

Hadoop & Complex (Systems) Research Thoughts About Large Scale Data Management

GridKa School 2014, Karlsruhe, Germany 2014-09-02 Mirko Kämpf | Cloudera

AGENDA

Complex Systems Research and Big Data: Challenges
The role of Metadata in the Big Data world
Data Management in Hadoop
Cluster Federation using the ETOSHA data catalog

COLLECT SOME INFORMATION...

- Do you work interdisciplinary?
- Do you collect your own data?
- Do you share data sets often?
- Do you contribute sub sets of results?
- Do you integrate multiple data sets?

I: CHALLENGES

- As research projects become more complex, larger groups **need scalable collaborative technology and methods**.
- As soon as datasets become to large it is impossible for them to be copied between partners. We need data federation for public and private (anonymous) data.
- Research results have to be transparent and reproducible. This requires not just access to research papers, but also access to comprehensive contextual information.



Characteristics of Complex Systems

Source: http://www.idiagram.com/examples/characterisitcs.gif

CHARACTERISTICS OF BIG DATA



Source: http://api.ning.com/files/tRHkwQN7s-Xz5cxyIXG004GLGJdjoPd6bVfVBwvgu*F5MwDDUCiHHdmBW-JTEz0cfJjGurJucBMTkIUNdL3jcZT8IPfNWfN9/dv1.jpg

IS IT A BIG DATA PROBLEM?

- Share and integrate measured data and simulation results from multiple sub projects on different scales (space and time) into one large meta model.
- Volume: For better statistical results large datasets are required.
- Velocity: Becomes more important in real world real time applications.
- Variety: New methods lead to evolution of procedures and schema updates.
- Veracity: Data collection tools using novel methods may not be perfect.

2: THE ROLE OF METADATA

- What can be expected from a data set or: "What is in it?"
 - Data set description is provided by metadata.
- How can I get access to a particular dataset?
 - legacy systems: FTP, NFS, and HTTP downloads
 - state of the art: cloud storage & distributed file systems (Amazon S3)
 - Location, version and schema is also metadata.

DEFNITION: DATA & METADATA

from: WIKIPEDIA

Data: Data as an abstract concept can be viewed as the **lowest level** of abstraction,

from which information and then knowledge are derived.

Metadata: Metadata (metacontent) is defined as the **data providing information** about

one or more aspects of the data, such as:

- Location on a computer network where the data were created
- Purpose of the data
- Time and date of creation
- Creator or author of the data
- Means of creation of the data
- Standards used ...

SUMMARY

according to Dr. J. Pomeranz in its Coursera course about metadata:

- (1) Metadata is a description about
 - natural or artificial
 - physical or digital things
- (2) Metadata can be:
 - descriptive
 - structural
 - administrative
- (3) Metadata exists on a per item or per collection level
- (4) Metadata can be **embedded** or **linked**

Metadata covers multiple aspects of a system.

PROCEDURAL & ENTITY METADATA

- Entity Metadata is about a document, a record, a file, a table, a collection.
- **Procedural Metadata** describe the collection, generation, filtering, and transformation procedures.



RESEARCH APPROACH & DATA REPRESENTATION





METADATA & CONTEXT



METADATA & CONTEXT



3: DATA MANAGEMENT IN HADOOP

- Hadoop is a distributed platform to store and process large data sets in parallel.
 - Hadoop offers multiple processing engines:
 - traditionally the Map-Reduce engine
 - Machine learning with Mahout
 - Graph processing with Giraph
 - Spark and GraphX for in memory processing
 - Hadoop offers data management capabilities:
 - HDFS for simple file storage
 - HBase for random access based on a row key
 - SOLR to search in indexed datasets

HCATALOG / HIVE

- Define table structure for data stored in
 - simple HDFS files
 - HBase tables
- Query such "Hive-Tables" also with Impala, Pig & MapReduce
- Access data sets using the Kite-SDK and within Spark:
 - 1 val f = classOf[org.my.package.spark.hcatalog.HCatInputFormat]
 - 2 val k = classOf[org.apache.spark.SerializableWritable[org.apache.hadoop.io.Writable]]
 - 3 val v = classOf[org.apache.hcatalog.data.HCatRecord]
 - 4 val conf = new org.apache.hadoop.conf.Configuration()
 - 5 org.apache.hcatalog.mapreduce.HCatInputFormat.setInput(conf, "db", "table")
 - 6 val data = sc.newAPIHadoopRDD(conf, f, k, v)

Source: https://gist.github.com/granturing/7201912

LIMITATIONS:

- Hive-Metastore contains primarily "technical" metadata.
- Meaning of a data set is not managed.
- Context and life cycle of a data set is not actively managed.
- Processing parameters are not available in core Hadoop.
- Hive-Metastore is shared, but only for users in one cluster.

4: WHAT IS ETOSHA?

- **Etosha** will be a distributed, decentralized, and fault tolerant metadata service.
- Its main purpose is **publishing and linking** of information about datasets.
- Etosha enables **cluster federation** for Hadoop.
- Etosha allows **knowledge management** while HCatalog and the Hive-Metastore are focused on low-level technical

ETOSHA ARCHITECTURE





heddood

Financial Data Analysis From Wikipedia DataSet.ABC.Test1

ClouderaTrainingCluster

Financial Data Analysis V3

Prep C5 Admin Training

MK.EC2 2014.02.03 1

DevEnv.master2

My Virtual Cluster

NC

VM3

ProjectA

ProjectB

Mirko

Nate

Training

Flume DB Import 01

Giraph Algorithm 01

MS-Excel-File Import 01

Mahout Algorithm 01

DataSet.CDEFG.Test2

Giraph Algorithm 01

Mahout Algorithm 01

Oozie Workflow 01

MR Job 01

Shakespeare Test Dataset

Wikipedia Click Count Data

Oozie Workflow 01 Sqoop DB Import 01

MR lob 01

ETOSHA CONTEXT BROWSER

Context browser works across clusters.



- merge multiple layers
- interconnect thing from multiple categories ...

NEXT STEPS

- Expose administrative, structural, and descriptive metadata about Hadoop data sets via SPARQL endpoint.
- Implement a data set discovery API to find relevant data.
- Enable subquery delegation between Hadoop clusters.
- Collect and publish metadata about public open datasets in the Etosha Data Catalog (EDC) project.

MANYTHANKS ...

- ... to the organizers of this great event!
- ... to my friend and contributor Eric Tessenow.
- ... to my colleagues from Cloudera.
- ... and most importantly, to the members of the open source community!