# Data Issues for WCRP Weather and Climate Modeling

James L. Kinter III[1] and Karl E. Taylor[2]

A White Paper
based on presentations at the

First Session of the WCRP Modeling Panel
Exeter, UK
6 October 2005

As described in *The World Climate Research Programme Strategic Framework 2005-2015: Coordinated Observation and Prediction of the Earth System (COPES)*[3], the WCRP is embarking on an ambitious, decade-long observing and modeling activity that is intended to improve understanding of the mechanisms that determine the mean climate and its variability, with the ultimate objective of providing the soundest possible scientific basis for a predictive capability for the total climate system. The framework calls for an integrated approach in which the roles of the atmosphere, ocean, land and cryosphere are considered in comprehensive models of the climate system. The continuum of prediction problems, from weather-to-climate and days-to-decades, will be addressed by a hierarchy of models that should become increasingly similar to one another, merging eventually into "unified models" with common infrastructure and interchangeable parameterizations, and variously configured to address a wide range of problems. This fact, when considered with the anticipated increasing international coordination of model development, integration and analysis, implies that the modeling and model output data management challenges will be very large.

The COPES framework calls for free and open access to data, with a movement toward a more unified system of data management and access across WCRP projects. Observational data are to be managed in such a way as to facilitate reprocessing and reanalysis, possibly repeated many times. Model output data must be easily accessible to experts in various disciplines and regions and managed in such a way as to facilitate in-depth analysis of multiple model data sets. The COPES data management plan is required to address many facets of the data handling challenge, including management, stewardship and access to data, and special issues concerning climate system data assimilation, synthesis and reanalysis, and model initialization. The development of the data management plan is to be done in coordination with other WMO activities, notably THORPEX[4], GEO/GEOSS[5], and GODAE[6].

This white paper raises several data management issues of relevance to COPES, with particular emphasis on how to facilitate analysis of model output data sets by a dispersed community

---

[1] Center for Ocean-Land-Atmosphere Studies, 4041 Powder Mill Road, Suite 302, Calverton, MD 20705 (http://www.iges.org; correspondence may be sent to kinter@cola.iges.org)

[2] Program for Climate Model Diagnosis and Intercomparison, Lawrence Livermore National Laboratory 7000 East Avenue, Livermore, CA 94550 (http://www-pcmdi.llnl.gov/; correspondence may be sent to taylor13@llnl.gov)

[3] WCRP-123, WMO/TD-No. 1291, August 2005; available from WCRP, World Meteorological Organization, 7 bis, avenue de la Paix, P.O. Box 2300, 1211 Geneva 2, Switzerland

[4] http://www.mmm.ucar.edu/uswrp/programs/thorpex.html

[5] http://www.epa.gov/geoss/index.html

[6] http://www.bom.gov.au/bmrc/ocean/GODAE/

of experts.  The model output data and metadata requirements anticipated in the next decade are summarized, current and emerging capabilities are discussed, and recommendations for how to proceed are presented.

**Model Output Volume in 2015**

The modeling systems likely to be used to realize the COPES vision will be highly sophisticated computer programs that represent at high resolution the coupled ocean-atmosphere-land-cryosphere system and all the dynamical, physical and biogeochemical processes that are relevant on a broad spectrum of time scales. As estimated in the Appendix, the potentially useful data produced by a single model for a single type of application will occupy $O(10^{17} - 10^{19})$ bytes of storage.  For coordinated applications across O(10) modeling centers, the COPES data management plan needs, therefore, to include provisions for data volumes of order $10^0 - 10^2$ exabytes,[7] distributed across multiple data centers worldwide.

**Relevant Technology Trends**

There are several technology trends that have a bearing on plans for COPES data management. The growth in data density in magnetic storage systems has accelerated in the past decade. Prior to the 1990s, disk density doubled every three years or so. With the introduction of magnetoresistive read heads in 1991, the doubling time was reduced to two years, and, since the giant magnetoresistive head (GMR) reached the market in 1997, density has been doubling every year. As an example of data volume growth, including archives on less volatile media, the Mass Storage System (MSS) at the U.S. National Center for Atmospheric Research acquired its first petabyte ($10^{15}$ bytes) over about 18 years of data accumulation. The MSS reached its second petabyte 18 months later. The measured rate of growth of the MSS at NCAR is about 50 bytes per sustained kiloflop (KF)[8].

The capacity (bandwidth) of wide area networks (WAN) that link computing centers worldwide has likewise experienced exponential growth over the past decade or more. In the 1990s, the typical highest speed within a data center was 100 Mb/s[9] (over fiber distributed data interface, FDDI, technology) and the WAN speed was typically up to 45 Mb/s. Today, the fastest data center network is about 1-10 Gb/s, and, while 10 Gb/s is possible in WAN, the practical limit is currently 1 Gb/s. Therefore, WAN bandwidth has barely kept pace with disk storage volumes, and, given the fact that large data transfers over long haul networks were unwieldy or impossible in the past, this situation has not improved.

With petaflop-class computing throughput and accelerating exponential online storage growth, exabyte ($10^{18}$ bytes) data volumes will be the norm by 2015. With relatively similar or slower WAN bandwidth growth, the networks will not be able to keep pace with the volume of weather and climate model output. Inevitably, such data will of necessity be widely distributed worldwide, and sophisticated subsetting and on-demand processing and visualization will be absolutely required. Similarly, the trend in data management systems and software has been away from the centrally designed, implemented and maintained systems that characterized data centers in the 1990s and earlier. The new generation of data management systems and software are integrated systems of independently designed, implemented and maintained system elements. One example of this is the Open-source Project for a Network Data Access Protocol (OPeNDAP)[10], which has been conceived

---

[7] An exabyte is $10^{18}$ bytes.

[8] Computer performance in modeling applications is typically measured in floating point results per second (flops); a kiloflop (KF) is $10^3$ flops.  Sustained supercomputer performance is expected to reach 100 teraflops (TF; $10^{12}$ flops) in 2010 and 1-10 petaflops (PF; $10^{15}$ flops) by 2015.

[9] Mb/s - megabits or $10^6$ bits per second; Gb/s - gigabits or $10^9$ bits per second

[10] http://www.opendap.org/

as an access/delivery element in this environment of distributed data system elements. OPeNDAP is a software framework used for data networking that makes local data accessible to remote locations.

**Subsetting, On-Demand Processing, and Multiple File Formats**
Of critical importance for future data management systems will be the capabilities to extract subsets of the very large data sets, stored in different formats, and to create derived data sets that represent the results of on-demand processing on the server side, delivered over the network. These capabilities significantly reduce or eliminate the bottleneck, expected to worsen over the next decade, caused by the growing group of scientists seeking data from large repositories of model output. The subsetting capability allows users to retrieve a specified temporal and/or spatial subdomain from a large data set, meeting a user's needs while minimizing the amount of data transferred. Generating derived data sets in response to user requests by processing data on the server side can further reduce the volume of data that must be transported over the WAN, and makes it possible for experts to share intermediate results and analysis methods with each other.

Model and observational data sets can currently be accessed through various web and internet-based interfaces (e.g., ftp, GDS[11], ESG[12], LAS[13]), each with different evolving capabilities. All can transfer complete files from a single site over the internet. Some provide a level of security, e.g., denying access except to approved users. Some can extract subsets of data and perform simple server-side calculations (e.g., obtain a single pressure level, a climatological mean, a zonal mean), and some can perform more complex server-side calculations. There is a rudimentary capability among some tools to transfer data from disperse sites, but make it look like a single aggregated site.

Another important aspect of distributed data management is the fact that data are served in a variety of formats (e.g., GRIB[14], binary, netCDF[15], HDF[16], BUFR[17], and GrADS station format), while the programs scientists use to analyze those data are often format-specific. Software can be written that unifies the variety of data formats into a single framework to simplify data analysis. One example of such software is the GrADS Data Server[18] (GDS), which provides subsetting and analysis services across a wide range of commonly used meteorological and oceanographic data formats.

The capability to analyze ensembles of predictions or simulations from multiple models, station observations, objective analyses and remote sensing data in a single analytic framework – is an essential component of the COPES framework. This capability can facilitate the assimilation of observational data and model output, and it can make it possible to perform analyses in a variety of ways, across models, across ensemble members, across real time and across forecast/simulation time.

**Metadata Standards**
There are a number of issues associated with so-called "metadata" that describes observations and model output data. Metadata useful in describing model data sets include, for example: model documentation, experimental design, fields stored, precise definition of the physical quantities represented, units (and calendar for time), space-time location, grid information (e.g., grid-cell area), and processing applied (interpolation, climatological averaging, zonal or other spatial averaging, etc.). A uniform representation (or translation) of all these and other metadata is required across models and experiments to facilitate automated analysis. Also, data "discovery" tools are needed to find and

---

[11] Grads Data Server - http://www.iges.org/grads/gds/gds.html

[12] Earth System Grid - https://www.earthsystemgrid.org/

[13] Live Access Server - http://ferret.pmel.noaa.gov/Ferret/LAS/ferret_LAS.html

[14] WMO gridded data format standard - http://www.wmo.ch/web/www/WDM/Guides/Guide-binary-2.html

[15] http://www.unidata.ucar.edu/software/netcdf/

[16] http://hdf.ncsa.uiuc.edu/

[17] WMO station report format standard - http://www.wmo.ch/web/www/WDM/Guides/Guide-binary-1A.html

[18] http://www.iges.org/grads/gds/gds.html

interpret these metadata and help users locate desired model output.

The current situation in metadata available for weather and climate model output data sets is somewhat disheartening but has been improving in recent years. Model documentation is, in general, uneven and difficult to find. One notable recent exception is the Intergovernmental Panel on Climate Change[19] (IPCC) data set for which standards were adopted and enforced for model documentation[20]. The description of experimental designs is likewise in a somewhat primitive state, and the metadata describing fields stored is often insufficient. A uniform representation of model output metadata, across models and experiments, exists only for a few well-coordinated multi-model projects (e.g., AMIP[21], PMIP[22], DEMETER[23] and IPCC). Data discovery tools are similarly available only for limited sets of model experiments.

In the cases of the coordinated multi-model projects, serving a wide variety of users, data have usually been sent to a central repository such as the Program for Climate Model Diagnosis and Intercomparison (PCMDI). In those cases, output has only been accepted with strict conformance to rigid metadata requirements. The quality control and data management in such cases is relatively easy. As data set sizes increase, centralized data archives will evolve along already planned pathways to become distributed repositories with centralized cataloging services.

One of the issues of metadata management is the (dis)similarity among metadata from diverse sources. To address this issue, some coordinated experiments have turned to metadata standards and conventions. One example is the CF standard[24], which provides specifications that govern the creation of fully self-describing netCDF files. The CF standard, which is becoming increasingly accepted by the climate modeling community, is an extension of the earlier, more limited COARDS[25] standard. The application of standards like CF encourages the storage of metadata that can be useful in model diagnosis and enables the development of common software that can "understand" model output from diverse sources. Software can be written to facilitate conformance with the standard.[26]

In the case of the IPCC data set, there was a major scientific payoff for the specification of standards and conventions. The output from 21 different models is being widely analyzed by over 400 registered users. Over 25 terabytes ($25 \times 10^{12}$ bytes; 60,000 files) were collected and more than 60 terabytes (290,000 files) have been disseminated to analysts. Over 200 manuscripts have been written, based (at least in part) on the IPCC data set, which will likely attract continued scientific interest for several years to document the current generation models and provide a baseline for assessing future models.

**Distributed Database Management, Cataloging, and Discovery**

As described above, future data set archives will become increasingly distributed, which poses several challenges.  It is likely that because of security issues or hardware limitations, some institutions may be unable to serve their own data and will need to send it to a publicly accessible repository.  In such cases a distributed approach to *management* of the data would be attractive,

---

[19] http://www.ipcc.ch/

[20] http://www-pcmdi.llnl.gov/ipcc/info_for_analysts.php

[21] Atmospheric Model Intercomparison Project - http://www-pcmdi.llnl.gov/projects/amip/index.php

[22] Paleoclimate Model Intercomparison Project - http://www-lsce.cea.fr/pmip2/

[23] http://www.ecmwf.int/research/demeter/

[24] http://www.cgd.ucar.edu/cms/eaton/cf-metadata/

[25] http://ferret.wrc.noaa.gov/noaa_coop/coop_cdf_profile.html

[26] As an example, for the IPCC data set, the Climate Model Output Rewrite (CMOR) program facilitated conformance with IPCC requirements (http://www-pcmdi.llnl.gov/softward-portal/). The CMOR "input tables" provided metadata information and quality control specifications, so that CMOR could trap mistakes in model output before it was released for analysis. The CMOR input tables can be tailored to the needs of future data-sharing activities.

allowing contributing groups to maintain their own data remotely.  This would obviously require remote access to the data beyond simply reading it and therefore raises additional security issues.

There is also a growing need to enable data discovery across many diverse data servers. One solution to this problem is to gather the metadata from various servers in a browsable, searchable catalog. The typical implementation includes a "crawler" that periodically updates a central metadata catalog from all known data servers, a browser interface for human exploration of the data being served, and a programmable interface to automate the use of the output.

As one example, under the Earth System Grid (ESG) project a cataloging system has been developed that spans data centers included in the grid.  One way the IPCC data set is being served is through such a catalog (but in this case limited to a single site).  As another example, a prototype developed at the Center for Ocean-Land-Atmosphere Studies (COLA) called *Greta*[27] "crawls" nightly over the entire data holdings at COLA and a subset of the data served at NCAR to "harvest" metadata automatically and produce a single catalog of all the data being served. Once the metadata have been collected, THREDDS[28] and Lucene[29] are used to create a searchable database. A search on keywords, metadata content, or space-time coordinates can be conducted from any web browser or, importantly, from a customizable code that can make use of the search output. The output is human-readable, as in any web browser, but also machine-readable: the output can be returned as plain text that permits parsing the output with customizable code, or it can be returned in XML[30]-form, which permits parsing the output with XML-enabled programs.

Inevitably, it is expected that a number of different solutions will be developed to assist in discovery and cataloging of distributed data sets.  In recognition of this, an international collaboration, called the Global Organization for Earth System Science Portal (GO-ESSP)[31], is discussing various approaches to providing distributed access to weather and climate data. The goal is to weave together various frameworks designed for data discovery, access, and analysis.

## Recommendations

In order to address the data management issues described above while working toward the goals of the COPES framework, we suggest that the following actions should be taken.

1. The WMP should endorse a distributed model output data management plan that:
    a. Takes into account the projected high data volume
    b. Avoids bandwidth and logistical problems associated with transmitting data to a single repository
    c. Minimizes single point of failure issues
    d. Accommodates site and data set requirements for various levels of access control and user authentication
    e. Makes optimal use of subsetting and server-side analysis capabilities by capitalizing on analysts' familiarity with existing tools and techniques
    f. Encourages further development of browsable, searchable cataloging capabilities with programmable interfaces that enable automated downstream processing
    g. Encourages emerging international frameworks designed to facilitate data discovery, access, and analysis
    h. Provides workable options for groups unable to serve their own data.

---

[27] http://www.iges.org/about_greta.html

[28] Thematic Realtime Environmental Distributed Data Services - http://www.unidata.ucar.edu/projects/THREDDS/

[29] An open source Java toolkit for text indexing and searching - http://sourceforge.net/projects/lucene/

[30] Extensible Markup Language - http://www.w3.org/TR/REC-xml/

[31] http://go-essp.gfdl.noaa.gov/

2. The WMP should consider the adoption of standards and conventions that facilitate the comparison and interoperability of metadata and model output data from diverse sources by:
   a. Encouraging the use of CF-conventions for netCDF files
   b. Establishing for coordinated experiments IPCC-like requirements. which, for example, specify required attributes for: modeling group, experiment, variable, units, calendar, location and area of grid cells, and other information useful to analysts[32]
   c. Capitalizing on modeling groups' familiarity with standards-conforming software

**Appendix: Estimate of data storage needs**

A rough estimate of data storage needs for a "single experiment" performed with the climate models of the future can be made as follows: Such models will have $O(10^2)$ levels representing the vertical structure in the system and $O(10^8)$ columns, subsampled before saving at a resolution perhaps a factor of 100 lower, or in some cases run only at a lower, $O(10$ km$)$, resolution, yielding $O(10^6)$ saved columns. The models will output $O(10^2)$ three dimension fields and $O(10^3)$ two-dimensional fields, representing the prognostic and diagnostic variables that characterize the physical, chemical and biological state of the system. Data will be saved $O(10^3)$ times per run, whether it is a relatively short weather prediction run or a longer climate simulation run – typically, sampled every half hour for weather prediction, four times per day for seasonal prediction, and monthly for climate simulation. The model integrations will be instantiated $O(10^1 - 10^2)$ times to represent ensembles that can be used to estimate uncertainty in each of $O(10^2 - 10^3)$ cases – e.g., three years of weather prediction cases or $O(10^3)$ choices of uncertain parameter values in climate prediction cases. Thus, $O(10^{10} - 10^{11})$ bytes will be stored for each of $O(10^3)$ save times in $O(10^4 - 10^5)$ runs per experiment suite, which means the global repository of COPES model output data sets will amount to $O(10^{17} - 0^{19})$ bytes or $O(0.1$ to $10)$ exabytes $(10^{18}$ bytes$)$ per model per suite of experiments for each of $O(10)$ modeling groups worldwide.

---

[32] e.g., http://www-pcmdi.llnl.gov/ipcc/IPCC_output_requirements.htm