# Performance and Productivity of Emerging Architectures

Presented by

## Jeremy Meredith

## Sadaf Alam

## Jeffrey Vetter
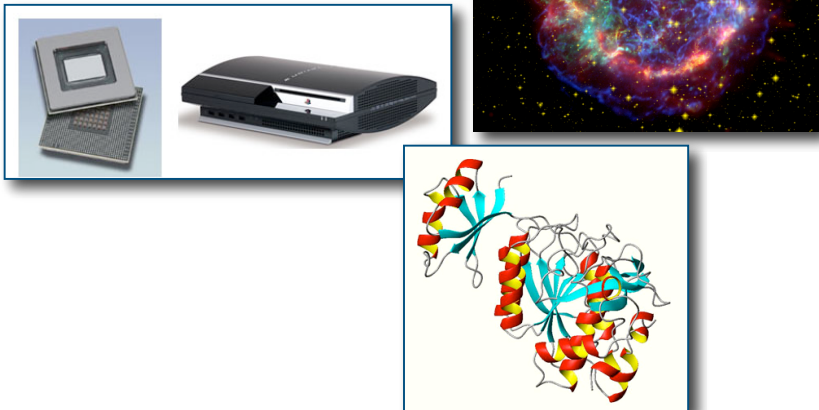
Future Technologies

SC07

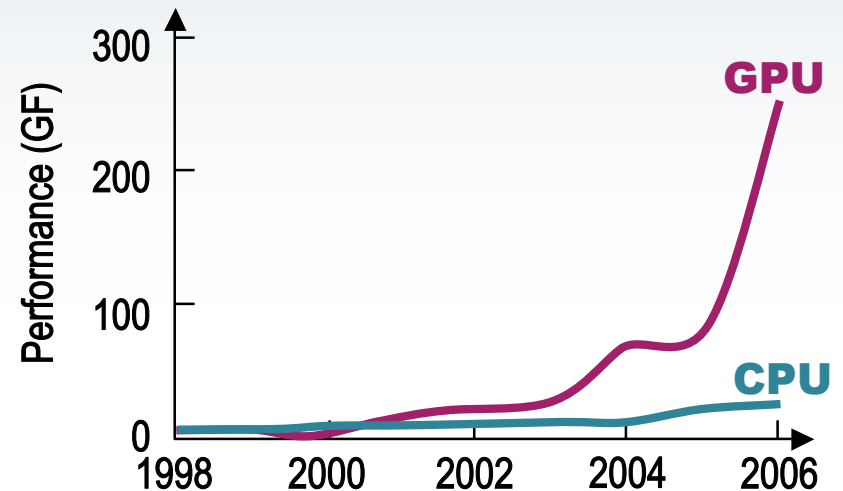OAK RIDGE
National Laboratory

# Overview

| Motivation | Goals |
|---|---|
| • Technologies on the horizon several orders of magnitude more powerful and power-efficient than microprocessors | • Gauge performance of emerging devices using high-level languages |
| • Grand challenge applications requiring several orders of magnitude more performance than now available | • Estimate productivity with respect to contemporary microprocessing devices |

# Intel Xeon multicore CPU



DIB
1066/1333 MHz
8.5/10.5 GB/s

Up to 21 GB/s

Dempsey
Woodcrest
Clovertown

ESB-2
I/O
bridge

ESI

Blackford
MCH

x8

FBD
FBD
FBD
FBD

x8    x8

Configurable set of PCIe ports

PCIe slot

Sunrise
Lake

PCI-X

10 GbE
I/OAT
iSCSI

SAS/SATA-2    PCI-X    10 GbE
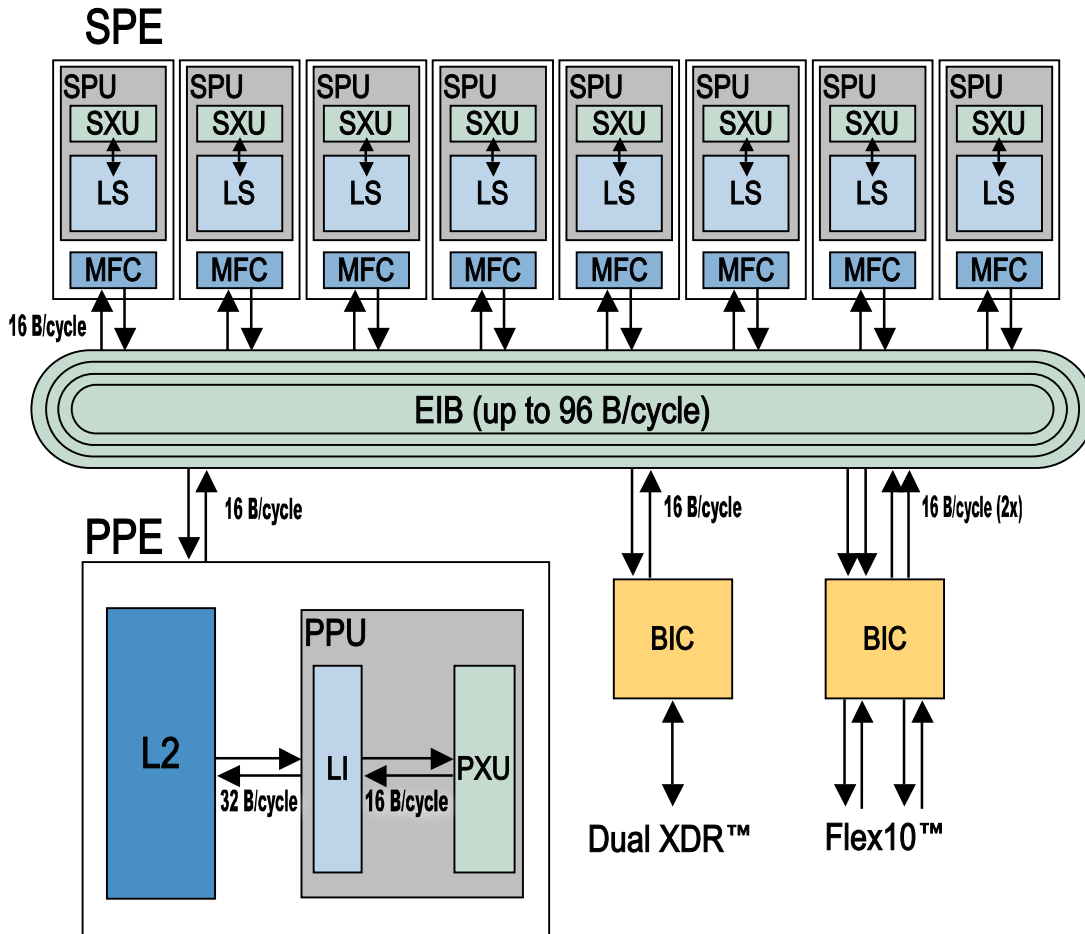
- Dual-core (Woodcrest)

- Quad-core (Clovertown)

- Programming models
  - MPI
  - OpenMP
  - Pthreads

- Recourse contention
  - Memory bandwidth
  - I/O bandwidth
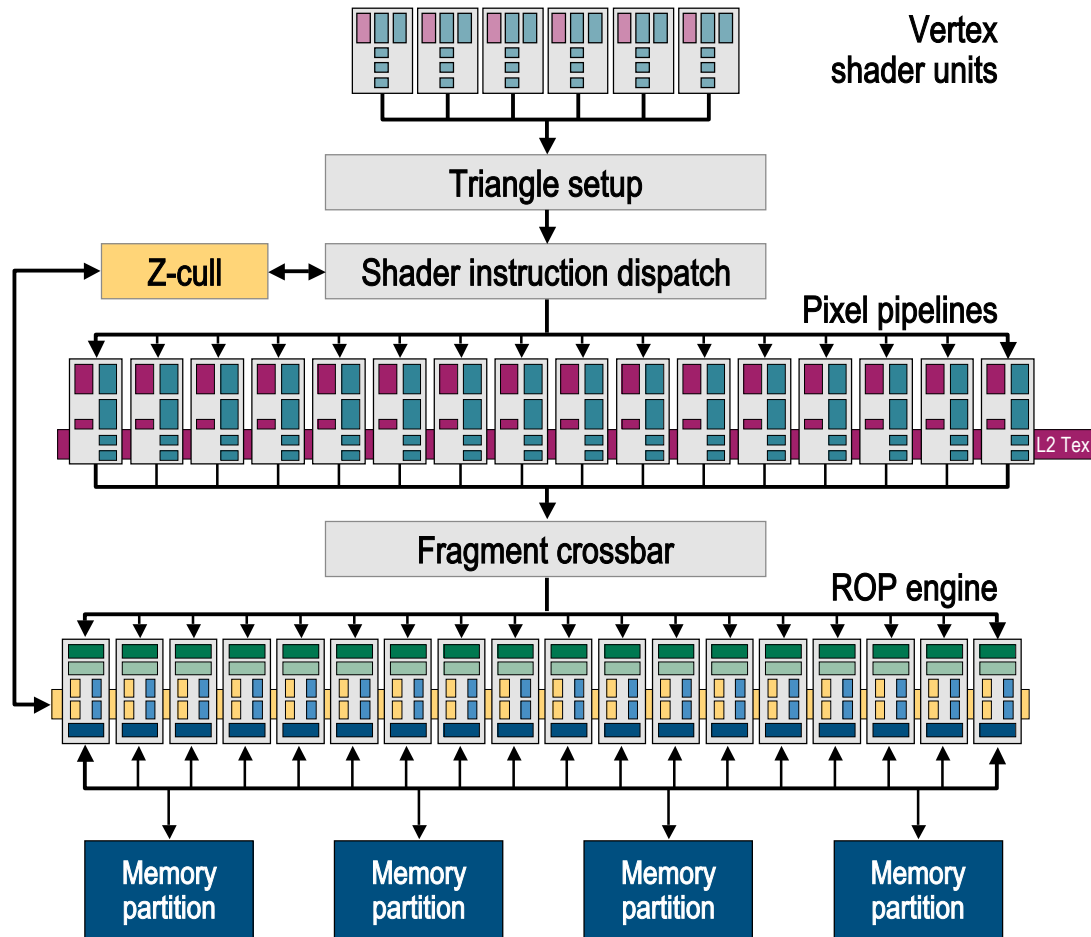
OAK
RIDGE
National Laboratory

# IBM Cell broadband engine



- Cell heterogeneous multicore processor
  - 8 synergistic processing element (SPE) cores
  - 200 GFLOPS at 3.2 GHz
  - 200-GB/s memory bandwidth

- Programming models
  - Multisource compilation
  - Threadlike launch semantics for SPEs

*Source: M. Gschwiind et al., Hot Chips-17, August 2005*

OAK RIDGE National Laboratory

# Graphics processing unit (GPU)



- **NVIDIA 7900GTX**
  - 24 SIMD pixel pipelines
  - 200 GFLOPS at 650 MHz
  - 50-GB/s memory bandwidth

- **Programming models**
  - Multiple-source compilation
  - Gather semantics define parallelism
  - Host CPU drives program setup and data transfer

OAK
RIDGE
National Laboratory

# Extreme multithreading with Cray MTA-II

Programs running in parallel

Concurrent threads of computation

i=n
i=3
i=2
i=1

Sub-problem A
Sub-problem B

i=n
i=1
i=0
Sub-problem A

Serial Code

Hardware streams (128)

Unused streams

Instruction ready pool

Pipeline of executing instructions

- Cray XMT system
  - MTA-I and MTA-II processor architecture
  - XT3/4 scalable infrastructure
  - AMD Torrenza technology

- Fine-grain multithreading (128 concurrent threads)

- Uniform memory hierarchy in MTA-I and MTA-II

OAK RIDGE National Laboratory

# Application kernels
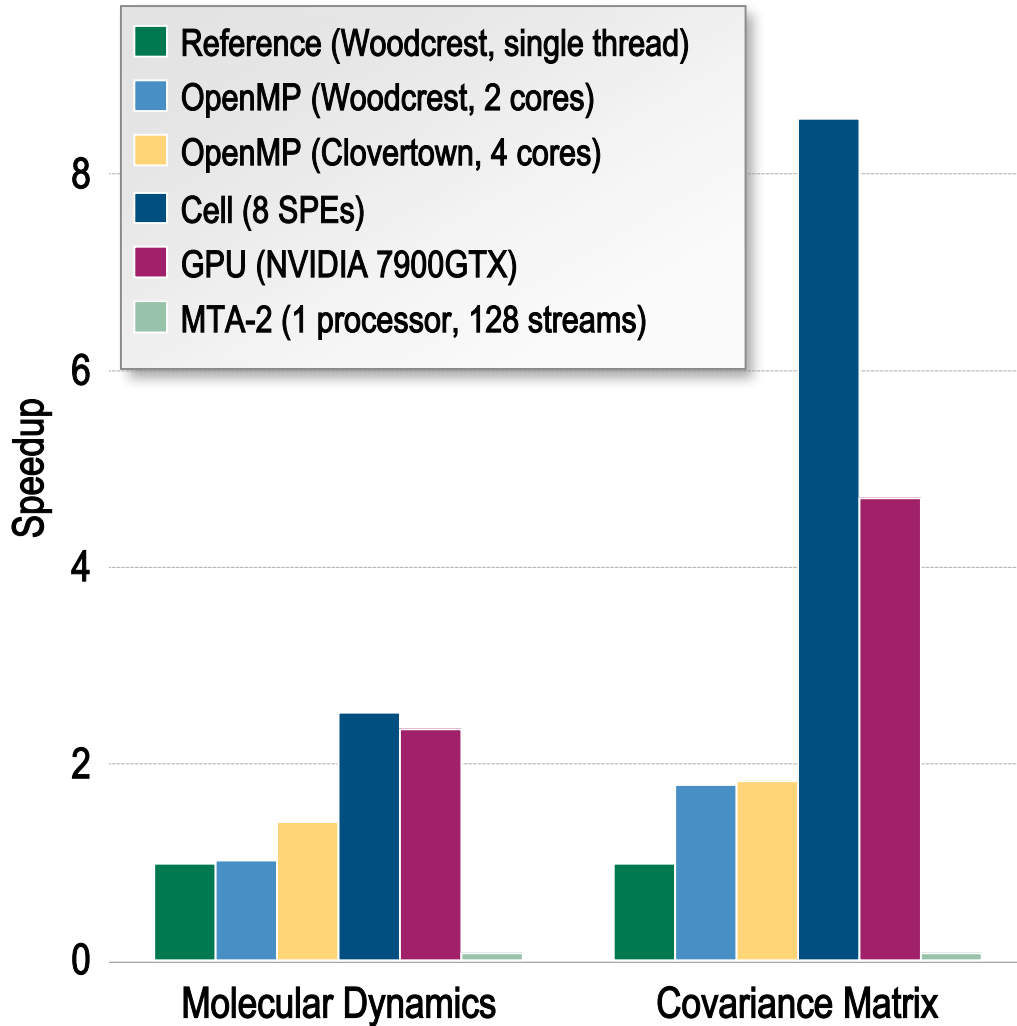
## Molecular dynamics (MD) calculations

- Applications in biology, chemistry, and materials
- Newton's second law of motion
- Bonded calculations
- Nonbonded electrostatic and van der Waals forces
- Workload characteristics
  - Floating-point intensive
  - Irregular memory access patterns



## Covariance matrix calculation

- Applications in hyperspectral imaging and machine learning
- Determines covariance among samples in data
- Workload characteristics
  - Memory-bandwidth and floating-point intensive
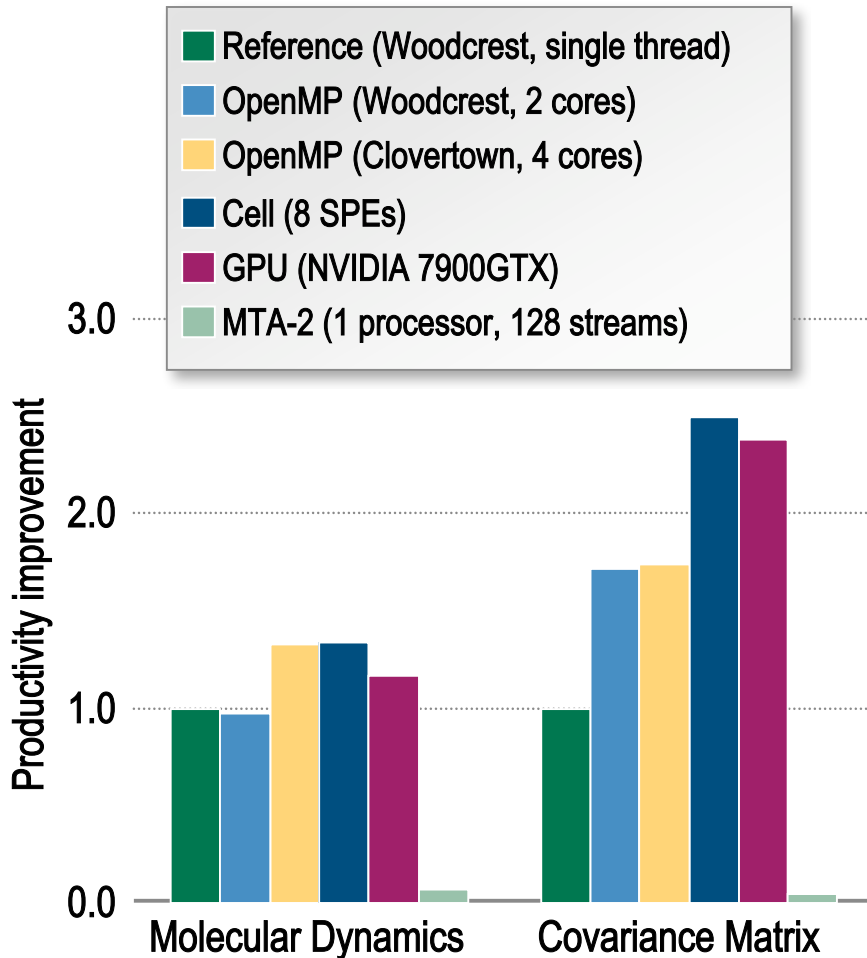  - Regular memory access patterns

OAK RIDGE
National Laboratory

# Performance



Legend:
- Reference (Woodcrest, single thread)
- OpenMP (Woodcrest, 2 cores)
- OpenMP (Clovertown, 4 cores)
- Cell (8 SPEs)
- GPU (NVIDIA 7900GTX)
- MTA-2 (1 processor, 128 streams)

Y-axis: Speedup (0, 2, 4, 6, 8)

X-axis categories: Molecular Dynamics, Covariance Matrix

- **OpenMP improves performance moderately on commodity CPUs**

- **High parallelism of Cell and GPU can improve performance, but more so when memory access is regular**
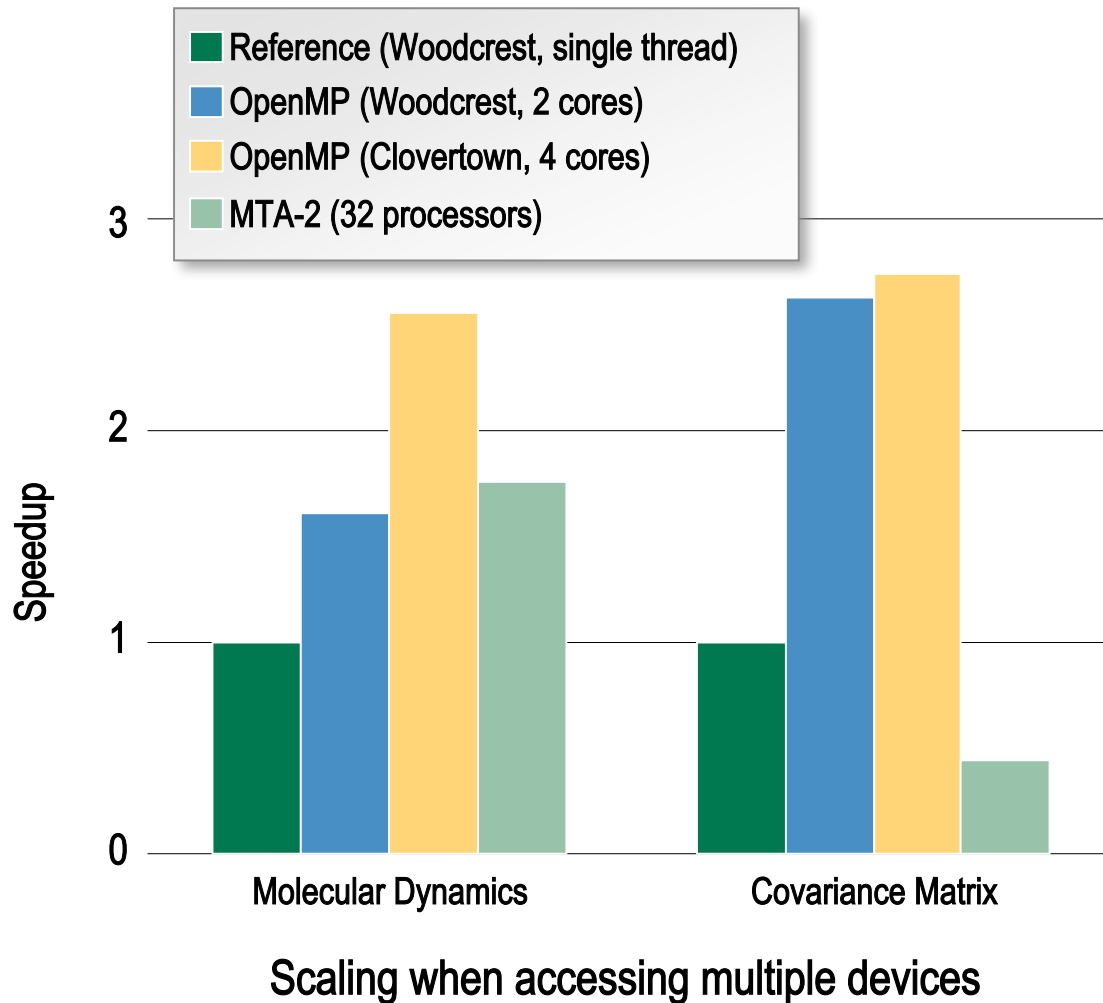
OAK RIDGE National Laboratory

# Productivity



Legend:
- Reference (Woodcrest, single thread)
- OpenMP (Woodcrest, 2 cores)
- OpenMP (Clovertown, 4 cores)
- Cell (8 SPEs)
- GPU (NVIDIA 7900GTX)
- MTA-2 (1 processor, 128 streams)

Y-axis: Productivity improvement (0.0, 1.0, 2.0, 3.0)

Categories: Molecular Dynamics, Covariance Matrix

- Despite small performance increases through OpenMP, the ease of using it means that productivity can still be increased

- Conversely, the high speed of Cell and GPU means that even substantial effort results in higher productivity

OAK RIDGE
National Laboratory

# Productivity scaling



Scaling when accessing multiple devices

Increased parallelism from all architectures with little increased effort results in higher productivity

OAK RIDGE
National Laboratory

# Contacts

**Jeremy Meredith**

Future Technologies Group
(865) 241-5842
jsmeredith@ornl.gov

**Sadaf Alam**

Future Technologies Group
(865) 241-1533
alamrs@ornl.gov

**Jeffrey Vetter**

Future Technologies Group
(865) 576-7115
vetter@ornl.gov

OAK RIDGE
National Laboratory