



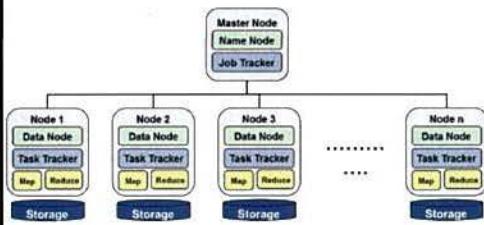
Development of U.S. NDC Performance Metrics Through Large Scale Analysis of System Log Files With Hadoop Distributed File System Based Tools

**W. N. Junek, C.A. Houchin, J. A. Wehlen III, A. Poffenberger, and R. C. Kemerait
Air Force Technical Applications Center, Patrick Air Force Base, Florida**

Abstract

Optimizing the performance of the United States National Data Center (U.S. NDC) geophysical data processing pipeline system is critical for efficiently monitoring international compliance to nuclear test ban treaties. The U.S. NDC software stores system performance information for each data processing event in a collection of semi-structured alphanumeric log files. On average, the system generates 140,000 log files per day, which are stored in different directories. Currently, acquisition of system specific performance information or isolation of error messages must be pursued from each log file individually. This manual parsing process is time consuming and often leads to incomplete log files in JavaScript Object Notation (JSON), which is a highly structured data format that can be easily parsed. Here, we show how U.S. NDC system performance information can be parsed from JSON files and analyzed using a Hadoop Distributed File System based tools such as HIVE, Zeppelin, and IP Spark. The Hadoop architecture is designed to store and process large alphanumeric datasets, which can be easily scaled to accommodate our continuously growing collection of performance information extracted from U.S. NDC log files.

Hadoop Distributed File System Architecture Overview

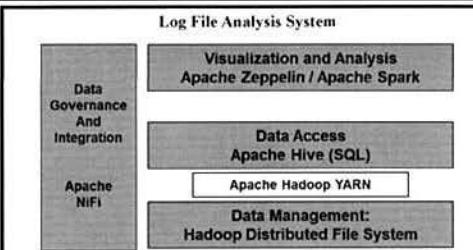


- **Hadoop Distributed File System (HDFS):** Java based file system that provides scalable and reliable data storage that is designed to span large clusters of commodity servers.
 - **MapReduce:** Framework for processing parallelizable problems across a large number of nodes. Takes advantage of data locality, processing it on or near the storage asset to reduce the distance it must be transmitted.

U.S. NDC HDFS System Specifications



- Nodes: 12 Dell PowerEdge Servers, 6x3Ghz Cores, 8x4TB SAS Drives, 128 GB RAM

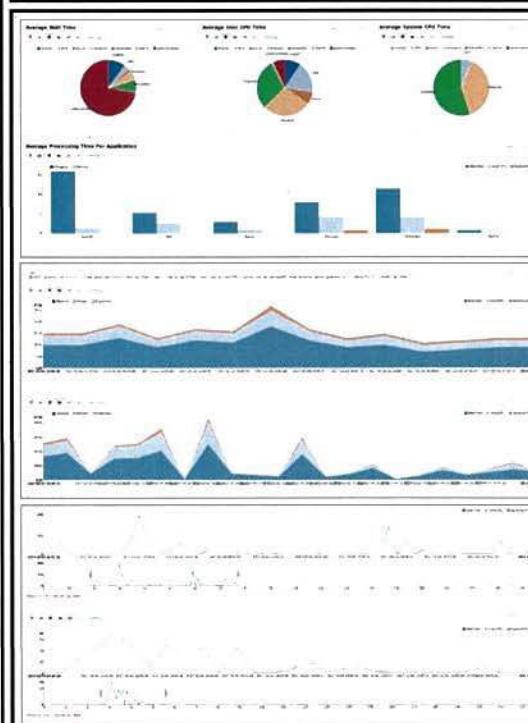
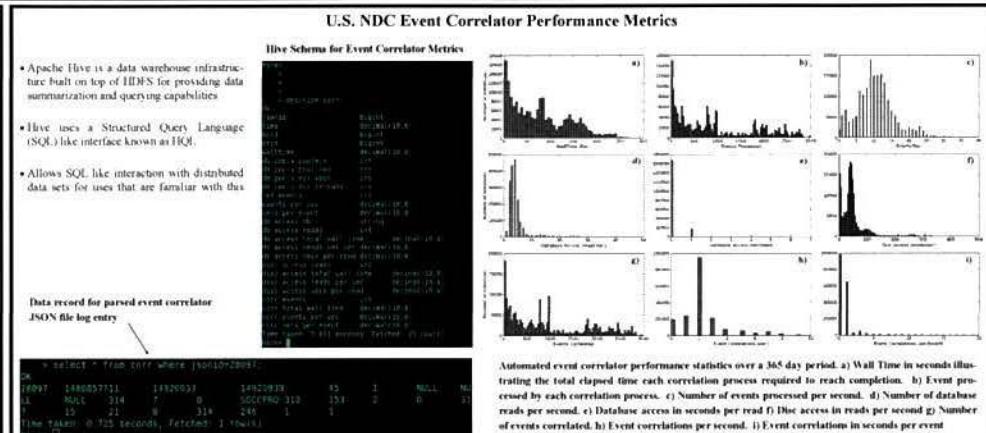
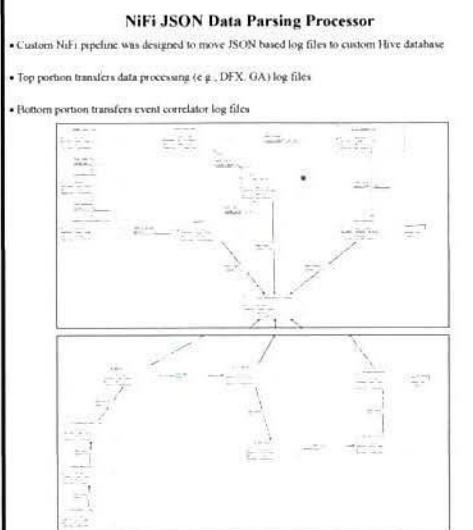


- US NDC System generates over 140,000 log files per day
 - Manual parsing of log files is time consuming, error prone, and leads to incomplete data sets
 - HDFS systems are ideal for ingesting, storing, and analyzing large scale alpha numeric datasets

JSON Data File Format

- U.S. NDC log files store system performance information in numerous unstructured text files
 - The U.S. NDC software now outputs system log files in JavaScript Object Notation (JSON)

U.S. NDC HDFS System Specifications



Web Browser Based Performance Metrics Dashboard

Apache Zeppelin

- Web-based notebook for interactive analytics
 - Multi-purpose toolbox for data ingestion, analytics, visualization
 - Supports multiple languages including Apache Spark and Python
 - Renders numerous chart styles to create interactive output
 - Share notebooks among collaborators

Apache Spark

- Engine for large-scale data processing
 - Runs on Hadoop and can access HDFS
 - Fast in-memory computing platform
 - Supports SQL, streaming and machine learning libraries
 - Application development in Java, Scala, Python, and R

Apache NiFi

- Flexible data collection, transformation, and routing platform
 - Supports numerous formats, schemas, and protocols
 - Web-based user interface for design, control, feedback and monitoring
 - Data Provenance for real-time visual chain of custody and metadata
 - Real-time streaming data: send, ignore, or store received data

- The U.S. NDC has acquired a HDFS cluster to study the large scale data acquisition and analysis
 - Modified US NDC system to output log files in JSON format
 - Developed HIVE based data model for archiving log file data
 - Analyzed log file data using Spark and Zeppelin based tools
 - Results of this work will influence U.S. NDC future data processing and hardware architecture design decisions for processing, storing, and analyzing large datasets