IBM

GREAT LAKES CONSORTIUM

Cache Level	Capacity	Array	Policy	Comment	
L1 Data	32K	Fast SRAM	Store-thru	Local thread storage update	
Private L2	256K	Fast SRAM	Store-In	De-coupled global storage update	
Fast L3 Region	Up to 4M	eDRAM	Partial Victim	Reduced power footprint (up to 4M)	
Shared L3	32M	eDRAM	Adaptive	aptive Large 32M shared footprint	



Memory Goals





I

NESA

IBM

GREAT LAKES CONSORTIUM



Memory Solutions





On-Chip Memory Controllers



IBM

GREAT LAKES CONSORTIUM



Memory



NESA IEM

GREAT LAKES CONSORTIUM



Memory Technology

- L1 32KB Instruction, 32KB Data / core
- L2 = 256KB / core
- L3 = 4MB eDRAM / core
- Memory Capacity: Up to 128 GB / p7
- 8 Channels of SuperNova buffered DIMMs / p7
 - 2 memory controllers / p7
 - 4 memory busses per memory controller

Each of the 4 busses are: 1B wide Write, 2B wide Read)



T

Dual SuperNova DIMM



Off-Chip Connectivity

Interface	Signal Type	Info Width	Frequency	Bandwidth
Off-chip Cache	none	none	none	none
Memory Channels	Differential	28 bytes	6.4 Ghz	180 GB/s
I/O Bridge	Single-ended	20 bytes	2.5 Ghz	50 GB/s
SMP Interconnect	Single-ended	120 bytes	3.0 Ghz	360 GB/s
Total Bandwidth				590 GB/s

(Note that bandwidths shown are raw, peak signal bandwidths)

I

GBEAT LAKES CUNS

- Better signaling technology
- On-chip L3 reduces BW requirements
- Advanced scheduling improves BW utilization





Shared Memory and I/O



NESA IEM

GREAT LAKES CONSORTIUM



High-End Server Resilience for Good MTBF



GREAT LAKES CONSORTIUM



Feeds and Speeds per QCM

- 32 cores
- 8 Flop/cycle per core
- 4 threads per core max
- 3.5 4 GHz
- 1 TF/s
- 32 MB L3
- 512 GB/s memory BW (0.5 Byte/flop)
- 800 W (0.8 W/flop)





4 chips, 32 cores

ONE MULTI-CHIP MODULE (aka QCM, Node, Octant)





Compute Intensive Quad-chip MCM

<u>QCM</u>

•4 x p7

•32 Cores (4 x 8 =32)

• Up to 512 GB memory





IBM

GREAT LAKES CONSORTIUM

NESA



Node (Octant) - Logical

•P7-1H Quad Block Diagram

•32 w SMP / 2 Tier SMP Fabric

•4 chip Processor MCM

•Hub SCM

On-Board Optics



I

TRM

GREAT LAKES CONSORTIUM



Hub Chip Module (Torrent)

- ✓ <u>Connects QCM to PCI-e</u> (two 16x and one 8x PCI-e slot)
- ✓ <u>Connects 8 QCM's Together</u> via low latency, high bandwidth, copper fabric.
 - Enables a single hypervisor to run across 8 QCM's
 - Allows I/O slots attached to the 8 hubs to be directed to the compute power of any of the 8 QCM's
 - Provides a message passing mechanism with very high bandwidth
 - Provides the lowest possible latency between 8 QCM's (7.6TF) of compute power
- <u>Connects four P7-IH planers Together</u> via the L Remote Optical connections (Super Node)
- Connects up to 512 Super Nodes Together via the D Optical Buses





IBM

GREAT LAKES CONSORTIUM

I

NESA



1.1 TB/s HUB

- 192 GB/s Host Connection
- 336 GB/s to 7 other local nodes
- 240 GB/s to local-remote nodes
- 320 GB/s to remote nodes
- 40 GB/s to general purpose I/O





8 MCMs, 32 chips, 256 cores

ONE DRAWER





First Level Interconnect

≻L-Local

≻HUB to HUB Copper Wiring

≻256 Cores



I

TRM

GREAT LAKES CONSORTIUM





TEM

GREAT LAKES CONSORTIUM

I

NESA



Every HUB Connected to Every Other HUB in CEC

TEM

GREAT LAKES CONSORTIUM

T



(CEC = Central Electronic Complex: IBM Jargon for drawer or packaging unit)





4 drawers, 32 MCMs, 128 chips, 1024 cores

ONE SUPERNODE







T

- **Second Level Interconnect**
- Optical 'L-Remote' Links from HUB
- Construct Super Node (4 CECs)
- 1,024 Cores
- Super Node



GREAT LAKES CONSORTIUM











Disk Storage

- Disk Enclosure is Connected to CEC via SAS PCIe Cards
- <u>Redundant DCA / DE</u>
- Up to 384 SFF DASD drives / DE
- 8 Storage Groups (STOR 1-8) with 48 drives each
 - ✓ Two Port Cards per STOR
 - ✓ 12 carriers per STOR
 - Carriers contains up to 4 drives each
- Data Rates:
- ✓ Serial Attach SCSI (SAS) SDR = 3.0 Gbps per lane (SEC to Drive)
- ✓ Serial Attach SCSI (SAS) DDR = 6.0 Gbps per lane (SAS Adapter in Node to SEC)









1-12 2U CECs (256 – 3073 cores) 0-6 4U drawers (0-2304 disks)

ONE RACK



<u>Rack</u> •990.6w x 1828.8d x 2108.2 •39"w x 72"d x 83"h •~2948kg (~6500lbs)

Data Center In a Rack

Compute Storage Switch 100% Cooling PDU Eliminated

Input: 8 Water Lines, 4 Power Cords Out: ~100TFLOPs / 24.6TB / 153.5TB 192 PCI-e 16x / 12 PCI-e 8x



BPA

200 to 480Vac
370 to 575Vdc
Redundant Power
Direct Site Power Feed
PDU Elimination

Storage Unit

4U
0-6 / Rack
Up To 384 SFF DASD / Unit
File System

<u>CECs</u>

2U 1-12 CECs/Rack 256 Cores 128 SN DIMM Slots / CEC **•**8,16, (32) GB DIMMs 17 PCI-e Slots Imbedded Switch Redundant DCA **NW Fabric** •Up to:3072 cores, 24.6TB (49.2TB) WCU Facility Water Input 100% Heat to Water Redundant Cooling CRAH Eliminated





IEM

GREAT LAKES CONSORTIUM

I

NESA

ONE SYSTEM: BLUE WATERS



BLUE WATERS SUSTAINED PETASCALE COMPUTING

Blue Waters

- Approximately 10 PF/s peak
- More than 300,000 cores
- More than 1 PetaByte memory
- More than 10 Petabyte disk storage
- More than 0.5 Exabyte archival storage
- More than 1 PF/s sustained on scientific applications





National Petascale Computing Facility

