## Open source software for Big Data : Experiences in indexing and browsing geo-archival records

## Authors

- **Jefferson Robert Heard**, RENCI, UNC Chapel Hill, *United States*

- **Richard Marciano,** UNC Chapel Hill, *United States*

### Introduction and Related Work

Archival records are not like records in a database. The original structure, attributes, and metadata of an archival record are as important to archival integrity as the record itself. These attributes include things like file structure, directory hierarchy, and file permissions. Organic, human structures do not always match up well with storage and retrieval patterns. Thus an archive of hundreds of millions or a billion records becomes a big data problem.

Additionally, certain kinds of records are difficult to browse and search. Much work has been done in making interfaces for searching text. Less work has been done on making geographic records searchable or browsable. In this talk, we present open source approaches to indexing and browsing geographic records in a large archival collection.

Some amount of work has been done in browsing collection of archival data. In particular, treemaps have been employed [1] along with metadata visualization to enable browsing of archival metadata and records. Large amounts of work have been done on how to effectively store archival data. In particular, IRODS [2], and Globus Online [4] have emerged as "Data Grid" software that provide rich capabilities to a user interested in preserving a dataset.

In this presentation we will detail our experiences with developing a system to scalably index and browse geographic records in an archival setting.

## Approach

Our archival collection is known as the CI-BER Testbed. It is a still-growing 70 million file, 41TB collection consisting of data from across 133 different agencies of the US government. The CI-BER Testbed was developed to create a system for testing the scalability of archival systems across many kinds of heterogenous data, from files consisting of gigantic chunks to directories containing hundreds of thousands of files, to deeply nested structures. The CI-BER Testbed is housed on an iRODS data grid that has been federated with the US National Archives and Records Administration (NARA)'s own data grid.
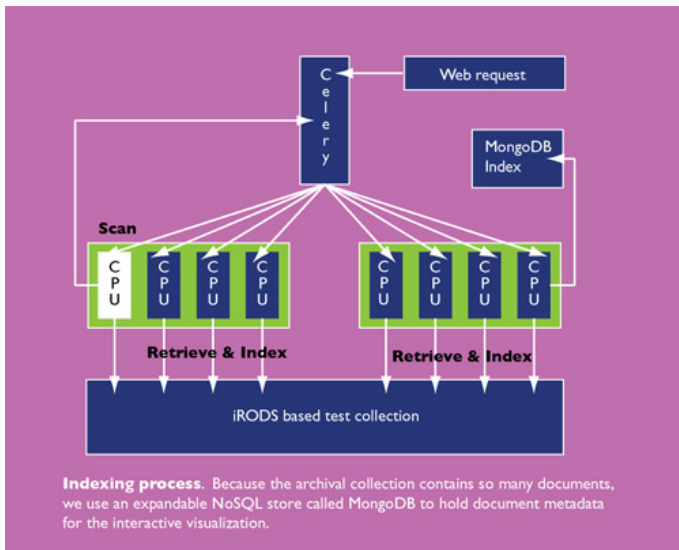
IRODS, the Integrated Rule Oriented Data System is an open source data grid software originally developed by the DICE Group. It provides a filesystem abstraction that allows for arbitrary file, directory, and file-system level metadata, a rich permissions system, user-defined "rules" that govern how data is collected into the filesystem, and "microservices" that can operate as independent agents on a collection. IRODS is used extensively in the academic and archiving world to manage large data collections.

We have built a distributed indexing system that authenticates with this data grid and crawls individual collections looking for geographic files. An administrator sends a web request to index a particular collection, then the indexer uses the user's IRODS credentials to mount the IRODS collection and crawls that collection for filenames. The filenames are then farmed out as individual tasks to each of the worker threads. Each worker thread attempts first to open the file with GDAL and then OGR [3]. If either of these works, then the worker thread passes the opened file to a number of metadata extraction functions and the metadata, along with the file, is added to the index.

For the index, we selected an open-source "NoSQL" database, MongoDB [5]. MongoDB was selected both for its scalability to large amounts of data and its

schema-free nature. We wanted the index to be extensible, so that alternative methods for metadata extraction could continue to append metadata to each individual record. These metadata appendices to each record could then be visualized independently for the records that have them. Core to our metadata, however, are the geographic boundaries of a particular dataset, translated into WGS84, the name of the dataset, and for coverages, the number of bands, or for feature sets, the name of layers and for each layer the name and type of record fields.
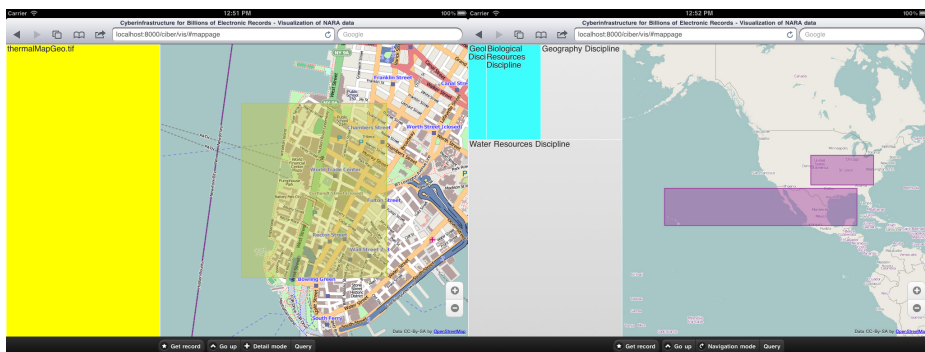
Directory paths within the archives have metadata from the individual records aggregated and appended. This results in each directory level having a boundary box associated with it as well as the number of archival records contained in that directory.



**FIGURE 1**
The indexer

This MongoDB index is used to power our visualizations. Our main visualization is a browsing interface for the records based on the open source JQueryMobile [6], D3 [7], and OpenLayers [8]. In the visualization, we extend the treemap for use specifically with geographic records by adding a geographic map to the interface and unifying the two.

The user is initially presented with an overview. The overview centers the map on the boundary of the entire collection set. The first level of the treemap contains one cell for each indexed collection. Treemap cells are sized by the number of actual record items contained within the collection. For levels containing leaf nodes in the treemap, each cell is sized by the physical area of the bounding box (in square meters) relative to its peers. Each collection has its own individual bounding box on the map. If the user touches a bounding box, it highlights the corresponding object in the treemap. If the user touches the treemap, the collection is descended into and the app retrieves the next level of the index.



**FIGURE 2**

Detail and overview showing the treemap's grouping of records and their location on the map

Additionally, there is a "detail mode." If the user touches on a treemap cell the app will center and expand on the cell's bounding box. Tapping "get record" retrieves the metadata record associated with the bounding box. For an individual record, this is the geographic metadata associated with the file. For a subcollection, this is the bounding box and the number of records in the subcollection. That metadata record is visualized and a link to the record itself is provided.

## Future Directions

In the near future, we intend to incorporate this visual browsing interface into larger, interactive mapping applications. These applications will serve to encourage community involvement in archives and a view of archives as a first-class Open Source medium.

## Acknowledgements

[1] ESTEVA, M., XU, W., JAIN, S.D., LEE, J.L., MARTIN, W.K., Assessing the preservation condition of large and heterogeneous electronic records collections with visualization, *International Journal of Digital Curation 6*(1):45-57 Feb. 2011.

[2] ZHU, B., MARCIANO, R., MOORE, R., HERR, L., SCHULTZ, J., Digital Repository: Preservation Environment and Policy Implementation, Springer Verlag 2012, *International Journal on Digital Libraries (IJDL).*

[3] THE OSGEO GROUP. GDAL and OGR, 2012. http://www.gdal.org/

[4] THE GLOBUS GROUP. *Globus Online*, 2012. http://www.globusonline.org/

[5] 10GEN, INC. *MongoDB*, 2012. http://www.mongodb.org/

[6] THE JQUERY GROUP. *JQueryMobile*, 2012. http://jquerymobile.org/

[7] BOSTOCK, M. D3.js. *Data Driven Documents,* 2012. http://www.d3js.org/

[8] OPENGEO, ENC. *OpenLayers*, 2012. http://www.openlayers.org/