

LA-UR- 02-6411

Approved for public release;
distribution is unlimited.

Title: Monte Carlo Radiation Transport & Parallelism

Author(s): Lawrence J. Cox, X-5
Susan E. Post, CCN-8

Submitted to: Los Alamos Computer Science Institute (LACSI) Symposium,
Santa Fe, NM
October 14-16, 2002



Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the University of California for the U.S. Department of Energy under contract W-7405-ENG-36. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this content, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

Monte Carlo Radiation Transport & Parallelism

Larry Cox, X-5

Susan Post, CCN-8

The Los Alamos National Laboratory logo, which includes a stylized blue and white graphic of a particle or atom to the left of the text 'Los Alamos' in a large, bold, sans-serif font. Below 'Los Alamos' is the text 'NATIONAL LABORATORY' in a smaller font.

Los Alamos

Abstract

This talk summarizes the main aspects of the LANL ASCI Eolus project and its major unclassified code project, MCNP. The MCNP code provide a state-of-the-art Monte Carlo radiation transport to approximately 3000 users world-wide. Almost all hardware platforms are supported because we strictly adhere to the FORTRAN-90/95 standard. For parallel processing, MCNP uses a mixture of OpenMP combined with either MPI or PVM (shared and distributed memory).

This talk summarizes our experiences on various platforms using MPI with and without OpenMP. These platforms include PC-Windows, Intel-LINUX, BlueMountain, Frost, ASCI-Q and others.



Eolus/MCNP Tradition at Los Alamos

- For decades, Monte Carlo radiation transport codes, from MCS to MCNP, have been developed and supported by the Monte Carlo team at LANL
- Concurrently, the extensive nuclear and atomic data libraries have also been under constant development
- This tradition continues in the Eolus ASCI Project and related efforts in X-5 of the WP Directorate
 - 14 MCNP team members
 - Physical Data team also in X-5
 - Two application teams (user groups) in X-5



Eolus/MCNP Tradition at Los Alamos

- Additional support is provided by other LANL groups
 - CCN-8, X-3
 - Porting, parallel efficiency analysis, testing
 - IM-8
 - Automated regression testing, software engineering improvements, modularity, quality assurance
 - CCS-4
 - Collaboration on Monte Carlo methods
 - T-16
 - Nuclear data and methods research

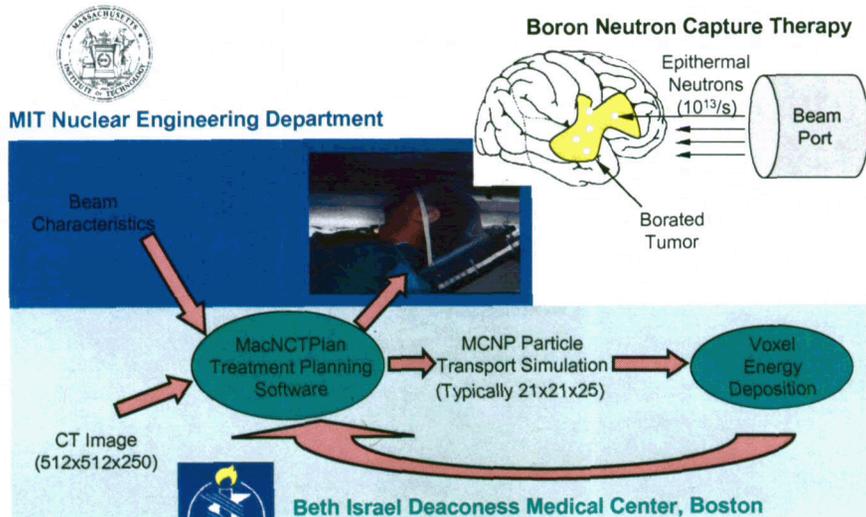


Some MCNP Applications

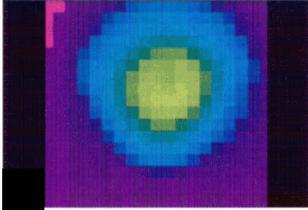
- Nuclear Criticality Safety
- Radiation Shielding
- Nuclear Safeguards
- Detector Design and Analysis
- Nuclear Well Logging
- Personnel Dosimetry
- Health Physics
- Accelerator Target Design
- Medical Physics and Radiotherapy
- Fission and Fusion Reactor Design
- Waste Storage/Disposal
- Radiography
- Aerospace Applications
- Decontamination and Decommissioning



Neutron Capture Therapy

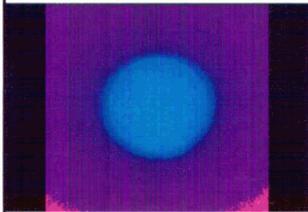


MCNP BNCT Simulation Results



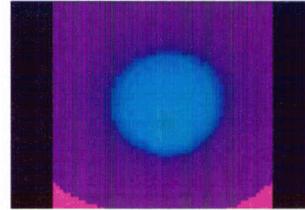
Typical MCNP BNCT simulation:

- 1 cm resolution (21x21x25)
- 1 million particles
- 1 hour on 200 MHz PC



ASCI Blue Mountain MCNP simulation:

- 4 mm resolution (64x64x62)
- 100 million particles
- 1/2 hour on 6048 CPUs



ASCI Blue Mountain MCNP simulation:

- 1 mm resolution (256x256x250)
- 100 million particles
- 1-2 hours on 3072 CPUs



US-Japan Joint Reassessment of Atomic Bomb Radiation Dosimetry in Hiroshima and Nagasaki

Steve White
Alexandra Heath
Paul Whalen

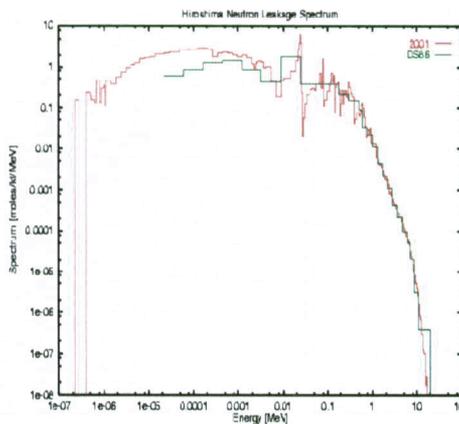


Hiroshima; August 6, 1945 ...




Los Alamos
 NATIONAL LABORATORY

Hiroshima Neutron Leakage Spectrum



LANL MCNP CALCULATION

	<i>New</i>	<i>DS86</i>	<i>new/DS86</i>
Total Neutrons (Moles/kt)	1.768E-01	1.773E-01	0.9972
Neutron energy (Average MeV)	0.3106	0.2976	1.0437
Yield (Kt)	16.1	(12 - 20)	


Los Alamos
 NATIONAL LABORATORY

Modernization of MCNP

- MCNP 4c3 frozen since April 2001
 - MCNP 4c3 is the last old-style MCNP version that will be released
- In the last 15 months, every line of code has been reworked
 - Conversion to ANSI-Standard Fortran-90
 - Completely new installation system
 - Modified patching method
- An ongoing process...



Eolus Parallelism

Analog Monte Carlo radiation transport can be made embarrassingly parallel because MC histories are independent

- Parallel efficiency is limited only by availability of memory to copy all information (geometry, cross sections, etc.) and by the overhead required to add up the results
- Providing statistical analysis and tracking of sequential results complicates the algorithm
- The actual algorithm chosen can also effect efficiency
 - Collective operations versus serial sweep of results
 - Reliance on shared memory locks or global synchronization versus designing in independence



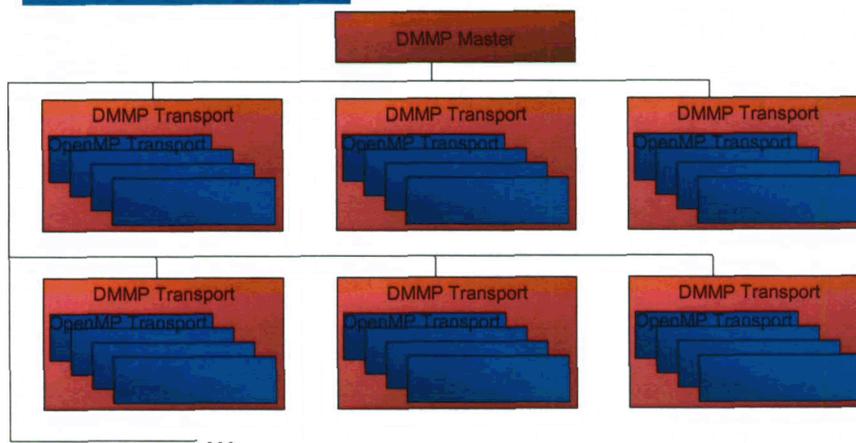
Enhanced Parallel Processing

- We use both distributed and shared memory parallelism separately or in combination
 - MPI is supported
(<http://www.amix.ncsl.gov/mpl>) 
 - PVM support is enhanced
(<http://www.epm.ornl.gov/pvm>) 
 - OpenMP support is added for shared-memory (threaded) parallel processing
(<http://www.openmp.org>) 

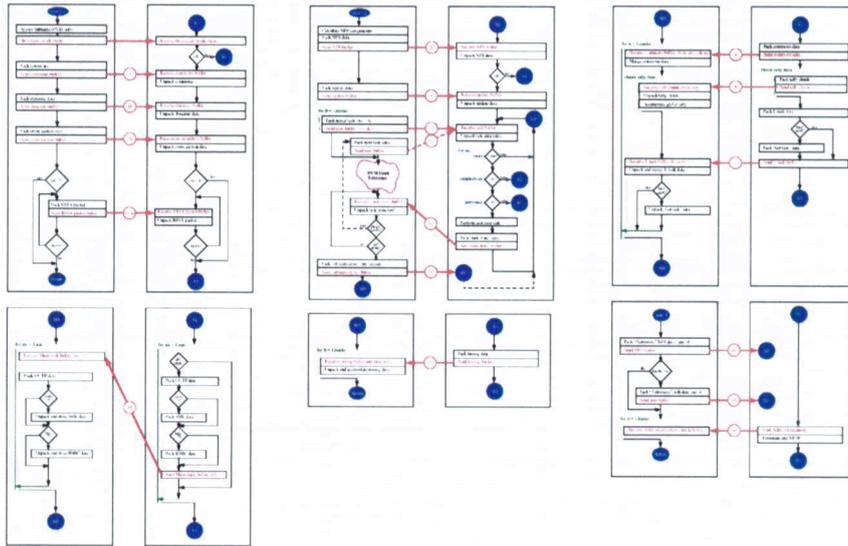
Either MPI or PVM can be used with OpenMP for optimal performance



Current Eolus Parallelism Algorithm



Distributed Communication



Mixed Mode Parallelism on ASCI Platforms

- **BlueMountain (IRIX64):** 128 PE boxes in a deep NUMA configuration
 - OpenMP efficiency scales approximately as $(N)^{1/2}$
 - Best efficiency at about 8-10 threads per DMMP transport task
 - 2% to 10% effect on sequential runs with OpenMP compiled in
- **White (AIX):** 16 PE nodes (non-NUMA)
 - Ideal mode: 1-4 MPI tasks per nodes with 4-16 threads per tasks
 - 20% to 80% (or greater) hit when compiled with OpenMP even for sequential runs
 - Is it the threading scheme (THREADPRIVATE data) or the libraries?
 - Can it be improved?
- **Q (OSF1):** 4 PE nodes (non-NUMA)
 - Ideal mode: One MPI task per node with 4 threads
 - 0% to 20% loss of efficiency when OpenMP is compiled in
 - Considering an MPI-only public version
 - Will be more severely limited by memory with 4 copies per node



Analysis & Efficiency Issues

- MPI runs sequentially as well as sequential executables
- Between 5-50% effect on sequential timing by compiling in OpenMP (system & compiler dependent)
 - Does not matter if locks are replaced with critical regions
 - AIX is the worst; SGI is the best; OSF1 is acceptable
- Runtime analysis involves use of special compilers that...
 - Do not work (insurmountable compiler errors on SGI);
 - Require removal of standard language features
 - Generate behavior that is so different from that from native OpenMP compilers that we cannot tell if it is the same code or if problems are the same problems



Analysis and Efficiency Issues

- A minor code change led to a 20-50% effect on AIX
 - We have not yet determined why
- Runtimes are erratic for a certain class of problems
 - Correlated to MPI startup speed on OSF1 (30% variation)
 - Correlated to processor count (20% at 32pe to 400% at 256pe)
- Use of Fortran pointers versus allocatables
 - 50% effect on IRIX64 ; 15% on OSF1; no effect on AIX



Criticality Performance

- Spreadsheet here from 9/12/02 (spost)



Parallel Support on PCs

- MPI Version – MPICH.NT.1.2.4
 - <http://www-unix.mcs.anl.gov/~ashton/mpich.nt/>
 - Uses GUI mpirun
 - Socket communication; allowed on LANL networks
 - Tested on Dual CPU Win 2000, Cluster of 2
- PVM Version – PVM 3.4.3
 - <http://www.csm.ornl.gov/~sscott/PVM/Software/>
 - Must start pvm on master, add other PCs to cluster
 - uses rsh; not allowed on LANL networks.
 - Tested on Dual CPU Win 2000, Cluster of 1



Parallel Pros/Cons (PC)

- MPI
 - Well suited for homogeneous cluster, less overhead
 - Cluster of Windows NT/2000/XP only
 - No control-c interrupt capability
 - If any task/Machine dies, does not recover gracefully
- PVM
 - Cluster of Linux/Unix/Windows Machines
 - Has control-c interrupt
 - If one task/Machine dies, recovers unless master died
 - Possible to spawn appropriately on single/dual CPU cluster

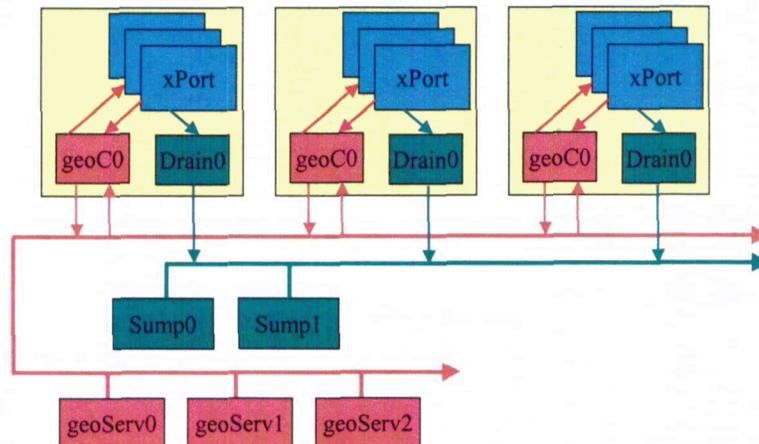


Domain Decomposition

- It may be necessary to implement various forms of domain-decomposition due to addressable memory limits
 - Already have encountered >32GB geometry files



One Possible Domain Decomposition Algorithm



Eolus Parallel Support Needs

- Efficiency in parallel methods
 - MPI, OpenMP
- Parallelism algorithm analysis
 - Improvements to internal algorithms, memory layout...
- Debugging in parallel modes
 - Totalview is pretty good
- Coverage analysis in parallel modes
- Efficiency analysis in parallel modes



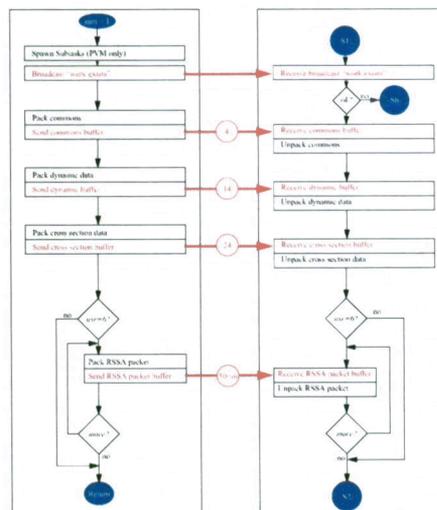
Contacts and Support

- **Email contacts**
 - User forum: mcnp-forum@lanl.gov
 - Direct Contact: ljcox@lanl.gov
- **Web sites**
 - <http://www-xdiv.lanl.gov/x5/MCNP/index.html>
 - <http://epicws.epm.ornl.gov/ENOTE.html>



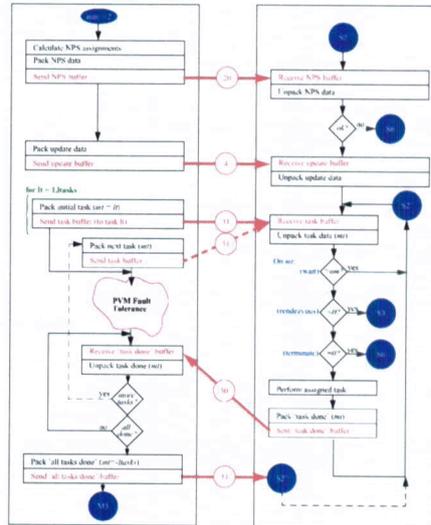
Initialization: Master to all Transport Tasks

- All data is broadcast from the master task to the transport tasks
 - Problem setup data
 - Mostly read-only



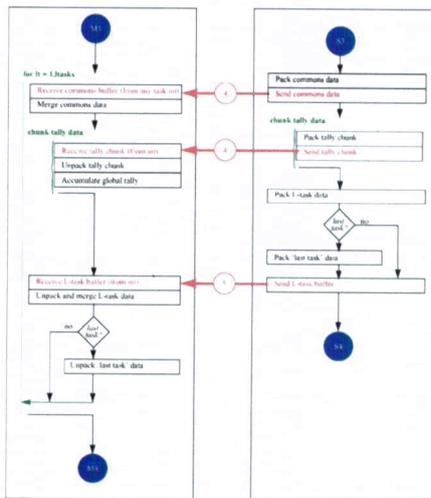
New assignments made and executed

- Involves broadcast and point to point messages



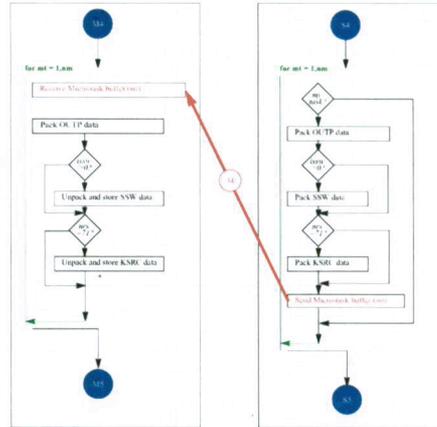
L-Task data accumulation

- Point to point from Transport tasks to Master
- First-come, First-served
 - Order of reception not important



M-Task data accumulation

- Micro-task data is taken in order to preserve some state info
 - Source data, output messages
 - Possibly more than one per L-task
 - Order is important



Timing Data Accumulation & Termination

- Timing data taken on first-come, first-served basis
- Termination is a broadcast work assignment signifying no more work
 - Acknowledgement required

