# MADOCA II Data Logging System Using NoSQL Database for SPring-8

**A.Yamashita and M.Kago**

**SPring-8/JASRI, Japan**

MADOCA

# NoSQL

# OR: How I Learned to Stop Worrying and Love Cassandra

# Outline

- **SPring-8 logging database**

- **Why NoSQL, why Cassandra**

- **Implementation**

- **Production Run**

MADOCA

# SPring-8 Logging database

- **Relational database system (RDBMS) has been used since 1997**

- **Grows from 871 signals to  27,626 signals (end of 2014)**

- **7,000 signal inserts per seconds**

- **4TB raw data at end of 2014**

- **SACLA ( X-ray FEL) is also using it**

MADOCA

# SPring-8 Logging database

- **What made the system live long?**
  - **Uniform data store**
  - **Simple access**

# SPring-8 Logging database

- **What made the system live long?**
  - **Uniform data store**
    - **Every data**
    - **Every time**
    - **in one database**
  - **Simple access**

# SPring-8 Logging database

- **What made the system live long?**

  - **Uniform data store**

  - **Simple access**

    - **Just Key + time range access**
      **get ("sr_mag_ps_b/current_adc",**
         **"2014/10/10 19:24:12",**
         **"2014/10/10 22:00:00")**

MADOCA

# RDBMS to NoSQL

- **For the next generation SPring-8-II**

- **We changed logging database for SPring-8 from RDBMS to a NoSQL database; Cassandra**
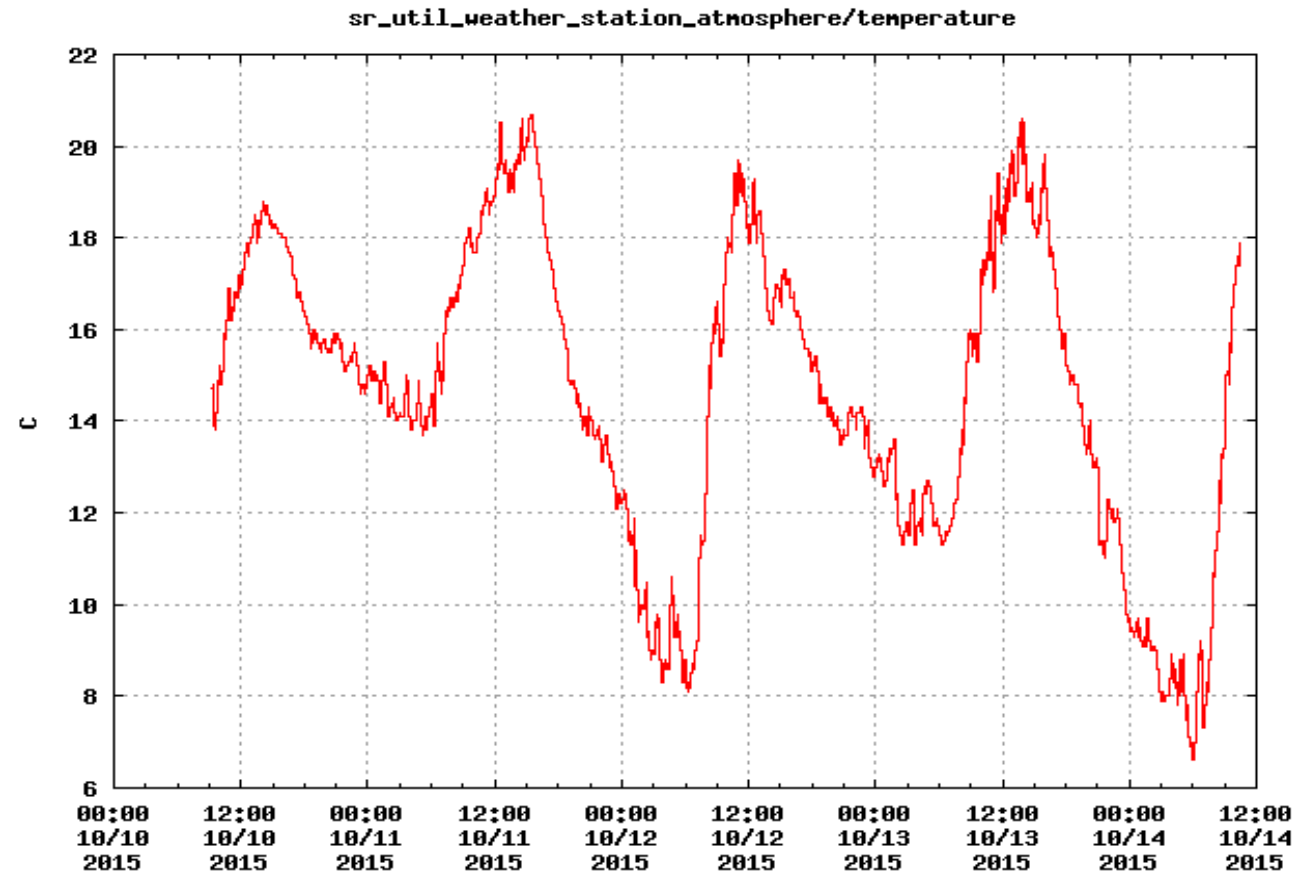
- **Why NoSQL, Why Cassandra ?**

# RDBMS is great

- **We are currently using RDBMS for**
  - **Configuration management**
  - **Parameter management**
  - **Alarm record**
  - **Etc**
- **But,**

MADOCA

# RDBMS limitation in logging

- **Performance**

- **Scalability**

- **Availability**

- **Flexibility**

# Logging in Accelerator Control

- **Time series data**

# Logging in Accelerator Control

- **Time series data**

- **Write many and rare read**

# Logging in Accelerator Control

- **Time series data**

- **Write many and rare read**

- **Is RDBMS is suitable for this task?**

MADOCA

# Logging in Accelerator Control

- **Time series data**

- **Write many and rare read**

- **Is RDBMS is suitable for this task?**

- **Looking for new database**

  - **Keeping advantage of the old system**

  - **Make up for its shortcomings**

MAD○CA

# NoSQL (Not only SQL)

- **Simplicity of design, simpler "horizontal" scaling to clusters of machines, which is a problem for relational databases, and finer control over availability. (Wikipedia)**

# NoSQL (Not only SQL)

- **Simplicity** of design, <u>simpler "horizontal" scaling</u> to clusters of machines, which is a problem for relational databases, and finer control over <u>availability</u>.  (Wikipedia)

# NoSQL variations

- - Key-value

- - Graph

- - Document

- **No solutions for time-series data in above NoSQL**

- **Wide-column**

MADOCA

# Wide-column database

- ## One type of NoSQL (Not only database)

**Row Key**

| |
|---|
| sig1:20130504 |

| |
|---|
| sig2:20130504 |

| |
|---|
| sig3:20130504 |

**Column key**

**Column value**

MADOCA

# Wide-column database

- **Columns are added when data added.**

| Row Key | Column |
|---|---|
| | t0 |
| sig1:20130504 | value0 |
| sig2:20130504 | |
| sig3:20130504 | |

**Column key**

**Column value**

# Wide-column database

- **Columns are added when data added**

| Row Key | Column | | |
|---|---|---|---|
| **sig1:20130504** | t0 | t1 | |
| | value0 | value1 | |
| **sig2:20130504** | T'0 | | |
| | Value'0 | | |
| **sig3:20130504** | T"0 | | |
| | Value"0 | | |

**Column key**

**Column value**

MADOCA

# Wide-column database

- **Row is added at any time**

| Row Key | Column | | |
|---|---|---|---|
| **sig1:20130504** | t0 | t1 | t2 |
| | value0 | value1 | value2 |
| **sig2:20130504** | T'0 | T'1 | |
| | Value'0 | Value'1 | |
| **sig3:20130504** | T"0 | | |
| | Value"0 | | |
| **sig4:20130510** | | | |

**Column key**

**Column value**

MADOCA

# Wide-column database

- **And it grows**

| Row Key | Column | | | |
|---|---|---|---|---|
| **sig1:20130504** | t0 | t1 | t2 | t3 |
| | value0 | value1 | value2 | value3 |
| **sig2:20130504** | T'0 | T'1 | T'2 | |
| | Value'0 | Value'1 | Value'2 | |
| **sig3:20130504** | T''0 | | | |
| | Value''0 | | | |
| **sig4:20130510** | T''0 | | | |
| | Value'''0 | | | |

**Column key**

**Column value**

MADOCA

# Wide-column database

- **Suitable for time-series data logging**
  - **Each data has its own time-stamp**
    - **cyclic data + event driven data in same place**
  - **Access**
    - **Key+ column range**
      - **same as current access method**

MADOCA

# Which Wide-column DB?

- **Major wide-column database**
  - **Apache Cassandra**
  - **Apache Hbase**
  - **Hypertable**

MADOCA

# Apache Cassandra

- **We select by its availability**

- **Every node has the same role**
  - **No master node**
    - **No single point of failure (SPOF)**
    - **HBase and Hypertable have masternode: SPOF**

# Apache Cassandra

- **Our criteria**
  - **Reliability**
  - **Scalability**
  - **Flexibility**
- **Consistency is covered by the other DB**

# Reliability

- **Most essential**

# Reliability

- **Most essential**

- **Cassandra**
  - **No master node, no single point of failure**
  - **Data redundancy**
    - **3 data replicas**

# Scalability

- **Just add nodes when you need more power**

  - **No cluster reboot is needed**

- **Apple is operating 100,000 node cluster for iTunes**

# Flexibility

- **Insert at any time**

# Flexibility

- **Insert at any time**
  - **Signal by signal**

# Flexibility

- **Insert at any time**

- **Data type**

# Flexibility

- **Insert at any time**

- **Data type**

  - **Store data using object serialization**

    - **Not using cassandra's data type**

    - **blob type column only**

MADOCA

# Flexibility

- **Insert at any time**

- **Data type**

  - **Store data using object serialization:MessagePack**

    - **Very fast**

    - **Low overhead  8 Byte float -> 9 Byte string**

    - **Self described**

      - **NO Interface Definition Language like Protocol Buffer**

# Consistency

- **Cassandra does not guarantee consistency**

- **In our cluster, it takes about 1 second after insert to obtain consistent value.**

  - **No real-time access**

MADOCA

# Redis

- **Covers Cassandra's inconsistency**

# Redis

- **Covers Cassandra's inconsistency**

- **Stores newest data only**

MADOCA

# Redis

- **Covers Cassandra's inconsistency**

- **Stores newest data only**

- **In-memory key-value database**

# Redis

- **Covers Cassandra's inconsistency**

- **Stores newest data only**

- **In-memory key-value database**

  - **Very fast by key access**

  - **Newest value+ meta data only**

  - **Data packed by MessagePack**

# Redis

- **Covers Cassandra's inconsistency**

- **Stores newest data only**

- **In-memory key-value database**

- **Two redis servers are running  in parallel for redundancy**

# Implementation

- **Data acquisition system**

- **Cassandra structure**

- **Performance**

# Entire system

# Entire system

# Entire system

# Message creation



| | |
|---|---|
| **Key** | **LGsr_mag_ps_b/current_adc:** |
| **Metadata** | **{"tm":1444814445053727,"tl":1296000000,"cy":1000}** |
| **Data** | **419.6774238114116** |

# Message structure

- ## 3 part message

| Key | LGsr_mag_ps_b/current_adc: |
|---|---|
| Metadata | {"tm":1444814445053727,"tl":1296000000,"cy":1000} |
| Data | 419.6774238114116 |

- ## Key: raw string

- ## Metadata and data are packed by MessagePack

MADOCA

# Convert to Cassandra command



insert('LGsr_mag_b/current_adc:201510',
{1444814445053727:419.6774238114116},
1296000000)

**Blue means MessagePacked string**

# Convert to Redis Command



insert('LGsr_mag_b/current_adc:',
'{"tm:1444814445053727}419.6774238114116')

**Blue means MessagePacked string**

# Writer

- **Converts messages to insert commands**

- **Plug-in structure using 0MQ's in-process pub/sub**

# Writer

- **Converts messages to insert commands**

- **Plug-in structure using 0MQ's in-process pub/sub**
  - **Other DB engine or anything may be added in the future**

# Structure of Cassandra

- **key : one key / one signal one day**
- **LGsr_mag_ps_b/current_adc:20151003**

**Keyspace: database**

**Column Family: Table**

| Row Key | Column | | | |
|---|---|---|---|---|
| **sig1:20130504** | t0 | t1 | t2 | t3 |
| | value0 | value1 | value2 | value3 |
| **sig2:20130504** | t0 | t1 | | |
| | value0 | value1 | | |
| **sig3:20130504** | t0 | t1 | t2 | t3 |
| | value0 | value1 | value2 | value3 |

Column key

Column Value

# Performance; write to Cassandra

## Values insert/sec (batch)

# Read from Cassandra

- One day data = 60s*60min*24hour
- Done during normal writing operation



Time to get one day data (sec)

| Mean | 0.76sec |
|------|---------|
| Sigma | 0.06sec |
| 95% | 0.86sec |

MADOCA

# Read from Redis



Online time to get one data (ms)

| Mean | 0.77ms |
|------|--------|
| Sigma | 0.47ms |
| 95% | 1.4ms |

# One year Test

- **Test performed about one year**

- **No major trouble**

  - **Test**

    - **Forced to shutdown a node and recovery**

- **Some modification are needed for the production run**

MADOCA

# For production run

- **Data migration from RDBMS**

- **Structure modification**

- **Monitoring tools**

- **Client libraries**

- **Node added**

# Data migration from RDBMS

- **Data  since 1997**

  - **4TB in RDBMS (logical file size, become larger in RAID disks)**

  - **0.75TB/node 9TB in total in 12 nodes (3 replicas)**

MADOCA

# Structure changed

- **One large column family was divided into small column families of each month**

  - **Cassandra's compaction operation**

    - **Batch operation**

    - **Columns that marked as "delete" are deleted at this time.**

MADOCA

# Temporary disk space for compaction

- One column family needs same size temporary disk space at compaction.

- One big column family cannot be larger than ½ disk space.

# Structure changed

- **One big column family was divided into small column families of  one month**

| Keyspace | |
|---|---|
| **Column Family** | |

➡️

| Keyspace | |
|---|---|
| **Column Family 201501** | |
| **Column Family 201502** | |
| **Column Family 201503** | |

- **Backup becomes easy by copying separate files**

MADOCA

# Monitoring tool

- **Server system monitoring by Zabbix.**
  - **Not only SNMP but also JXM**
    - **Cassandra is written by JAVA**
    - **JAVA VM monitoring**

- **DAQ system monitoring tool**
  - **For experts**
  - **For operators**

MADOCA

# Zabbix screen

# DAQ system monitoring for experts

# DAQ System monitoring for shift operators

# Client libraries

- **Mainly C and C++**

  - **For applications written in C**

  - **Same interfaces, no modification to source code of application**

    - **Just re-link**

- **Python modules**

  - **for Web applications**

- **We used**

  - **ZeroMQ, Messagepack, Cassandra CQL/Thrift , Redis**

MADOCA

# Cassandra Cluster

| | |
|---|---|
| Number of nodes | 12 |
| Server | Dell PowerEdge R420 |
| CPU | Intel Xeon E5-2420 2.2GHz |
| Memory | 16GB |
| System disk | 600GB 15k rpm |
| Data disks | 3TB 7200 rpm x3 |
| Cassandra | 2.0.10 |
| JavaVM | JRE1.7.0-67-b01 |

MADOCA

# Summary

- **We implemented new data acquisition and store system with new technologies**

- **Apache Cassandra provides high-performance, reliable, scalable and flexible data store that was impossible by RDBMS**

- **We build supporting infrastructures for healthy operations**

- **The system is stably running more than one year including test run**

MADOCA