

Volume Management of Hadoop DataNode

Xu Wang, Leitao Guo, and **Zhiguo Luo**
China Mobile Research Institute(CMRI)

June 1, 2011

Agenda

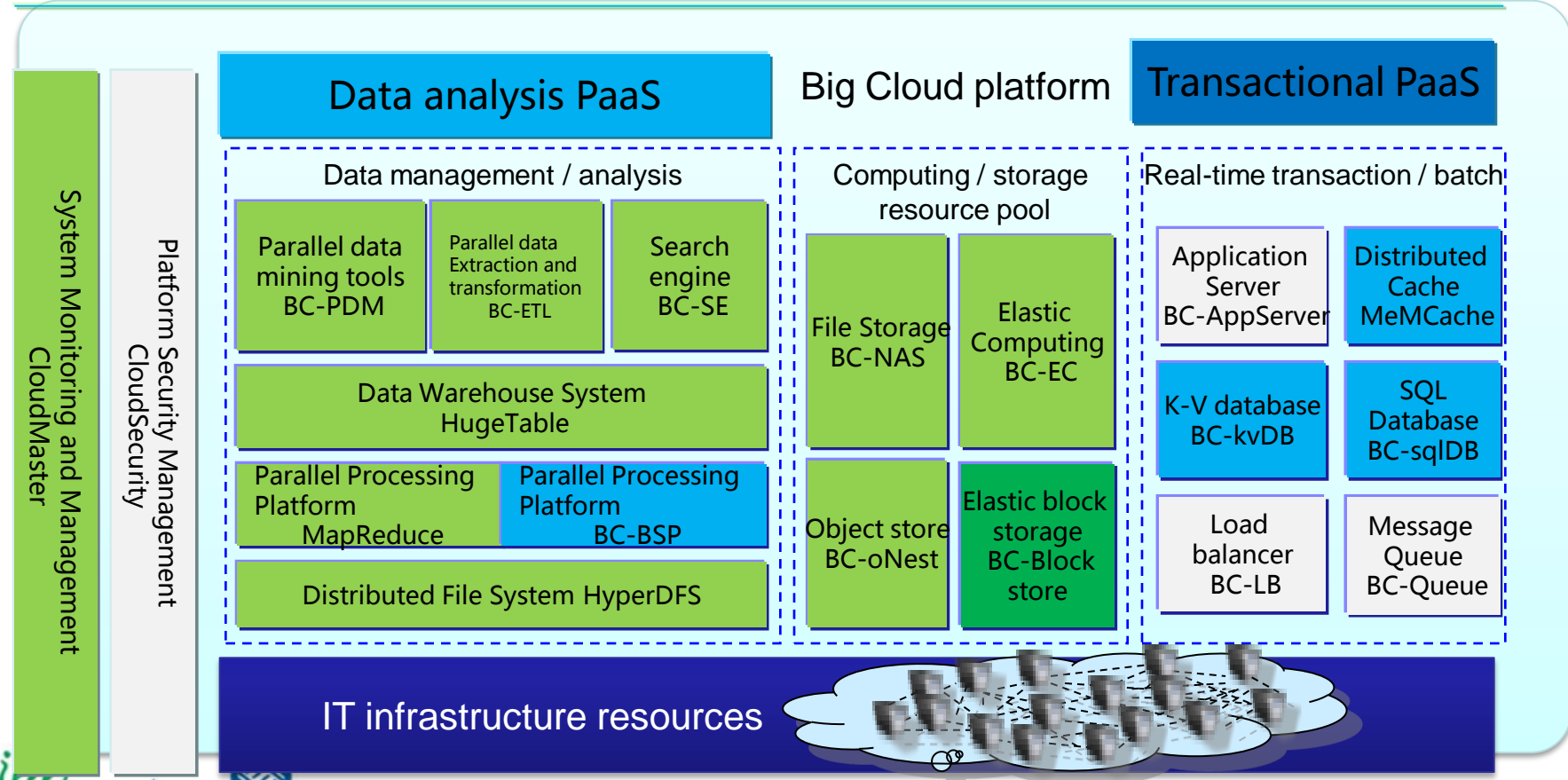
- ④ Cloud R&D in China Mobile
- ④ Volume Management of Hadoop DataNode

CMRI “Big Cloud” Project

- ④ What is “Big Cloud”
 - ④ A CMRI project: research of key technologies and prototype system development for China Mobile Cloud Computing Infrastructure
- ④ Goal
 - ④ To meet the growing demands for high performance, low cost, high scalability, high reliability of China Mobile IT Infrastructure (computing, storage)
 - ④ To meet the demands of China Mobile to deliver Internet business and services

"Big Cloud " Product Evolution Direction

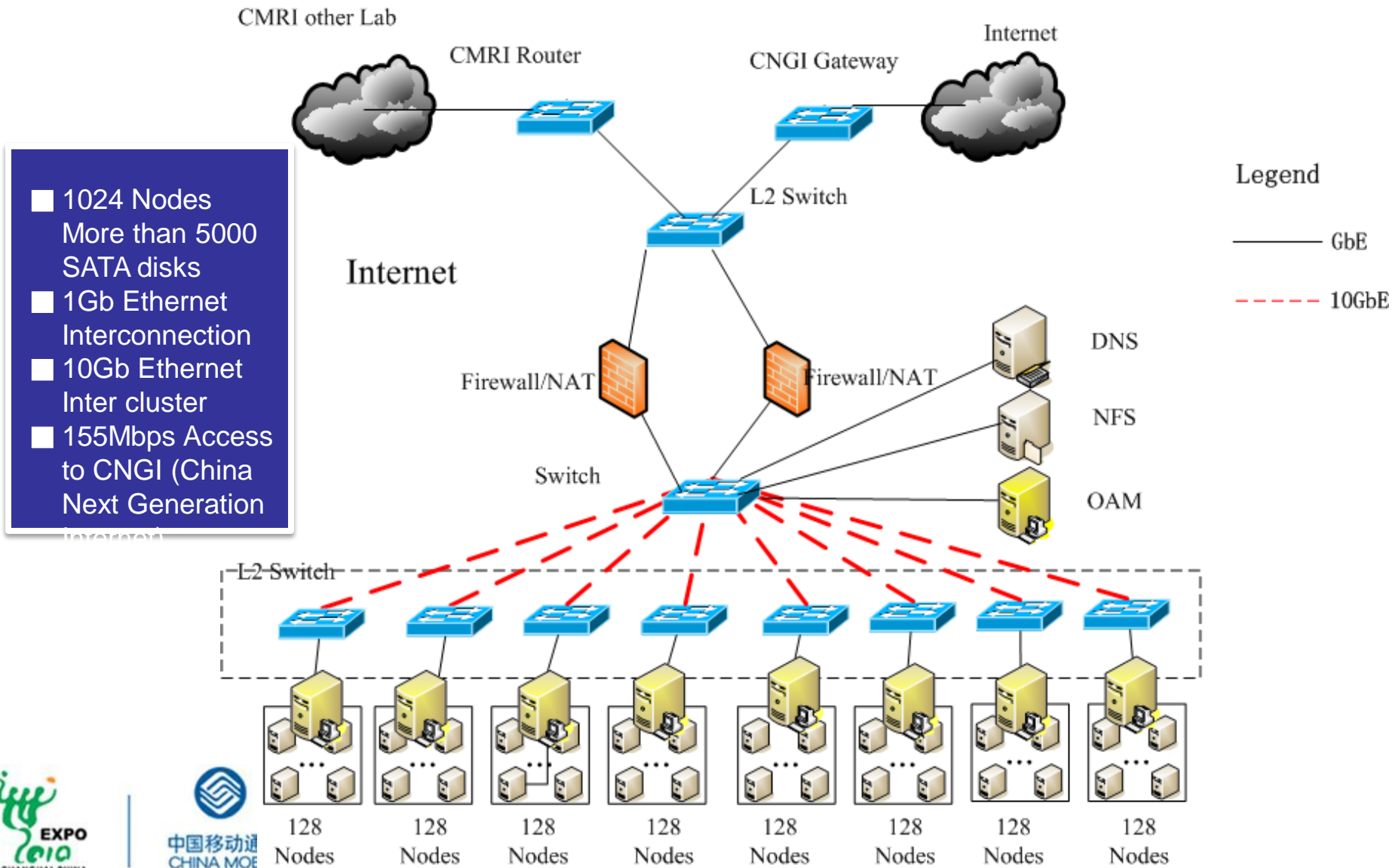
BASS KPI Centralized computing
BASS ETL/DM
Settlement System
Signaling System
Resource pool system of Cloud Computing
IoT Application
E Mail
IDC Services



BC 1.0
BC 1.5
2.0
3.0
integrate



Testbed in CMRI Lab

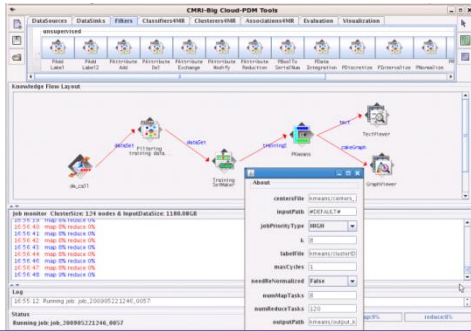


Cloud based Parallel Data Mining tool——BC-PDM/ETL

- ④ A data mining solution for large-scale data analysis
 - ④ Massive scalability - based on Hadoop
 - ④ Low cost - commodity machines and open source software
 - ④ Customization – facing to application requirements
 - ④ Easy to use - similar user interface to commercial tools
- ④ Parallel Data Mining algorithms based on MapReduce
 - ④ ETL operations – parallel data extract, clean/transfer, load operations
 - ④ Data mining algorithms – parallel classification, clustering, association rules algorithms
 - ④ Data Exploration – parallel basic statistics descriptor, univariable analysis, multivariable analysis algorithms
 - ④ Social Network Analysis – parallel network structure analysis, social network discovery, social network development algorithms

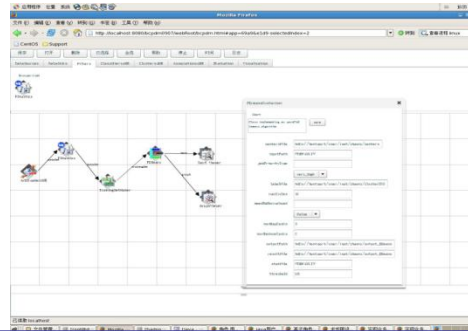
BC-PDM/ETL Evolution

Version I (2008)



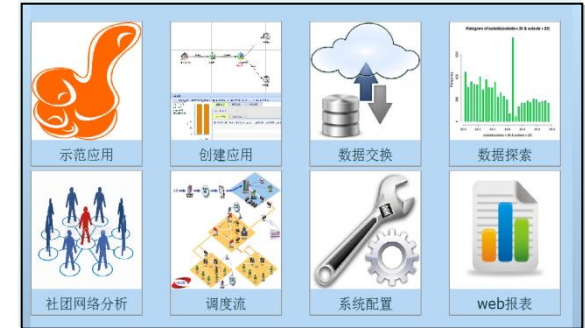
- Workflow management
 - GUI - Drag Operation for application modeling design
 - Job Monitoring
 - Flow Configuration
- ETL (14 different ETL operations from 6 categories)
 - Statistic, attribute processing, data sampling, query, data processing, redundancy data processing
- Data mining Algorithm (9 algorithms from 3 categories)
 - Clustering, Classifier, Association Analysis
- Visualization
 - Text, decision tree, cake graph, and histogram

Version II (2009-2010)



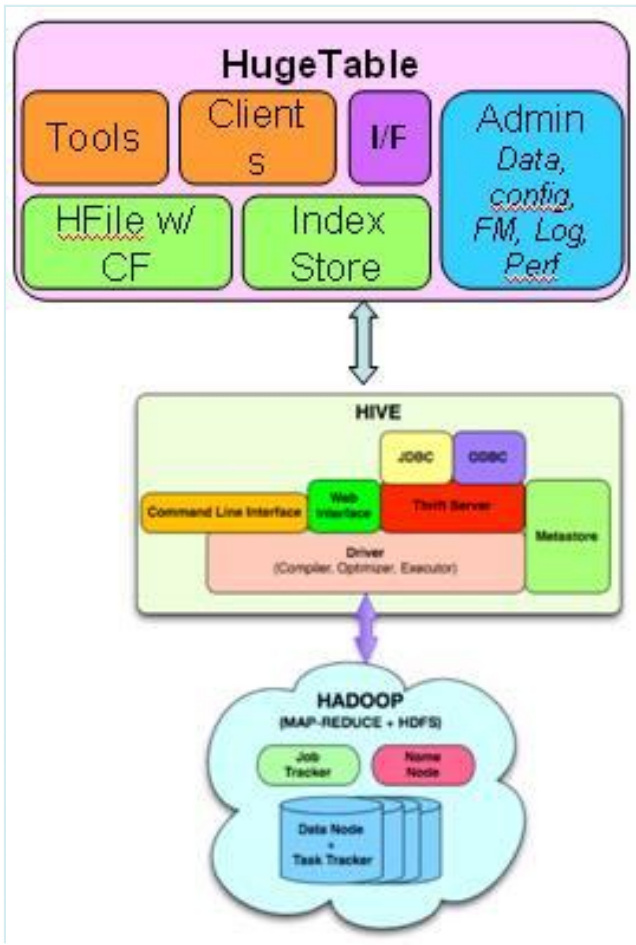
- Web based GUI: SaaS mode for users
- DE(Data Exploration): Simple data analysis and preview
- ETL (25 more): Join, Group by, Expression, case when, Update, and etc.
- Data mining Algorithm (4 more): Classifier, Sequence Association Analysis
- Data Transfer Tool: Provide data upload and download tools for SaaS
- Security: Multi-talent and user group for branch, ACL for data access

Version III (2010-2011)



- Web based GUI
 - SaaS mode for users
- UI: Workflow GUI, CLI, SQL
- Data Exchange: between DB and HDFS
- Integrate Hive: SQL-like scripts
- Parallel Social Network Analysis
- ETL(45 from 7 categories) DM(14 from 3 categories) Data exploration (3 categories)
- application schedule system
- Application templates
- API: java Api, web service

Hugetable: Structured Data Management System and Applications



- Build on top of HDFS、MapReduce、Hive、HBase and Zookeeper
- Provide index, join, fast load/export
- Provide JDBC/SQL, MapReduce, native API under access control
- Open Management Interface for system management
- Some system integrator companies are integrating the Hugetable into their applications for CMCC

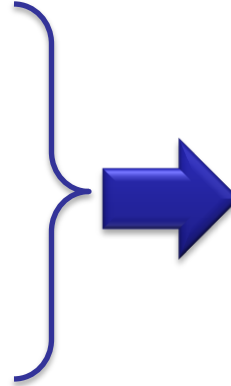
Agenda

- ④ Cloud R&D in China Mobile
- ④ Volume Management of Hadoop DataNode

Background

SATA Disk failure happens much more frequently than other failure

Servers tend to equip with more disks



Servers should be able to survive through disk failures

-
- ⊙ HDFS 0.20 and before: disk fail is a fatal error of Hadoop Datanode
 - ⊙ HDFS-457: disk fail will lose a disk, but the server can still work.

Hard disk failure observation

```
smartd: Device: /dev/sdb, 2  
Currently unreadable (pending)  
sectors
```

After about 10
days or 1 month

```
smartd: Device: /dev/sdb,  
FAILED SMART self-check. BACK  
UP DATA NOW!
```

sub-health status

After several
days

```
fatal: qmgr_active_feed:  
DFF5D10411: rename from  
deferred to active: Read-only  
file system
```

We can guess the disks' health status from smartd information, find the sub-health disks and then take action ahead of time.

Objectives

- ④ To forecast the disk failure automatically, and take actions ahead of time to improve the data availability, such as
 - ④ to notify the NameNode to replicate the blocks stored in the failure disk by DataNode block-report. [similar with HDFS-457]
 - ④ or to move the blocks to other healthy disks in the same DataNode and then invalidate the sub-healthy disk. [invalidate a disk online]

Implementation status

- ④ Hadoop community don't care about forecasting the disk failure too much, but regard it's useful to replace the failure disk online.
- ④ So, we implemented the functionality of changing disks online (HDFS-1362) firstly.
- ④ We also tried to implement a functionality of invalidating a disk online but not open source until now.
- ④ In our testbed, we automatically observe and guess the disk status, and trigger the action externally.

Volume Management of Datanode

FSVolumeSet

Maintain volumes:

- Scan volumes while startup
- Remove failed volume (feature of HDFS-457)

ReplicaMap

Maintain blocks

- Map block to on-disk file
 - Responsible for all kinds of block related operation
- You need traverse the ReplicaMap to find blocks on a specified volume if you want to remove a volume.

Reconfiguration of Hadoop

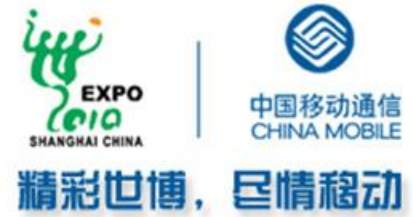
- ⊗ HADOOP-7001 introduced a framework for reconfiguration of Hadoop
 - ⊗ Triggered by Servlet
 - ⊗ Reload specified configuration from config file
 - ⊗ Keep the consistent between runtime configuration and the persistent configuration

Volume Refreshment

- Based on the framework provided by HADOOP-7001
- Implement a refresh method
 - Got the disks configured in files
 - Check which disk is in configuration file but not in service
 - Check and load the above volumes
 - Load blocks if there are some in the volumes

Future Works

- ④ A disk rebalance function
 - ④ Disk might be un-balanced, because:
 - ④ Change failed disk
 - ④ Mass delete some data
 - ④ Consumed by other system, such as mapReduce local storage.
 - ④ This function may relate to another...
- ④ Monitor disk usage status of Datanode
 - ④ Check data distribution among disks
 - ④ Triggered by servlet
- ④ Improve the functionality of invalidating a disk online and push it to mainline.
 - ④ Challenges: when we configure the RAID card, it may conceal some details of disk details, which we observe by smartd. We need to do more verifications later.



Thank You!