# 10 Million Smart Meter Data with Apache HBase

5/31/2017

OSS Solution Center

Hitachi, Ltd.

## Masahiro Ito

**Open Source Summit Japan 2017**

HITACHI
Inspire the Next

- Masahiro Ito
  - ➢ Software Engineer at Hitachi, Ltd.
  - ➢ Focus on development of Big Data Solution with Apache Hadoop and its related OSS.
  - ➢ Mail: masahiro.ito.ph@hitachi.com

  - ➢ Book and Web-articles (in Japanese)
    - Apache Spark ビッグデータ性能検証（Think IT Books）

    - **ユースケースで徹底検証！** HBaseでIoT時代のビッグデータ管理機能を試す
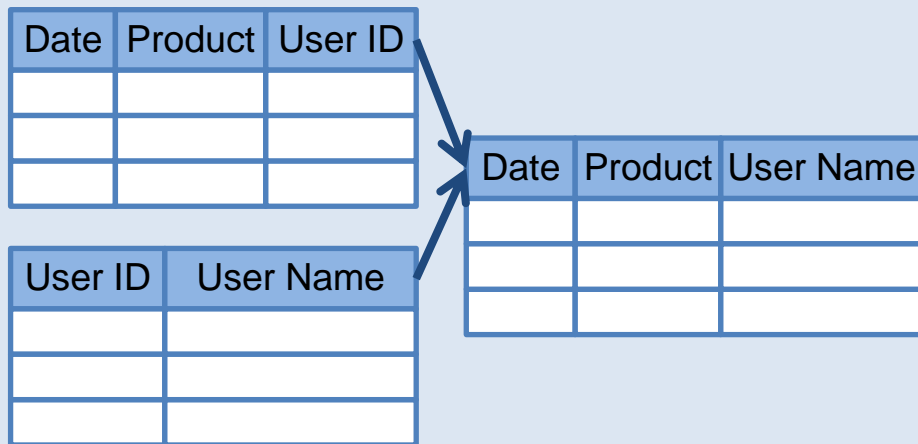      - – https://thinkit.co.jp/series/6465

# 1. Motivation

- The internet of things (IoT) and NoSQL
  - ➢ Various sensor devices generate large amounts of data.
  - ➢ NoSQL has higher performance and scalability than RDB.
  - ➢ HBase is one of NoSQL.

- Is HBase suitable for sensor data management?
  - ➢ HBase seems to be suitable for managing time series data such as sensor data.
  - ➢ I will introduce the result of performance evaluation of HBase with 10 million smart meter data.

# 2. What is NoSQL?
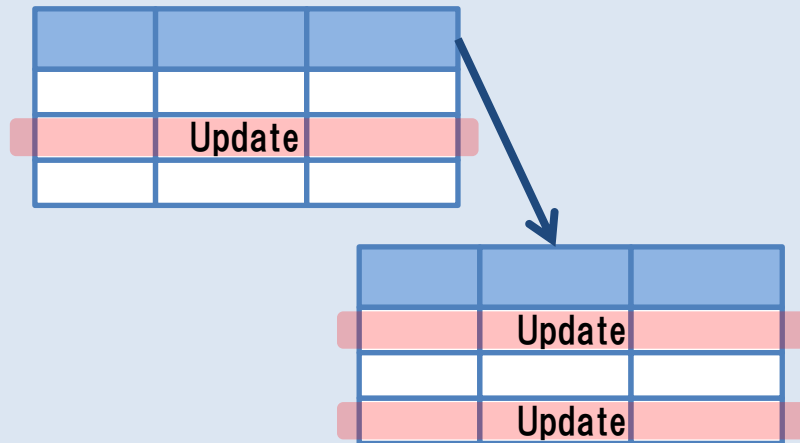
# NoSQL（Not only SQL）

- NoSQL refers to databases other than RDB (Relational DataBase).

- Motivations of NoSQL include:
  - ➢ More flexible data model (not tabular relations).
  - ➢ High performance and large disk capacity.
    - With simpler "horizontal" scaling to clusters of machines.
  - ➢ etc.

- NoSQL databases are increasingly used in big data and real-time web applications.

Relational model

| Date | Product | User ID |
|------|---------|---------|
|      |         |         |
|      |         |         |
|      |         |         |

| User ID | User Name |
|---------|-----------|
|         |           |
|         |           |
|         |           |

| Date | Product | User Name |
|------|---------|-----------|
|      |         |           |
|      |         |           |
|      |         |           |

ACID Transaction

- Table format (tabular relations)
- SQL interface
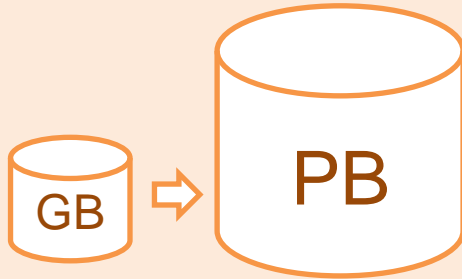  - ➢ Supports complex queries

- **A**tomicity
- **C**onsistency
- **I**solation
- **D**urability

# 3 Vs of Big Data: Challenges of RDB for big data
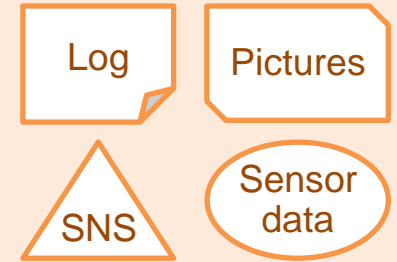
## Volume
Need to manage large amount of distributed data.

GB ⇒ PB

## Velocity
Need to process large number of requests in real time.

## Variety
Need to manage data of various structures.

Log  Pictures  SNS  Sensor data

**RDB**

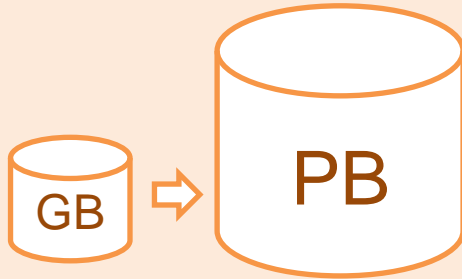| | | |
|---|---|---|
| Transaction control over distributed data is difficult. | Exclusive control of transaction is overhead. | It is incompatible with the predefined table. |

8

# 3 Vs of Big Data: Challenges of RDB for big data

## Volume
Need to manage large amount of distributed data.

GB ⇒ PB

## Velocity
Need to process large number of requests in real time.

## Variety
Need to manage data of various structures.

Log | Pictures
SNS | Sensor data

**RDB**

Transaction control over distributed data is difficult.

Exclusive control of transaction is overhead.

It is incompatible with the predefined table.

**NoSQL**

Limiting the scope of transaction control makes it possible to improve performance and disk capacity with scale out.

9

# 3 Vs of Big Data: Challenges of RDB for big data
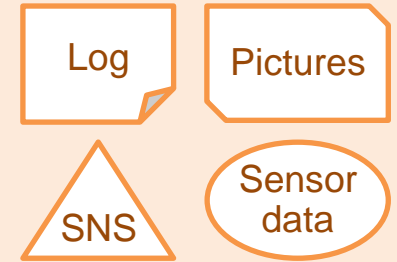
## Volume
Need to manage large amount of distributed data.

GB ⇒ PB

## Velocity
Need to process large number of requests in real time.

## Variety
Need to manage data of various structures.

Log | Pictures
SNS | Sensor data

**RDB**

Transaction control over distributed data is difficult.

Exclusive control of transaction is overhead.
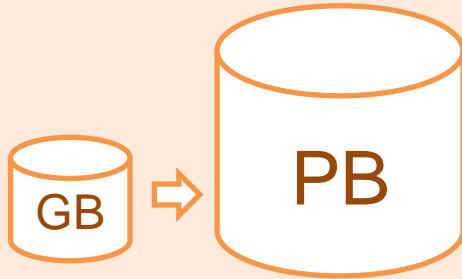
It is incompatible with the predefined table.

**NoSQL**

Limiting the scope of transaction control makes it possible to improve performance and disk capacity with scale out.

Adopted flexible data structure other than table.

10

# There are lots of NoSQL in the world (many others)

Riak

Cassandra

Redis

HBase

MongoDB

Neo4j

Couchbase

TITAN

# NoSQL is generally classified by data model



Key value store

Riak

Redis

Wide column store

Cassandra

HBase

Document store

MongoDB

Couchbase

Graph database

Neo4j

TITAN

12

# NoSQL is generally classified by data model

## Key value store

Low latency access with simple data structure.

| Key | Value |
|-----|-------|
|     |       |
|     |       |
|     |       |

## Wide column store

Each row has different number of columns.

| Key | Value | Value | Value |
|-----|-------|-------|-------|
|     |       |       |       |
|     |       |       |       |
|     |       |       |       |

Store structure data such as JSON.

| Key | Document |
|-----|----------|
| 001 | {<br>  ID: 001<br>  User: {<br>    Name: "Engineer"<br>  }<br>} |

## Document store

Represent relationship between data as graph structure.



## Graph database

13

# 3. Overview of HBase architecture

- HBase is distributed, scalable, versioned, and non-relational (wide column type) big data store.

- A Google Bigtable clone.
  - ➢ Implemented in Java based on the paper of Bigtable.

- One of the OSS in Apache Hadoop eco-system.

15

# Relationship between HBase and Hadoop（HDFS）

- HBase build on HDFS (Hadoop Distributed File System).

**Hadoop**

**MapReduce**
**[Parallel processing framework]**

**YARN (Yet Another Resource Negotiator)**
**[Cluster resource management framework]**

HBase can read/write many small data
with low latency.
⇒ HBase is a complement to HDFS.

**HBase**
**[Distributed database]**

**HDFS (Hadoop Distributed File System)**
**[Distributed File System]**

**Commodity servers**

- HDFS can read/write large files with high throughput.
- However, it is not suitable for read/write small data.

- HBase processes the request and HDFS saves the data.



Managing RegionServers

**Master Node**
- HBase Master
- HDFS NameNode
- ZooKeeper

**Client Node**
- HBase Client

Data    Data    Data

Managing data

**Slave Node**
- HBase RegionServer
- HDFS DataNode
- Disk ..... Disk

**Slave Node**
- HBase RegionServer
- HDFS DataNode
- Disk ..... Disk

**Slave Node**
- HBase RegionServer
- HDFS DataNode
- Disk ..... Disk

Data is stored in HDFS and data is replicated between nodes.

19

# Data model: Conceptual view

Rows in a table are sorted by RowKey

Each row can have a different number of columns.

Value is stored in Cell. The past values are stored together with Timestamp.

| Timestamp | Value |
|-----------|-------|
| 20170310 | CCC |
| 20170124 | BBB |
| 20160930 | AAA |

➢ This table looks like a RDB's table.

20

- Data is stored as key value.
  - The keys are sorted in the order of **RowKey**, **Column** (ColumnFamily:qualifier), **Timestamp.**
  - It is a "multi-dimensional sorted map".
    - SortedMap<**RowKey**, SortedMap<**Column**, SortedMap<**Timestamp**, **Value**>>>

**Conceptual view of Table**

| RowKey | fam1 | | fam2 | |
|---|---|---|---|---|
| | Col1 | Col2 | Col3 | Col4 |
| Row 1 | - | | Val_03 | Val_04 |
| Row 2 | Val_05 | Val_06 | - | |

**Physical view of Table**

Key — Value

| RowKey | Column (**ColumnFamily**:qualifier) | Timestamp | Type | Value |
|---|---|---|---|---|
| Row 1 | fam1:Col1 | 20170310 | Delete | - |
| Row 1 | fam1:Col1 | 20170310 | Put | Val_01 |
| Row 1 | fam2:Col3 | 20170215 | Put | Val_03 |
| Row 1 | fam2:Col4 | 20170309 | Put | Val_04 |
| Row 2 | fam1:Col1 | 20170310 | Put | Val_05 |
| Row 2 | fam1:Col2 | 20160104 | Put | Val_06 |
| Row 2 | fam2:Col3 | 20170221 | Delete | - |
| Row 2 | fam2:Col3 | 20170204 | Put | Val_07 |

# Operations and functions

- **Operations**
  - ➢ Put, Get, Scan, Delete, etc.

- **Functions**
  - ➢ Index
    - • Only be set to RowKey and Column.
  - ➢ Transaction
    - • Only within one Row.

**Put** a row

**Get** a row with random access

**Scan** multiple rows with sequential access

**Delete** a value by adding tombstones

| RowKey | Column | Timestamp | Type | Value |
|--------|--------|-----------|------|-------|
| Row 1 | fam1:Col1 | 20170310 | Delete | - |
| Row 1 | fam1:Col1 | 20170310 | Put | Val_01 |
| Row 2 | fam2:Col3 | 20170215 | Put | Val_03 |
| Row 2 | fam2:Col4 | 20170309 | Put | Val_04 |
| Row 3 | fam1:Col1 | 20170310 | Put | Val_05 |
| Row 3 | fam1:Col2 | 20160104 | Put | Val_06 |
| Row 4 | fam2:Col3 | 20170221 | Delete | - |
| Row 4 | fam2:Col3 | 20170204 | Put | Val_07 |

- How is a table physically divided?

**Table**

| RowKey | Column | ••• | Value |
|--------|--------|-----|-------|
| Row 1 | fam1:Col1 | ••• | Val_01 |
| Row 1 | fam1:Col2 | ••• | Val_02 |
| Row 1 | fam1:Col3 | ••• | Val_03 |
| Row 1 | fam2:Col1 | ••• | Val_04 |
| Row 2 | fam1:Col1 | ••• | Val_05 |
| Row 2 | fam2:Col2 | ••• | Val_06 |
| Row 2 | fam2:Col3 | ••• | Val_07 |
| Row 3 | fam1:Col1 | ••• | Val_08 |
| Row 3 | fam2:Col1 | ••• | Val_09 |
| Row 4 | fam1:Col2 | ••• | Val_10 |
| Row 4 | fam1:Col4 | ••• | Val_11 |
| Row 4 | fam2:Col3 | ••• | Val_12 |
| Row 4 | fam2:Col5 | ••• | Val_13 |

# Table is divided into **Region** with the range of RowKey

**Table**

**Region (Row1-2)**

| RowKey | Column | ••• | Value |
|--------|--------|-----|-------|
| Row 1 | fam1:Col1 | ••• | Val_01 |
| Row 1 | fam1:Col2 | ••• | Val_02 |
| Row 1 | fam1:Col3 | ••• | Val_03 |
| Row 1 | fam2:Col1 | ••• | Val_04 |
| Row 2 | fam1:Col1 | ••• | Val_05 |
| Row 2 | fam2:Col2 | ••• | Val_06 |
| Row 2 | fam2:Col3 | ••• | Val_07 |

**Region (Row3-4)**

| Row 3 | fam1:Col1 | ••• | Val_08 |
|-------|-----------|-----|--------|
| Row 3 | fam2:Col1 | ••• | Val_09 |
| Row 4 | fam1:Col2 | ••• | Val_10 |
| Row 4 | fam1:Col4 | ••• | Val_11 |
| Row 4 | fam2:Col3 | ••• | Val_12 |
| Row 4 | fam2:Col5 | ••• | Val_13 |

24

# Data is distributed on the cluster via Regions

- Automatic sharding
  - ➢ Regions are automatically split and re-distributed as data grows.

- Simple horizontal scaling
  - ➢ Adding slave nodes improves performance and expands disk capacity.



Region holds data across HBase (as cache in memory) and HDFS (as file in disk).

# Summary of HBase architecture

- ## Simple horizontal scaling
  - ➢ Adding slave nodes improves performance and expands disk capacity

- ## Data is stored as sorted key value
  - ➢ Like multi-dimensional sorted map.
  - ➢ By designing RowKey carefully, data that are accessed together are physically co-located.

- ## Limited the index and transaction
  - ➢ Index : Only be set to RowKey and Column.
  - ➢ Transaction: Only within one Row.

# 4. Performance evaluation with 10 million smart meter data

# i. Evaluation scenario

# Smart meter data management

- We assumed the Meter Data Management System for 10 million smart meters.
  - Smart meters collect consumption of electric energy from customers.
    - Send the collected data to the Meter Data Management System every 30 minutes.
  - The collected data is used for power charge calculation and demand forecast analysis, etc.



Power plants

Power Grid

Total 10 million meters

Data Analysis System

Meter Data Management System

Data from smart meters (every 30min.)

- Write 10 million records every 30 minutes in HBase.
- Read to analyze records stored in HBase.

# Contents of performance evaluation

① **Write performance**
Measure write time and throughput of 10 million records.

③ **Read performance**
Measure read time and throughput in two kinds of analysis use cases.

Analyst

10 million smart meters

Gateway servers (with HBase clients)

HBase Cluster

Analysis server (with HBase client)

② **Data compression performance**
Measure data compression ratio and compression / decompression time.

# Evaluation environment

HITACHI
Inspire the Next

1 Client Node
1 Master Node
(Virtual Machine)

1Gbps LAN

10Gbps SW

10Gbps LAN

4 Slave Nodes
(Physical Machines)

disk    disk    ...    disk

|  | Client Node | Master Node |
|---|---|---|
| CPU Core | 16 | 2 |
| Memory | 12 GB | 16 GB |
| # of disk | 1 | 1 |
| Capacity of disk | 80 GB | 160 GB |

|  | Per slave node | Total |
|---|---|---|
| CPU Core | 32 | 128 |
| Memory | 128 GB | 512 GB |
| # of disk | 6 | 24 |
| Capacity of disk | 900 GB | - |
| Total capacity of disks | 5.4 TB (5,400 GB) | 21.6 TB (21,600 GB) |

32

# Table design

- Divided the table into 400 Regions in advance.
  - ➢ 100 Regions per RegionServer
  - ➢ Region split key: 0001, 0002, …, 0399

| Region (〜0001) | Region (0001〜0002) | Region (0002〜0003) | ○ ○ ○ | Region (0399〜) |
|---|---|---|---|---|

To distribute data among Regions, add 0000 to 0399 (meter ID modulo 400) to the head of RowKey. This technique is called "Salt".

| RowKey (&lt;Salt&gt;-&lt;Meter ID&gt;-&lt;Date&gt;-&lt;Time&gt;) | Column (ColumnFamily:qualifier) | Timestamp | Type | Value |
|---|---|---|---|---|
| 0000-0000000001-20170310-1100 | CF: | | Put | 3.241 |
| 0000-0000000001-20170310-1030 | CF: | | Put | 0.863 |
| … | … | | Put | 0.430 |
| 0000-0000000001-20160910-1100 | CF: | | Put | 0.044 |
| 0001-0000000002-20170310-1100 | CF: | | Put | 2.390 |
| … | … | | Put | 1.432 |

ii.　Evaluation of write performance

# Evaluation of write performance

- Generate 10 million records with HBase clients.
- Send put request using multi clients.
- Measured the write time and throughput of 10 million records.



10 million
smart meters

Gateway servers
(with HBase clients)

HBase Cluster
(RegionServers)

HBase client

**Tuning parameters**
① # of clients
② # of send records per request

**Tuning parameters**
③ # of Regions

35

- Write time and throughput of 10 million records.



- Stored multiple records by one request:
  - ➢ Records per request: 1 to 10,000 ⇒ Throughput: 526 to 46,729 records/sec (89x)
- Increased the number of clients:
  - ➢ # of Clients: 1 to 64 ⇒ Throughput: 46,729 to 327,869 records/sec (7x)

36

# iii. Evaluation of Compression performance

# Compressor and data block encoding

- HBase tends to increase data size for the following reasons.
  - ➢ The number of records increases because data is stored in key value format.
  - ➢ Each record length is long because a key is composed of many fields.

- Compress data with a combination of compressor and data block encoding.

| PREFIX | | SNAPPY |
| PREFIX_TREE | ✕ | GZIP |
| DIFF | | |
| FAST_DIFF | | LZ4 |

**Data Block Encoding**
Limit duplication of information in keys.

**Compressors**
Compress block of HFiles.

- Measured the file size, write time, and read time of 10 million records.

38

# Data block encoding performance with 10 million records

## HFile size

Encoding

| | |
|---|---|
| NONE | 586 MB |
| PREFIX | 425 MB |
| PREFIX_TREE | 404 MB |
| FAST_DIFF | 311 MB |
| DIFF | 311 MB |

0 MB  200 MB  400 MB  600 MB  800 MB

**HFile size**

Reduced to 53%
by DIFF encoding

## Write time

| | |
|---|---|
| NONE | 31 sec |
| PREFIX | 55 sec |
| PREFIX_TREE | 47 sec |
| FAST_DIFF | 50 sec |
| DIFF | 46 sec |

0 sec  20 sec  40 sec  60 sec

**Write time**

Increased 48%
by DIFF encoding

## Read time

| | |
|---|---|
| NONE | 45 sec |
| PREFIX | 45 sec |
| PREFIX_TREE | 50 sec |
| FAST_DIFF | 46 sec |
| DIFF | 43 sec |

0 sec  20 sec  40 sec  60 sec

**Read time**

Reduced 4%
by DIFF encoding

# Compressor performance with 10 million records

# Compressor and data block encoding performance with 10 million records



## HFile size

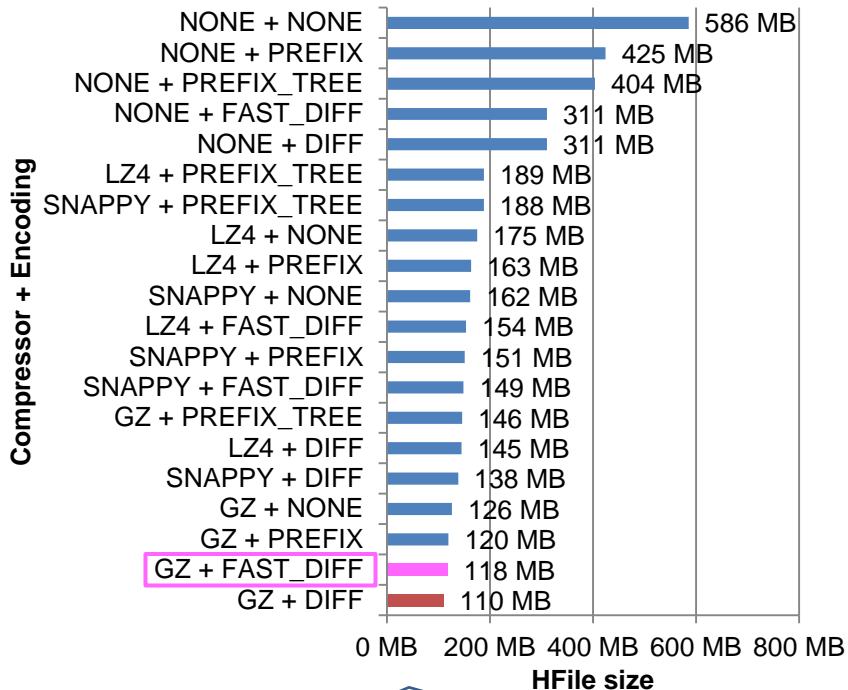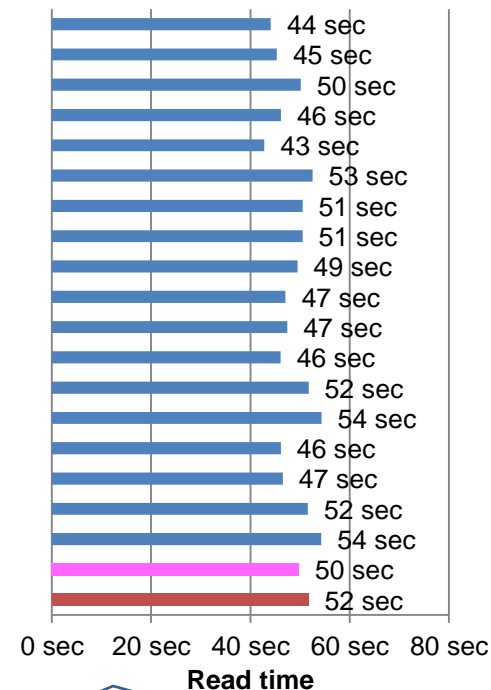| Compressor + Encoding | HFile size |
|---|---|
| NONE + NONE | 586 MB |
| NONE + PREFIX | 425 MB |
| NONE + PREFIX_TREE | 404 MB |
| NONE + FAST_DIFF | 311 MB |
| NONE + DIFF | 311 MB |
| LZ4 + PREFIX_TREE | 189 MB |
| SNAPPY + PREFIX_TREE | 188 MB |
| LZ4 + NONE | 175 MB |
| LZ4 + PREFIX | 163 MB |
| SNAPPY + NONE | 162 MB |
| LZ4 + FAST_DIFF | 154 MB |
| SNAPPY + PREFIX | 151 MB |
| SNAPPY + FAST_DIFF | 149 MB |
| GZ + PREFIX_TREE | 146 MB |
| LZ4 + DIFF | 145 MB |
| SNAPPY + DIFF | 138 MB |
| GZ + NONE | 126 MB |
| GZ + PREFIX | 120 MB |
| GZ + FAST_DIFF | 118 MB |
| GZ + DIFF | 110 MB |

## Write time

| Compressor + Encoding | Write time |
|---|---|
| NONE + NONE | 31 sec |
| NONE + PREFIX | 55 sec |
| NONE + PREFIX_TREE | 47 sec |
| NONE + FAST_DIFF | 50 sec |
| NONE + DIFF | 46 sec |
| LZ4 + PREFIX_TREE | 51 sec |
| SNAPPY + PREFIX_TREE | 41 sec |
| LZ4 + NONE | 63 sec |
| LZ4 + PREFIX | 50 sec |
| SNAPPY + NONE | 51 sec |
| LZ4 + FAST_DIFF | 49 sec |
| SNAPPY + PREFIX | 42 sec |
| SNAPPY + FAST_DIFF | 41 sec |
| GZ + PREFIX_TREE | 46 sec |
| LZ4 + DIFF | 47 sec |
| SNAPPY + DIFF | 52 sec |
| GZ + NONE | 45 sec |
| GZ + PREFIX | 46 sec |
| GZ + FAST_DIFF | 41 sec |
| GZ + DIFF | 51 sec |

## Read time

| Compressor + Encoding | Read time |
|---|---|
| NONE + NONE | 44 sec |
| NONE + PREFIX | 45 sec |
| NONE + PREFIX_TREE | 50 sec |
| NONE + FAST_DIFF | 46 sec |
| NONE + DIFF | 43 sec |
| LZ4 + PREFIX_TREE | 53 sec |
| SNAPPY + PREFIX_TREE | 51 sec |
| LZ4 + NONE | 51 sec |
| LZ4 + PREFIX | 49 sec |
| SNAPPY + NONE | 47 sec |
| LZ4 + FAST_DIFF | 47 sec |
| SNAPPY + PREFIX | 46 sec |
| SNAPPY + FAST_DIFF | 52 sec |
| GZ + PREFIX_TREE | 54 sec |
| LZ4 + DIFF | 46 sec |
| SNAPPY + DIFF | 47 sec |
| GZ + NONE | 52 sec |
| GZ + PREFIX | 54 sec |
| GZ + FAST_DIFF | 50 sec |
| GZ + DIFF | 52 sec |

Reduced to 19%
by GZip + FAST_DIFF

Increased 33%
by GZip + FAST_DIFF

Increased 14%
by GZip + FAST_DIFF

iv. Evaluation of read performance
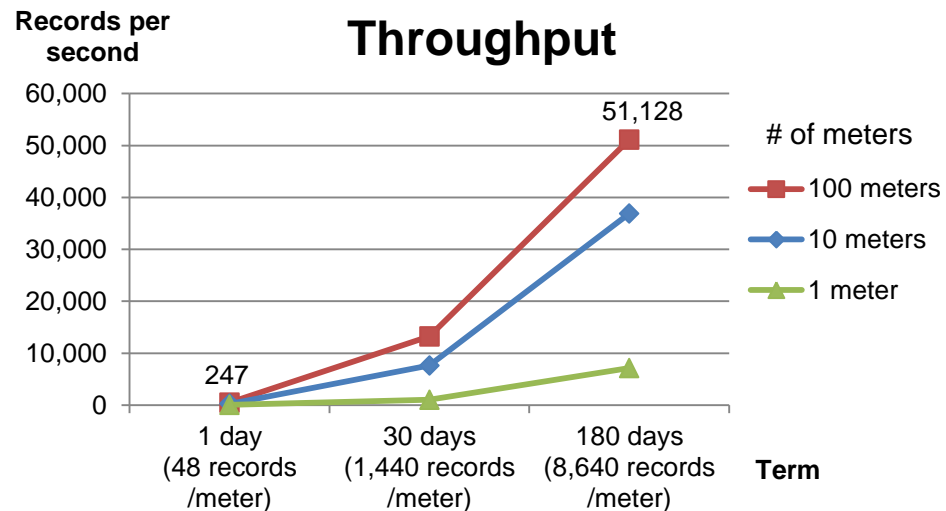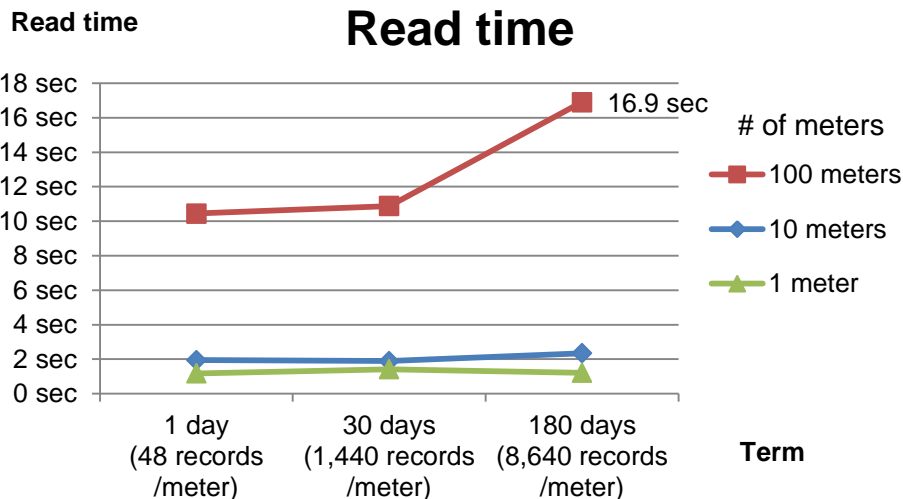
# Evaluation of read performance

- Measure the read time and throughput in two kinds of analysis use cases.
  - ➢ **Use case A**: Scan time series data of a few meters.
    - To display the transition of power consumption per meter in the line chart.

  - ➢ **Use case B**: Get the latest data of many meters.
    - To calculate the average and total value of the latest power consumption.

  - ➢ Evaluation settings
    - Dataset: 10 million meter * 180 days records (Compressed by FAST_DIFF + GZ )
    - Disabled caches and make sure to read data from disk.



HBase Cluster
(RegionServers)

Analysis server
(with HBase client)

Analyst

**Read**

HBase
client

**Tuning parameters**
① # of request threads

# Use case A: Scan time series data of a few meters

- ## Scan meter data for **1-180 days** of **1-100 meters**.
  - ➤ Scan time series data of one meter by one scan.
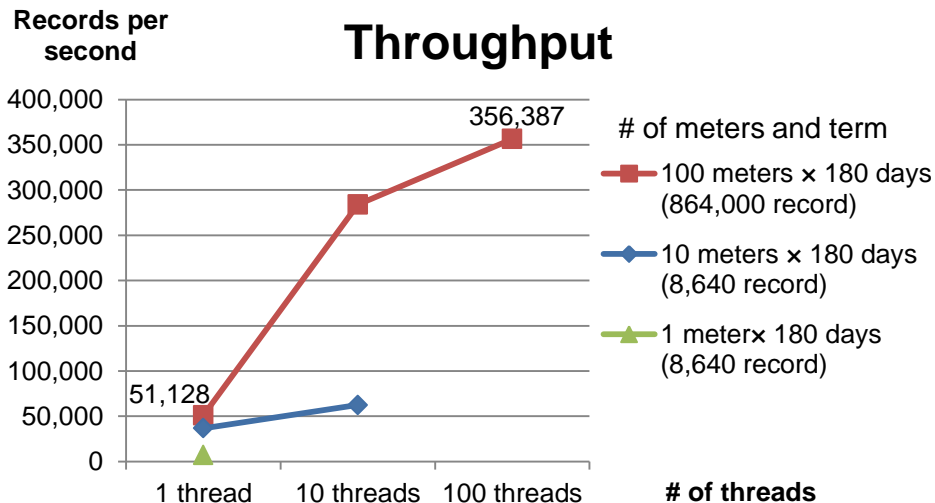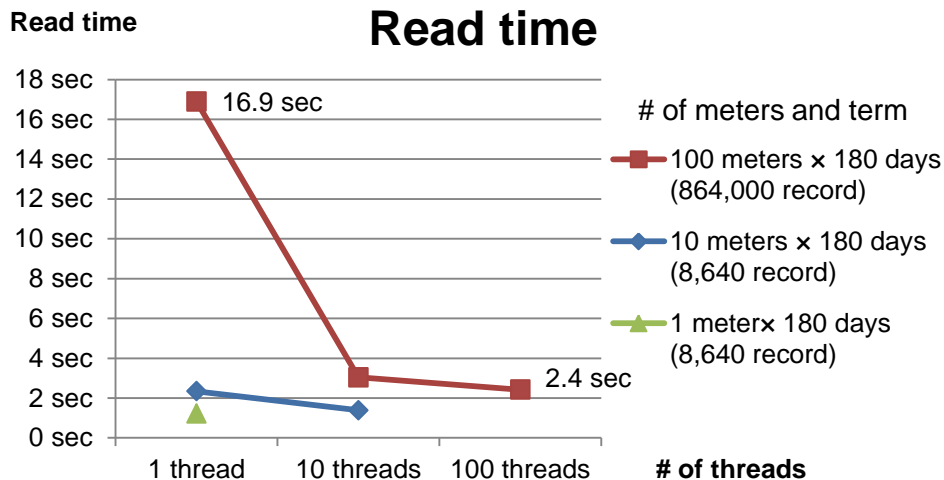


**Since read multiple data with one Scan, the throughput improves as the term was longer.**
➤ Term: 1 to 180 days ⇒ Throughput: 247 to 51,128 records/sec (207x)

44

# Use case A: Scan time series data of a few meters（with multi thread）

- ## Scan meter data for **180 days** of **1-100 meters**.
  - ➢ Scan request was executed in multi thread. (Maximum 1 Scan 1 thread)

**Read time**

### Read time



18 sec
16 sec — 16.9 sec
14 sec
12 sec
10 sec
8 sec
6 sec
4 sec
2 sec — 2.4 sec
0 sec

1 thread | 10 threads | 100 threads

**# of threads**

# of meters and term

■ 100 meters × 180 days
(864,000 record)

◆ 10 meters × 180 days
(8,640 record)

▲ 1 meter× 180 days
(8,640 record)

**Records per second**

### Throughput



400,000
350,000 — 356,387
300,000
250,000
200,000
150,000
100,000
51,128 — 50,000
0

1 thread | 10 threads | 100 threads

**# of threads**

# of meters and term

■ 100 meters × 180 days
(864,000 record)

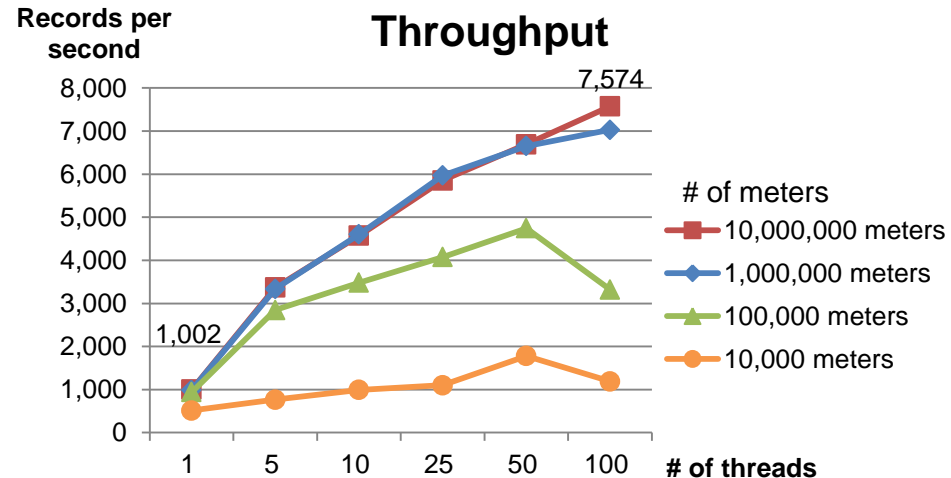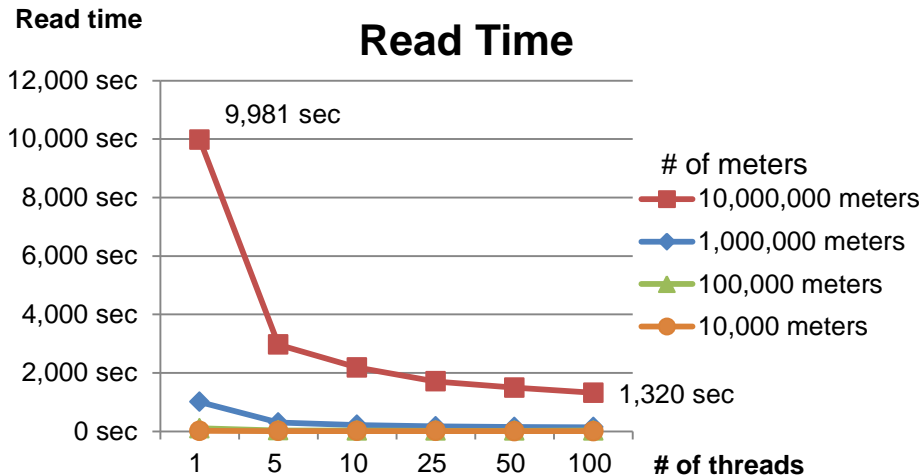◆ 10 meters × 180 days
(8,640 record)

▲ 1 meter× 180 days
(8,640 record)

**Throughput was improved by running Scan requests in parallel.**
➢ **# of threads: 1 to 100 ⇒ Throughput: 51,128 to 356,387 records/sec（7x）**

45

# Use case B: Get the latest data of many meters (with multi thread)

- Get the latest time (30 minutes) data of **10,000 to 10 million meters**.
  - ➢ Scan request can not be applied to these data.
  - ➢ Requests are executed in multi thread.
  - ➢ Batch execution of multiple "Get" request by one "batch" request.



**Read Time**

Read time

9,981 sec

# of meters
- 10,000,000 meters
- 1,000,000 meters
- 100,000 meters
- 10,000 meters

1,320 sec

# of threads

**Throughput**

Records per second

7,574

# of meters
- 10,000,000 meters
- 1,000,000 meters
- 100,000 meters
- 10,000 meters

1,002

# of threads

**Throughput was improved by running Get requests in parallel.**
  ➢ # of threads: 1 to 100  ⇒  Throughput: 1,002 to 7,574 records/sec (7.5x)

46

# Comparison of Scan request with Get request
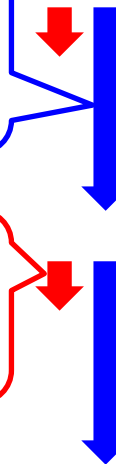
**Use case A:**
**Scan** 180 days time series data of 100 meters with 100 thread.
= **Throughput 356,387 records/second**

**Use case B:**
**Get** the latest 30 min. data of 10,000,000 meters with 100 thread.
= **Throughput 7,574 records/second**

| RowKey (<Salt>-<Meter ID>-<Date>-<Time>) | ... | Value |
|---|---|---|
| 0000-0000000001-20170310-1100 | | 3.241 |
| 0000-0000000001-20170310-1030 | | 0.863 |
| ... | | ... |
| 0000-0000000001-20160910-1100 | | 0.044 |
| ... | | ... |
| 0200-0000000201-20170310-1100 | | 10.390 |
| 0200-0000000201-20170310-1030 | | 14.325 |
| ... | | ... |
| 0200-0000000201-20160910-1100 | | 9.32 |
| ... | | ... |

- Scan request's throughput was about 47x higher than the Get request.
- Careful RowKey design is important.
  - Place the data that are accessed together physically co-located.

# 5. Summary

- HBase is suitable for storing time series data generated by sensor devices.

- Lessons from performance evaluation:
  - ➢ Careful RowKey design to be able to scan data is important.
    - Scan request's throughput was more than 47x that of Get request.
  - ➢ HBase has high multi-client / multi-thread concurrency.
    - Throughput of the Put / Scan / Get request with multi-client / multi-thread is 7x faster than single-client / single-thread.
  - ➢ Choosing the appropriate compression setting.
    - The storage size of time series data could be reduced to 19%.

# Trademarks

- Apache HBase and Apache Hadoop are either a registered trademark or a trademark of Apache Software Foundation in the United States and/or other countries.
- Other company and product names mentioned in this document may be the trademarks of their respective owners.