

## **The impact of SciDAC on US climate change research and the IPCC AR4.**

Michael Wehner  
Scientific Computing Group, MS-50F  
Ernest Orlando Lawrence Berkeley National Laboratory  
1 Cyclotron Rd. Berkeley, CA 94720 USA  
mfwehner@lbl.gov

### **Abstract**

SciDAC has invested heavily in climate change research. We offer a candid opinion as to the impact of the DOE laboratories' SciDAC projects on the upcoming Fourth Assessment Report of the Intergovernmental Panel on Climate Change.

As a result of the direct importance of climate change to society, climate change research is highly coordinated at the international level. The Intergovernmental Panel on Climate Change (IPCC) is charged with providing regular reports on the state of climate change research to government policymakers. These reports are the product of thousands of scientists' efforts. A series of reviews involving both scientists and policymakers make them among the most reviewed documents produced in any scientific field. The high profile of these reports acts a driver to many researchers in the climate sciences. The Fourth Assessment Report (AR4) is scheduled to be released in 2007. SciDAC sponsored research has enabled the United States climate modeling community to make significant contributions to this report.

Two large multi-Laboratory SciDAC projects are directly relevant to the activities of the IPCC. The first, entitled "Collaborative Design and Development of the Community Climate System Model for Terascale Computers", has made important software contributions to the recently released third version of the Community Climate System Model (CCSM3.0) developed at the National Center for Atmospheric Research. This is a multi-institutional project involving Los Alamos National Laboratory, Oak Ridge National Laboratory, Lawrence Berkeley National Laboratory, Pacific Northwest National Laboratory, Argonne National Laboratory, Lawrence Livermore National Laboratory and the National Center for Atmospheric Research. The original principal investigators were Robert Malone and John B. Drake. The current principal investigators are Phil Jones and John B. Drake. The second project, entitled "Earth System Grid II: Turning Climate Datasets into Community Resources" aims to facilitate the distribution of the copious amounts of data produced by coupled climate model integrations to the general scientific community. This is also a multi-institutional project involving Argonne National Laboratory, Oak Ridge National Laboratory, Lawrence Berkeley National Laboratory, Lawrence Livermore National Laboratory and the National Center for Atmospheric Research. The principal investigators are Ian Foster, Don Middleton and Dean Williams.

Perhaps most significant among the activities of the "Collaborative Design" project was the development of an efficient multi-processor coupling package. CCSM3.0 is an extraordinarily complicated physics code. The fully coupled model consists of separate submodels of the atmosphere, ocean, sea ice and land. In addition, comprehensive biogeochemistry and atmospheric chemistry submodels are under intensive current development. Each of these submodels is a large and sophisticated program in its own right. Furthermore, in the coupled model, each of the submodels, including the coupler, is a separate multiprocessor executable program. The coupler package must efficiently coordinate the communication as well as interpolate or aggregate information between these programs. This regridding function is necessary because each major subsystem (air, water or surface) is allowed to have its own independent grid.

Equally important to the success of the CCSM3.0 in contributing to the IPCC process was the large effort of the “Collaborative Design” project directed to insuring that the code was computationally efficient on the computing platforms available to the climate community. The predecessor to CCSM3.0, the Climate System Model (CSM), was developed for the Cray SMP architecture of the mid 1990s and was highly vectorized. The publicly released versions of CSM did not run on large distributed memory parallel computer systems. CCSM3.0, on the other hand, was expressly designed from the outset for such systems. It is a mixed mode parallel code with both Message Passing Interface (MPI) and OpenMP constructs. CCSM3.0 development benefited from the significant DOE investment in another NCAR coupled climate code, the Parallel Climate Model (PCM). In PCM, all component submodels are separate parallel codes as in CCSM3.0. However, in PCM the submodels are executed sequentially on the same set of processors rather than being executed simultaneously on separate sets of processors as is done in CCSM3.0. This constraint limits the parallelism to the least scalable of the submodels.

The land and ocean submodels are entirely different codes than in CSM (the ocean submodel is a descendant of that used in PCM). Furthermore, although the atmospheric component is a direct descendant of its counterparts in CSM and PCM, large portions of the code are either entirely new or were rewritten for better efficiency on cache-based multiprocessors. At the beginning of the CCSM3.0 design cycle, vectorization was not considered, since this was not relevant to the available machines. Later in the development of CCSM3.0, but prior to its public release, a significant allocation on the large NEC SX based vector machine at the Earth Simulator Center (ESC) in Yokohama, Japan was made available by CRIEPI. The challenge then became not only to vectorize the entire program, but also to preserve both efficiency on cache-based machines and the readability of the code. This latter point is especially key as the number of source code contributors to CCSM3.0 is very large. Scientists from the multiple institutions funded by the “Collaborative Design” project played an integral role in meeting this and the many other requirements of this model development and the public release of the CCSM3.0 source code was made in June of 2004. With three levels of parallelism, MPI, OpenMP and vectorization, and hundreds of lines of coding, the resulting multi-executable program is arguably the most computationally sophisticated large scale simulation code in history.

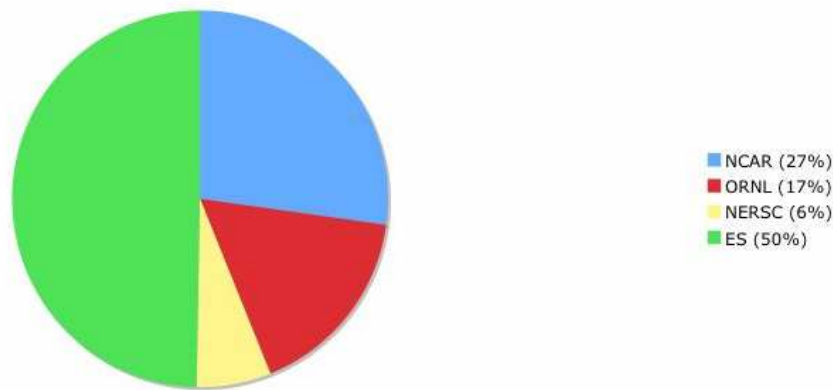
To be sure, there is more computing science to be done in the context of successors to CCSM3.0. The load balance between the submodels as well as within the submodels is not well characterized. Also, as resolution is increased, ever more restrictive time step conditions result in poor scalability with increasing numbers of processors. Better efficiency, obtained either algorithmically or through alternate parallelization strategies would appear to be necessary to permit increases in model fidelity. Finally, future physics parameterizations should be expected to increase in complexity, further taxing computational resources.

The fruit of these labors is substantial. Enabled by sizable allocations at NCAR, ORNL, NERSC and the ESC, the CCSM3.0 is the largest single contributor to the IPCC AR4 database. Not only is the model represented in more transient scenarios than any other model, more statistically independent realizations of each of these scenarios have been integrated than any other model. Also, the resolution of the atmosphere is exceeded by only one other model (which has not been ensemble integrated) and the ocean is at the highest resolution of any model. Differences between the CSM contribution to the IPCC Third Assessment Report (TAR) and the CCSM3.0 contribution to the IPCC AR4 are shown in the following table:

	Atmospheric resolution	Oceanic Resolution	Length of control run(s) in years	# of historical runs (# of scenarios)	# of future runs (# of scenarios)
CSM	T42L17	2°x2°x45L	300	3(3)	4(4)
CCSM3.0	T85L28	~0.3°x0.3°x40L	~1000 600	10(1)	55(6)

For the IPCC AR4 database, the vast majority of NCAR and DOE compute cycles initially available to the community was on IBM Power3 and Power 4 systems. Additionally, a large allocation of computing time at the Earth Simulator Center in Yokohama, Japan was made available. The distribution of compute resources used is shown in the pie chart below. It was obtained from a timeline published at NCAR prior to the completion of the integrations and may not be precisely what was used. It shows the relative fraction of simulated years integrated at each of the major supercomputer centers.

**Source of cycles for production T85 CCSM**



Source: <http://www.cgd.ucar.edu/ccr/ipcc/>

The integrations performed at US supercomputing centers were done solely on IBM Power 3 and Power 4 systems. The CRAY X-1 at ORNL was not used for the IPCC AR4 production runs. Note that half of the cycles that were planned to be used were at the Earth Simulator Center.

This increase in production is scientifically important in many ways. The need for large numbers of integrations is important in almost any climate change study. This is especially true of studies that attempt to characterize the uncertainty of climate change. For instance, in studies to detect recent climate change, larger numbers of statistically independent simulations of the twentieth century allow for a better definition of the model produced fingerprint of climate change. Studies considering future climate change require multiple realizations of a given emissions scenario to quantify the range of possible outcomes. Considering the differences between alternative emissions scenarios furthers the need for large numbers of integrations. Long control runs are also important to climate change studies in order to better quantify internal climate noise. While drift-free control runs of a few hundred years may be adequate to characterize interdecadal variability, studies considering climate change of the entire twentieth century period require control runs of at least a millennia to avoid artificial skill in the interpretation of the model. In summary, large ensembles and long control runs are necessary to better understand the signal to noise properties of climate change.

The Earth System Grid II (ESG) project, the other SciDAC funded laboratory project related to the IPCC AR4, attempts to facilitate the distribution of data produced by climate models for analysis. Because of the success of the CCSM3.0 in exploiting the current generation of either cache-based or vector machines, the volume of data generated by that model exceeds 100TB for the IPCC AR4 studies alone. Next generation architectures will enable even higher resolutions furthering the demand to rapidly transfer large amounts of data between institutions. Multi-model studies, often better describing the climate system than a single model, require yet more data and present the need for data standards. The ESG project aids in the transfer of data by creating data catalogs describing what is available. To accomplish this, metadata standards appropriate to climate model data have been defined. “Middleware” tools to exploit advances in hardware technology relevant to data movement have been developed. While still in an early stage, the ESG project has delivered a web-based portal that enables climate researchers to locate and acquire climate model data that is distributed across DOE and NSF supercomputing sites.

Climate models produce copious amounts of data. Because climate models are so computationally intensive, the inclination is to err on the side of caution and save many of the intermediate fields. In a large database of many different simulation scenarios, most of these will never be retrieved by an analyst because there are no observational counterparts to compare to. However, it must be recognized that model developers do need these fields in at least a limited number of simulations to diagnose the internal behavior of their models.

As a climate model integration proceeds, a large number of quantities are either accumulated for temporal averages or output instantaneously at scheduled intervals. Most models collect a large number of monthly averaged fields. These can be either at a surface or throughout the entire physical domain. Longer time averages are generally calculated afterwards from the monthly dataset. This portion of the output generally consists of the largest number of fields but may not be the largest portion in terms of bits. Higher frequency data can be of any interval but daily and six hourly frequencies are commonly specified. For many analysis purposes, surface fields suffice for the high frequency datasets. However, this is not always the case as illustrated by the needs of mesoscale atmospheric modelers to save prognostic variables from global models serving as input boundary conditions that resolve the diurnal cycle.

The very nature of how model data is produced leads to distribution problems. As the model is advanced in time, output must be saved. Because of the length of climate simulations, output data files must regularly be closed and successor output files opened. At the end of the simulation, the user is left with a series of files, each of which contains many fields but over a relatively brief time span. The climate model analyst, on the other hand is usually interested in only a few fields but over a long period of time. Extracting a single field out of the database is a time consuming process as the entire database must be processed to get at this very small portion of it. Typically, bandwidth, disk space and memory limitations conspire against making this a straightforward task. Therefore, serving raw model output data is an inefficient means to distribute data to a community of external data users.

ESG has developed and deployed two linked web portals where scientists can register, browse data catalogs, search for specific datasets, check system status, and download data to their own institutions. The first portal, [www.earthsystemgrid.org](http://www.earthsystemgrid.org), installed at NCAR serves 70TB of data in 600,000 files from the NCAR climate models, CCSM3.0 and PCM. It supplies both “raw” history files along with post-processed files containing a single variable. The user may choose to select a collection of files or request an arbitrary subset of an online dataset by specifying which variables, times, and geographic region they’d like to acquire. It is also the point of distribution for the CCSM3.0 source code, model initialization datasets, and tools for analysis. The second

portal, linked from [www.pcmdi.llnl.gov](http://www.pcmdi.llnl.gov), serves the multi-model IPCC AR4 database. This dataset consists of entirely post-processed data and uses an older version of the ESG environment than does the NCAR site. Significant differences including searching and subsetting functions exist in the two implementations of the webportal. The IPCC AR4 database has been prepared using data from 23 modeling groups around the world using the PCMDI routine, Climate Model Output Rewriter (CMOR). In addition to standard naming and directory conventions, CMOR enforces compliance with the CF standard and provides a comprehensive metadata description. Without these standards, distribution of the IPCC AR4 data to the community would have been difficult if not impossible. With these standards, 42 TB of distributed data has resulted in at least 220 papers submitted for review between February and June 2005.

The climate community now has a new enabling data resource to tap for their research efforts. Building upon this, future directions for ESG center around expanding the available data holdings, enhancing the organization of the datasets, and streamlining data download capabilities. For datasets that are online, the system provides a high-level interface that allows users to flexibly specify a spatiotemporal constraint along with which set of fields that they want. However, a large fraction of the climate database is distributed across several sites on archival storage systems. While the files are available, they are not formatted in a manner most optimal for typical user access patterns. Strategies for addressing this include expanding available online storage for popular datasets, optimizing the formatting of the individual files, and developing new interfaces to make it easier to specify and download a large number of individual files. Additional enhancements to the metadata and dataset browsing interface will elevate the user's ability to rapidly locate experiments and fields of interest.

### **Conclusion**

The success of the "Collaborative Design and Development of the Community Climate System Model for Terascale Computers" and the "Earth System Grid II: Turning Climate Datasets into Community Resources" SciDAC projects are intimately connected. For it is as important for the first to accelerate the production of climate model output data, as it is for the second to accelerate the distribution of this data to the analysis community. The impact of these two projects on the IPCC AR4 is significant. The CCSM3.0 is the largest contributor (at about 30% in June 2005) to the IPCC AR4 database of climate model integrations. The distribution of this database to the community has resulted in a large number of studies in a very brief period of time. Further distribution of this database to a yet larger community is both desirable and practicable with the planned ESG web portal enhancements.

Work supported by the Office of Biological and Environmental Research of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231

### **References**

Intergovernmental Panel on Climate Change website, <http://www.ipcc.ch/>  
Jerry Meehl, presentation at 10<sup>th</sup> Annual CCSM workshop, Breckenridge, CO June 2005  
Dean Williams and Don Middleton, personal communication, June 2005  
Program for Climate Model Diagnosis and Intercomparison website, <http://www-llnl.pcmdi.gov>  
Earth System Grid website, <http://earthsystemgrid.org>