# DOE Evaluation of the Cray X1

**Mark R. Fahey**

**James B. White III (Trey)**

**Center for Computational Sciences**

**Oak Ridge National Laboratory**

# Acknowledgement

**Research sponsored by the Mathematical, Information, and Computational Sciences Division, Office of Advanced Scientific Computing Research, U.S. Department of Energy, under Contract No. DE-AC05-00OR22725 with UT-Battelle, LLC.**

# Outline

$\sum$ **CCS X1**

$\sum$ **Evaluation Overview**

$\sum$ **Applications**

   **Climate**

   **Fusion**

   **Materials**

   **Biology**

# Phase 1 – March 2003

$\Sigma$ **32 Vector Processors**

    **8 nodes, each with 4 processors**

$\Sigma$ **128 GB shared memory**

$\Sigma$ **8 TB of disk space**

**400  GigaFLOP/s**

# Phase 2 – September 2003



∑ **256 Vector Processors**

**64 nodes**

∑ **1 TB shared memory**

∑ **20 TB of disk space**

**3.2 TeraFLOP/s**

# X1 evaluation

$\sum$ **Compare performance with other systems**

    **Applications Performance Matrix**

$\sum$ **Determine most-effective usage**

$\sum$ **Evaluate system-software reliability and performance**

$\sum$ **Predict scalability**

$\sum$ **Collaborate with Cray on next generation**

# Hierarchical approach

$\Sigma$ **System software**

$\Sigma$ **Microbenchmarks**

$\Sigma$ **Parallel-paradigm evaluation**

$\Sigma$ **Full applications**

$\Sigma$ **Scalability evaluation**

# System-software evaluation

$\Sigma$ **Job-management systems**

$\Sigma$ **Mean time between failure**

$\Sigma$ **Mean time to repair**

$\Sigma$ **All problems, with Cray responses**

$\Sigma$ **Scalability and fault tolerance of OS**

$\Sigma$ **Filesystem performance & scalability**

$\Sigma$ **Tuning for HPSS, NFS, and wide-area high-bandwidth networks**

$\Sigma$ **See Buddy Bland's Talk**

> **"Early Operations Experience with the Cray X1"**
> **Thursday at 11:00**

# Microbenchmarking

$\sum$ **Results of standard benchmarks**

    **http://www.csm.ornl.gov/~dunigan/cray**

    **See Pat Worley's talk today at 4:45**

$\sum$ **Performance metrics of components**

    **Vector & scalar arithmetic**

    **Memory hierarchy**

    **Message passing**

    **Process & thread management**

    **I/O primitives**

$\sum$ **Models of component performance**

# Parallel-paradigm evaluation

$\sum$ **MPI-1, MPI-2 one-sided, SHMEM, Global Arrays, Co-Array Fortran, UPC, OpenMP, MLP, …**

$\sum$ **Identify best techniques for X1**

$\sum$ **Develop optimization strategies for applications**

# Scalability evaluation

$\Sigma$ **Hot-spot analysis**

- **Inter- and intra-node communication**
- **Memory contention**
- **Parallel I/O**

$\Sigma$ **Trend analysis for selected communication and I/O patterns**

$\Sigma$ **Trend analysis for kernel benchmarks**

$\Sigma$ **Scalability predictions from performance models and bounds**

# Full applications

$\Sigma$ **Full applications of interest to DOE Office of Science**

   **Scientific goals require multi-tera-scale resources**

$\Sigma$ **Evaluation of performance, scaling, and efficiency**

$\Sigma$ **Evaluation of ease/effectiveness of targeted tuning**

# Identifying applications

$\sum$ **Draft evaluation plan**

$\sum$ **Prototype Workshop at ORNL Nov. 5-6**

$\sum$ **Feb 3-5, 2003: Fusion**

$\sum$ **Feb 6, 2003: Climate**

$\sum$ **March 2, 2003: Materials**

$\sum$ **May 9, 2003: Biology**

$\sum$ **Future DOE-wide workshops**

# Workshop Goals

$\Sigma$ **Set priorities**

    **Potential performance payoff**

    **Potential science payoff**

$\Sigma$ **Schedule the pipeline**

    **porting/development**

    **processor tuning**

    **scalability tuning**

    **production runs - *science!***

    ***small* number of applications in each stage**

# Identifying applications

$\sum$ **Potential application**

   **Important to DOE Office of Science**

   **Scientific goals require multi-terascale resources**

$\sum$ **Potential user**

   **Knows the application**

   **Willing and able to learn the X1**

   **Motivated to tune application, not just recompile**

# Climate

$\sum$ **3 codes**

  **CAM, CLM, POP**

$\sum$ **Participants from**

  **NCAR, LANL, LBNL, ORNL, NASA-Goddard, CRIEPI, Cray, NEC**

$\sum$ **Want to optimize for NEC and Cray**

  **May require different optimizations**

# Climate: CAM

$\sum$ **People involved**

   **Cray(1), NEC(2), NCAR(1), ORNL(2)**

$\sum$ **Porting, profiling ongoing at Cray**

$\sum$ **NEC expects single node optimizations for SX-6 complete by early Fall**

   **Coordination between NEC and Cray?**

$\sum$ **Radiation and Cloud models are focus of most work**

# Climate: CLM

$\Sigma$ **Land component of the Community Climate System Model**

$\Sigma$ **Undergoing changes to data structures to make easier to extend and maintain**

    **Fortran user-defined types with pointers**

$\Sigma$ **Vectorization involvement**

    **NCAR(1), ORNL(2), Cray(1), NEC(1)**

    **Coordination with NEC to be worked out**

$\Sigma$ **See Trey White's presentation**

    **Wednesday at 8:45**

OAK RIDGE NATIONAL LABORATORY
U. S. DEPARTMENT OF ENERGY

# Climate: POP

- $\Sigma$ **Organization involvement**
  - LANL(1), Cray(1), NCAR(2), CRIEPI(2)
- $\Sigma$ **Need to coordinate between CRIEPI and Cray**
- $\Sigma$ **Significant optimizations already implemented, successful**
  - Vectorization and Co-Array Fortran
- $\Sigma$ **Remaining issues**
  - Parallel algorithm issues
  - I/O issues
- $\Sigma$ **See Pat Worley's presentation**
  - "Early Performance Evaluation of the Cray X1"
  - Today at 4:45

# Fusion

- $\Sigma$ **Workshop held Feb 3-5 @ ORNL**
- $\Sigma$ **Participants from**
  - **General Atomics, Princeton Plasma Physics Lab, University of Wisconsin, University of Iowa, Cray, ORNL**
- $\Sigma$ **6 codes**
  - **M3D and NIMROD (extended MHD)**
  - **GYRO and GTC (micro turbulence)**
  - **AORSA and TORIC (RF plasma interactions)**
- $\Sigma$ **Concurrent work by different teams**
- $\Sigma$ **Too many codes?**
  - **Provides flexibility when impediments encountered**

# Fusion: NIMROD

$\Sigma$ **CCS Teaming with developer and Cray to port and optimize**

    **Cray has actively participated**

$\Sigma$ **Uses F90 reshape quite a bit**

    **Exploits a known weakness in the compiler**

    **Cray filed SPR**

$\Sigma$ **Uses F90 sums extensively inside loops that should be vectorizable**

    **Compiler cannot vectorize, arrays are actually pointers**

    **Cray filed SPR**

$\Sigma$ **Dramatic effect on performance**

    **Cannot predict how fast will be when compiler fixed**

# Fusion: NIMROD (cont.)

$\Sigma$ **Data structures are derived types of pointers with allocatable attribute**

$\Sigma$ **Pointers vs allocatable arrays**

How much performance can be gained by replacing pointers?

$\Sigma$ Would benefit other architectures too

Analysis needed before code rewrite can even be discussed

Climate Land Model success is important here

$\Sigma$ See Trey White's presentation

$\Sigma$ Wednesday at 8:45

**OAK RIDGE NATIONAL LABORATORY**
**U. S. DEPARTMENT OF ENERGY**

# Fusion: GYRO

$\sum$ **Developer and CCS teaming**

$\sum$ **Implemented in F90, no derived data-types**

$\sum$ **Hand-coded transpose operations using loops over MPI_Alltoall**

    **Expected to scale to ES class machine**

    **Scaling is ultimately limited by this**

$\sum$ **UMFPACK library for field solves**

# Fusion: GYRO (cont.)

$\Sigma$ **Functional port complete**

> **Has identified a couple bugs in code**

$\Sigma$ **Several routines easily vectorized by manual loop interchange, and directives**

$\Sigma$ **Vectorized sin/cos calls by rewriting code**

> **Numerical integration routine**
>
> **Bisection search**

$\Sigma$ **Hand optimizations have yielded 5X speedup so far (more work to do)**

> **About 35% faster than PWR4 (not enough!)**

# Fusion: GTC

$\sum$ **Developer ran GTC on SX6**

$\sum$ **Cray had previously looked at parts of GTC**

$\sum$ **Result:**

   **Developer is directly working with Cray**

$\sum$ **GTC has been ported to the X1**

   **Some optimizations introduced**

   **Work ongoing**

# Fusion: AORSA

$\sum$ **Uses ScaLAPACK**

$\sum$ **Cray has ScaLAPACK implementation**

    **Not tuned**

    **Cray pursuing ScaLAPACK optimizations**

$\sum$ **Ported**

    **Performance worse than expected**

    **Culprit is matrix scaling routine**

      $\sum$**Fix implemented, tests underway**

$\sum$ **With Cray Benchmarking group**

# Fusion: M3D

$\sum$ **M3D uses PETSc**

 - **Parallel data layout done within this framework**

 - **Uses the iterative solvers**

 - **Accounts for 90% of time**

$\sum$ **Need to port PETSc to X1**

 - **Estimate of 6 man-months**

 - **Require significant changes**

# Materials

$\sum$ **Primary codes**

    **Dynamic Cluster Algorithm, FLAPW, LSMS, Socorro**

$\sum$ **Secondary codes**

    **LAMMPS, GP, FEFF/TD-DFT, a M-C code**

$\sum$ **Majority are C++ and use MPI and/or OpenMP**

$\sum$ **Contacts for each code identified**

# Materials: Dynamic Cluster Alg.

$\Sigma$ **MPI, OpenMP, PBLAS, BLAS**

  **significant amount of time spent in dger**

  **significant amount of time spent in cgemm**

  **On the IBM Power4, the blas2 calls dominate**

$\Sigma$ **A quick port was performed**

  **Optimizations targeted a couple routines adding a few directives**

  **Took a couple days**

  **Showed dramatic speedup over IBM Power4**

  **For the small problem that was solved**

  $\Sigma$ **time doing calculations became nearly negligible**

  $\Sigma$ **formatted I/O became dominant**

# Materials: LSMS

$\sum$ **Locally Self-consistent Multiple Scattering**

$\sum$ **Code spends most of its time**

      **matrix-matrix mulitplications**

      **Computing partial inverse**

$\sum$ **Communication involves exchanges of smaller matrices with neighbors**

$\sum$ **Expected to vectorize well**

$\sum$ **Developers are moving to sparse-matrix formulations to scale to larger problems**

$\sum$ **With Cray Benchmarking group**

# Materials: FLAPW

$\Sigma$ **Full Potential Linearized Augmented Plane Wave (FLAPW) method**

    **All-electron method**

    **Considered to be most precise electron structure method in solid state physics**

$\Sigma$ **Validation code and as such is important to a large percentage of the materials community**

$\Sigma$ **Cray has started porting it**

# Biology

$\sum$ **Workshop May 9 (few days ago!)**

$\sum$ **Bioinformatics plan to exploit special features of X1**

  Not one or two codes that are more important

  Probably can use primitives in BioLib

  Collaborate with Cray on adding more primitives to BioLib

$\sum$ **Molecular dynamics based biology codes expected to vectorize**

  AMBER is used a lot

  Cray working on AMBER port already, nearly done

# Conclusions

$\sum$ **Future:**

  **Chemistry and Astrophysics workshops**

$\sum$ **Early results are promising**

$\sum$ **Evaluation continues, much to do**

$\sum$ **Workshops very productive**

  **Focused set of codes to port is important**

  **Identify users/teams**

# References

$\Sigma$ **System software evaluation**

    **Buddy Bland's talk Thursday at 11:00**

$\Sigma$ **Results of standard benchmarks**

    **http://www.csm.ornl.gov/~dunigan/cray**

    **Pat Worley's talk today at 4:45**

$\Sigma$ **Optimization Experiment with CLM**

    **Trey White's talk Wednesday at 8:45**

OAK RIDGE NATIONAL LABORATORY
U. S. DEPARTMENT OF ENERGY

# Contacts

$\sum$ **Trey White**

     **whitejbiii@ornl.gov**

$\sum$ **Mark Fahey**

     **faheymr@ornl.gov**

$\sum$ **Pat Worley**

     **worleyph@ornl.gov**

$\sum$ **Buddy Bland**

     **blandas@ornl.gov**

$\sum$ **consult@ccs.ornl.gov**

**OAK RIDGE NATIONAL LABORATORY**
**U. S. DEPARTMENT OF ENERGY**