

## Atmospheric Sciences and Climate Applications Using HDF and HDF5

MuQun Yang  
Robert E. McGrath  
Mike Folk

National Center for Supercomputing Applications  
University of Illinois, Urbana-Champaign

### 1. Introduction to HDF and HDF5

The Hierarchical Data Format (HDF) developed at the National Center for Supercomputing Application (NCSA) at University of Illinois at Urbana-Champaign has been used by many atmospheric science applications since it is first released in 1988.

Since 1999, NCSA has developed a more general and robust data format, called HDF5, which was aimed to support the future demands of Earth Science such as large data storage, performance and flexibility. NCSA HDF group encourages applications to use HDF5 for better performance and maintenance. You may read extend abstracts of P2.42 for more information about HDF5. In summary, HDF5 have the following features to attract many important atmospheric Sciences applications to use HDF5.

- 1) Flexible data model
- 2) Open-source, free software
- 3) Portability
- 4) Emphasize on performance
  - Support MPI- IO
  - Better sub-setting
  - In-memory compression
  - Alternative storage

In this paper, we will demonstrate how HDF and HDF5 provide the flexibility and efficiency in six well-known Earth Science applications.

### 2. NASA EOSDIS project

The NASA Earth Observing System The Earth Observing System Data and Information System (EOSDIS) is a large archive of systematic measurements of the Earth's climate collected by a series of satellites [1]. The EOSDIS archives store and distribute data using HDF and HDF5. Data from these missions will be available for many decades, to study climate and climate change.

HDF provides a rich data model that supports most of the types of data required by scientists. To support the specific needs of remote sensed data, the EOS project defined a storage profile, called HDF-EOS, and metadata standards [2]. HDF-EOS defines a standard way to store georeferenced data using HDF and HDF5.

The EOSDIS archives contain several petabytes of data stored in HDF files. From 1999, data from Terra, Aqua, Landsat 7, and other satellites has been distributed in HDF-EOS4 and HDF4. Starting in 2004, data from the Aura satellite will be delivered in HDF-EOS5 and HDF5 [3].

### 3. NPOESS project

The National Polar-orbiting Operational Environmental Satellite System (NPOESS) is a system of polar orbiting weather satellites and ground equipment used for the collection, analysis and distribution of weather data to government and civilian users, and will also provide long term climate records [4].

NPOESS Interface Data Segment will deliver data in HDF5. HDF5 was chosen for use on NPOESS because it has the capability to operate well in high performance, data intensive environments. HDF5 can store data in a variety of ways. NPOESS has chosen to standardize their organization of the HDF5 files so data can be easily and consistently accessed and shared amongst the community [4].

The NPOESS Preparatory Program (NPP) will be used as a bridge between the existing Earth Observing System (EOS) program and the NPOESS Program, to assure a continuous record of the Earth's climate [5]. The NPP will provide an opportunity to utilize new instruments, algorithms and data delivery packages prior to utilize new instruments, algorithms and data delivery packages prior to NPOESS. Data from the NPP satellite will be delivered through the NPOESS Interface Data Segment using HDF5.

### 4. WRF-HDF5 IO modules

WRF (Weather and Research Forecasting Model) [6] has been mainly developed and maintained at NCAR. It is a regional weather model that has become to be intensively used for both weather research and prediction. Due to the heavy computational volume, WRF is using multi-layer domain decomposition method to run the model in parallel supercomputing environments.

#### 4.1 Sequential IO issues to address

WRF officially supported sequential NetCDF IO module. There are two problems with the current WRF IO module.

The first one is that the IO module may become IO bottleneck and the IO module may exceed the memory capacity in some applications.

As the model resolution becomes higher, huge volume of disk storage for large WRF applications may need to use in-memory data compression technique to store model data without taking extra computer time.

#### 4.2: How HDF5 can help

- 1) HDF5 fully support MPI-IO in many platforms, which is the standard parallel IO library. With parallel-IO, the large WRF application can avoid or reduce the IO bottleneck and the potential problem of exceeding the memory capacity.
- 2) HDF5 supports two in-memory compression packages: deflate compression and szip compression[7].

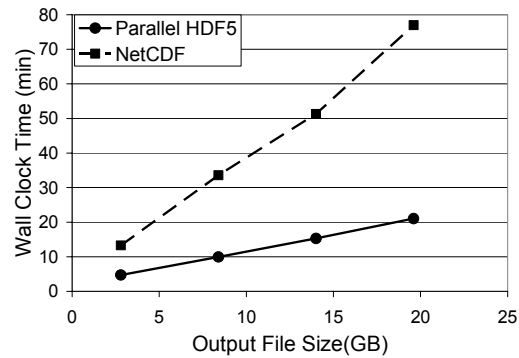
#### 4.3: WRF-HDF5 IO modules

NCSA implements WRF-HDF5 sequential IO module and WRF-parallel HDF5 IO module.

- 1) Both the sequential and the parallel HDF5 IO module should work with the current WRF 2.0 release. Applications can do szip and gzip compression within sequential IO module. Applications can do parallel IO through parallel HDF5 IO module. Parallel HDF5 IO is incorporated into WRF model package in WRF 2.0 release[6]. You can also download the source code from [8].

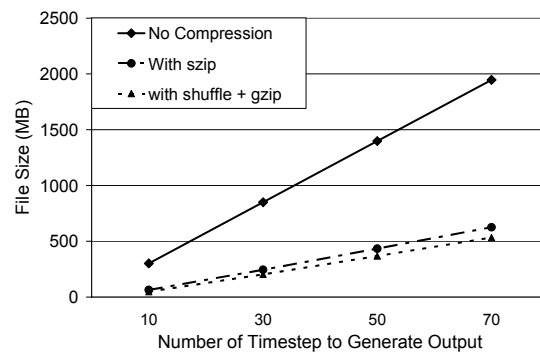
#### 4.4: Performance results

Figure 1 shows a comparison of the wall clock time between WRF-parallel HDF5 IO module and WRF-sequential netCDF IO module for a real WRF application at NCAR IBM with 256 processors. The maximum hyperslab size per timestep is 17 MB. The WRF-parallel HDF5 IO module outperforms WRF-netCDF IO module as the output file size increases. When the output file size reached around 20GB, the model run took about 80 minutes wall clock time when using WRF-sequential netCDF IO module; compared to 20 minutes when using parallel WRF-Parallel HDF5 IO module.



**Figure 1: Performance of Parallel HDF5 to sequential NetCDF for Different Output File Size IBM Power 3 (256 Processors)**

Figure 2 shows file size of different runs for each option in another application. As expected, the overall HDF5 file size in either SZIP compression or shuffle with GZIP compression is much less than the same data without compression. For example, for the data at timestep 70, the file size is about 2000 MB without compression, compared to less than 600MB for either SZIP or GZIP compression.



**Figure 5: Performance of Sequential HDF5 with Different Compression Methods IBM Power 4 (16 processors)**

Compared with sequential NetCDF IO module:

- Parallel HDF5-WRF IO module can greatly improve WRF IO performance for **some** WRF applications.
- Sequential HDF5-WRF IO module can greatly reduce the WRF model output data size.

## 5. OPENDAP/DODS HDF4 and HDF5 servers

The Distributed Oceanographic Data System (DODS) [9] is a distributed data management system that was designed for scientists to remotely use their application programs over the Internet in a consistent way [10]. DODS uses a Web client-server model. A client sends a data request as across the Internet to a web server. The

server checks the request and answers with the requested data if the request is legal. DODS has been widely used by many atmospheric and climate applications.

DODS maintains HDF4 and HDF5 servers [9] to help Earth scientists to obtain data in HDF4 and HDF5.

Many applications use DODS to do data subsetting, so data access time for data subsetting through internet is extremely important. Furthermore, data compression technique may also be used to enlarge the relative bandwidth of the data transmitted through the internet. An evaluation of DODS-HDF4 server with AVHRR Oceans Pathfinder data and found that using HDF4 chunking storage and compression techniques may greatly improve the data access time [11].

NCSA also implemented a DODS-HDF5 server prototype associated with a DODS-HDF5 white paper [12] and a DODS-HDF5 data mapping paper [13]. NCSA also demonstrates how DODS-HDF5 can work with real NASA data through DODS Ferret client [14]. Since there are no real HDF5 dataset in 2000, we used the NCSA HDF4 to HDF5 conversion utility to convert a real NASA data from HDF4 to HDF5 and then used DODS ferret client to demonstrate the server work.

## 6. NetCDF4

NCSA and Unidata are collaborating to merge netCDF and HDF5. In version 4.0 the netCDF API will be extended and implemented on top of the HDF5 data format [15]. Users of netCDF in numerical models will benefit from support for packed data, larger datasets, data compression, and parallel I/O, all of which are available with HDF5. HDF5 users will benefit from the availability of a simpler high-level interface suitable for array-oriented scientific data; wider use of the HDF5 data format; the wealth of netCDF software for data management, analysis and visualization; and the body of experience that has evolved in the years since netCDF began.

## 7. HL-HDF5

National meteorological services in Scandinavia countries have used HDF5 to store their radar data in order to exchange radar data among Sweden, Norway and Finland with quality-related information. HDF5 was selected because it is open-source, well-designed and supports in-memory compression. Using HDF5 with deflate compression to store weather data performed better than compression through BUFR [16]. They created an information model to use HDF5 to store radar data. Their weather radar information model can be used to store individual scans, images and products. They also created a high-level interface to HDF5 called HL-HDF5 to better facilitate the data management [17].

## 8. Conclusions

We made a brief introduction of six atmospheric and climate applications using HDF and HDF5. An incomplete application lists can be found at [18] and more general information of HDF5 can be found at <http://hdf.ncsa.uiuc.edu>.

## References

1. Savtchenko A., Ouzounov D., Ahmad S., Acker A., Leptoukh G., Koziana J. and Nickless D.. Terra and Aqua MODIS products available from NASA GES DAAC. *Advances in Space Research, Volume 34, Issue 4, 2004*, 710-714
2. Leptoukh, G.; Ouzounov, D.; Savtchenko, A.; Ahmad, S.; Li Lu; Pollack, N.; Zhong Liu; Johnson, J.; Jianchun Qin; Sunmi Cho; Li, J.Y.; Kempner, S. and Bill Teng. HDF/HDF-EOS data access, visualization and processing tools at the GES DAAC, *Geoscience and Remote Sensing Symposium, 2003. IGARSS '03. Volume: 6*, 3571 - 3573
3. Schoeberl, M.R.; Douglass, A.R.; Hilsenrath, E.; Luce, M.; Barnett, J.; Beer, R.; Waters, J.; Gille, J.; Levelt, P.F. and DeCola, P. The EOS Aura Mission, *International Geoscience and Remote Sensing Symposium (IGARSS), Volume 1, 2001*, 227-232.
4. Goldberg, A.M. Delivering Earth's shapes & colors in near-real time: NPOESS products and their characteristics for users, *Geoscience and Remote Sensing Symposium, 2003. IGARSS '03, Volume: 5*, 3035.
5. Murphy, R.E.; Henegar, J.; Wharton, S.; Guenther, B. and Kealy, P.M.; Extending climate data records from the EOS era into the NPOESS era, *Geoscience and Remote Sensing Symposium, 2003. IGARSS '03, Volume: 2*, 1332 - 1334.
6. WRF: <http://wrf-model.org/>
7. SZIP: [http://hdf.ncsa.uiuc.edu/doc\\_resource/SZIP/](http://hdf.ncsa.uiuc.edu/doc_resource/SZIP/)
8. WRF-HDF5 IO module: <http://www.ncsa.uiuc.edu/apps/WRF-ROMS>
9. DODS: <http://www.unidata.ucar.edu/packages/dods/index.html>
10. Davis E. and J. Gallagher, Using DODS to Access and Deliver Remote Data, *The 15th International Conference on Interactive Information and Processing Systems(IIPS) for Meteorology, Oceanography, and Hydrology*, 1999.
11. Kapadia A. and N. Yeager. Performance Evaluation of HDF Version 4 Chunking Facility when used for Subsetting Pathfinder Data, 1999 <http://hdf.ncsa.uiuc.edu/apps/dods/perfeval/report.html>
12. Yang M. and R. E. McGrath. The HDF5-DODS Server Prototype, 2001 <http://hdf.ncsa.uiuc.edu/apps/dods/DODS-White-paper.pdf>
13. Yang M. and R. E. McGrath. HDF5-DODS Data Model and Mapping, 2001 <http://hdf.ncsa.uiuc.edu/apps/dods/HDF5-DODS-Mapping.pdf>

14. Yang M. and R. E. McGrath. Demonstration of HDF5-DODS Server Prototype with the DODS-Ferret Client, 2001  
<http://hdf.ncsa.uiuc.edu/apps/dods/DODS-Ferret-demo.pdf>
15. Russ R. and E. Hartnett. "Merging netCDF and HDF5," *Bulletin of the American Meteorological Society, Combined Preprints: 84th American Meteorological Society (AMS) Annual Meeting*, 2004, 1457-1460
16. Michelson D. B., Holleman I., Hohti H., and Salomonsen M., HDF5 information model and implementation specification for weather radar data, 2003,  
[http://www.smhi.se/cost717/doc/WDF\\_02\\_200204\\_1.pdf](http://www.smhi.se/cost717/doc/WDF_02_200204_1.pdf)
17. Henja A, and D. B. Michelson, A High Level Interface to the HDF5 File Format,  
<ftp://ftp.ncsa.uiuc.edu/HDF/HDF5/contrib/hl-hdf5/README.html>
18. HDF5 Users, <http://hdf.ncsa.uiuc.edu/users5.html>

#### **Acknowledgements:**

Authors want to thank Daniel Michelson for providing information about how HDF5 is used to store weather radar data.

The netCDF4 project is funded by NASA award AIST-02-0071.

The WRF IO Module was done as part of NSF-funded Modeling Environment for Atmospheric Discovery Expedition (MEAD) (<http://www.ncsa.uiuc.edu/AboutUs/FocusAreas/MEADExpedition.html>)

This report is based upon work supported in part by a Cooperative Agreement with NASA under NASA grant NAG 5-2040. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Aeronautics and Space Administration.