# Application Performance on Dual Processor Cluster Nodes

by

Kent Milfeld

milfeld@tacc.utexas.edu

**Avijit Purkayastha, Kent Milfeld, Chona Guiang, Jay Boisseau**
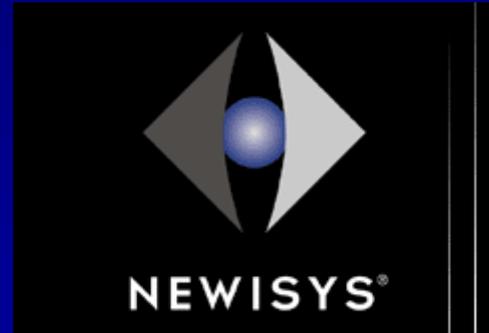
TACC

TEXAS ADVANCED COMPUTING CENTER

The University of Texas At Austin

# Thanks

- **Newisys (Austin, TX)**
  AMD Opteron System

- **Dell (Austin, TX) & Cray**
  Intel Xeon System

# OUTLINE

- HPC needs for Single- & Dual-processor Commodity system Nodes

- The architecture of Intel Xeon & AMD Opteron Systems

- Single & Dual Processor Xeon & Opteron Performance Comparison

  - Measured Memory Characteristics

  - Parallel vs Serial Execution of Codes on a Node

    - Kernels
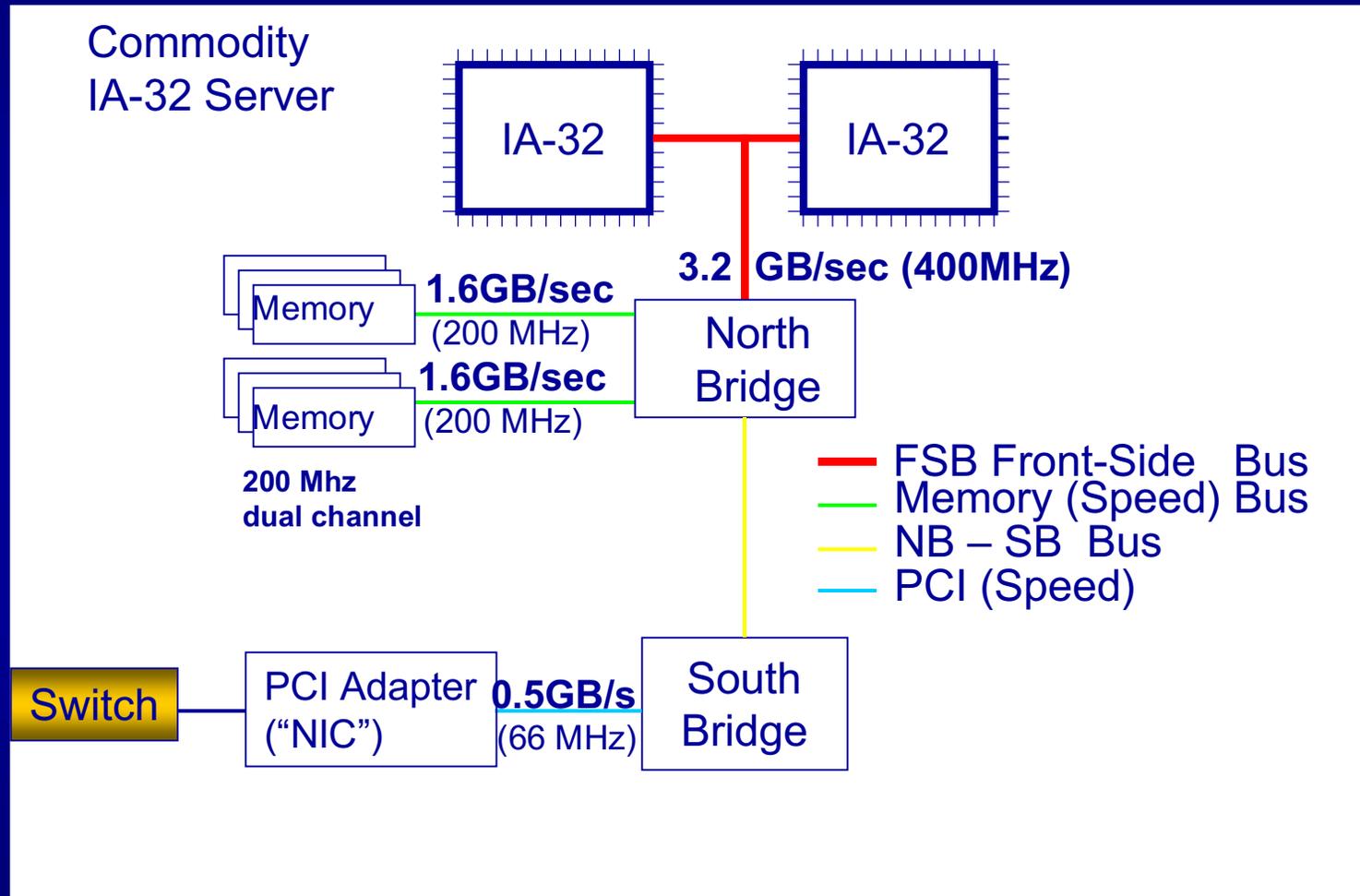
    - Applications

**TACC**

3

# Motivation

- 1995-2000 Commodity Massively Parallel Systems used uni-processor nodes:
  - Beowulf Systems
  - SP2(SC)
  - T3E
- Today the e-commerce market has driven the price of SMP servers down.
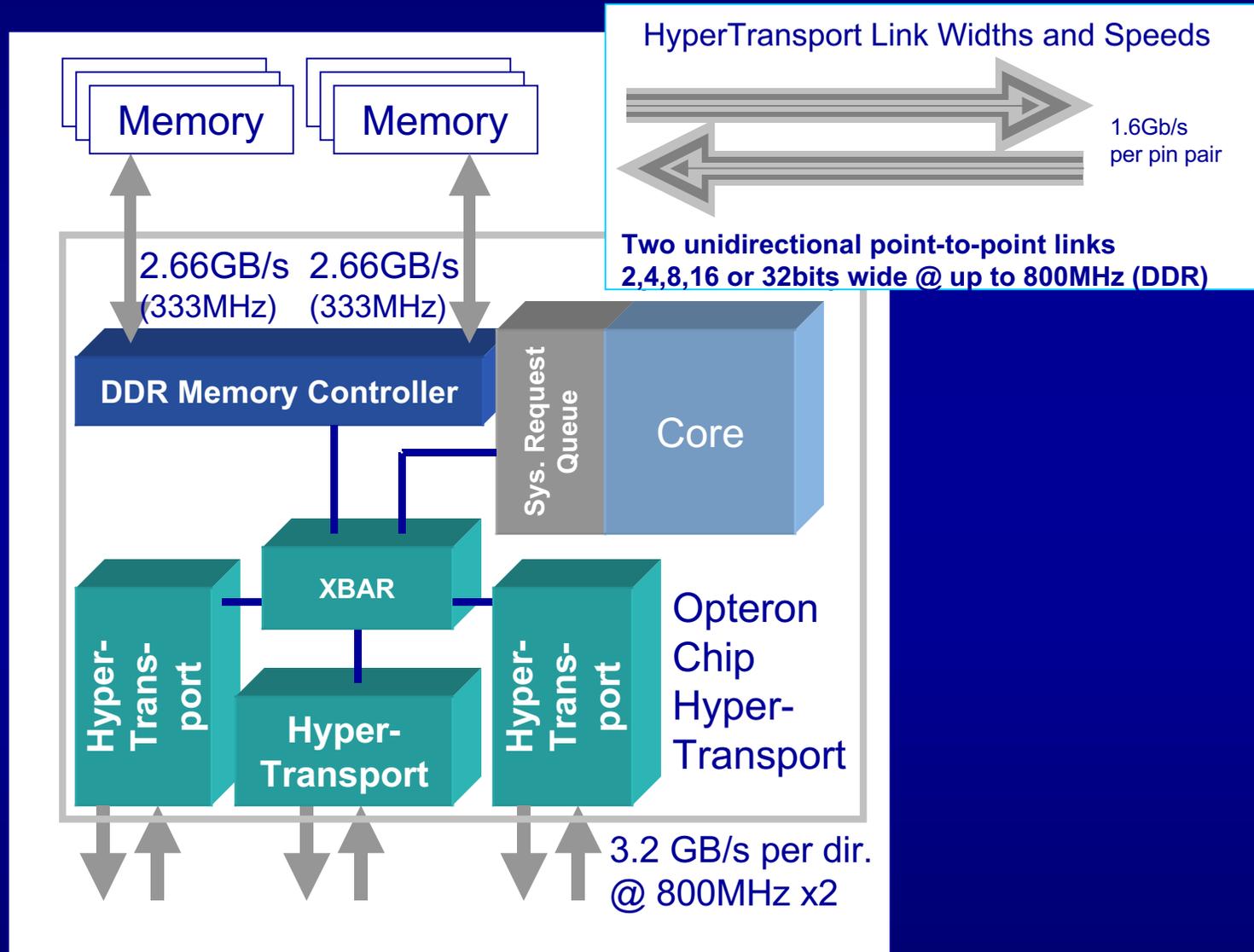  - Dell, Gateway, HP/Compaq, … compete for this market.

TACC

4

The University of Texas At Austin

# Motivation

## Dual processor scoreboard for HPC Applications:

Single    Dual

| Single | Dual | |
|:---:|:---:|:---|
| ☐ | ☒ | Peak performance (TFLOP) |
| ☐ | ☒ | Cost Per Processor |
| | | Memory Subsystem |
| ☒ | ☐ |     No shared bus system |
| ☒ | ☐ |     No Coherence in Caches (processor and "northbridge" & OS) |
| ☒ | ☐ |     No False Sharing |
| ☒ | ☐ | Memory Size |
| | | Message Passing |
| ☒ | ☐ |     No Shared interconnect adapters |
| ☐ | ☒ |     On-node MPI performance |
| ☐ | ☐ | I/O Performance |
| ☒ | ☐ |     Local |
| ☒ | ☐ |     Parallel |

TACC

The University of Texas
At Austin

# Intel Architecture

Commodity
IA-32 Server

IA-32 —— IA-32

**3.2 GB/sec (400MHz)**

Memory —**1.6GB/sec** (200 MHz)— North Bridge

Memory —**1.6GB/sec** (200 MHz)—

**200 Mhz dual channel**

— FSB Front-Side Bus
— Memory (Speed) Bus
— NB – SB Bus
— PCI (Speed)

Switch — PCI Adapter ("NIC") **0.5GB/s** (66 MHz) — South Bridge

# Intel Architecture

Memory Memory

## HyperTransport Link Widths and Speeds

1.6Gb/s per pin pair

**Two unidirectional point-to-point links
2,4,8,16 or 32bits wide @ up to 800MHz (DDR)**

2.66GB/s (333MHz)  2.66GB/s (333MHz)

**DDR Memory Controller**

Sys. Request Queue

Core

XBAR

Hyper-Trans-port

Hyper-Transport

Hyper-Trans-port

Opteron Chip Hyper-Transport

3.2 GB/s per dir. @ 800MHz x2

The University of Texas At Austin

# AMD Architecture

AMD Opteron

AMD Opteron

**6.4GB/s**
**Coherent**
**HT**

**6.4GB/s**
**HT**

2.1/2.7
GB/sec

AMD-8151
HT
AGP Tunnel

**6.4GB/s**
**HT**

Dual Channel
266/333 MHz
(PC2100/2700)

AMD-8131
HT
PCI-X Tunnel
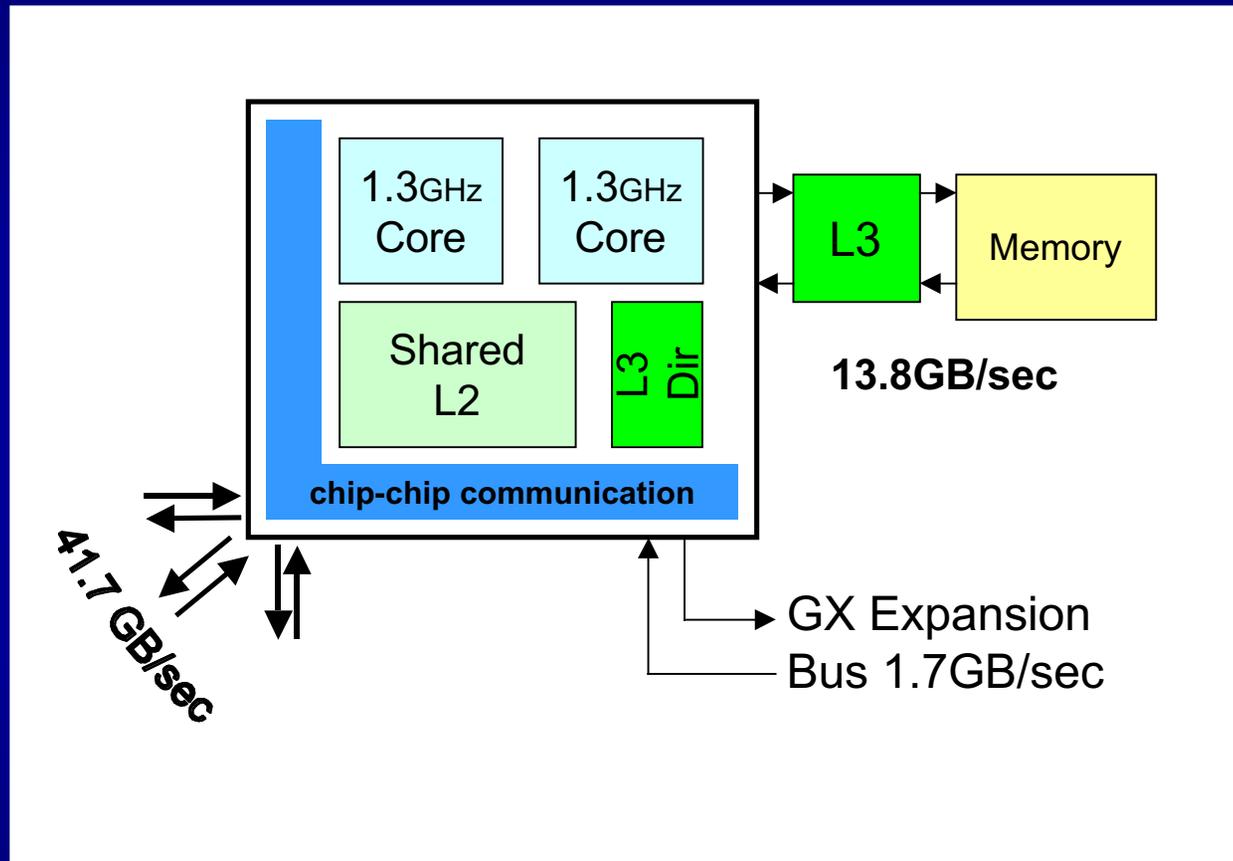
# IBM Power4

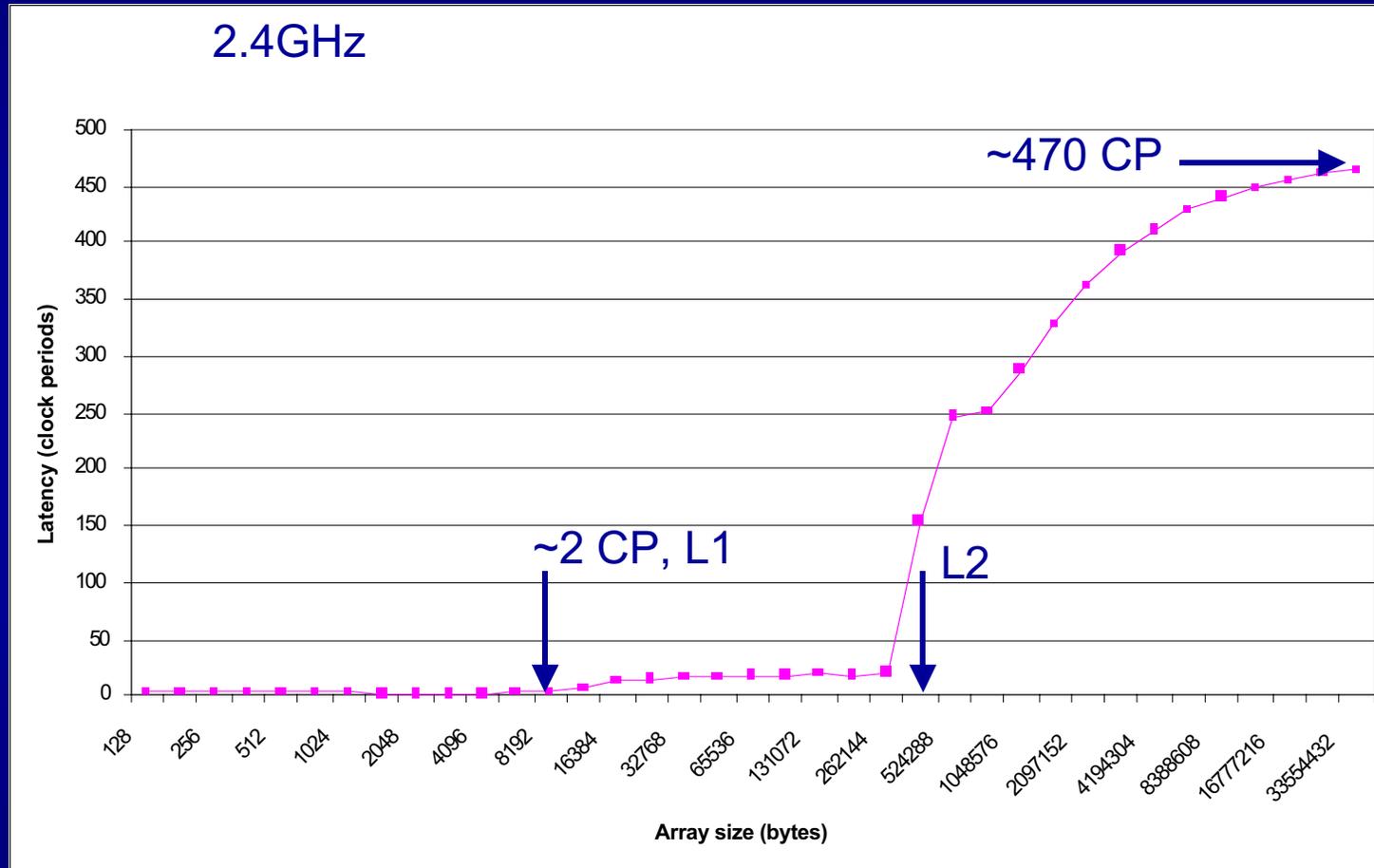# Memory Latency
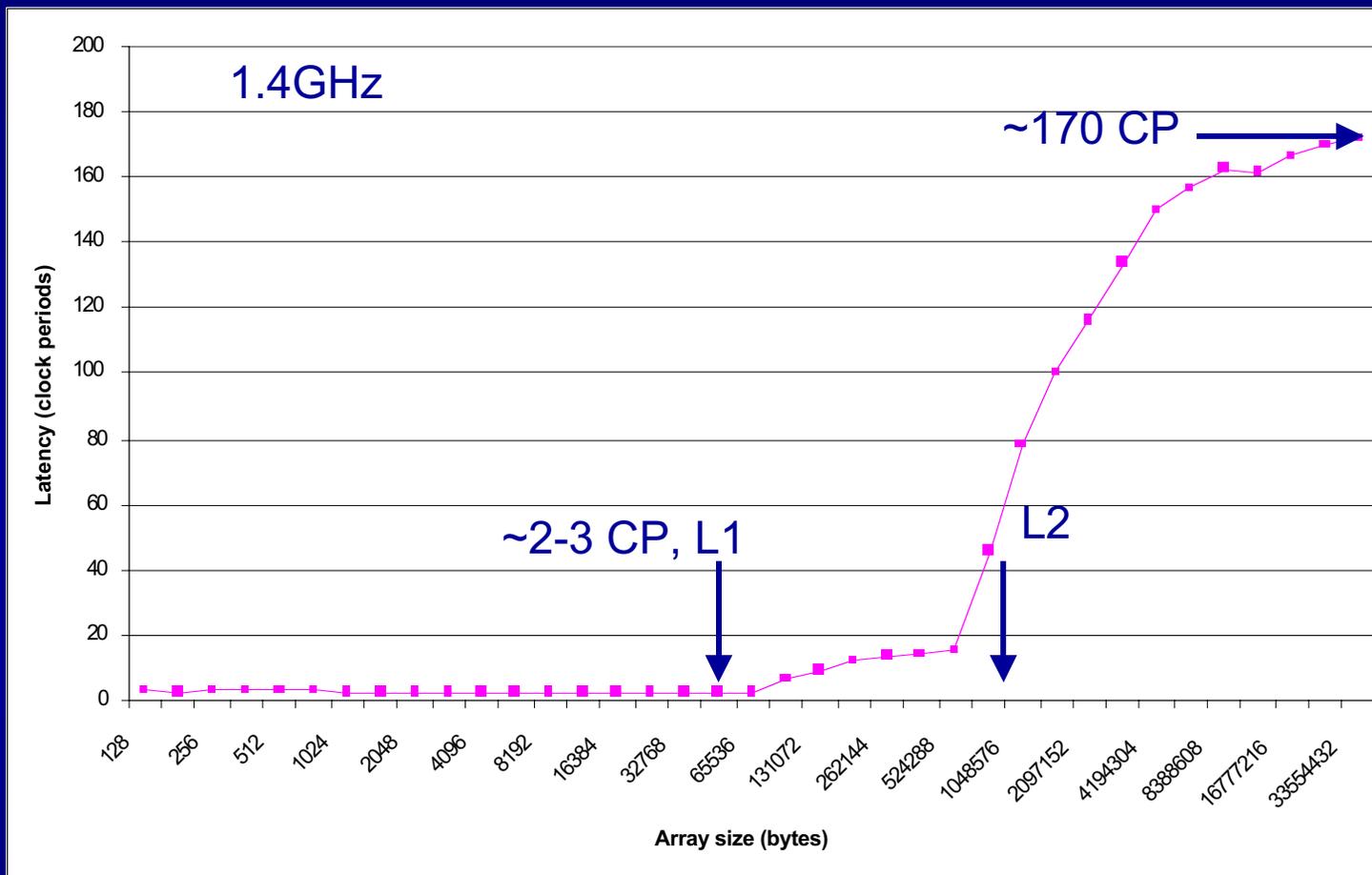
```
I1  =  IA(1)
DO  I  =  2,N
   I2  =  IA(I1)
   I1  =  I2
END  DO
```

1.) Load IA with sequence 1→N.
2.) Randomize IA entries.
3.) Measure Clock Periods of loop. (CPs/N = single memory access time = latency)
4.) Loop does not optimizes: no prefetching or streams

# Memory Latency  Xeon



2.4GHz

Latency (clock periods)

~470 CP

~2 CP, L1

L2

Array size (bytes)

128  256  512  1024  2048  4096  8192  16384  32768  65536  131072  262144  524288  1048576  2097152  4194304  8388608  16777216  33554432

TACC

11

The University of Texas
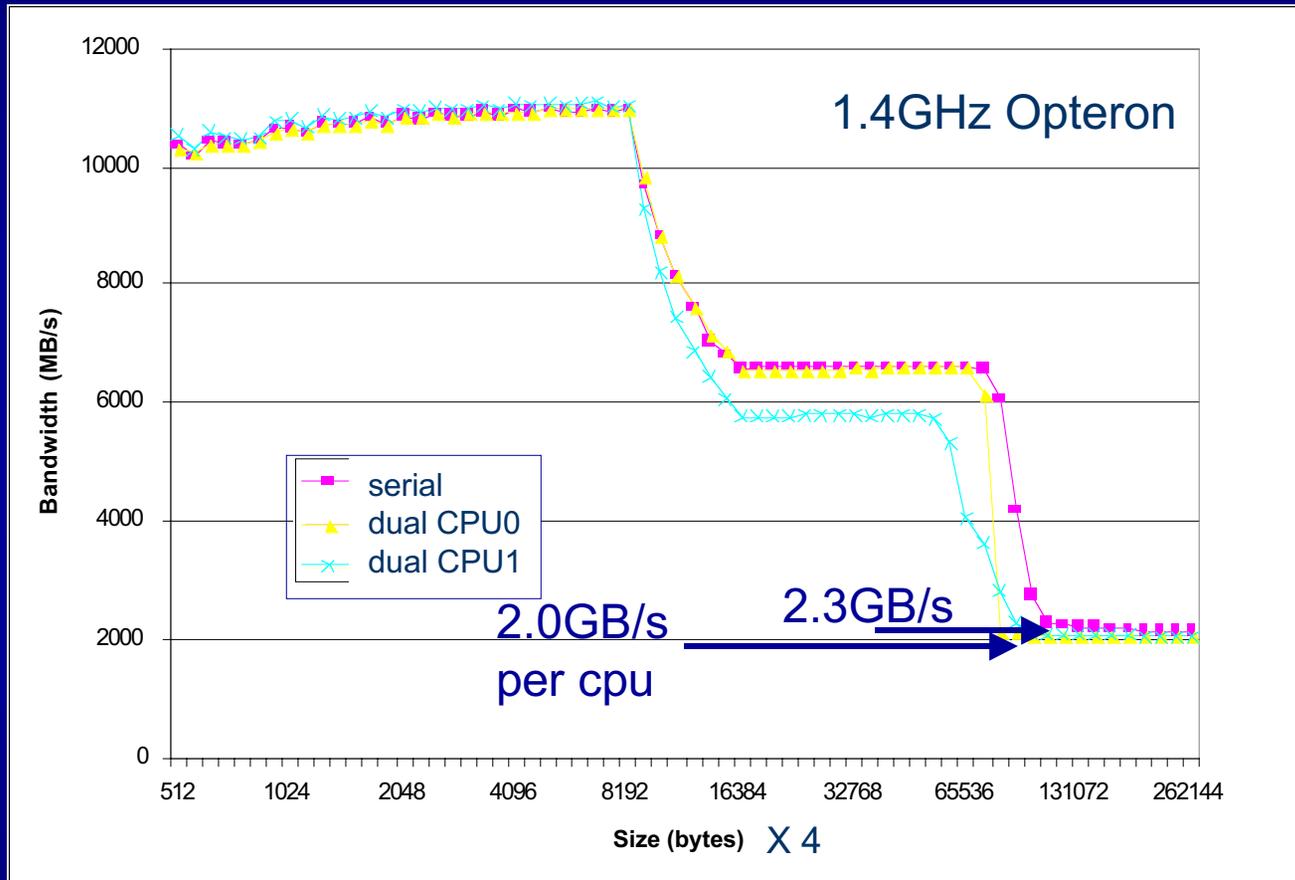At Austin
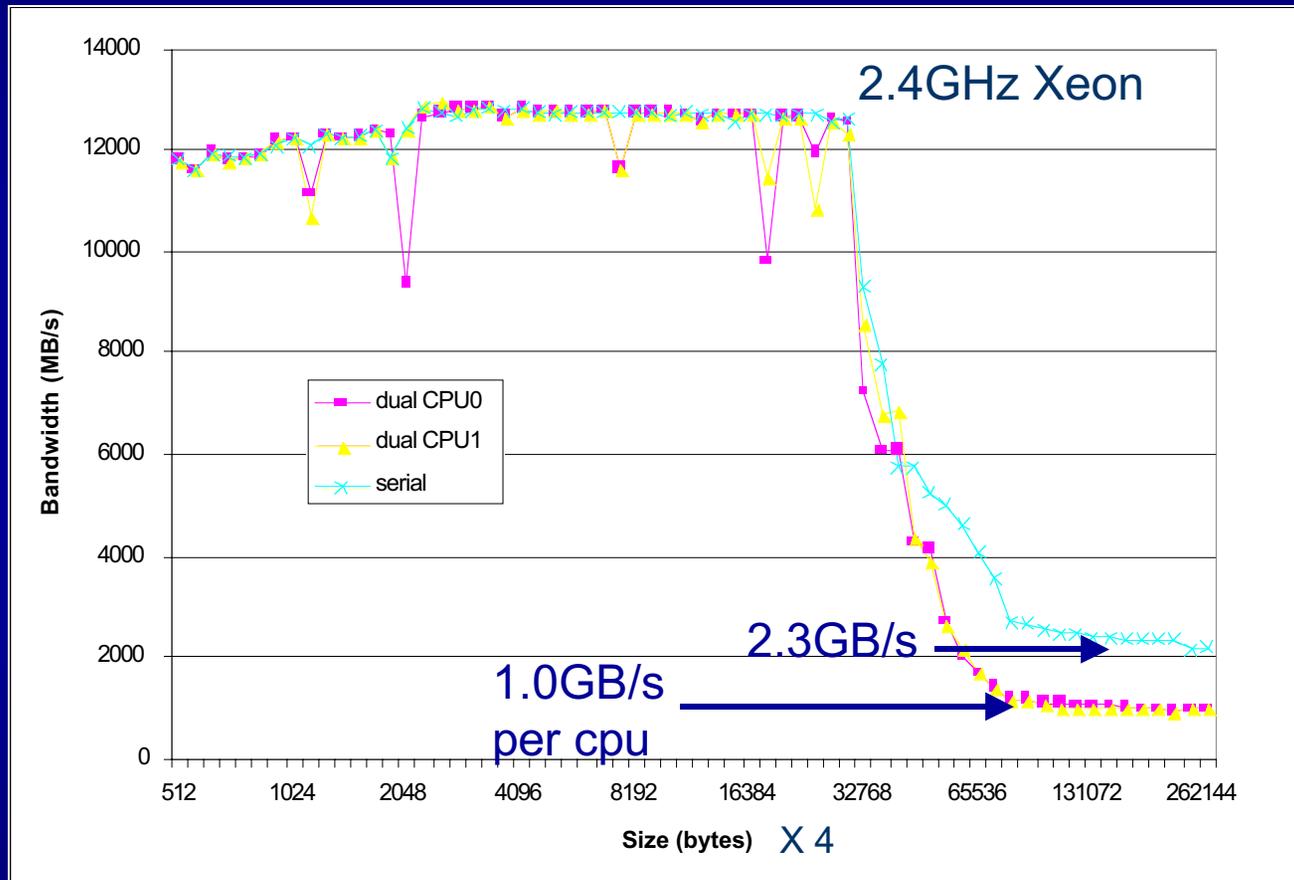
# Memory Latency AMD

# Memory Bandwidth

```
DO I = 1,N
   S = S + A(I)
   T = T + B(I)
END DO
```

1.) -O3, unrolling = 2
2.) Two streams—
    gives high, reasonable
    bandwidths expected
    across memory & caches

# AMD SP/DP memory bandwidth

# Xeon SP/DP memory bandwidth

# STREAM Results

| Kernel | Intel Xeon | AMD Opteron |
|--------|-----------|-------------|
| Copy   | 1213      | 2162        |
| Scale  | 1206      | 2093        |
| Add    | 1381      | 2341        |
| Triad  | 1375      | 2411        |

**Serial** Execution, (MB/sec).

| Kernel | Intel Xeon | AMD Opteron |
|--------|-----------|-------------|
| Copy   | 1167      | 3934        |
| Scale  | 1162      | 4087        |
| Add    | 1273      | 4561        |
| Triad  | 1281      | 4529        |

**Parallel** Execution, two threads (MB/sec).

TACC

The University of Texas
At Austin

# MPI On-Node Bandwidth

**It should be faster than node-to-node.**

(MB/sec)

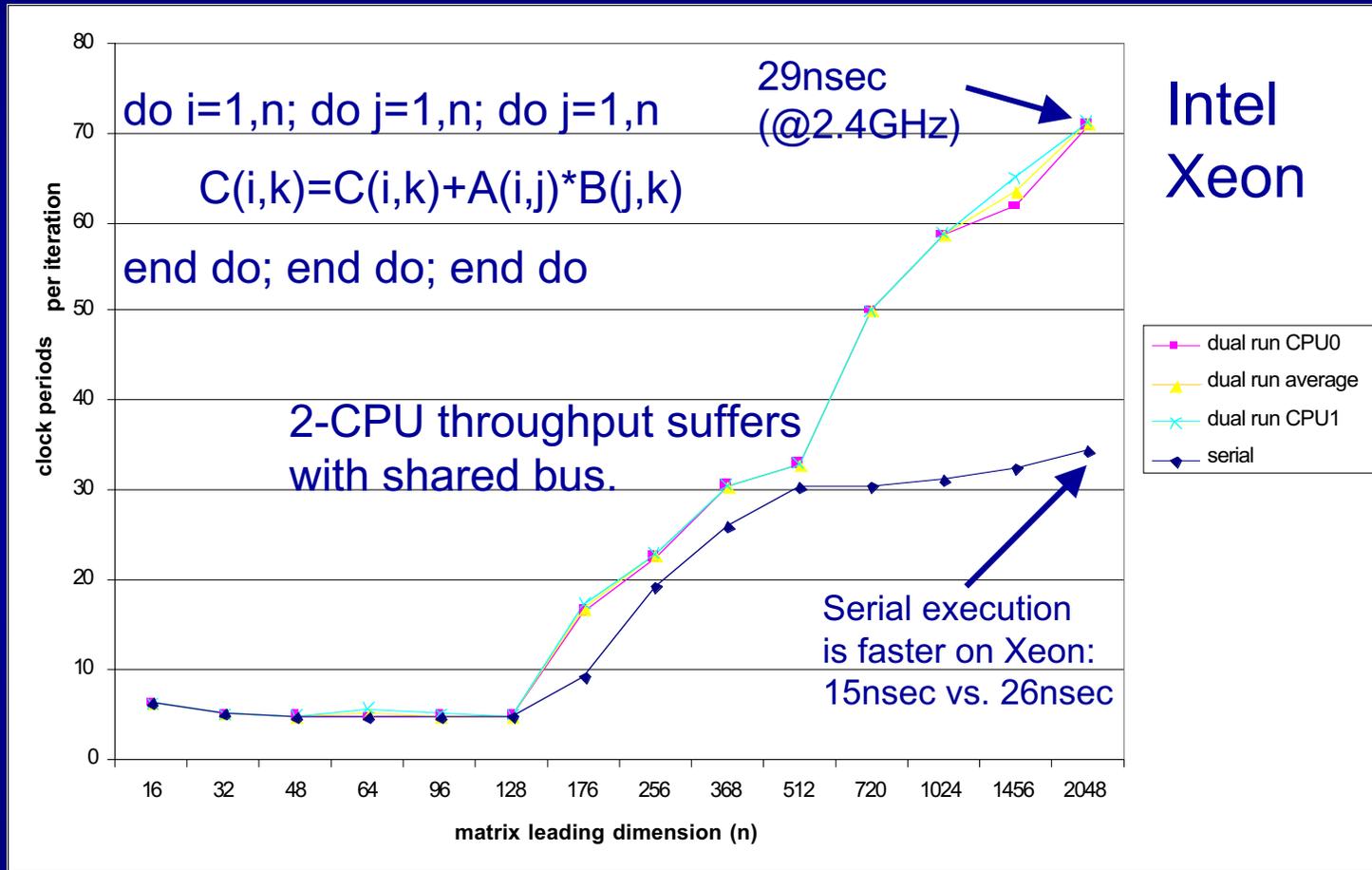| | |
|---|---|
| DELL 2650 | 295 @ 2MB |
| Opteron Suse-64 ch_p4 | 172 @ 2MB |
| Opteron Suse-64 ch_shmem | 404 @ 2MB |
| IBM P690 HPC | 1324 @ 2MB |
| IBM P690 Turbo | 1398 @ 2MB |
| IBM P655 HPC | 1684 @ 2MB |

Different implementations
of MPI will vary with On-Node
Performance.

TACC

17

# Hand Coded Matrix-Matrix Multiply

Accesses Memory with 1 stream and 1 strided pattern.
(Don't do this at home in your optimized code ☺.)

do i=1,n; do j=1,n; do j=1,n

   C(i,k)=C(i,k)+A(i,j)*B(j,k)

end do; end do; end do

26nsec
(@1.4GHz)

AMD
Opteron

2x throughput when
run on two CPUs.

- ■— dual run CPU0
- ▲— dual run average
- ✕— dual run CPU1
- ✳— serial

clock periods per iteration (y-axis: 0, 5, 10, 15, 20, 25, 30, 35, 40)

matrix leading dimension (n): 16, 32, 48, 64, 96, 128, 176, 256, 368, 512, 720, 1024, 1456, 2048

# Hand Coded Matrix-Matrix Multiply

do i=1,n; do j=1,n; do j=1,n

   C(i,k)=C(i,k)+A(i,j)*B(j,k)

end do; end do; end do

29nsec
(@2.4GHz)

Intel
Xeon

2-CPU throughput suffers
with shared bus.

Serial execution
is faster on Xeon:
15nsec vs. 26nsec

Legend:
- dual run CPU0
- dual run average
- dual run CPU1
- serial

Y-axis: clock periods    per iteration (0, 10, 20, 30, 40, 50, 60, 70, 80)

X-axis: matrix leading dimension (n): 16, 32, 48, 64, 96, 128, 176, 256, 368, 512, 720, 1024, 1456, 2048

# Library Matrix-Matrix Multiply (DGEMM)
## MKL 5.1 Library



**May be much higher with Opteron-optimized Libs – (e.g., NAG Lib.)**

# Remote & Local Memory Read/Write



Swap the 2 j columns

$A(i,j) = time \cdot A(i,j)$

AMD Opteron

Remote Access

Local Access

**y-axis:** clock cycles   per iteration

**x-axis:** matrix leading dimension (n)

400  700  1000  1400  2500  4000  5500  7000  8500  10000  11500  13000

1.) Each processor writes a column to local memory.

2.) Each processor reads/writes to same column. (Local Access)

3.) Each processor "swaps" column index and reads/writes to remote memory. (Remote Access)

# Applications

- **SM**: Stommel model of ocean "circulation" ; solves 2-D partial differential equation.
  - Uses <u>Finite Difference</u> approx for derivatives on discretized domain, (timed for a constant number of Jacobi iterations).

  **Memory Intensive**

- **MD**: Molecular Dynamics of argon lattice.
  - Uses Verlet algorithm for propagation (displacement & velocities).

  **Compute Intensive**

| Platform | Serial SM (sec) | Parallel SM (sec) |
|---|---|---|
| AMD Opteron 2P | 68.0 | 43.4 |
| Intel Xeon 2P | 57.5 | 63.4 |

| Platform | Serial MD (sec) | Parallel MD (sec) |
|---|---|---|
| AMD Opteron 2P | 9.48 | 5.3 |
| Intel Xeon 2P | 7.14 | 5.4 |

TACC

22

# Summary

| | SERIAL | | Parallel | |
|---|---|---|---|---|
| | Opteron | Xeon | Opteron | Xeon |
| Latency | Low | High | Overlapped | Overlapped |
| Band-width | ~2GB/s | ~2GB/s | 2x | 1x |
| MxM (per CP) | Lower | Higher | 2x mem | 1x mem |
| MXM (time) | Higher | Lower | slightly lower | slightly higher |
| DGEMM | Low | High | Scale: 1.9x (MKL 5.1 not optimized for AMD) | Scale:1.8x 2x Opteron performance |

# Summary

- Performance of dual-processor systems varies with memory architecture and processor speed.
  - AMD memory bandwidth scales by 2x when second processor is used– (using "local" memory).
  - Xeon memory bandwidth is shared by second processor.
  - Xeon outperforms Opteron on serial compute-intensive codes (due to speed: 2.4GHz Xeon vs. 1.4GHz Opteron); but lead can be eliminated with dual-processor execution of (parallel) programs when memory bandwidths & synchronizations are involved.

TACC

24