

# VECTORIZING THE COMMUNITY LAND MODEL

**Forrest M. Hoffman<sup>1</sup>**  
**Mariana Vertenstein<sup>2</sup>**  
**Hideyuki Kitabata<sup>3</sup>**  
**James B. White III<sup>1</sup>**

## Abstract

In this paper we describe our extensive efforts to rewrite the Community Land Model (CLM) so that it provides good vector performance on the Earth Simulator in Japan and the Cray X1 at Oak Ridge National Laboratory. We present the technical details of the old and new internal data structures, the required code reorganization, and the resulting performance improvements. We describe and compare the performance and scaling of the final CLM Version 3.0 (CLM3.0) on the IBM Power4, the Earth Simulator, and the Cray X1. At 64 processors, the performance of the model is similar on the IBM Power4, the Earth Simulator, and the Cray X1. However, the Cray X1 offers the best performance of all three platforms tested from 4 to 64 processors when OpenMP is used. Moreover, at low processor counts (16 or fewer), the model performs significantly better on the Cray X1 than on the other platforms. The vectorized version of CLM was publicly released by the National Center for Atmospheric Research as the standalone CLM3.0, as a part of the new Community Atmosphere Model Version 3.0 (CAM3.0), and as a component of the Community Climate System Model Version 3.0 (CCSM3.0) on June 23, 2004.

Key words: Community Land Model, CLM, Community Climate System Model, CCSM, vectorization, Cray X1, Earth Simulator

## 1 Introduction

### 1.1 THE MODEL

The Community Land Model (CLM) is a single-column (snow–soil–vegetation) biogeophysical model of the land surface. Written in Fortran 90, CLM can be run off-line (i.e. run in isolation using stored atmospheric forcing data), coupled to an atmospheric model (e.g. the Community Atmosphere Model (CAM)), or coupled to a climate system model (e.g. the Community Climate System Model (CCSM)) through a flux coupler (e.g. Coupler 6). When coupled, CLM exchanges fluxes of energy, water, and momentum with the atmosphere.

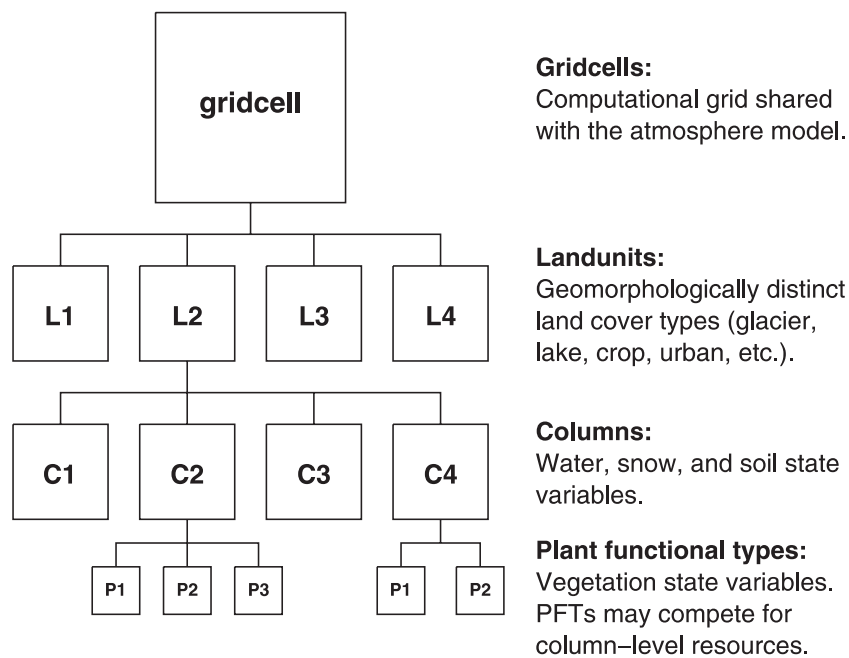
Originally developed as the standalone Common Land Model (Dai et al., 2003), the code was significantly modified for integration with the CAM (Bonan et al., 2002) and the CCSM (Kiehl and Gent, 2004) in 2001 and 2002. This prior development work targeted cache-based scalar multiprocessor computer platforms. The availability of the Earth Simulator in Japan and the Cray X1 at Oak Ridge National Laboratory (ORNL) spawned renewed interest in running Earth systems models (including the CLM) on vector architectures. However, the prior development of the CLM for multiprocessor architectures had inadvertently resulted in code which would not vectorize. This led to the need to rewrite the model to provide good performance on both vector and scalar architectures in 2003.

After the substantial development and testing described below, the vectorized version of the model was publicly released by the National Center for Atmospheric Research (NCAR) as the standalone Community Land Model Version 3.0 (CLM3.0), as a part of the new Community Atmosphere Model Version 3.0 (CAM3.0), and as a component of the Community Climate System Model Version 3.0 (CCSM3.0) on June 23, 2004. The CLM3.0 Developer's Guide (Hoffman et al., 2004), the CLM3.0 User's Guide (Vertenstein et al., 2004), the Technical Description of the CLM (Oleson et al., 2004), and the CLM's Dynamic Global Vegetation Model (DGVM): Technical Description and User's Guide (Levis et al., 2004) provide the developer, user, or researcher with details of implementation, instructions for using the model, a scientific description of the model, and a scientific description of the DGVM integrated with the CLM, respectively.

<sup>1</sup>OAK RIDGE NATIONAL LABORATORY CLIMATE AND CARBON RESEARCH INSTITUTE OAK RIDGE, TENNESSEE 37831–6008, USA (FORREST@CLIMATE.ORN.LGOV)

<sup>2</sup>NATIONAL CENTER FOR ATMOSPHERIC RESEARCH BOULDER, CO, USA

<sup>3</sup>CENTRAL RESEARCH INSTITUTE OF ELECTRIC POWER INDUSTRY, JAPAN



**Fig. 1** The CLM subgrid hierarchy.

Several multidecadal simulations have been carried out with CCSM3.0, and the system will be a major contributor to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change (IPCC). Simulation output is available from the Earth System Grid at <http://www.earthsystemgrid.org/>. An analysis of the results of these simulations emphasizing the contribution of CLM3.0 will appear with related papers in a special issue of the *Journal of Climate* (Dickinson et al., 2005).

## 1.2 MODEL GRID HIERARCHY

The horizontal land surface heterogeneity in the CLM is represented by a nested subgrid hierarchy composed of gridcells, landunits, columns, and plant functional types (PFTs) as shown in Figure 1. Each gridcell can have a different number of landunits, each landunit can have a different number of columns, and each column can have multiple PFTs. Gridcells represent the computational grid which is shared with the atmospheric physics.

The landunit, the first subgrid level, is intended to capture the broadest spatial pattern of subgrid heterogeneity. It serves primarily to distinguish physical soil properties. Specific landunits include glacier, lake, wetland, urban,

and vegetated. The column captures variability in soil and snow state variables within a landunit. Water and energy states and fluxes are tracked at the column level. The PFT, the third subgrid level, captures the characteristic biophysical and biogeochemical functions of broad categories of vegetation and bare soil. Up to four out of 15 possible PFTs differing in physiology and structure may be contained within a single column.

Biophysical processes are simulated for each subgrid unit independently, and prognostic variables are maintained for each subgrid unit. Processes related to soil and snow require PFT-level properties to be aggregated up to the column level. Aggregation is usually accomplished by computing a weighted sum for each quantity over all PFTs within a column. Similarly, different PFTs compete for the resources tracked at the column level. A complete description of the biophysical processes simulated by the CLM is available in Oleson et al. (2004).

## 2 Model Vectorization

The primary obstacles to vectorization were the layout of the internal data structures and the organization of the processing loops in CLM2.1. The hierarchical subgrid

organization is reflected in the data structures used in the model code. In CLM2.1, the hierarchy was implemented as arrays of pointers to derived data types at each subgrid level. The flux and state variables were implemented as scalars in every instance of a derived data type. While this object-oriented layout is reasonable for cache-based scalar platforms, it is not conducive to vector processing since it results in large, unpredictable strides among variable elements. In addition, the model had many high-level loops over various grid and subgrid units. In each loop, science subroutines were called separately and repeatedly for each grid or subgrid unit member. The loops within the science subroutines were all short and contained negligible work.

Early vectorization experiments, carried out separately by coauthors Kitabata at CRIEPI on the Earth Simulator and White at ORNL on the Cray X1, demonstrated significant performance improvement by eliminating the hierarchy of derived data types in favor of traditional Fortran arrays and by moving long grid and subgrid unit loops “down” into science subroutines where they were then typically interchanged with short loops (White, 2003). The Earth Simulator has no cache but has very long vector pipelines (eight 2304-element pipelines per arithmetic processor); the Cray X1 has a cache and shorter vector pipelines (four 256-element pipelines per MSP). As a result, memory contiguity and long vector lengths are more important on the Earth Simulator than on the Cray X1.

While these experiments were being carried out, a strategy for rewriting the CLM code was being developed to meet the requirements put forth by the research community. Researchers required a single CLM source code that would run well on both vector and scalar architectures while maintaining the hierarchical nature of the internal data structures. The use of conditionally compiled code blocks (using `#ifdef` preprocessor macros) was to be minimized, and the code modifications could not reduce the performance on the existing scalar platforms.

## 2.1 DATA STRUCTURES

After considerable discussion and experimentation on the Cray X1, a new data structure organization was developed which would retain the subgrid hierarchy while allowing for loop vectorization. New data structures were prototyped in the model by adding them in as vector versions of individual code branches were added to the code. A testing methodology (developed by White as a part of his vectorization experiments) was adopted for validating simulation results at every time-step. During developmental test runs, new data structures were initialized from old data structures, the original subroutine was called to update the old data structures, the new vector-friendly subroutine was called to update the new data

structures, and the results in the two sets of data structures were then compared to ensure only round-off differences resulted from the new code branch. The old data structures and subroutines were removed from the code as vector development neared completion.

To begin, a vectorized version of the most costly code branch, the `Biogeophysics1()` routine, which calls the computationally intensive `CanopyFluxes()` subroutine, was added to the code. In this experiment, data were copied from the original data structures into the new data structures, and then the original `Biogeophysics1()` subroutine was called, followed by the vectorized version of the subroutine, which used the new data structures. For each time-step in the model run, after the non-vector and vector subroutines were called, the results contained in the old and new data structures were compared to ensure that only round-off differences existed between the two sets of results.

Timing utilities included in the CLM code were used to measure performance differences between the original and the new vector versions of the code branch. Using the proposed hierarchical data structures, the performance improvement matched that of prior vectorization experiments (White, 2003). The strategy for further vectorization was to use the same testing methodology for each code branch in turn to ensure that model results were equivalent at each time-step and to monitor performance improvement. As vectorization progressed, the code was tested by Kitabata on the Earth Simulator and by others on NEC platforms to ensure similar performance gains on these systems. The new data structures and vectorization strategy were subsequently presented to and approved by the CCSM Land Model Working Group and NCAR researchers in the summer of 2003.

The new data structures now in CLM3.0 consist of derived data types for each subgrid unit as shown in Figure 2. Each grid and subgrid unit data type in the hierarchy contains a number of physical and chemical state and flux data types. These data types typically contain real arrays for the state and flux variables which were previously represented as scalars in multiple instances of data types.

Vertical heterogeneity is represented by a single vegetation layer, 10 layers for soil, and up to five layers for snow, depending on the snow depth. Multilayer quantities are stored as two-dimensional real arrays. Using real arrays at every level in the hierarchy maximizes opportunities for contiguous memory access, thereby providing significantly better performance on vector architectures. All arrays contained in derived data types are implemented as Fortran 90 pointers. Memory for these arrays is allocated dynamically during model initialization.

Each of these grid and subgrid data types also contains arrays of integers, which serve as array indices to the higher subgrid levels or as initial and final bounds on the

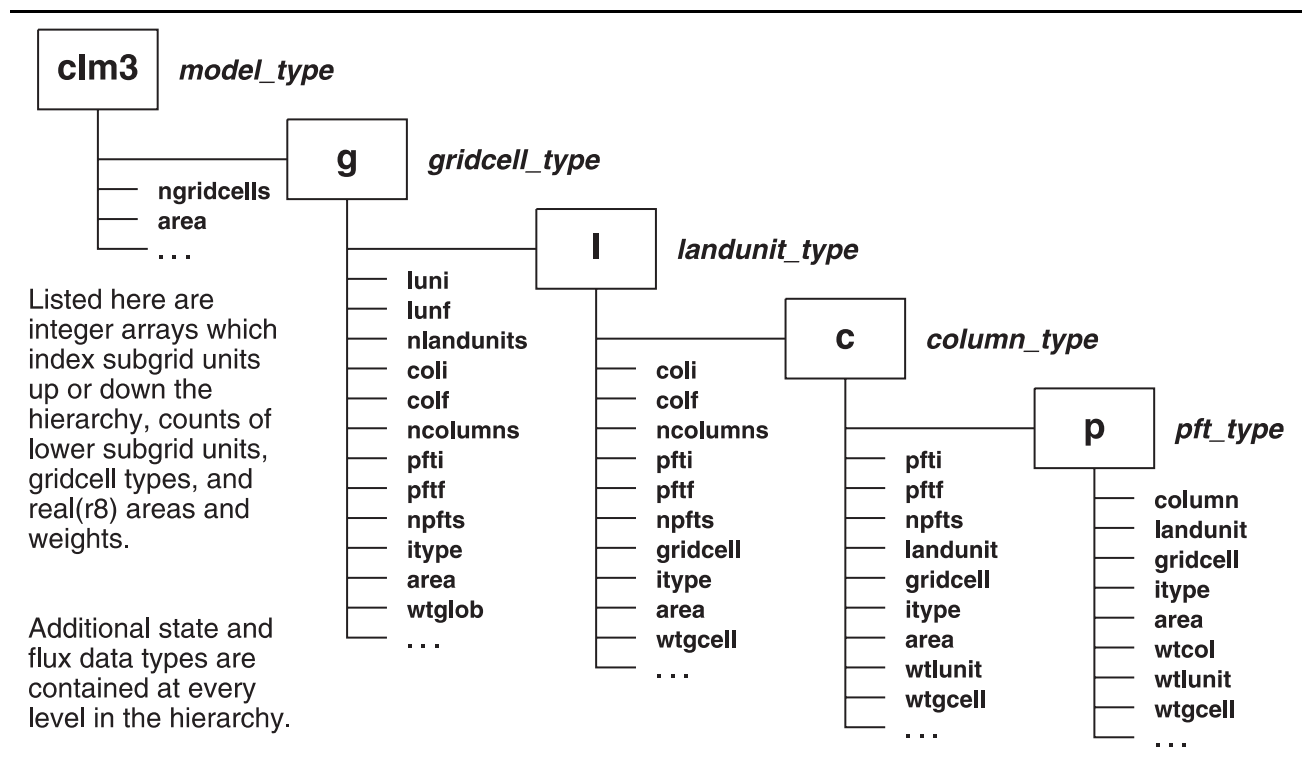


Fig. 2 The new CLM3.0 data structure hierarchy.

lower subgrid levels. For example, as shown in Figure 2, the column-level data type contains `landunit` and `gridcell` integer arrays, which refer to the appropriate `landunit` and `gridcell` for a given column. The column-level data type also contains integer arrays called `pfti`, `pftf`, and `npfts`, which refer to the first and last PFTs and the total number of PFTs, respectively, for a given column. Each subgrid data type also contains real arrays for surface areas and area weights at all higher grid levels.

**2.1.1 Clumps for Decomposition** CLM uses the Message Passing Interface (MPI) for distributed memory parallelism and uses OpenMP for shared memory parallelism. On the Cray X1, the model can also stream across processing units (called Single Streaming Processors, or SSPs) within a Multistreaming Processor (MSP) based on Cray Streaming Directives (CSDs) included in the new code. In order to provide performance portability across a wide range of current and future computer architectures, a decomposition strategy partially implemented in CLM2.1 was fully developed in CLM3.0.

When compiled for distributed memory parallelism (with the preprocessor macro `SPMD` defined), each MPI process will create an instance of the data structures shown in Figure 2 containing only the subset of data assigned to that process. A cache-friendly blocking structure is superimposed on the data structure hierarchy for improved computational efficiency. This blocking structure implicitly controls the vector length of most computations. Gridcells are grouped into blocks (called clumps) of nearly equal computational cost, and these clumps are subsequently assigned to MPI processes.

The computational cost of a `gridcell` is approximately proportional to the number of PFTs contained within it. However, since computational cost for some PFTs is higher than for others and since similar PFTs tend to cluster geographically, balancing the workload across MPI processes requires a more complex scheme than simply assigning contiguous blocks of `gridcells` to clumps. To minimize the potential for load imbalance, `gridcells` are assigned in cyclic (or round-robin) fashion to a predetermined number of clumps. The clumps are then assigned

```

nclumps = get_proc_clumps()
!$OMP PARALLEL DO PRIVATE (nc,begg,endg,begl,endl,begc,endc,begp,endp)
!CSD$ PARALLEL DO PRIVATE (nc,begg,endg,begl,endl,begc,endc,begp,endp)
do nc = 1,nclumps
  call get_clump_bounds(nc, begg, endg, begl, endl, begc, endc, begp, endp)
  .
  .
  .
  call Hydrology1(begc, endc, begp, endp, &
    filter(nc)%num_nolakec, filter(nc)%nolakec, &
    filter(nc)%num_nolakep, filter(nc)%nolakep)
  .
  .
  .
end do
!$OMP END PARALLEL DO
!CSD$ END PARALLEL DO

```

**Fig. 3** Example high-level loop in the `driver` routine.

in cyclic fashion to available MPI processes. This scheme has proved to sufficiently distribute gridcells of various computational costs among MPI processes, yielding very good parallel load balancing characteristics for most process counts and surface data sets. A small penalty is incurred during history output since the data in memory must be remapped to the standard south-to-north ordering for output.

Clumps not only define the workload for an MPI process, they also serve to block data for shared memory parallelism when using OpenMP or streaming on the Cray X1. The number of clumps per MPI process is determined by the parallel configuration of the model at run-time, but it may be set explicitly by setting the `clump_pproc` namelist variable to the desired number of clumps per process. When run serially or with MPI-only parallelism, one clump per process is used. When OpenMP is enabled, the number of clumps per process is set to the maximum number of OpenMP threads available. On the Cray X1 when OpenMP is disabled, CSDs are interpreted by the compiler in place of OpenMP directives, and the number of clumps per process is set to four to take maximum advantage of the four SSP units on an MSP.

**2.1.2 Filters for Vector Performance** In addition to clumps, another set of structures, called “filters”, was added to better support vectorized processing of columns and PFTs. Filters group like columns or PFTs based on their process-specific categorization and are used for indirect addressing into the main data structure hierarchy.

Filters are created for snow, non-snow, lake, non-lake, and bare soil columns and PFTs for each clump of gridcells. Many filters are initialized once, but the snow and non-snow filters must be reconstructed as snowfall and melting occur, and the vegetated PFT filter must be reconstructed in DGVM as die-off and establishment occur.

## 2.2 CODE REORGANIZATION

Loops over columns and PFTs previously located in the top-level `driver` routine were moved down into science subroutines to provide opportunities for vectorization. In CLM3.0, the highest level loops in the `driver` routine run over clumps for each MPI process and provide for OpenMP or Cray Streaming parallelism. Science subroutines, called within these loops, are passed local clump bounds for gridcells, landunits, columns, and PFTs as needed. Relevant filters, in the form of counts and vectors of array indices, are also passed as needed to science subroutines.

Figure 3 shows a portion of a high-level loop from the `driver` routine. First, the number of clumps assigned to the process is obtained and stored in `nclumps`. The subsequent loop over all clumps is wrapped with OpenMP and Cray Streaming directives to support shared memory parallelism. Within the loop, the bounds for gridcells, landunits, columns, and PFTs are obtained for the clump being processed by calling `get_clump_bounds()`. Then a science subroutine, `Hydrology1()`, is called and passed the column and PFT bounds as well as the

```

! Assign local pointers to derived type members (landunit-level)
clandunit => clm3%g%l%c%landunit
itype     => clm3%g%l%itype
! Assign local pointers to derived type members (column-level)
cgridcell => clm3%g%l%c%gridcell
t_grnd    => clm3%g%l%c%ces%t_grnd
h2osno    => clm3%g%l%c%cws%h2osno
snowdp    => clm3%g%l%c%cps%snowdp
snowage   => clm3%g%l%c%cps%snowage
!dir$ concurrent
!cdir nodep
do f = 1, num_nolakec
  c = filter_nolakec(f)
  l = clandunit(c)
  g = cgridcell(c)
  .
  .
  .
  if (itype(l) == istwet .and. t_grnd(c) > tfrz) then
    h2osno(c) = 0._r8
    snowdp(c) = 0._r8
    snowage(c) = 0._r8
  end if
  .
  .
  .
end do

```

**Fig. 4** Example filter loop in a science subroutine.

non-lake filters for columns and PFTs for the clump being processed. Additional science subroutines are subsequently called within the same loop. The `driver` routine consists primarily of two such high-level loops, which call most of the science subroutines used by the model.

Within science subroutines, vector loops run over grid or subgrid units. Other short loops run over snow and soil levels within a column or PFT. In most cases, the vector loops were inserted into the short loops for vectorization. Many loops were split into multiple loops and temporary local arrays (vectors) were used to “carry” data between them. Many of the vector loops use filters for indirect addressing of relevant columns and PFTs. Since arrays in data structures are implemented as pointers, compilers cannot determine if vector dependences exist. As a result, compiler directives are required in order to obtain loop vectorization in most cases on both the Cray X1 and the Earth Simulator.

Figure 4 shows an example of a filter loop within a science subroutine. First, local pointers are created to shorten the notation used in equations. The subsequent

loop over all non-lake columns is preceded by Cray X1 and Earth Simulator compiler directives. The first directive tells the Cray X1 compiler that the loop is concurrent, meaning it may be streamed and vectorized. The second directive tells the Earth Simulator compiler that no vector dependences exist in the loop. Within the loop, the column index is obtained from the non-lake column filter vector, the appropriate landunit index is obtained from the column’s landunit vector, and the appropriate gridcell index is obtained from the column’s gridcell vector. Next, the landunit type and the ground temperature of the column are checked. If the landunit contains water and the ground temperature is above freezing, three variables are initialized to zero. Other computations are usually performed within such loops.

### 3 Vector Performance

Preliminary vectorization of the CLM was completed in October 2003. The new vectorized model has a smaller memory footprint than CLM2.1. The new data structures

simplify history updates and reduce the complexity and number of MPI gathers and scatters. The new vectorized model runs 25.8 times faster than the CLM2.1 code on the Cray X1 and even runs 1.8 times faster on the IBM Power4 in standalone mode. To gauge overall performance of the new CLM3.0 and determine optimum run-time configurations, timing tests were performed on the IBM Power4, the Earth Simulator, and the Cray X1. The timing tests consisted of a series of 30-day off-line runs at T85 resolution (approximately 1.5 degrees) varying both processor counts (with and without OpenMP) and the clumps-per-process tuning parameter. The source code used for the tests was that tagged `clm2_deva_51` in the NCAR code repository.

CLM3.0 may be built with a variety of options depending on the computer architecture and processor configuration to be used. Both MPI and OpenMP may be independently enabled or disabled. With both MPI and OpenMP disabled, the model can run serially on a single processor. With MPI enabled and OpenMP disabled, the model runs in distributed memory mode across a number of nodes. With MPI disabled and OpenMP enabled, the model can run on a single shared memory symmetric multiprocessor (SMP) node. With both MPI and OpenMP enabled (hybrid mode), CLM3.0 runs in hybrid distributed/shared memory mode across a number of SMP nodes. The optimum configuration for performance depends on the computer architecture. In general, a hybrid mode configuration provides the best performance on parallel systems with shared memory processors.

### 3.1 CRAY X1

On the Cray X1, CLM3.0 can be built with CSDs enabled; however, most CSDs are used only around loops which also use OpenMP directives. As a result, OpenMP can be enabled only when CSDs are not. When compiling with OpenMP (ignoring CSDs), the compiler will still multistream concurrent loops in science subroutines where possible. Calls to these science subroutines are typically contained within the high-level loops in the `driver` routine where the OpenMP directives and CSDs are located. Slight modification of these loops makes possible simultaneous use of OpenMP and CSDs. The performance impacts of CSDs versus OpenMP on the Cray X1 are described below.

For the Cray X1 with OpenMP disabled, the number of MPI processes is equal to the number of MSPs used. Four threads were used for all tests with OpenMP enabled on the Cray X1. The curves in the performance graphs presented here are all colored by the total number of processors used in the run. For the Cray X1, a processor refers to one MSP. The clumps-per-process tuning parameter was varied between 1 and 32 for all performance tests.

Figure 5 shows total run-times for CLM3.0 on the Cray X1 with both OpenMP and CSDs disabled for 1–128 MSPs. The number of MPI processes is equal to the number of MSPs, except that MPI was disabled for the 1 MSP case. For this test, streaming across SSPs occurs only in concurrent loops within science subroutines. The best performance is obtained when clumps per process is set to one, since this setting maximizes the vector length. Scaling beyond 16 MSPs is poor, but run-time continues to drop through 64 MSPs. Performance for 128 MSPs is no better than for 64 MSPs.

Figure 6 shows total run-times for CLM3.0 on the Cray X1 with OpenMP disabled and CSDs enabled for 1–128 MSPs. The number of MPI processes is equal to the number of MSPs, except that MPI was disabled for the 1 MSP case. For this test, streaming across SSPs occurs at the high-level loops in the `driver` routine. The best performance is obtained when clumps per process is set to four since this setting maximizes the vector length for each of the four SSPs. Scaling beyond 32 MSPs is poor, but run-time continues to drop through 64 MSPs. Performance for 128 MSPs is just slightly worse than for 64 MSPs.

Figure 7 shows total run-times for CLM3.0 on the Cray X1 with OpenMP enabled and CSDs disabled for 4–256 MSPs. Four threads were used per MPI process, and the number of MSPs shown is equal to the product of the number of MPI processes and the number of threads. MPI was disabled for the 4 MSP case. For this test, streaming across SSPs occurs only in concurrent loops within science subroutines. The best performance is obtained when clumps per process is set to four, since this setting maximizes the vector length for each of the four shared memory MSPs. Scaling beyond 64 MSPs is poor, but run-time continues to drop through 128 MSPs.

In all three cases, the best performance is obtained when vector lengths are maximized. On average, performance is improved by 20% using clumping and enabling CSDs around high-level loops to explicitly utilize the four SSPs on each MSP. However, better scaling is obtained by using OpenMP instead of CSDs for the same high-level loops in the `driver` routine. MPI performance and problem size limit distributed memory scaling to 64 MSPs, but adding shared memory processors via OpenMP improves scaling to 128 MSPs (i.e. 32 MPI processes and four OpenMP threads). It appears that the compiler does a good job of multistreaming loops in science subroutines (thereby keeping the SSPs sufficiently busy). This fact, combined with the reduced MPI communications, explains why using OpenMP to parallelize high-level loops results in better performance beyond eight MSPs.

In an effort to achieve additional performance on the Cray X1, the high-level loops in the `driver` routine were slightly modified so that OpenMP and Cray Streaming could be used simultaneously. The code modification

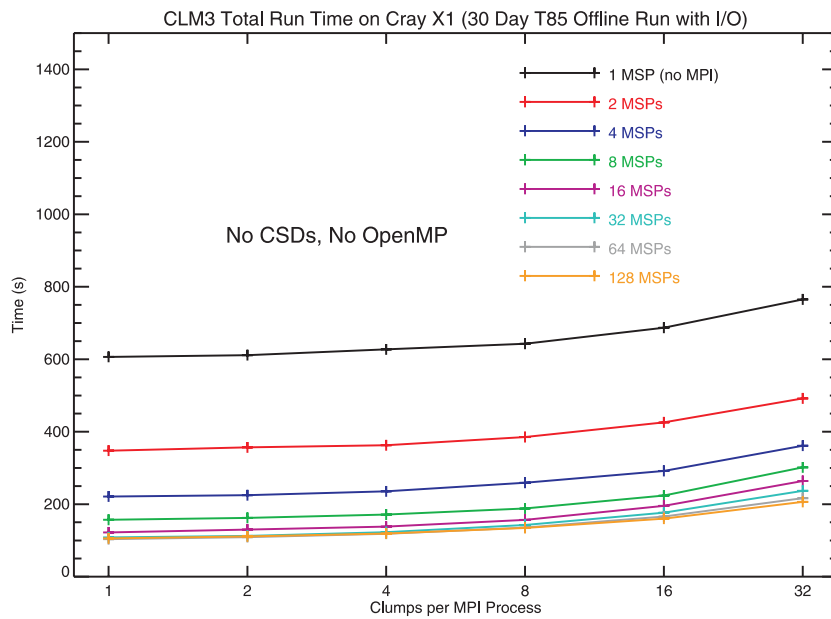


Fig. 5 Curves show total run-times for CLM3.0 on the Cray X1 with OpenMP and CSDs both disabled for 1–128 MSPs.

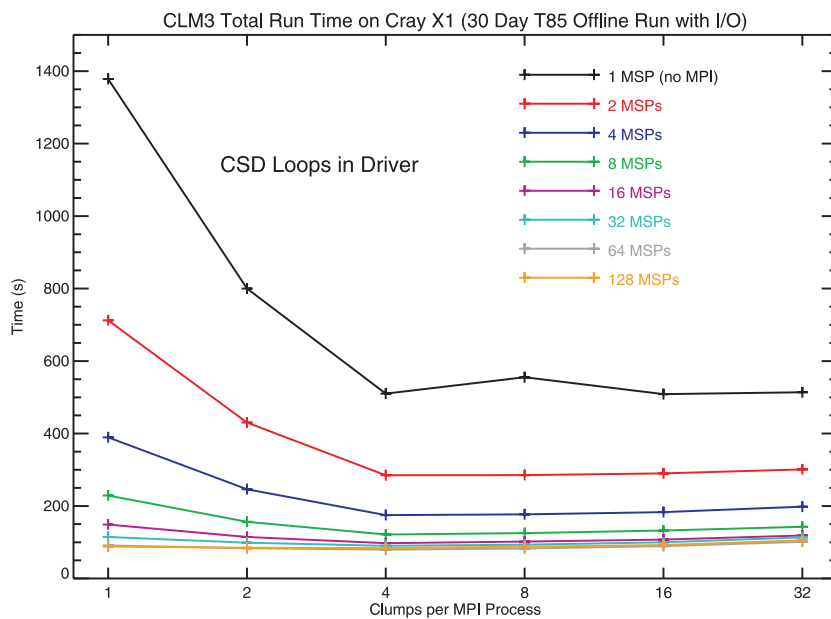
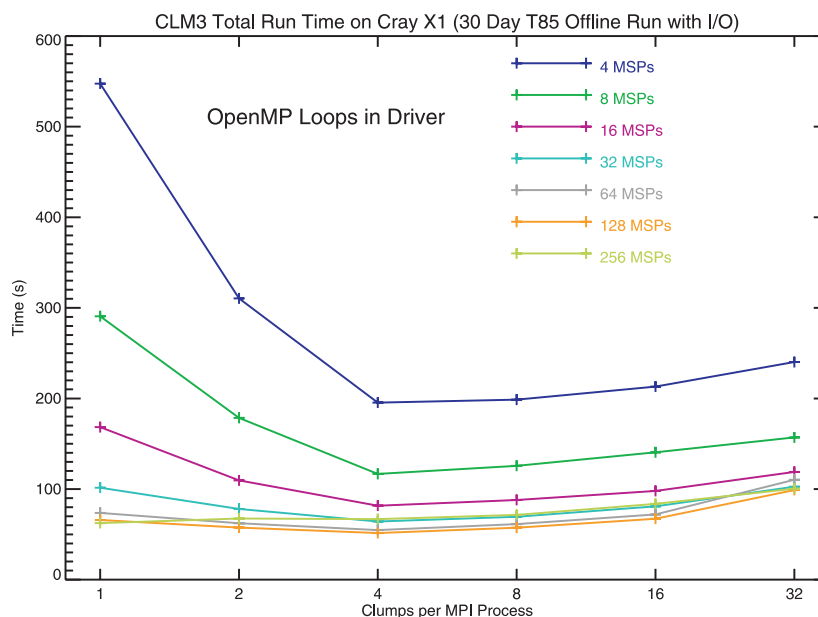


Fig. 6 Curves show total run-times for CLM3.0 on the Cray X1 using CSDs around high-level loops in the driver routine for 1–128 MSPs.





**Fig. 7** Curves show total run-times for CLM3.0 on the Cray X1 using OpenMP directives around high-level loops in the driver routine for 4–256 MSPs. Four threads were used per MPI process, so the number of processors shown is equal to the product of the number of MPI processes and the number of threads.

is shown in the listing in Figure 8. For this case, the optimum number of clumps per process is 16: four for the four OpenMP threads times four for the four SSPs.

As shown in Figure 9, this change slightly improves the performance of the OpenMP-only case up through 64 MSPs. Beyond 64 MSPs, the addition of CSD loops within OpenMP loops does not further improve model performance. The overall performance improvement from using OpenMP and CSDs simultaneously was smaller than expected because the vector lengths become too short (i.e. the problem size is too small) and the overhead of MPI communications limits further scaling. More importantly, this result is seen as confirmation that the automatic streaming of concurrent loops within science subroutines is very good. The SSPs are kept busy most of the time.

Profiling experiments were performed to further investigate performance characteristics of CLM3.0 on the Cray X1. The most time-consuming routines are `gettimeofday` (which is used by the timing utilities), an internal vector memory copy routine, and `canopyfluxes`. This routine is known to represent the majority of the calculations in the present version of CLM, and it has been extensively optimized. Additional routines near the top of

the profile listing are the `updateinput`, `areaovr_point`, `mkmxovr`, and `areaave` routines, which provide data exchange and interpolation functionality for the River Transport Model (RTM). The data exchange and interpolation functions, which have not been vectorized, are used repeatedly when running CLM3.0 in off-line mode; however, they are called only during initialization when run in fully coupled CCSM mode, so effort has not been spent optimizing these routines.

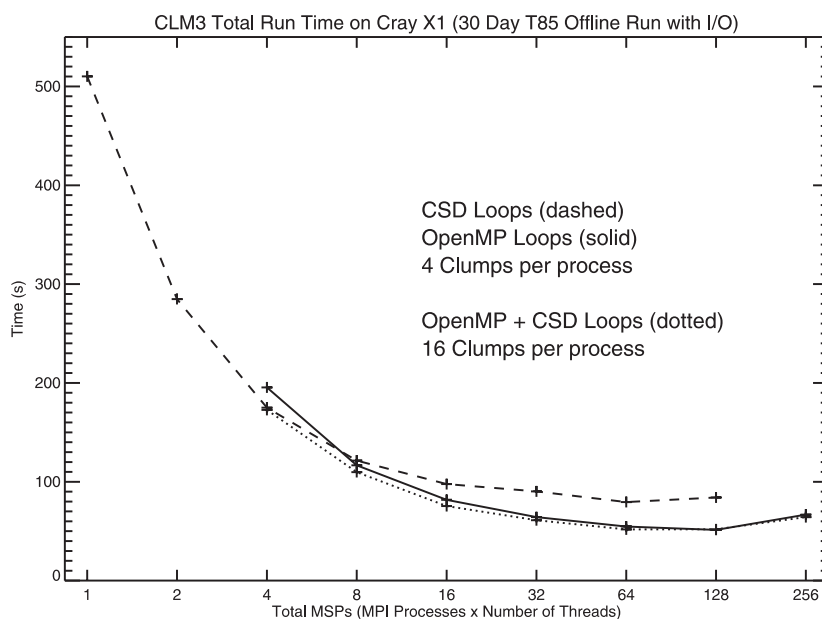
When run using two MSPs with CSDs enabled, MPI is used for communication, and `MPI_Gatherv` and `MPI_Bcast` show up at the seventh and eighth positions in the profiling report. A small portion of this time is attributable to synchronization (i.e. waiting for processes to catch up to the point in the program where the communication occurs); however, it is likely that MPI performance of the Cray X1 could be improved. CLM3.0 has very good load balancing, so synchronization should take very little time. `MPI_Gatherv` is used each time-step to generate a temperature diagnostic, and it is used to accumulate data for monthly history output. `MPI_Bcast` is used to distribute atmospheric forcing data read from disk by the master process to all processes on a monthly basis.

```

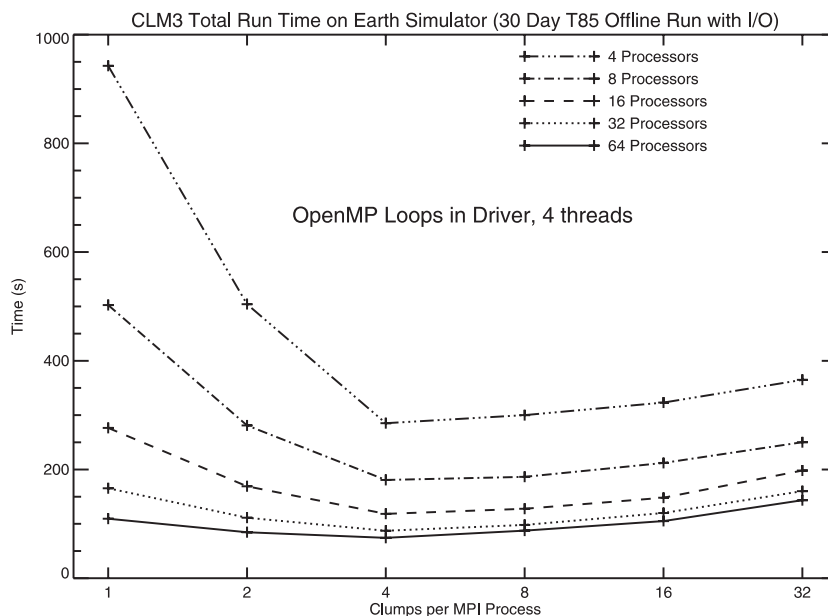
!$OMP PARALLEL DO PRIVATE (gnc,nc_beg,nc_end,nc,begg,endg,begl,endl, &
!$OMP & begc,endc,begp,endp)
do gnc = 1,nclumps/4
  nc_beg = (gnc - 1) * 4 + 1
  nc_end = min(nc_beg+3,nclumps)
!CSD$ PARALLEL DO PRIVATE (nc,begg,endg,begl,endl,begc,endc,begp,endp)
  do nc=nc_beg,nc_end
    call get_clump_bounds(nc, begg, endg, begl, endl, begc, endc, &
      begp, endp)
    .
    .
    call Hydrology1(begc, endc, begp, endp, &
      filter(nc)%num_nolakeec, filter(nc)%nolakeec, &
      filter(nc)%num_nolakeep, filter(nc)%nolakeep)
    .
    .
  end do
!CSD$ END PARALLEL DO
end do
!$OMP END PARALLEL DO

```

**Fig. 8** Modified clump loop from `driver` providing both OpenMP and Cray Streaming directives.



**Fig. 9** CLM3.0 total run-time on the Cray X1 using CSD loops (dashed), OpenMP loops (solid), and OpenMP with CSDs simultaneously (dotted).



**Fig. 10** Curves show total run-times for CLM3.0 on the Earth Simulator using OpenMP directives around high-level loops in the `driver` routine for 4–64 processors. Four threads were used per MPI process, so the number of processors shown is equal to the product of the number of MPI processes and the number of threads.

When run on 32 MSPs, `MPI_Bcast` and `MPI_Gatherv` are the second and third most expensive routines. Only `gettimeofday` represents more time samples. `MPI_Allgatherv`, used after each run of RTM, appears at number 10. With 32 MPI processes, MPI routines and interpolation for RTM dominate run-time. While some additional time may be lost to a slightly larger load imbalance than when using two processes, at 32 MPI processes communication time appears to exceed calculation time.

Because of the inefficiency of `gettimeofday` on the Cray X1, it was replaced by calls to the real-time clock routine (`rtc`) in later versions of the timing utilities used by CLM3.0. Nevertheless, disabling model timers on the Cray X1 results in a 10–15% performance improvement in overall run-time. Since they are not necessary, it is recommended that the timers be disabled for normal production runs.

### 3.2 EARTH SIMULATOR

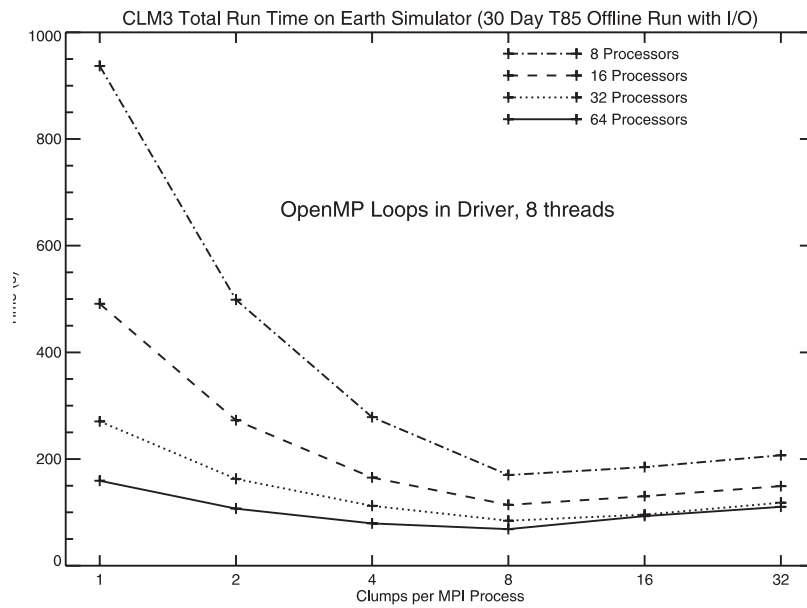
Similar timing experiments were carried out on the Earth Simulator. Since the Earth Simulator machine does not

use multistreaming processors, only OpenMP configurations were tested. One set of tests used four OpenMP threads, while the other used eight OpenMP threads. Like the timing tests on the Cray X1, the clumps-per-process parameter was varied from 1 to 32. The results from these timing tests are shown in Figures 10 and 11. The number of processors shown is equal to the product of the number of MPI processes and the number of threads.

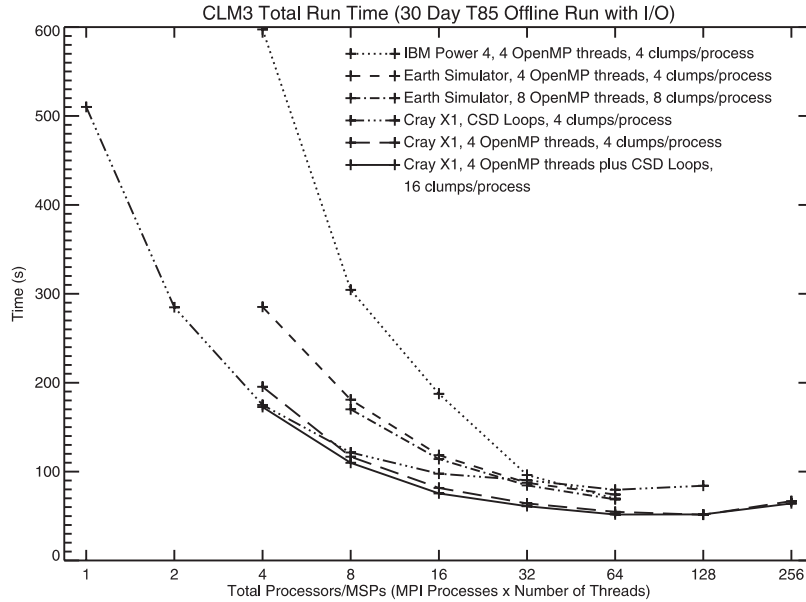
As expected, choosing four clumps per process provides the best performance for the four-thread case while eight clumps per process provides the best performance for the eight-thread case. These choices maximize vector lengths for each configuration. On the Earth Simulator, as on the Cray X1, the model scales only to about 64 processors. Using eight OpenMP threads yields slightly better performance than using four threads.

### 3.3 PERFORMANCE COMPARISON

Figure 12 shows a comparison of the best CLM3.0 total run-times for the IBM Power4 using four OpenMP threads; the Earth Simulator using both four and eight OpenMP threads; and the Cray X1 using CSDs, OpenMP,



**Fig. 11** Curves show total run-times for CLM3.0 on the Earth Simulator using OpenMP directives around high-level loops in the driver routine for 8–64 processors. Eight threads were used per MPI process, so the number of processors shown is equal to the product of the number of MPI processes and the number of threads.



**Fig. 12** Curves show performance of the CLM3.0 run in off-line mode on the IBM Power4, the Earth Simulator, and the Cray X1.

and both CSDs and OpenMP simultaneously. The curve for the IBM Power4 shows very good scaling to 32 processors (within a single node) and reasonable scaling to 64 processors. The IBM matches the best Earth Simulator performance and outperforms the Cray X1 when using only CSDs at 64 processors.

The shapes of the Earth Simulator and Cray X1 curves are similar, so the model appears to scale in a similar fashion on both machines. However, the Cray X1 performance numbers are usually better than those from the Earth Simulator. When using only CSDs around the driver loops, the Cray X1 outperforms the Earth Simulator out to 32 processors. At 64 processors, the Earth Simulator does slightly better. On the other hand, when OpenMP is used on the Cray X1 (with or without CSDs), it always outperforms the Earth Simulator on a per processor basis. Moreover, at low processor counts (16 or fewer), the model performs significantly better on the Cray X1 than on the other two platforms.

#### **4 Conclusions**

Vectorizing the CLM required significant changes to the internal data structures and every science subroutine in the model. The process of rewriting the code took about six months (or one man-year) to complete, and the planning and related development activities impacted model developers for an entire year. However, the resulting model code performs better on both vector and scalar architectures. These performance improvements will result in more and better research into land-surface processes and feedbacks as access to vector platforms widens.

The vectorized version of the CLM was publicly released by the NCAR as the standalone CLM3.0, as a part of the new CAM3.0, and as a component of the CCSM3.0 on June 23, 2004. Model development has continued since the release with frequent testing on the Cray X1 platform.

In general, vectorization should be considered during all model development to ensure performance portability across myriad computer platforms. As demonstrated here, internal data structures can often be designed to meet the goals of researchers while still providing acceptable good performance on both vector and scalar architectures. Writing loops so that they will vectorize is worth the up-front effort since it ensures good performance on a wide variety of systems. Even today's personal computers have small vector units which can be harnessed for high performance computational science. Guidelines for writing vector-efficient code can be found in various software development manuals including Cray's manual for optimizing applications (Cray Inc., 2004).

As development continues on the CLM, vectorization should remain a goal. Dynamic land use capabilities

recently added to an experimental version of the model required that history buffers be updated every time-step instead of whenever history output is written (usually monthly). This model change exposed some accumulation routines which were only partially vectorized. Since these routines were not previously called every time-step, effort had not been invested in vectorizing them. Model run-time on the Cray X1 more than doubled. Vectorizing these routines reclaimed almost the entire factor of 2 in performance. Frequent performance testing on vector platforms, in addition to the standard numerical testing, should help guard against frequent use or introduction of non-vectorizable code in the model.

The vector performance of the CLM3.0 is reasonable even when run in off-line mode. When coupled to the CAM and the CCSM, better performance is realized since atmospheric data are passed directly to the land model instead of being read from disk and broadcast to all MPI processes. At 64 processors, the performance of the model is similar on the IBM Power4, the Earth Simulator, and the Cray X1. However, the Cray X1 offers the best performance of all three platforms tested from 4 to 64 processors when OpenMP is used. Moreover, at low processor counts (16 or fewer), the model performs significantly better on the Cray X1 than on the other platforms.

#### **ACKNOWLEDGMENTS**

The authors wish to gratefully acknowledge the assistance of Matthew Cordery at Cray Inc., and Dave Parks and John Snyder at NEC Corp. This research was supported by the Office of Biological and Environmental Research of the U.S. Department of Energy (DOE) and by the National Science Foundation (NSF). This research used the resources of the Center for Computational Sciences at ORNL, which is managed by UT-Battelle, LLC, and is supported by the Office of Science of the U.S. DOE under Contract No. DE-AC05-00OR22725. The timings on the Earth Simulator were performed with support from the Ministry of Education, Culture, Sports, Science and Technology of the Japanese Government through the Project for Sustainable Coexistence of Human, Nature and the Earth. Additional resources were provided by the NCAR, the Central Research Institute of Electric Power Industry (CRIEPI) in Japan, Cray Inc., and NEC Corp.

#### **AUTHOR BIOGRAPHIES**

*Forrest Hoffman* is a researcher at ORNL where he holds joint appointments in the Computer Science and Mathematics and the Environmental Sciences Divisions. He received Masters and Bachelors of Science degrees in physics from the University of Tennessee. At ORNL,

Forrest established the first World Wide Web site at the Laboratory in 1995 and built ORNL's first Beowulf-style parallel computer, called The Stone SouperComputer, in 1997. A member of Sigma Pi Sigma, the American Geophysical Union, the IEEE Computer Society, and the Smoky Mountain Chapter of the American Meteorological Society, Forrest performs computational science research in global climate, landscape ecology, and terrestrial biogeochemistry on Linux clusters as well as some of the world's largest supercomputers in the National Center for Computational Sciences at ORNL. Forrest writes a monthly column for Linux Magazine called "Extreme Linux", and presently serves on the Advisory Committee for Advanced Research Computing (ARC) at Georgetown University.

*Mariana Vertenstein* has worked at NCAR since 1989 and is currently the head of the CCSM Software Engineering Group. She coordinates CCSM model development, testing, quality control and software releases. She also coordinates the interaction of CCSM, from a software perspective, with external collaborators, such as the DOE SciDAC project, the NASA Earth System Modeling Framework, the National Oceanic and Atmospheric Administration (NOAA) WRF group, the software engineers at the Earth Simulator Center and the ORNL. She is a co-chair of the CCSM Software Engineering Working Group. She is also the CCSM software engineering liaison to the CCSM Land Model Working Group. She is responsible for all software development of the CLM and for its release. She received a Ph.D. in chemical physics from Harvard University in 1987.

*Hideyuki Kitabata* is a visiting researcher in the Climate Group at the Central Research Institute of Electric Power Industry in Chiba-ken Japan where he performs climate modeling on the Earth Simulator. Hideyuki received his Ph.D. in fusion science in 1996 from the National Institute for Fusion Science at the Graduate University for Advanced Studies in Japan.

*James B. White III* (Trey) is an R&D staff member in the National Center for Computational Sciences at

ORNL. His research interests include enabling new science through high-end computation and improving software development for computational science. White has a Masters of Science degree in physics from Ohio State University.

## References

- Bonan, G. B., Oleson, K. W., Vertenstein, M., Levis, S., Zeng, X., Dai, Y., Dickinson, R. E., and Yang, Z.-L. 2002. The land surface climatology of the Community Land Model coupled to the NCAR Community Climate Model. *Journal of Climate* 15:3123–3149.
- Cray Inc. 2004. *Optimizing Applications on the Cray X1 System*. Cray Inc., S–2315–52 edition.
- Dai, Y. et al. 2003. The Common Land Model. *Bulletin of the American Meteorological Society* 84(8):1013–1023.
- Dickinson, R. E., Oleson, K., Bonan, G., Hoffman, F., Thornton, P., Vertenstein, M., Yang, Z.-L., and Zeng, X. 2005. The Community Land Model and its climate statistics as a component of the Community Climate System Model. *Journal of Climate*, in press.
- Hoffman, F., Vertenstein, M., Thornton, P., Oleson, K., and Levis, S. 2004. Community Land Model Version 3.0 (CLM3.0) Developer's Guide. Technical Memorandum ORNL/TM–2004/119, ORNL, Oak Ridge, TN, (<http://climate.ornl.gov/~forrest/pubs/CLM3DevGuideAndReference.pdf>)
- Kiehl, J. T. and Gent, P. R. 2004. The Community Climate System Model, version two. *Journal of Climate* 17:3666–3682.
- Levis, S., Bonan, G. B., Vertenstein, M., and Oleson, K. W. 2004. The Community Land Model's Dynamic Global Vegetation Model (CLM-DGVM): Technical description and user's guide. Technical Note NCAR/TN–459+IA, NCAR, Boulder, CO.
- Oleson, K. et al. 2004. Technical Description of the Community Land Model (CLM). Technical Note NCAR/TN–461+STR, NCAR, Boulder, CO.
- Vertenstein, M., Hoffman, F., Oleson, K., and Levis, S. 2004. Community Land Model version 3.0 (CLM3.0) user's guide. Technical report, NCAR, Boulder, CO.
- White, J. B. 2003. An optimization experiment with the Community Land Model on the Cray X1. *Proceedings of the 2003 Cray Users Group (CUG) Meeting*, Columbus, OH.