**NSRL** 中国科学技术大学
国家同步辐射实验室
NATIONAL SYNCHROTRON RADIATION LABORATORY

**PCaPAC 2018**
12th International Workshop on
Emerging Technologies and Scientific Facilities Controls
October 16-19, 2018    Hsinchu, Taiwan

# Design and Construction of the Data Warehouse Based on Hadoop Ecosystem at HLS-II

Yifan Song

Doctoral student, Controls Group

National Synchrotron Radiation Laboratory

University of Science and Technology of China

# Outline

- Motivation
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
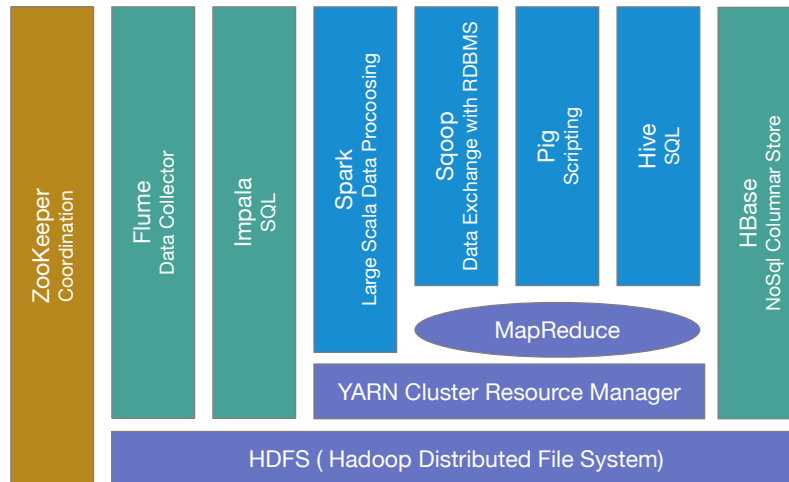- Current progress and next plan
- Summary

# Outline

- **Motivation**
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
- Current progress and next plan
- Summary

# Motivation

- The Hefei Light Source II (HLS-II) at National Synchrotron Radiation Laboratory is the first dedicated synchrotron radiation facility in China.

- As more and more data have been stored, the exciting data archiving tools can not meet our requirements of data query and processing.

- In order to deal with these problems, we are designing and constructing the data warehouse based on Apache Hadoop ecosystem.

- In addition, our laboratory is conducting pre-research on Hefei Advanced Light Source (HALS) project. Compared to HLS-II, HALS is a larger and more complex synchrotron facility. The work described in this paper also provides technical verification for the future construction of HALS.

# Outline

- Motivation
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
- Current progress and next plan
- Summary

# Hadoop Ecosystem

- Apache Hadoop is a collection of open-source software.
- Hadoop core components include HDFS, Yarn and MapReduce.
- Some open source projects based on Hadoop, including Spark and Impala, which provide necessary support for the whole life cycle of big data processing.
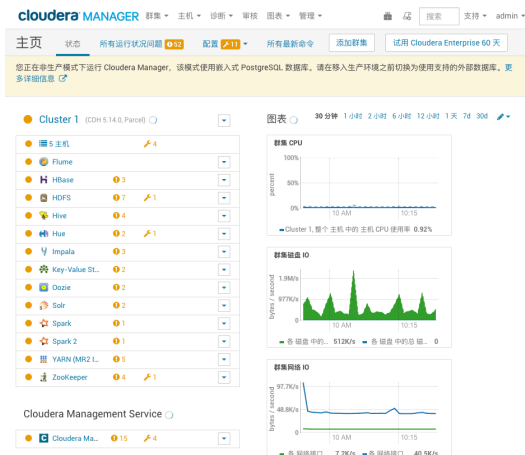
# Hadoop Ecosystem

- HDFS
  - a distributed file-system to store data.
- Spark
  - has the advantages of Hadoop MapReduce
  - but unlike MapReduce, its intermediate output can be stored in memory, instead of the HDFS.
  - Tests have shown that Spark is more efficient than MapReduce for the same operation. Spark is widely used in the field of ETL (Extract-Transform-Load) and machine learning.
- Impala
  - an open source new query system that can query PB-level big data stored in Hadoop's HDFS and HBase. Compared to MapReduce, Impala provides data query function more efficient and convenient.
- Sqoop
  - command-line interface application for transferring data between relational databases and Hadoop.
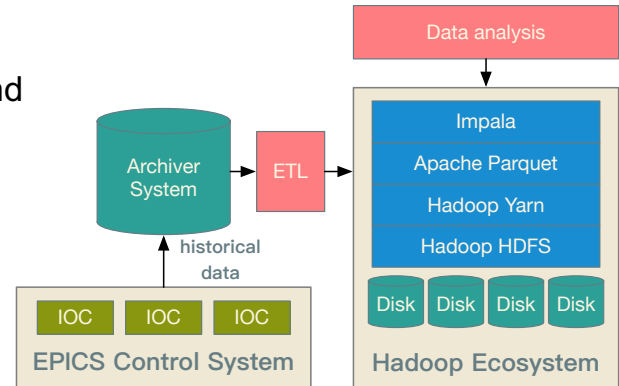
# Hadoop Configuration

- Cloudera Distribution Hadoop (CDH) is used as the tool for cluster construction.
  – handle version compatibility and configuration files
  – manage and monitor the state of the entire Hadoop cluster
- We built the test system with 5 nodes on the VMware vSphere virtualization platform based on CDH.



- adding or deleting services
- adding computer nodes
- viewing the status of the cluster
  – Many services are in the alarm state due to the limitations of the hardware performance of the test system.

# Outline

- Motivation
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
- Current progress and next plan
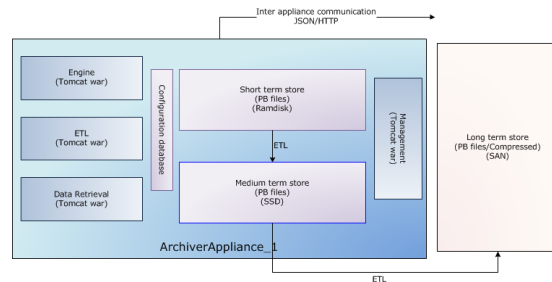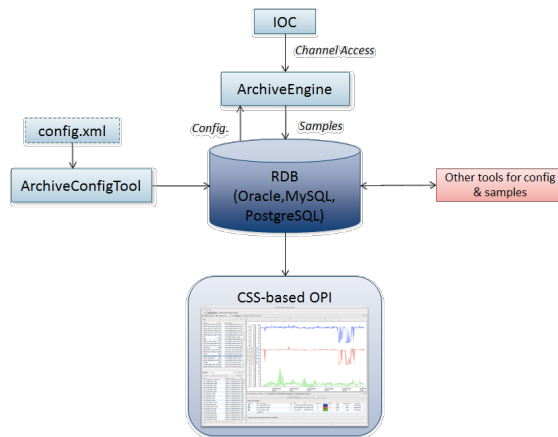- Summary

# Software Architecture

- Control system based on EPICS
  - the Records in EPICS IOC hold the real-time state of the controlled devices
- Archiver system and ETL program
  - didn't redevelop the data acquisition and archiving software
  - still use data archiving tools already available in the EPICS community
    - RDB Channel Archiver
    - Archiver Appliance
- Hadoop big data platform
  - migrate data from archiving tools to HDFS
  - store them in Parquet format
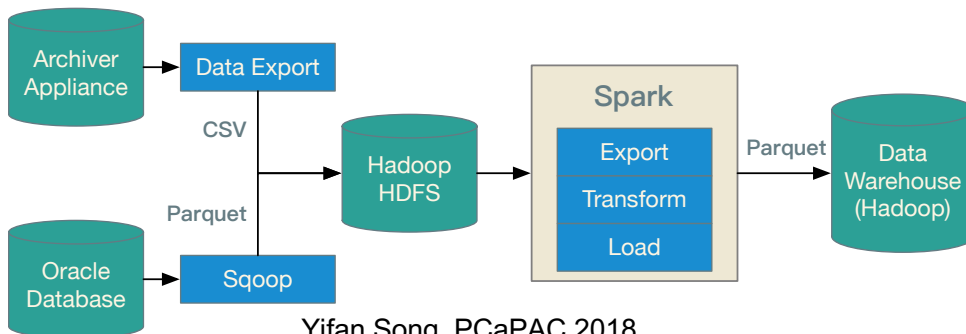  - Impala provides a convenient data analysis interface

# Outline

- Motivation
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
- Current progress and next plan
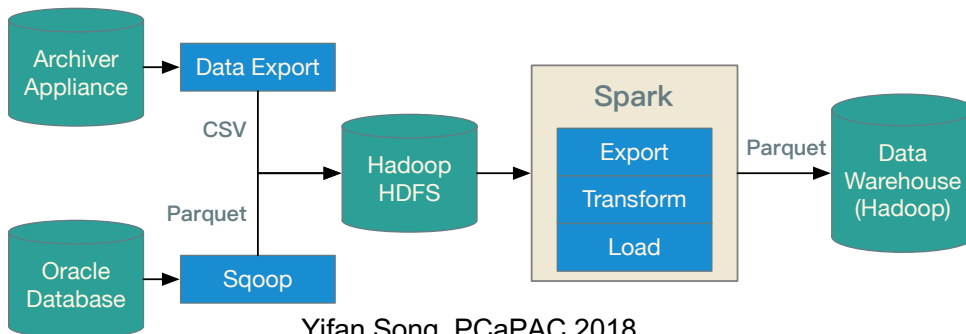- Summary

# Archiving tools at HLS-II

- There are two data archiving systems running at the same time
- EPICS Archiver Appliance
  - collects data from IOCs
  - stores them on the hard disk in Protocal Buffer (PB) format
- RDB Channel Archiver
  - collects data from IOCs
  - stores them in Oracle database
- Both types of data from different data sources need to be migrated to HDFS.

# ETL program

- ETL: Extract-Transform-Load
- First step
  - the original data in the archiving tools are exported to the distributed file system HDFS.
  - Archiver Appliance
    - the Python program is developed to get data from the Archiver Appliance via HTTP in CSV format
    - upload the CSV files to the HDFS file system.
    - This process is performed periodically in batch mode.
  - RDB Channel Archiver
    - Sqoop is used to migrate the data from Oracle database directly to the HDFS system.

# ETL program

- Second step
  - The ETL program written based on Spark is triggered after the data are uploaded.
  - It completes the ETL process, and finally export the Parquet format file into the data warehouse with a predefined data model.
  - Performs a series of data cleaning tasks, including the elimination of abnormal data points and the alignment of sampling time.
- workflow management program based on Airflow
  - ensure the stable operation of the ETL process, including the workflow task scheduling, metadata acquisition, logging and notification function.

# Data analysis

- Query data through the Impala command line.
- The user queries the first 10 records of the sample table in the css database through SQL statement and the result is returned.

```
[root@node01 ~]# impala-shell --impalad=192.168.113.42
Starting Impala Shell without Kerberos authentication
Connected to 192.168.113.42:21000
Server version: impalad version 2.11.0-cdh5.14.0 RELEASE (build d68206561bce6b26762d62c01a78e
6cd27aa7690)
********************************************************************************
Welcome to the Impala shell.
(Impala Shell v2.11.0-cdh5.14.0 (d682065) built on Sat Jan  6 13:27:16 PST 2018)

When pretty-printing is disabled, you can use the '--output_delimiter' flag to set
the delimiter for fields in the same row. The default is ','.
********************************************************************************
[192.168.113.42:21000] > use css;
Query: use css
[192.168.113.42:21000] > select channel_id, smpl_time, float_val from sample limit 10;
Query: select channel_id, smpl_time, float_val from sample limit 10
Query submitted at: 2018-10-08 09:42:13 (Coordinator: http://node02.cdh.nsrl:25000)
Query progress can be monitored at: http://node02.cdh.nsrl:25000/query_plan?query_id=924247bf
c18d2302:380a2e2f00000000
+------------+---------------+-------------------+
| channel_id | smpl_time     | float_val         |
+------------+---------------+-------------------+
| 1596       | 1525418774966 | 29.74092864990234 |
| 1540       | 1525418773003 | 36.93181228637695 |
| 1088       | 1525418773014 | 30.86883926391602 |
| 984        | 1525418774958 | 29.26275825500488 |
| 1100       | 1525418773027 | 25.61181831359863 |
| 458        | 1525418775031 | 23.49545860290527 |
| 1671       | 1525418774013 | 27.05657958984375 |
| 826        | 1525418799040 | 86.629150390625   |
| 1502       | 1525418834977 | 35.76647186279297 |
| 1596       | 1525418834965 | 29.72908020019531 |
+------------+---------------+-------------------+
Fetched 10 row(s) in 0.03s
```

Open impala shell

Use database

Query data

- Motivation
- Hadoop Ecosystem and Data Warehouse
  - Hadoop Ecosystem
  - Hadoop Configuration
- Software Architecture
- Implementation of data warehouse
  - Archiving tools at HLS-II
  - ETL program
  - Data analysis
- Current progress and next plan
- Summary

- Test system
  - Built the test system on the VMware vSphere virtualization platform
  - Developed the ETL program
- System in the production environment
  - Finish server hardware installation
  - 8 nodes
  - Will deploy the software on the server
  - Will carry out performance testing

# Current progress and next plan

- Will provide data reporting functions.
  - develop a series of programs for calculating data reports based on Impala, such as operating status statistics and integral current calculations.
  - control the response time of the report calculation task within a few seconds.
- Will provide more convenient data query entry.
  - Currently, query data by SQL statements on the Impala command line.
  - In the future, provide a web-based data query interface and develop the web applications to facilitate scientists and engineers to query data.
- Will carry out research on data mining and Artificial intelligence (AI) algorithm applications.
  - As the Spark's machine learning library, MLlib consists of some common learning algorithms and tools, including classification, regression, recommendation, etc..
  - We plan to use the spark MLlib to mine the characteristics of the data stored in the data warehouse, carry out research on AI algorithm, and improve the intelligence of the entire control system.

# Summary

- A data warehouse based on Hadoop ecosystem is designed and constructed for Hefei Light Source II (HLS-II).
- Cloudera CDH is very helpful for the cluster construction.
- The ETL program based on Spark and Sqoop migrates data to HDFS from RDB Channel Archiver and the EPICS Archiver Appliance continuously.
- The test System has been built and will be transferred to the production environment.
- Such a technology plan combines various open sources tools under Hadoop framework and the archiving tools in EPICS community. It provides a solution to solve the big data problem for the large-scale scientific facility.

# Q&A