



# Evaluation of Big Data Containers for Popular Storage, Retrieval, and Computation Primitives in Earth Science Analysis

Thomas Clune<sup>§</sup>, Kamalika Das<sup>§</sup>, Daniel Duffy<sup>§</sup>, Ted Habermann<sup>\*</sup>, Thomas Huang<sup>§§</sup>, Kwo-Sen Kuo<sup>§</sup>, Chris Mattman<sup>§§</sup>, Chaowei Phil Yang<sup>§§</sup>

AGU Fall Meeting, San Francisco, December 2015

<sup>§</sup>NASA Goddard Space Flight Center, <sup>§§</sup>NASA Ames Research Center, <sup>\*</sup>The HDF Group, <sup>§§</sup>NASA JPL, <sup>§§</sup>George Mason University, Contact: Kamalika.Das@nasa.gov

## Abstract

- Data containers are infrastructures that facilitate storage, retrieval, and analysis of data sets. Big data applications in Earth Science require a mix of processing techniques, data sources and storage formats that are supported by different data containers. The data containers compared in this study are
  - AsterixDB,
  - RasDaMan,
  - SciDB
  - Hadoop
  - HDF
- These infrastructures optimize different aspects of the data processing pipeline and are, therefore, suitable for different types of applications. These containers are also all undergoing rapid evolution and the ability to re-test, as they evolve, is very important to our handling of the large volumes of observational data and model output. We have identified a selection of steps that are relevant to most data processing exercises in Earth Science applications and we evaluate these systems for optimal performance for each of these steps in the data processing pipeline. The steps evaluated in this study:
  - Hardware/software dependencies
  - Data ingestion
  - Data preparation/processing
  - Data analysis
  - Result reporting

## Software and Hardware Dependency

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Version</b> AsterixDB 0.8.7-Snapshot	<b>Version</b> 9.1	<b>Version</b> 14.12	<b>Version</b> Cloudera 5.0	<b>Version</b> HDF5 v1.8.15
<b>Hardware dependency</b> Share Nothing Architecture, uses controller nodes	<b>Hardware dependency</b> 3GHz, 8GB RAM, 400MB HDD for installation	<b>Hardware dependency</b> Distributed file system for data (Optional) shared file system for software	<b>Hardware dependency</b> 12-24 1-4TB hard disks, 2 quad-/hex-/octo-core CPUs, running at least 2-2.5GHz	<b>Hardware dependency</b> Xeon, Ethernet, Ephemeral file system, S3
<b>Software dependency</b> JDK 1.7 Password-less SSH configuration	<b>Software dependency</b> Git, lib, Tomcat (or another suitable servlet container), Java Runtime Environment (JRE) 1.6 or higher, PostgreSQL 8.x	<b>Software dependency</b> PostgreSQL, Apache Maven, Apache log4cxx, Fedora mock, Google protobuf, ScaLAPACK, Shim, SciDB-Py, SciDB-R, SciDB cluster	<b>Software dependency</b> CentOS OpenVZ for RHEL 6 – LX C version 1.1.3 Infiniband JDK 1.7.0_67, python, perl	<b>Software dependency</b> HDF5 library v1.8.15, h5dump, h5repack, Python 3, h5py, numpy, ipyparallel

## Earth Science Application and Data

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Application:</b> Dynamic data subsetting and statistics aggregation using selected oceanographic data	<b>Application:</b> Dust storm analysis framework consisting of dust storm feature identification, attribute calculation, and object tracking.	<b>Application:</b> Identify grid cells meeting blizzard conditions using (imprecise) NWS definitions. Identify blizzard events using spatio-temporal CCL and appropriate statistics. Compare results with observed data.	<b>Application:</b> Climatology research to enable simple canonical operations including subsetting, averaging, searching for minimum and maximum values, etc.	<b>Application:</b> Supporting multiple applications and various data sets
<b>Data:</b> GHRSSST Level 4 CMC 0.2° Global Foundation Sea Surface Temperature Analysis. Grid size: 1800x901	<b>Data:</b> Non-hydrostatic Mesoscale Dust Model (NMM-dust) from NCEP, simulating dust event in Phoenix, Arizona during 3rd and 4th of July 2014.	<b>Data:</b> Modern Era Retrospective Analysis for Research and Analysis (MERRA)	<b>Data:</b> Modern Era Retrospective Analysis for Research and Analysis (MERRA)	<b>Data:</b> • NCEP/DOE Reanalysis II, for GSSTF, Daily Grid, v3 Spatial: 0.25°x0.25°, global Temporal: 1987-08, daily • NOAA Coral Reef Temperature Anomaly Database Spatial: ~4km global Temporal: 1982-12, weekly
<b>Subset Used:</b> Spatial span: 50 x 50 grid Temporal span: 4 months Size: 2.43 GB	<b>Subset Used:</b> Horizontal resolution: 3 km with 45 vertical levels in the vertical and the Vertical Resolution: Between 2.5 KM and ~5 KM Time Resolution: 3 hours.	<b>Subset Used:</b> Spatial resolution: 3/5°x 1/2° Hourly resolution: Hourly	<b>Subset Used:</b> MERRA data for northern India/Pakistan, North China Plain, California Central Valley, and Nile Valley Size: ~ 132 GB	<b>Subset Used:</b> Full data set Size of NCEP/DOE Reanalysis2 ~ 17GB NOAA Coral Reef temperature data ~ 24MB

## Data Ingestion and Workflow

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Time to convert to .adm file:</b> ~7541 sec ~ 32 sec/MB	<b>Time to convert Geotiff file:</b> ~1 sec/MB Can be parallelized	<b>Time to convert file:</b> Pull data on the fly (OpenDAP); Write (1-D) binary data; Load 1-D and redimension): ~0.3 sec/MB Can be parallelized	<b>Time to convert to HDFS file:</b> Involves sequencing, mapping, and then using Bloom filter for reducer ~0.4 sec/MB Can be parallelized	<b>Time to convert HDF file:</b> Data sets are re-chunked and compressed ~ 0.2 sec/MB Can be parallelized
<b>Disk space required:</b> Raw data ~ 235 MB AsterixDB format ~2.43 GB 10 fold increase in disk space requirement	<b>Disk space required:</b> Raw data ~ 45 MB AsterixDB format ~77 MB Less than 2 fold increase in disk space requirement	<b>Disk space required:</b> ~2.5 fold increase in disk space requirement	<b>Disk space required:</b> ~2.5 fold increase in disk space requirement	<b>Disk space required:</b> 63% reduction in file size
<b>Workflow:</b> 	<b>Workflow:</b> 	<b>Workflow:</b> 	<b>Workflow:</b> 	<b>Workflow:</b> 

## Data Preparation/Preprocessing

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Operations</b> • Subsetting (50x50 chunks) • Sorting (**bugs in current version)	<b>Operations</b> None: Data ingested in Geotiff format and entire data set used.	<b>Operations</b> • Subsetting • Table join • Constructing alternative representation	<b>Operations</b> • Write custom NetCDF to Hadoop convertor to keep files as sequence files • These files sent to Hadoop for storage • Hadoop splits and distributed the sequence files across HDFS; builds index for Hadoop access • Maintains NetCDF metadata for each file	<b>Operations</b> • HDF5 dataset chunks with all-missing data not stored during the data ingest stage. • No subsetting; entire temporal and spatial extent of data is used after ingestion • Sorting only in the temporal domain, if required, to ensure monotonic order of the temporal axis. • Data are indexed by calculating descriptive statistics for each HDF5 dataset chunk. • Initial data files are collated into a single file with optimized HDF5 dataset chunking/compression.
Operations are parallelizable	<b>Bottlenecks</b> • Physical memory and disk I/O are main performance bottlenecks • Performance can be really slow when subsetting portions of source images	Operations are parallelizable <b>Bottlenecks</b> • Operations are local, little or no data exchange • Performance bottlenecks are being investigated	Operations are parallelizable <b>Bottlenecks</b> • Processing is offline mode – not useful for adhoc queries • Very large and skewed data causes memory issues both at mappers and reducers	Operations are parallelizable <b>Bottlenecks</b> • File granularity (inefficient to copy same file to multiple nodes if number of nodes > number of files). • One processor performs aggregation of results, could result in bottleneck depending on data.

Acknowledgement: This research is supported by funding from the NASA ESTO-AIST Program.

## Data Analysis

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Primitives Tested</b> Standard statistics computation (mean, std deviation) • Calculations cannot be performed on 50x50 chunks due to unresolved bugs in software. Computing similarity or distance between every pair of records • Supports edit distance (on strings) and Jaccard coefficient (on sets) Largest connected subgraph search • Can be integrated with Pregel for graph computation.	<b>Primitives Tested</b> Extract individual dust storm object (region-growing based algorithm) Mean computation • Done using queries  Example queries -Select a single pixel from all images ~66msec -Select a subset from all images ~1sec Select mean value of each band of a single image ~0.3sec Select mean value of each band across all images ~4.5sec	<b>Primitives Tested</b> Use of Bloom filter to speed up Hadoop jobs by leveraging the probabilistic search capability. Speed up by 30-80% obtained Example performances: • 83.9% efficient for reading a single parameter ("T") from a single sequenced monthly means file • 29% efficient for single MR job across 4 months of data seeking "T" (period = 2)	<b>Primitives Tested</b> Use of Bloom filter to speed up Hadoop jobs by leveraging the probabilistic search capability. Speed up by 30-80% obtained Example performances: • 83.9% efficient for reading a single parameter ("T") from a single sequenced monthly means file • 29% efficient for single MR job across 4 months of data seeking "T" (period = 2)	<b>Primitives Tested</b> Standard statistics computation (mean, std deviation) • Calculation performed on original data (as obtained from the archive) ~Single node: 5.4 sec/GB ~50 nodes: 0.05 sec/GB Clustering • Searching for points/regions based on a set of temporal, spatial, data value conditions Data subsetting • Slicing: selection based on temporal, spatial, data value criteria

## Result Reporting

AsterixDB	Rasdaman	SciDB	Hadoop	HDF
<b>Graph Plotting:</b> • Pregel supports parallel graph computations using the Pregel programming model • Query results in json format can be used as input to visualizations (software or web visualizations) • Cannot be used for plotting figures with overlay for showing results on the Earth's grid • Cannot be automated	<b>Graph Plotting:</b> • Provides several Open Geospatial Consortium (OGC) standard interfaces through its web services wrapper, Petascope. • Can be used to plot figures with geographic overlays • Plotting can be automated but using spatial/temporal indexing which would require Petascope to store the temporal and spatial metadata	<b>Graph Plotting:</b> • Plotting is possible using "shim" SciDB client to interact with external tools like SciDB-py and SciDB-R • Exporting data also possible (for other external tools) • Figures with overlay can also be plotted using SciDB-py or SciDB-R. • Automation can be done via scripting around SciDB-Py or SciDB-R	<b>Graph Plotting:</b> • Plotting is possible by exporting data to standard formats and using external plotting software • Visualization tool IDL can be used to visualize and diagnose data stored in the native Hadoop file format, HDFS • Process can be made faster by using parallel reader for data ingestion before visualization	<b>Graph Plotting:</b> • Hdf file formats allows storage of meta information that can be used for plotting results using overlays.

## Observations

- AsterixDB is inefficient for big data applications because its storage format requires 10x more disk space than raw archive format. Current version has many bugs.
- Hadoop requires significant parameter tuning for optimal performance and has high bandwidth requirements.
- Most containers suffer from parallelization bottlenecks due to aggregation/merging of results at a single node
- HDF files can cause issues during concurrent read/copy in multicore architectures
- Rasdaman can be slow for large I/O operations and inefficient for big data applications. Also development support for Rasdaman is also low compared to some other containers
- SciDB data format is not compatible with other common big data processing frameworks thereby requiring duplicate data storage.

### Additional Contributors

John Thompson, NASA Goddard Namrata Malarout, NASA JPL	Fei Hu, George Mason University John Ready, HDF Group	Aleksander Jelenak, HDF Group Amidu Oloso, NASA Goddard
---	--	--