# The ESRI Geodatabase Framework

Future developments at ArcGIS 10.2 and 11

MSc Marco Boeringa

2013

*** **DISCLAIMER** ***: This is **not** an ESRI document. It has been created by the author solely, and any errors herein are his only. The document contains a diagram that tries to represent the most significant parts of what the author calls "The ESRI Geodatabase Framework", with a focus on showing the different options to access spatial data in an enterprise database. As this framework is rather involved, and has evolved through time, the author makes no assumption of completeness nor correctness by representing this figure. There are, and have been, other components and options (E.g. CAD Client, ArcSDE C/JAVA API, MapObjects, ArcView 3.x etc...), and there may be issues in the way the figure represents reality. Feel free to contact the author if you spot an obvious error or have suggestions for improvement. The author takes no responsibility for problems or damage caused by errors in, or misinterpretation of, the contents of this document.

## Author

MSc Marco Boeringa, May 2013

The Netherlands

Document version: V1.0

## Introduction

This document and the three diagrams it contains supplement the main document titled:

"**The ESRI Geodatabase Framework**"

that I wrote and published as a PDF on the ESRI User Forums.

I made this supplement to account for the recently announced deprecation plans for ArcGIS 10.2, that contain some important notifications of ESRI regarding the future support for components of the ESRI software product line. As these announcements will impact the ESRI Geodatabase Framework, I thought a supplement that visualized these changes might be of help to others in understanding the upcoming changes. I visualized the changes by creating three separate diagrams, displaying the changes and features in 10.1, 10.2 and a presumed next generation ArcGIS 11. I have added comments to the diagrams (e.g. "Future") to explain changes.

In addition, I have included the new **Spatial Framework for Hadoop** in the diagram for ArcGIS 11, as it more or less extents the ESRI Geodatabase Framework with an entirely new option: to store, handle and process **Big Data** via **Apache Hadoop**. Developed by ESRI, this framework extents Hadoop with spatial functions based on OGC ST_Geometry. Although it is included in the diagram of ArcGIS 11 only, it is already available at ArcGIS 10.1 and downloadable as an OpenSource project from GitHub.

Please note this is currently particularly "developer oriented", although the Spatial Framework for Hadoop does include a more user friendly Geoprocessing Toolset, but defining actual Job tasks for this may need Java programming, or at least Hive HQL, knowledge. Additionally, setting up and managing a Hadoop cluster is *not* a trivial task, and requires Linux/Unix expertise. There are enterprise grade platforms for Hadoop that supposedly ease the tasks of managing a Hadoop cluster considerably. Two examples are MapR and WANdisco, a comprehensive overview of Hadoop platforms can be found on this Hadoop Wiki page: http://wiki.apache.org/hadoop/Distributions andCommercial Support

Since many of the concepts, terms, and software product names regarding Hadoop and "Big Data" may be unfamiliar to traditional GIS users, I have included a small glossary encompassing the most important concepts of Big Data and this new addition to the ESRI Geodatabase Framework. You can find this glossary at the end of the document.

Please note that the inclusion of this new Hadoop option, does not mean ESRI will replace, or intents to replace, the ArcSDE Technology based Enterprise Geodatabase Framework by Hadoop in the future. Nor is it a replacement for the ArcSDE *Application Server*, at which former location in the ArcGIS 11 diagram it is drawn simply for practical reasons. It is a *new* and *different* option, that may offer benefits or new opportunities for analysis in case of very specific "Big Data" questions in the scientific research community, or for business analytics on giant consumer datasets.
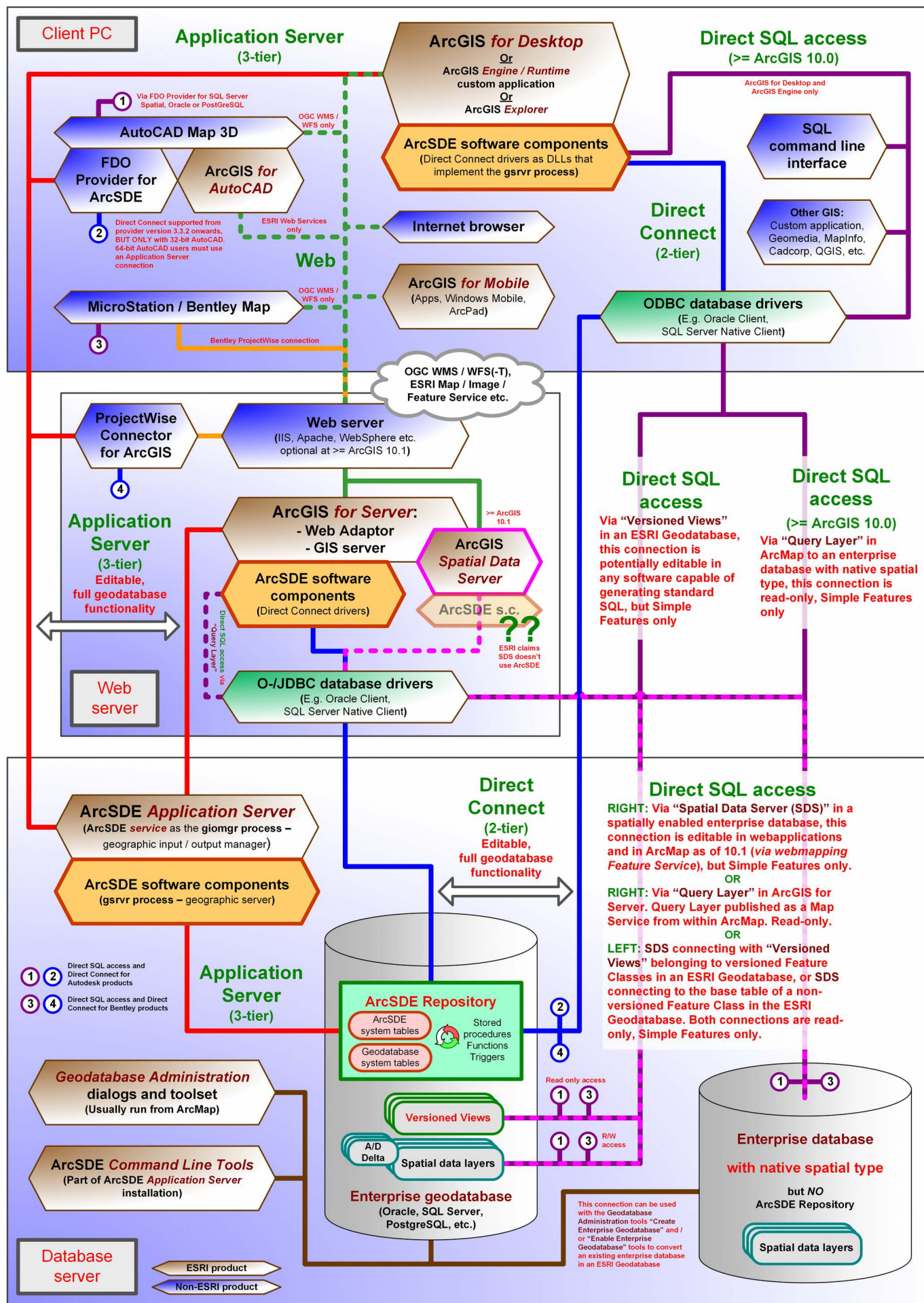
Some resources from ESRI concerning the new option:

**Video:** Big Data: Using ArcGIS with Apache Hadoop

**Resources:** http://www.esri.com/technology-topics/big-data/resources
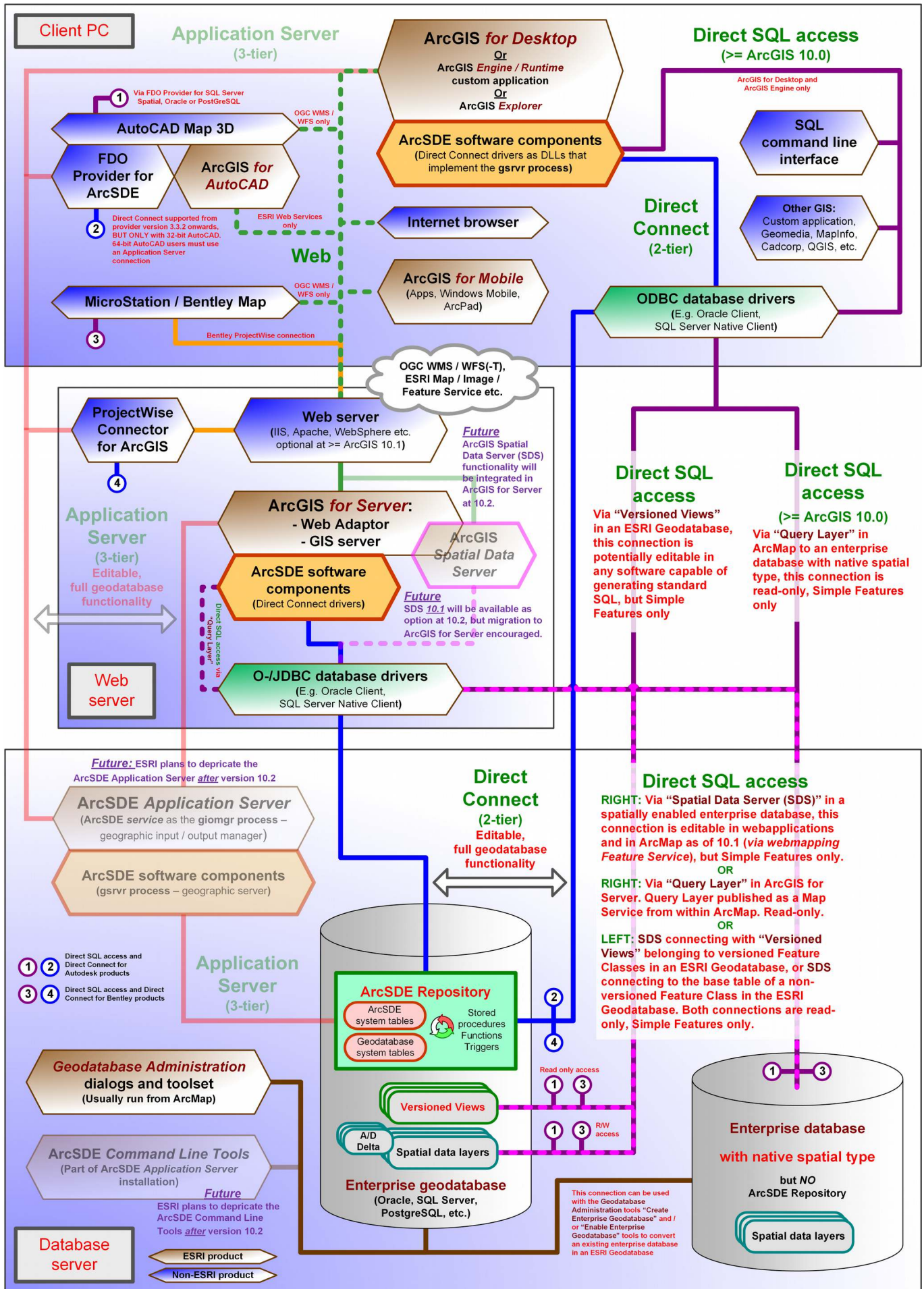
**Blog:** by one of the ESRI developers involved in the Spatial Framework for Hadoop

**GitHub:** Spatial Framework for Hadoop

The ESRI Geodatabase Framework at ArcGIS 10.1, © Marco Boeringa, 2013

The ESRI Geodatabase Framework at ArcGIS 10.2, © Marco Boeringa, 2013

4

The ESRI Geodatabase Framework at ArcGIS 11 and beyond,  © Marco Boeringa, 2013

5

## "Big Data" glossary

### Big Data

Any (geographic) dataset big enough to choke an ordinary enterprise database server. Think of "Facebook", "Google" type size "databases" as an extreme example. With the scaling up of the capacity of traditional enterprise databases and enhanced storage concepts for them, the point at which data can no longer be handled by an ordinary database is shifting. Still, there may be good reasons to use a distributed file storage system like Hadoop to not only store, but also take full advantage of parallel processing, of Big Data.

### The Cloud

A number of server computers, also called a **cluster**, providing services over a LAN (Local Area Network) or WAN (Wide Area Network), essentially operating ***as one single unit*** to the outside world. "The Cloud" usually refers to internet based implementations of a cluster, like Amazon Elastic Compute Cloud (Amazon EC2), but it might just as well be a big server room in the cellar of your office building, or even just a dozen desktop computers hooked up in a test configuration in your section of the office. Individual servers can be ordinary commodity desktop computers, or professional rack server hardware. A service could be a distributed file storage system like Hadoop, a webmapping service, or a geoprocessing service processing data, to give some arbitrary examples.

### *Distributed storage* and *parallel processing*

*Distributed storage* and *parallel (geo-)processing* using services spread (distributed) across many, potentially thousands, of "servers" on a LAN or WAN, like in "The Cloud". This set of servers is also called a **cluster**. Since the data is processed by many different computers, distributed storage and parallel processing can potentially be much faster than traditional processing on a single server or desktop computer.

### Elastic

A "buzzword" (fashionable expression) often employed when referencing The Cloud or applications running in a cloud / cluster of computers like Hadoop. "Elastic" means that individual servers / computers can be added or removed at will depending on need. This process is transparent and shouldn't affect running applications.

### File System

A means to logically store and retrieve data and documents on (hard) drives. Data safety, in the sense of preventing data loss, is a key feature of many file systems, meaning some form of redundancy is build-in to cope with hardware or software failures. **HDFS** is a file system, but more commonly known are **NTFS** from Windows, or **FAT** used on digital photo cameras' memory cards. File systems can be locally, on one drive, or distributed like HDFS.

### Geoprocessing Tools for Hadoop

Set of geoprocessing tools developed by ESRI that can be added to ArcToolbox to allow you to convert geographic data (e.g. a Feature Class) *to* and *from* JSON, and store the data to the HDFS file system and Hadoop. The functionality currently only supports Simple Features vector data. In addition, the Geoprocessing Tools for Hadoop allow you to start a MapReduce Job to extract data from Hadoop.

**Hadoop (Apache Hadoop)**

A distributed file system *and software framework* for storing and parallel processing of Big Data. Hadoop is based on **Java** and can easily run on commodity hardware. Hadoop and its underlying file system HDFS were not only designed to be able to store Big Data, but also to handle very high read/write throughput of the data from and to the distributed data store, allowing efficient processing of the data to answer "Big" questions. Hadoop is an OpenSource Apache project. Hadoop's HDFS and MapReduce are based on earlier Google work on these systems.

http://hadoop.apache.org/

http://en.wikipedia.org/wiki/Apache_Hadoop

**Hive (HQL or HiveQL)**

**Hive, HQL** (Hive Query Language)  or **HiveQL**, is a query language similar to SQL for NoSQL databases. It allows querying and filtering of data to select and process data from Hadoop.

http://hive.apache.org/

Hive queries are build and submitted using command line tools, or could be defined in an ODBC client, in case one uses for an Hadoop platform like MapR, which includes an OBDC driver for Hadoop. If using Hive, there is no need for writing Java MapReduce Jobs, as extraction of data from Hadoop is guided by the HQL query statements. Hive is an extension to Apache Hadoop, and must be installed separately if not part of a pre-configured installation of Hadoop.

**HDFS**

HDFS (Hadoop Distributed File System) is the file system used by Hadoop. It is a distributed system, meaning data is stored on many different computers / servers, also called a **cluster**. HDFS has build-in redundancy to prevent data loss. Redundancy means that multiple copies of data (usually three copies in case of Hadoop) are stored to avoid data loss, this is also called **replication**. In combination with *automated failure detection* and *recovery*, this ensures the health of the distributed file system, as these are vital elements to make such a distributed file system viable.

HDFS can not be accessed like an ordinary drive in Windows Explorer, except through custom extensions of Hadoop. MapR seems to offer, as one of the few Hadoop implementations, a supposedly reliable NFS mounting facility for HDFS compatible with Windows, meaning you can access the HDFS file system as an ordinary drive from Windows Explorer and drag & drop data to it. A more comprehensive overview of Hadoop projects related to NFS mounting can be found here: http://wiki.apache.org/hadoop/MountableHDFS, and Google turns up more results of ongoing work to provide robust NFS mounting, e.g. Hortonworks' "Simplifying data management: NFS access to HDFS"

**Job**

A small program or set of instructions that will be executed against the geographic data stored in Hadoop. Jobs can be written in for example Java to extract certain data from Hadoop using the ESRI's **Spatial Framework for Hadoop** geometry functionality. Or a job could be defined using ESRI's **Hive Spatial**. There are a couple of ongoing projects attempting to provide compatibility between Python and Java, allowing Hadoop Jobs to be written in Python. Examples of this are Hadoop Streaming, Jython etc. The following pages may be of use:

A Guide to Python Frameworks for Hadoop

Jython: Python for the Java Platform

Making Hadoop Simpler with Python Joins / Keys

**JSON / BSON**

JavaScript Object Notation (**JSON**) is a simple text file format used to store structured and unstructured content in a hierarchical manner using key-value pairs. Just like XML, it is a human readable text format (but JSON isn't XML based).

http://www.json.org/

**BSON** is binary stored JSON.

**Key/Value pairs (also called name/value pairs)**

A "key" is essentially a unique "name" for a value to be extracted from the JSON file, or any other structured data format using key/value pairs (e.g. dictionaries in programming). By submitting the key, you get back the associated value. E.g. you could store a list of values of your friends' names using a key called "MyFriends".

An example might look like this (formatting is arbitrary here and doesn't represent JSON):

("MyFriends",{"John","Peter","Patrick","Eric"})

**MapReduce**

MapReduce is a programming model in the Hadoop system allowing **distributed storage** and **parallel processing** of Big Data via the Map and Reduce functions (MapReduce).

Unfortunately for GIS users, the word "Map" in **Map**Reduce has nothing to do with the traditional sense of a "map" in GIS or cartography. "Map" in MapReduce is a mechanism to cut up a big dataset stored in Hadoop in manageable smaller chunks for parallel processing and distribute and manage them in the Hadoop cluster, including performing extraction of data selections using special filtering techniques. Since Hadoop is a distributed system, instead of "tables" stored in a traditional RDBMS, selecting or filtering data is not as straight forward as going through a single table collecting records. In Hadoop and HDFS, data has to be indexed and extracted / filtered from many different servers. The Map procedure and its underlying server implementation orchestrates this process.

"Reduce" in Map**Reduce** is the possibility to parallel process or do something with the data selected by the "Map" procedure, e.g. summarize or count values. Think of calculating the, min, max or average number of households on a state or township basis.

Both the Map and Reduce function are, and must be, defined by the user in a Hadoop **Job**.

**NoSQL database**

A NoSQL database is a type of database that fundamentally differs from traditional relational databases (**RDBMS**). There are no "database tables" in the traditional sense, with fixed field names or content type for each field (e.g. string, date, number). A NoSQL database doesn't enforce strict relational integrity like traditional databases, but rather stores simple key-value pairs like the ones used in JSON to allow storage and extraction of data based on keys. NoSQL databases can have "SQL like" languages for extraction of values by keys. The advantages are more flexibility in the content stored, but comes at the price of less integrity assurances.

Examples of NoSQL databases are **MongoDB, HBase** etc.

**Oozie**

A Job tracker software or workflow scheduler allowing you to manage and submit Jobs to Hadoop in a controlled and planned manner. Oozie is Java based.

**Spatial Framework for Hadoop**

A "spatial extension" to Hadoop developed by ESRI based on standard **OGC ST_Geometry** allowing you to store and retrieve spatial data on Hadoop.  This extension enables several spatial options in Hadoop:

- The Spatial Framework extents **Hive** with User Defined Functions (UDF), also called **Hive Spatial** by ESRI. These UDFs enable the use of standard OGC style ST_Geometry geometric functions in **Hive (HQL) queries** (e.g. ST_Polygon, ST_Within, ST_Intersection). So you can start writing HQL queries containing for example ST_Intersects, to determine if two geometries intersect.

- Addition of the ESRI Geometry API for Java. This allows you to write Java based MapReduce Job files containing spatial operators on geometries and other functions available in the ESRI Geometry API for Java (e.g. spatial indexing capability).

- A number of helper utilities or classes (**JSON Utilities**), like an **EsriFeatureClass** class object, allowing you to directly construct an ESRI Feature Class from JSON files.

The **Spatial Framework for Hadoop** can be found on this ESRI managed GitHub page:

https://github.com/Esri/spatial-framework-for-hadoop

A competitive system to the Spatial Framework for Hadoop by ESRI, is **SpatialHadoop** by the University of Minnesota. This framework is *not* an OGC ST_Geometry implementation though, like the one from ESRI, and as a consequence seems more limited in its spatial functions and types. See this page:

http://spatialhadoop.cs.umn.edu/index.html