# REAL TIME DATA WAREHOUSING AND ON LINE ANALYTICAL PROCESSING AT ABERDEEN TEST CENTER

# REAL TIME DATA WAREHOUSING AND ON LINE ANALYTICAL PROCESSING AT ABERDEEN TEST CENTER

**Michael J. Reil**
SFA Inc. / Aberdeen Test Center


**T. George Bartlett**
Aberdeen Test Center


**Kevin Henry**
Sverdrup Technology Inc. / Aberdeen Test Center

## ABSTRACT

This paper is a follow on to a paper presented at the 2005 International Telemetry Conference by Dr. Samuel Harley et. al., titled *Data, Information, and Knowledge Management.* This paper will describe new techniques and provide further detail into the inner workings of the VISION (Versatile Information System – Integrated, Online) Engineering Performance Data Mart.

## KEY WORDS

VISION, HDF5, data warehouse, data mart, On Line Analytical Processing (OLAP)

## INTRODUCTION

VISION embodies a comprehensive, holistic, top down approach to information collection, management, and ultimate transformation into knowledge. The Aberdeen Test Center has developed VISION to meet the increasing information demands that arise from the testing of complex military combat systems. One of the data marts managed under VISION is the Engineering Performance Data Mart, which was designed to manage developmental and operational engineering test data. This data mart will be the topic of this paper. Figure 1 illustrates the overall VISION data management scheme, with the Engineering Performance Data Mart circled.

## VISION Data Management
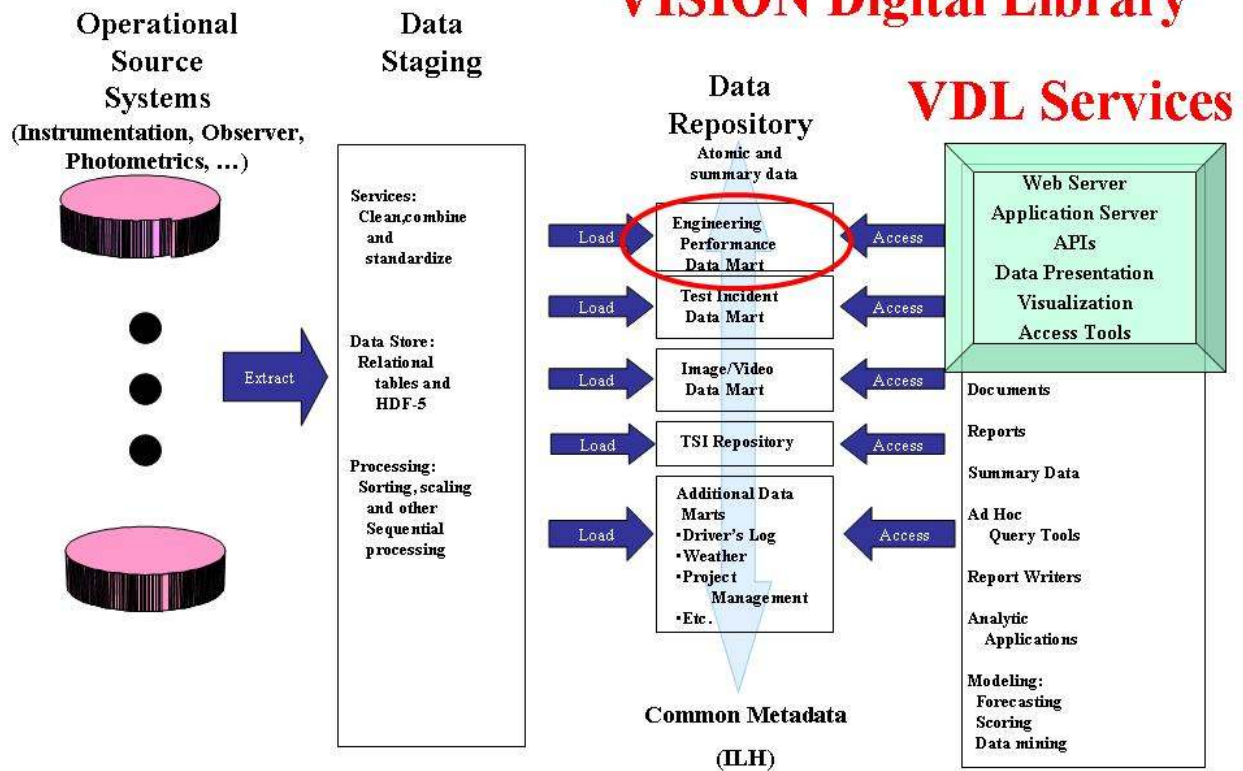
### VISION Digital Library



**Figure 1 VISION Data Management**

## PRE TEST SETUP

To gain the maximum ROI (Return on Investment) from test data, it must be tagged with all applicable metadata. Data consumers, which include analysts, modelers, simulation event planners and as yet unknown future users, will be far disconnected from the actual test that produced the data. Their only reliable means of placing the data into context is via the associated metadata. The VISION instrumentation suite of Advanced Distributed Modular Acquisition System (ADMAS) enforces this requirement by means of an XML based configuration file. The configuration file serves multiple purposes:

- Program the ADMAS for data collection.
- Physically couple the metadata with the data (if you have one, you have the other).
- Control the post-processing of the data for entry into the data mart.

A 'Configuration File Builder' application has been developed to assist the test engineer with this task. This application pulls what metadata it can from the project administrative databases, making it easier for the test engineer to enter correct values. There are a few required pieces of

metadata, including Project ID, System ID, Test Item ID and Test Location.  Figure 2 shows a screenshot of the 'Configuration File Builder' application.
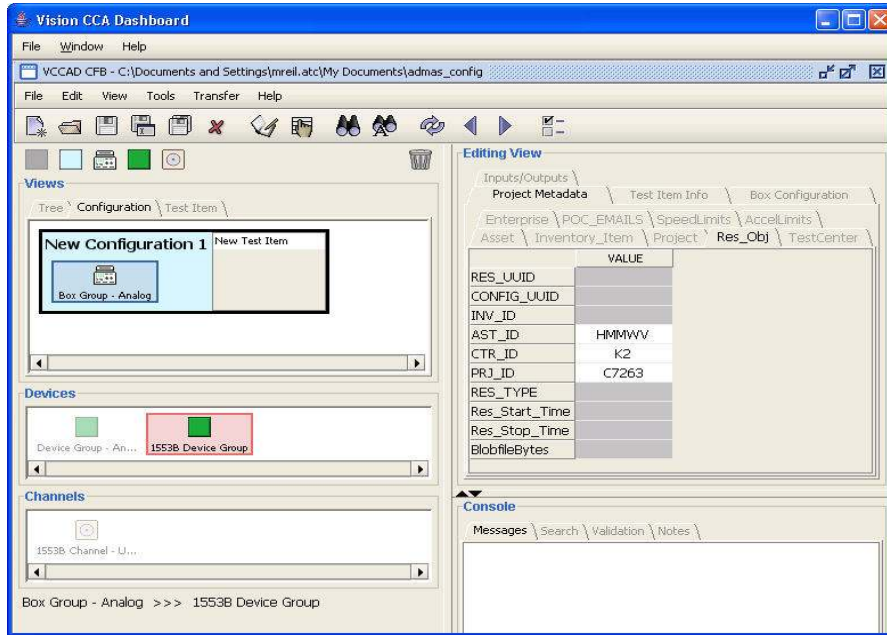


**Figure 2 Configuration File Builder**


## RAW DATA FILE COLLECTION


During the test event, the ADMAS records applicable 'Run Time' metadata such as date/time, GPS position, operator's name etc.   The pre-test entered metadata, the engineering data itself, and any run time collected metadata are stored in a single file (a BLOb – Binary Large Object) by the ADMAS.  The format of the raw data is optimized for the real-time requirements of the ADMAS.  The metadata however, is in the well known XML format at the head of every ADMAS BLOb file.  Figure 3 shows an ADMAS and the layout of all ADMAS BLOb files.
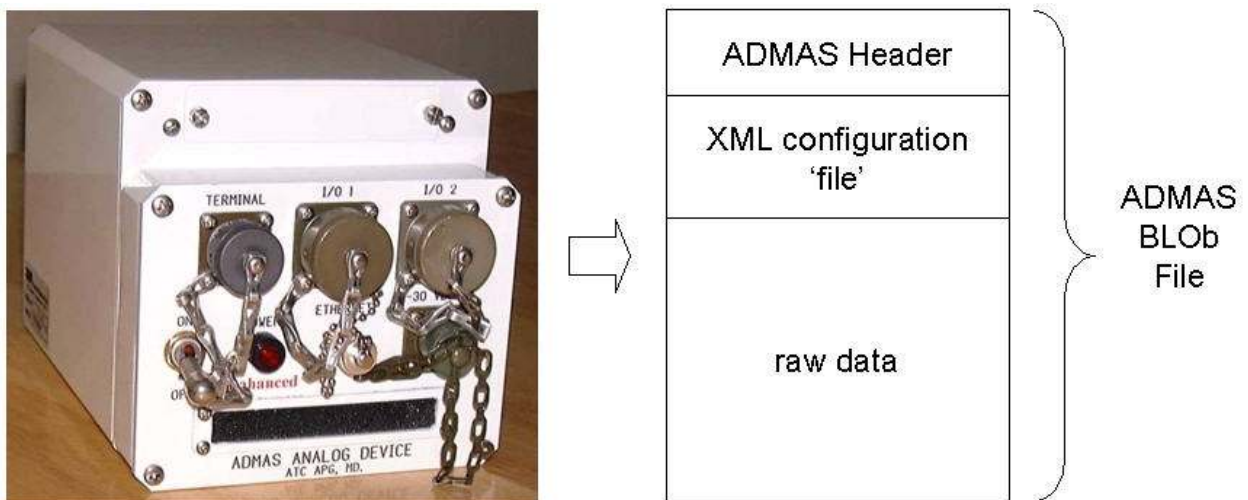


**Figure 3 ADMAS and BLOb file layout**

## DATA HARVESTING

ADMAS' record BLOb files on PC card flash drives.   BLOb files are harvested from the ADMAS via either wireless LAN, hard wire LAN or by physically removing the flash cards from the ADMAS and inserting them into a networked PC.  The raw BLOb files are transferred via SSL to the data mart for processing.  The method of transport is not really important, and can be customized to suit each range's unique networking requirements.  The only real requirement is that the files somehow make it to the data mart.   Figure 4 shows the generic process.
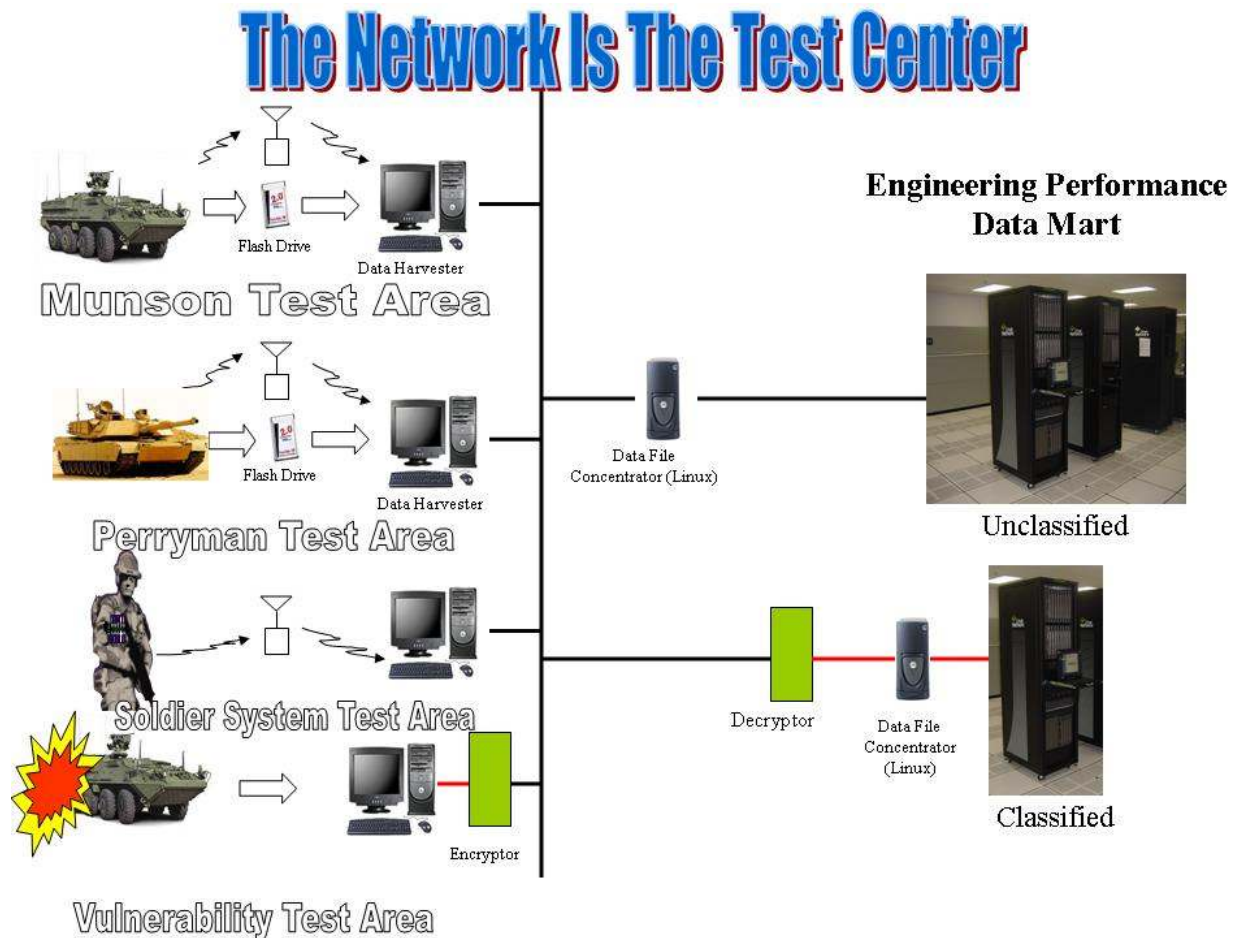


**Figure 4 Data Harvesting**

## DATA TRANSFORMATION AND HDF5

  Unattended, triggered processes run on the data mart to convert the raw data to processed data, which is stored in HDF5 files.  HDF5 (http://hdf.ncsa.uiuc.edu/HDF5/) is developed and maintained by NCSA, and is ideally suited for storing large scientific data sets.  Generally, one HDF5 file is produced for each raw data file.  The processing instructions in the raw file's XML

configuration section control what additional processing steps are done, such as validation, reporting etc.

The processes are triggered by the arrival of new raw data files using 'AntFlow' (http://onionnetworks.com/products/antflow/), which is an open source framework for workflow automation based on hot folders.  Ant (http://ant.apache.org/) is a java based open source framework for creating tasks.  These tasks are run by AntFlow when new files appear in the hot folders.  The tasks run the java applications, which do the converting, validating and reporting.  A typical data flow process is shown in figure 5.  This follows the 'Extract, Transform and Load' data warehousing technique (http://en.wikipedia.org/wiki/Extract,_transform,_load).
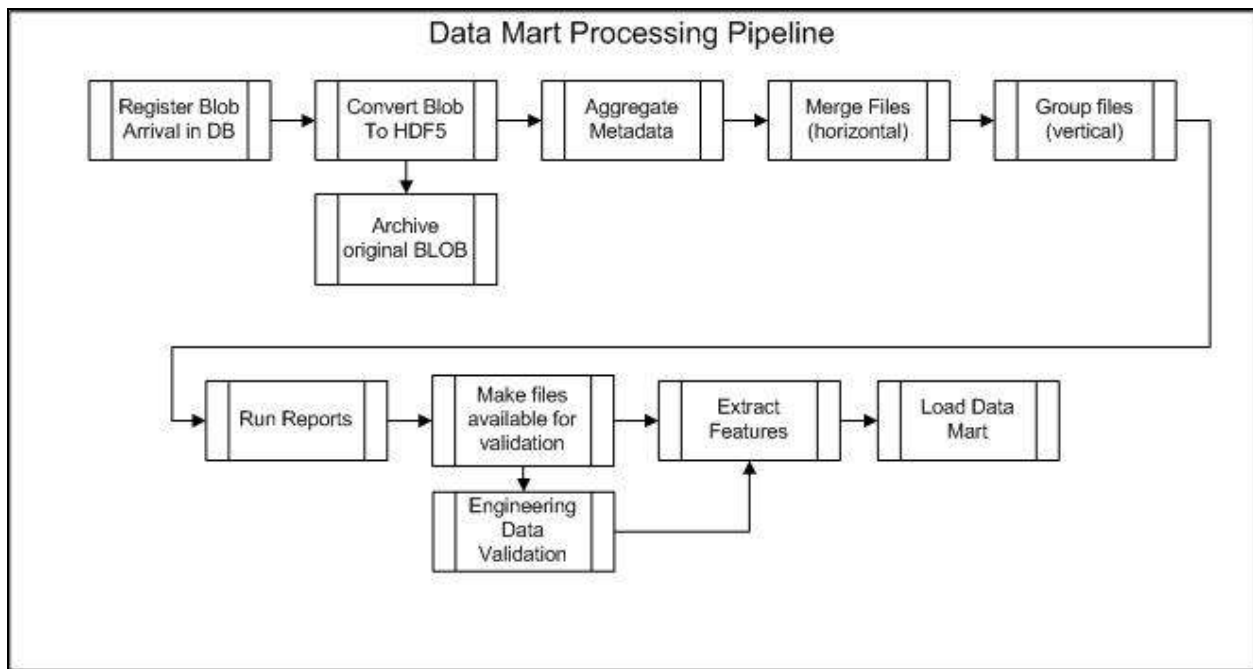


**Figure 5**


## DATA MART POPULATION

The Engineering Performance Data Mart consists of both a relational database and a collection of HDF5 files.  All of the metadata and none of the data are stored in the relational database.  The data is kept in the HDF5 files, with a pointer from the relational database to the individual HDF5 files.  This technique keeps the database small, and therefore queries against it run very fast.  Loading of the data mart is done via a java application that reads the metadata from the HDF5 files, connects to the relational database, and issues the appropriate SQL "Insert …"statements to populate the tables.  The HDF5 files are then moved to the distributed storage pool, where the OLAP tools can find them.

**DATA MART PHYSICAL MAKEUP**

The Engineering Performance Data Mart is comprised of two Linux servers, a Linux cluster, and a large distributed storage pool.  The data mart components consist of the following:

- Hardware
    - One Opteron dual processor CentOS-4 database server
    - One Opteron dual processor CentOS-4 application server
    - Distributed File System
        - One Opteron dual processor CentOS-4 Lustre metadata server
        - Six Opteron dual processor CentOS-4 Lustre object storage servers with 6.4 terabyte storage each (38.4 TB total).
    - One Linux cluster comprised of one head node and 32 slave nodes, each Opteron dual processors.  Gigabit Ethernet interconnect.
- Software
    - In house developed java classes
    - PostgreSQL 8.1 (http://www.postgresql.org/)
    - Apache Tomcat 5.5 application server (http://tomcat.apache.org/)
    - Lustre 1.4.6 distributed file system (http://www.lustre.org/).

Note that the software stack of the data mart is made up of open source software and the hardware stack is using 64 bit processors. The cluster is used to populate and query the data mart using distributed parallel processing. The Lustre file system stores the HDF5 files.


**ON LINE ANALYTICAL PROCESSING TOOLS**

A web based OLAP toolbox has been developed to allow querying of the data mart.  Access to the Engineering Performance Data Mart – which is a web based application – is via the VISION Digital Library (VDL) - as is access to all VISION data marts.  This single sign-on allows the project coordinators to control and track access to the data marts.  By logging into VDL and navigating to the project of interest, the user would click on the data mart URL (assuming permission), thus launching the OLAP toolbox in his browser.  Tools have been developed which allow the user to view aggregate composites, or to drill down to individual data files and parameters.  The user builds queries to submit to the data mart via the java GUI, by selecting metadata values that are presented as pick lists.  The more metadata that was collected before and during the test, the better the OLAP user will be able to define his/her query.  Some of the things that users can do include:

- View all of the collected metadata tags, and all of the different values used for these tags
- View how many files are available, and what time period they cover, for any combination of metadata
- View a GPS map of the location of the item under test for any combination of metadata
- Apply metadata filters to the query results
- View time series traces of all parameters
- Create custom plots of parameters from a pick list

- Save plots as JPEG files on users computer
- Download raw files, HDF5 files or Excel spreadsheets of selected data.
- Create composite histograms, time in limits, threshold crossing or trend reports across multiple files.
- Download Google Earth (http://earth.google.com/) files for displaying GPS locations during the entire test, or just select portions of the test.

Using the OLAP GUI, and the power of the Linux cluster, all of these tasks can be initiated by the user, run on the cluster, and the results sent back to the user in seconds. The vast majority of queries return to the user in under 5 seconds. To test the performance of the data mart, we ran a composite histogram query of an acceleration parameter across all files in an 8 month duration test (38000 data files, with 5.5 billion samples). The histogram plot was displayed in less than 2 minutes. Screenshots of the OLAP GUI are shown in figures 6 through 9.
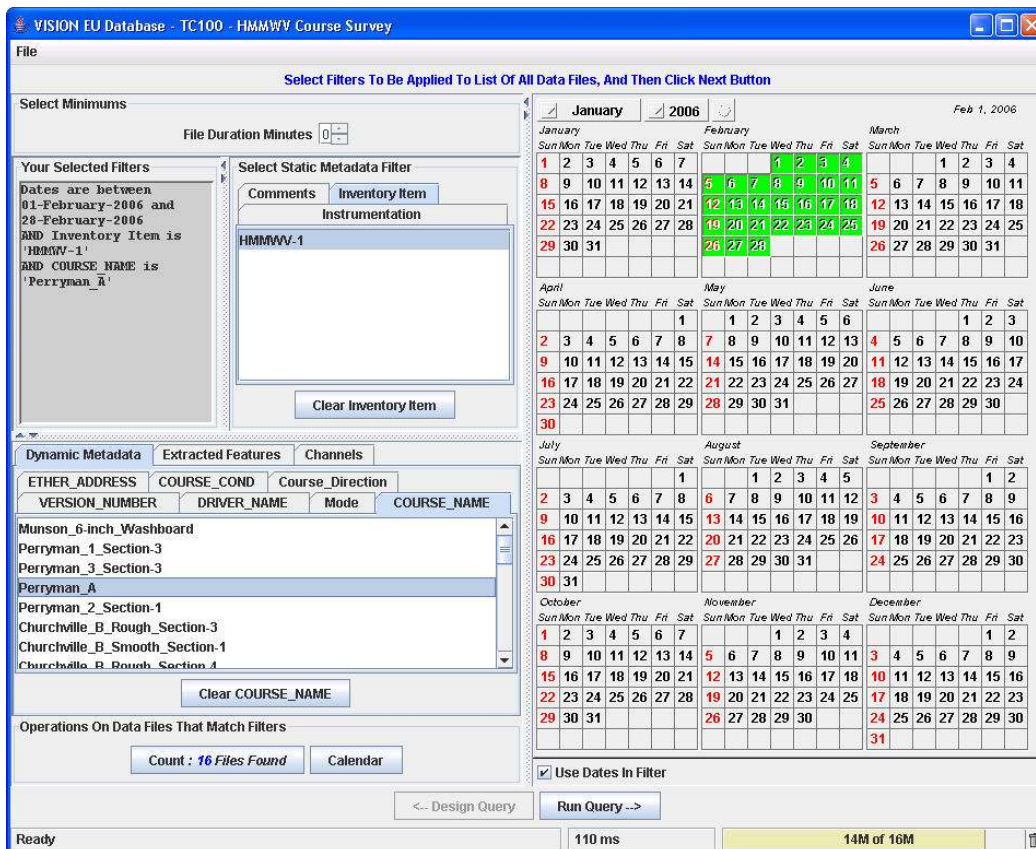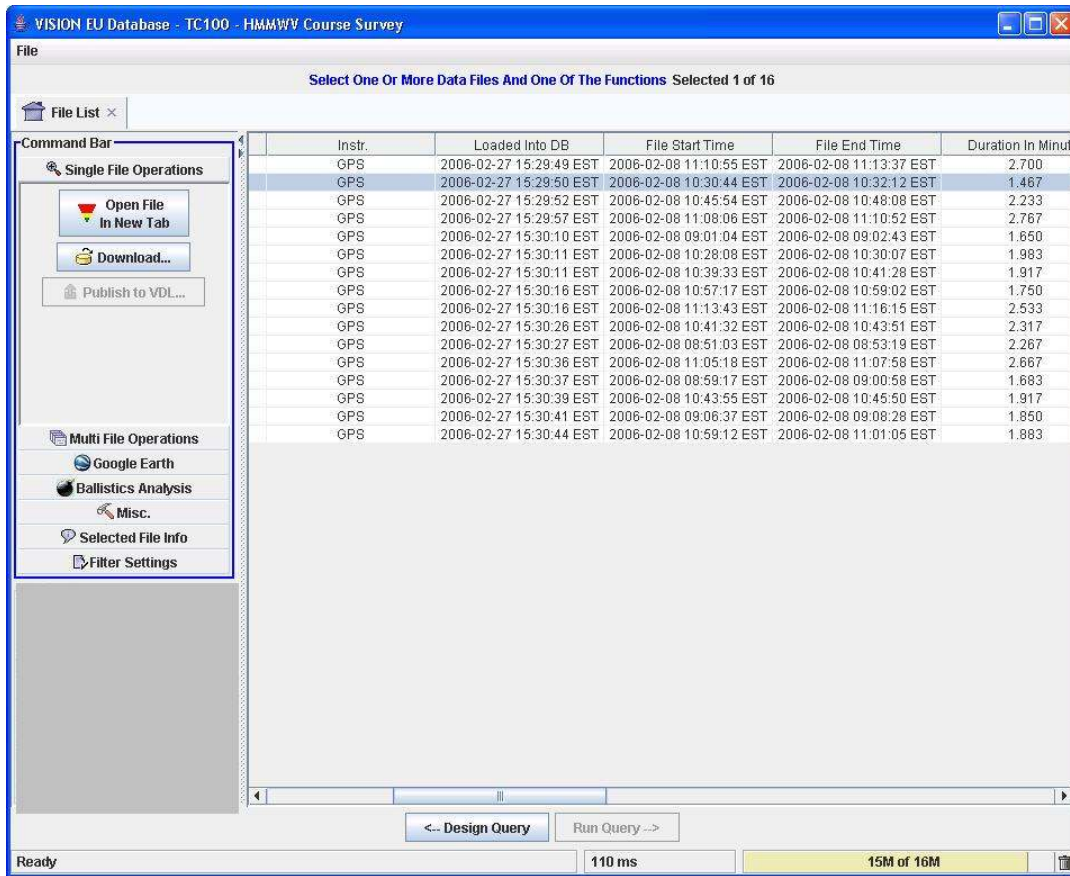


**Figure 6 – Metadata filter page**

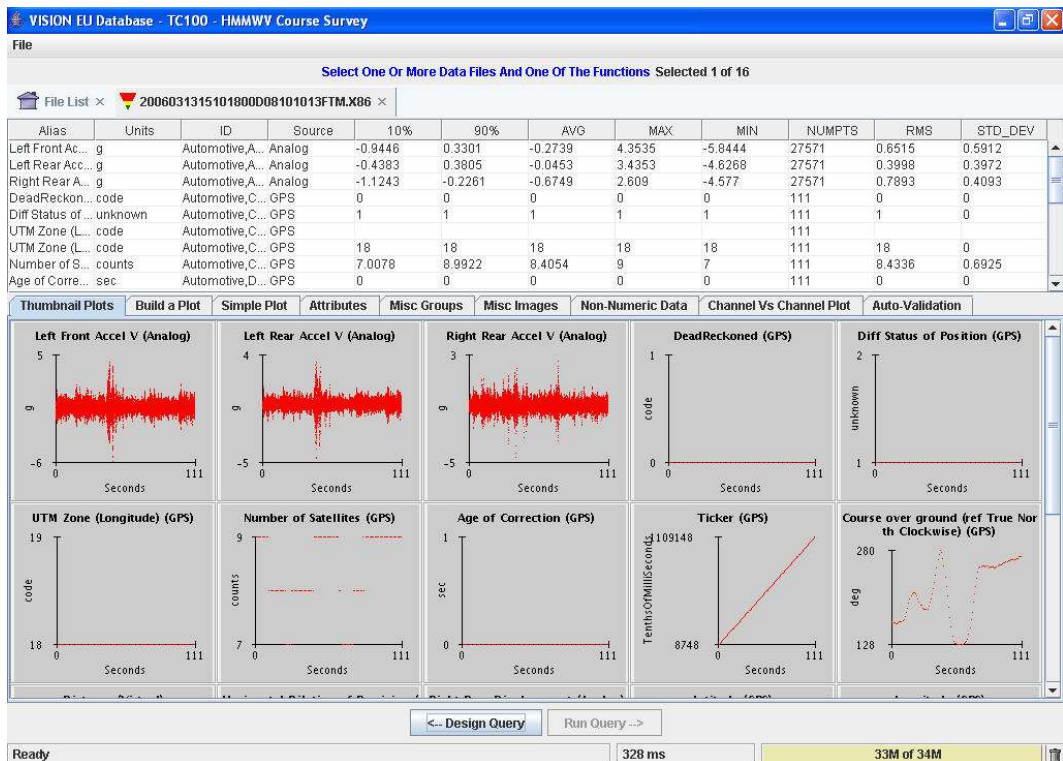**Figure 7 – List of files returned from query**

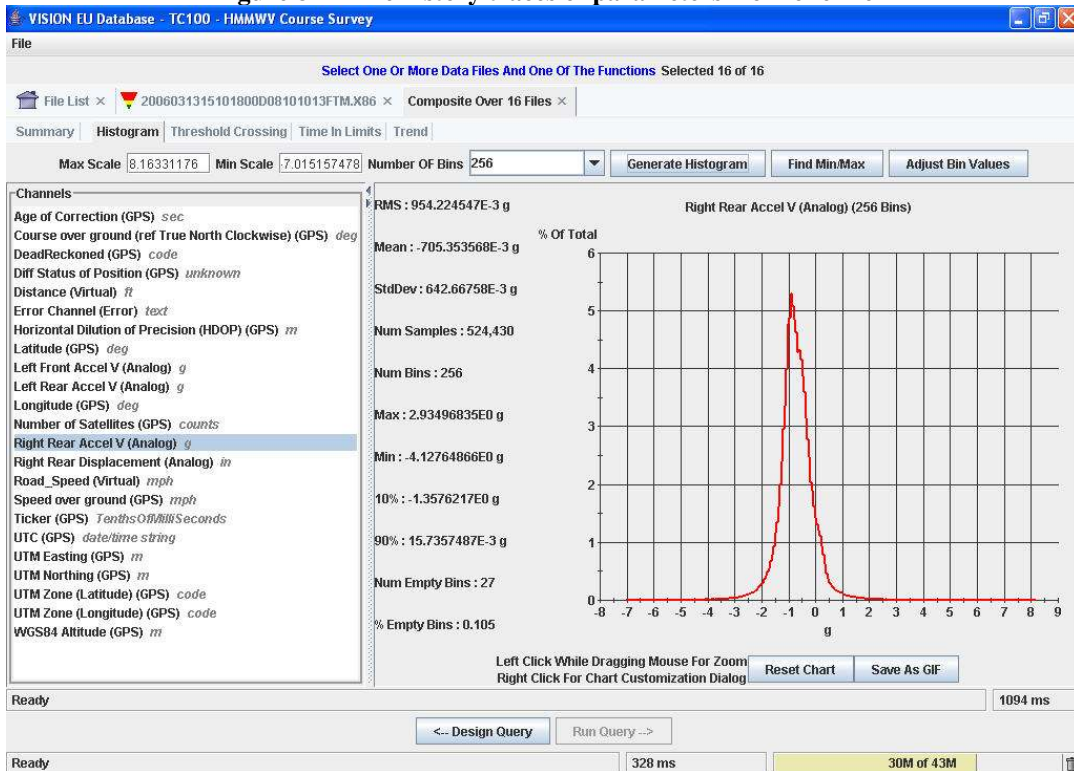**Figure 8 – Time history traces of parameters from one file**



**Figure 9 – Composite histogram of one parameter across 16 files**

**CONCLUSION**

The VISION Engineering Units Data Mart at the Aberdeen Test Center contains both developmental and operational test data for over 150 Army, Navy and Marine Corp projects, dating from 1999 to the present.  The metadata totals 100 gigabytes and the HDF5 data totals 1.5 terabytes.  The data mart has a flexible architecture, allowing for configuration and deployment in a number of different ways.   As an example, a remote data mart was set up at "internet challenged" sites at Camp Pendleton, CA and Marine Corp Base Hawaii to support testing of the Expeditionary Fighting Vehicle (EFV) in 2004 through 2005.  Because of the flexible configuration and deployment options of the system, the entire data mart was able to run on a single Dell workstation with one 2TB external network attached storage device.   The use of well defined interfaces at all of the software tiers allows for additional features to be developed and plugged in with little or no rework to the infrastructure.   As an example of this, in 2001 the Department of Transportation wrote a custom graphical analysis component (in java) for one of their projects that conformed to the analysis interface that we published.  They gave us the component, we plugged it into the OLAP GUI, and they were able to use it immediately.  This allowed them to leverage the data mart's metadata querying capability into their own analysis tool.  Through the use of open source software, notably HDF5, PostgreSQL and Linux along with in house developed cross platform software components (java based), the VISION Engineering Units Data Mart is in an excellent position to meet the unknown challenges of future generations of testing.

## ACKNOWLEDGMENTS

## REFERENCES

Harley, Samuel F., Data, Information, and Knowledge Management, International Telemetry Conference, 2005