

---

# Leveraging smart meter data for electric utilities:

## Comparison of Spark SQL with Hive

5/16/2017

Hitachi, Ltd. OSS Solution Center

Yusuke Furuyama

Shogo Kinoshita

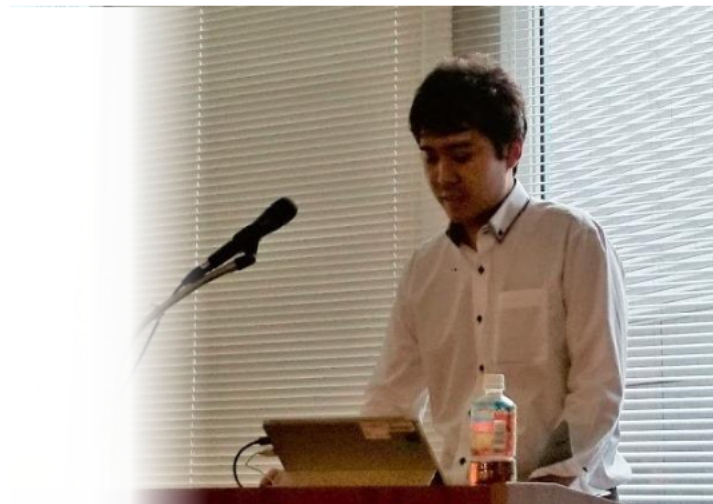
## ◆ Yusuke Furuyama

- Solutions engineer at Hitachi, Ltd.
- Offering and co-creating progressive Hadoop solutions to customers who are going to build enterprise system.



## ◆ Shogo Kinoshita

- Solutions Engineer at Hitachi, Ltd.
- Focusing on Hadoop eco-system (including Spark, Hive, Impala) and write web-articles, make presentations about evaluation of Hadoop-related OSS.



## Contents

---

1. Leveraging smart meter data [Sample use case for electric utilities]
2. Performance evaluation of MapReduce and Spark 1.6 (using Hive and Spark SQL)
3. Additional evaluation with Spark 2.1
4. Summary

---

**1. Leveraging smart meter data**  
**[Sample use case for electric utilities]**



Hitachi, Ltd.  
President & CEO  
Toshiaki Higashihara

**Established**

**February 1, 1920**

**Capital**

**458.7 billion yen**  
(as of end of Mar. 2017)

**Number of Employees**

**303,887**  
(as of end of Mar. 2017)

**Revenues**

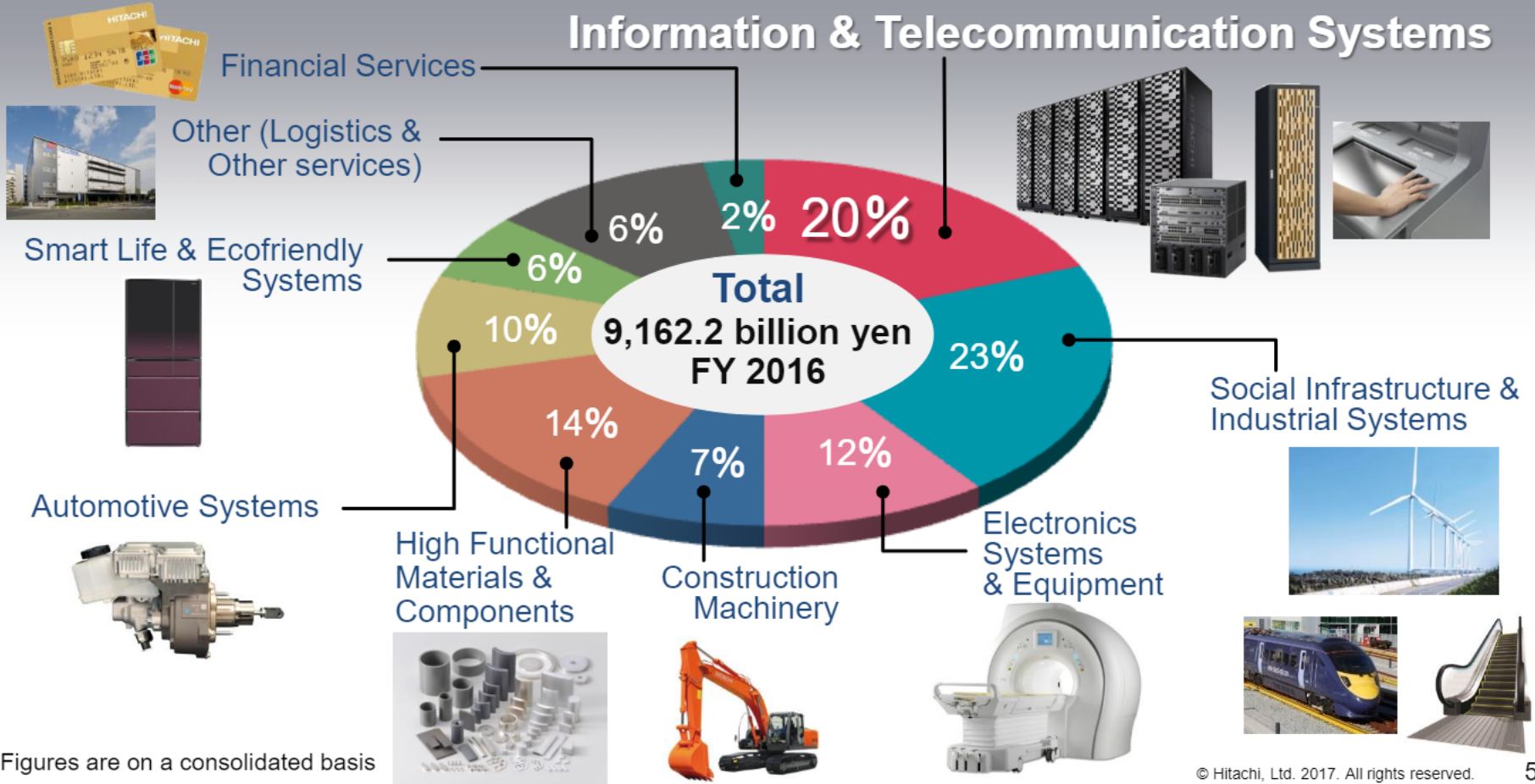
**9,162.2 billion yen**  
(FY2016 Consolidated)

**Operating Income**

**587.3 billion yen**  
(FY2016 Consolidated)

# 1-2 Revenue by Segments FY2016

## Information & Telecommunication Systems



\*Figures are on a consolidated basis

Solutions to social issues and social innovation  
by collaborative creation and open innovation

OT×IT Integrated Services / Cross-industrial Business / Business Ecosystem

O&M Service

Datahealth

Fintech

Energy Management

Total SCM

OT

Railway/machinery

×  
IT

Power

×  
IT

Healthcare

×  
IT

Manufacturing

×  
IT

Finance

×  
IT

IT

# 1-4 Situation of electric utilities in Japan and their needs

## ◆ Liberalization of the retail Electric Power Market

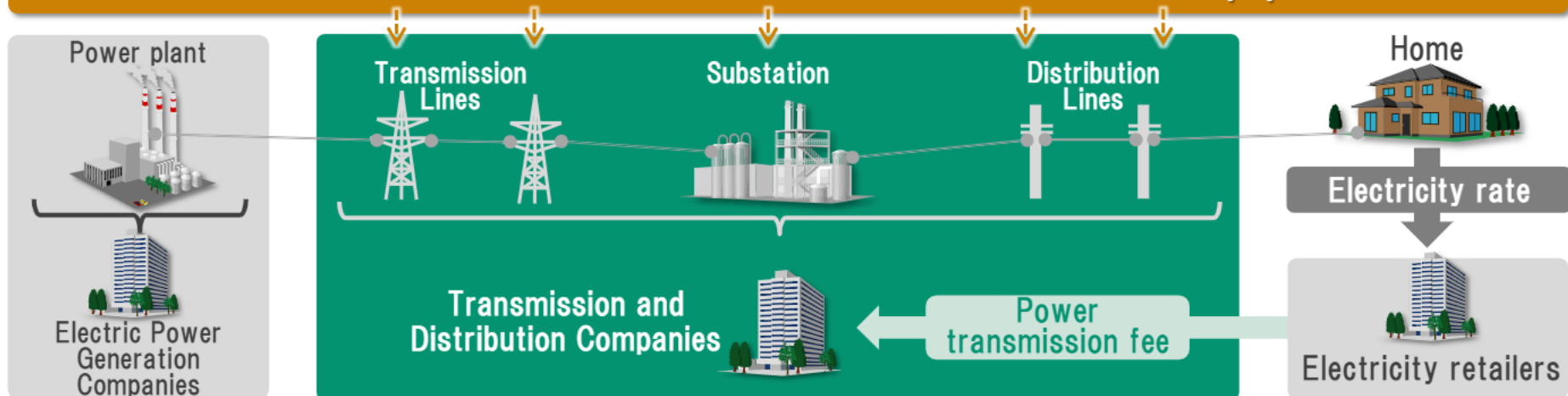
- Electric utilities in Japan have to adapt to competitive free market
- Request for price cut of power transmission fee from government of Japan

In Japan, the retail sale of electric power was fully liberalized in April 2016

## ◆ Needs to cut the cost for transmission and distribution equipment

- Transmission and distribution equipment have been replaced periodically
- Decide the timing of replacement by the condition of equipment

Maintenance team needs: Obtain the load status of each equipment

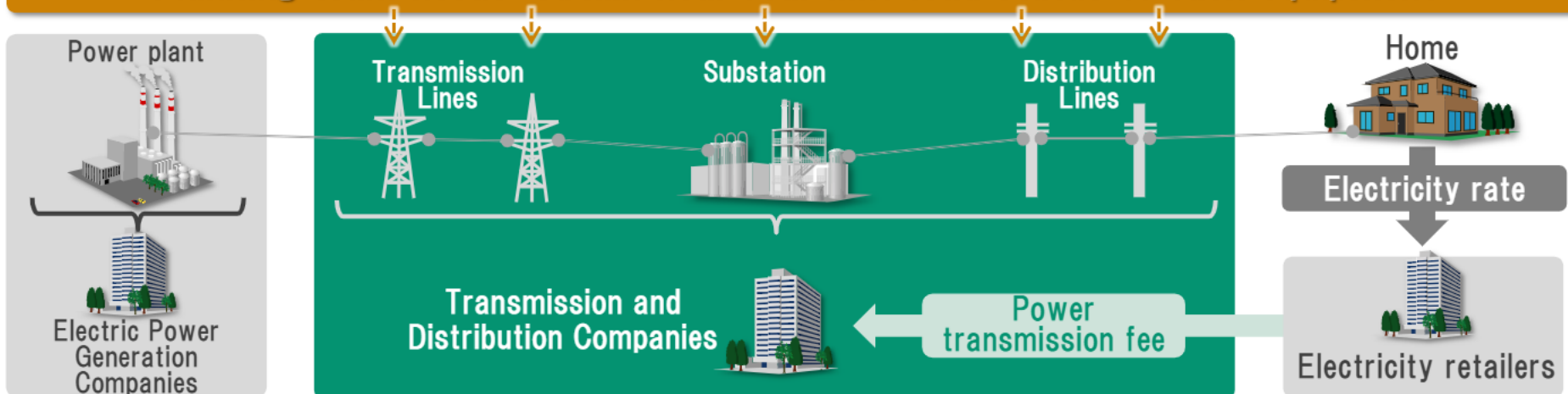




# 1-5 Situation of electric utilities in Japan and their needs (future)

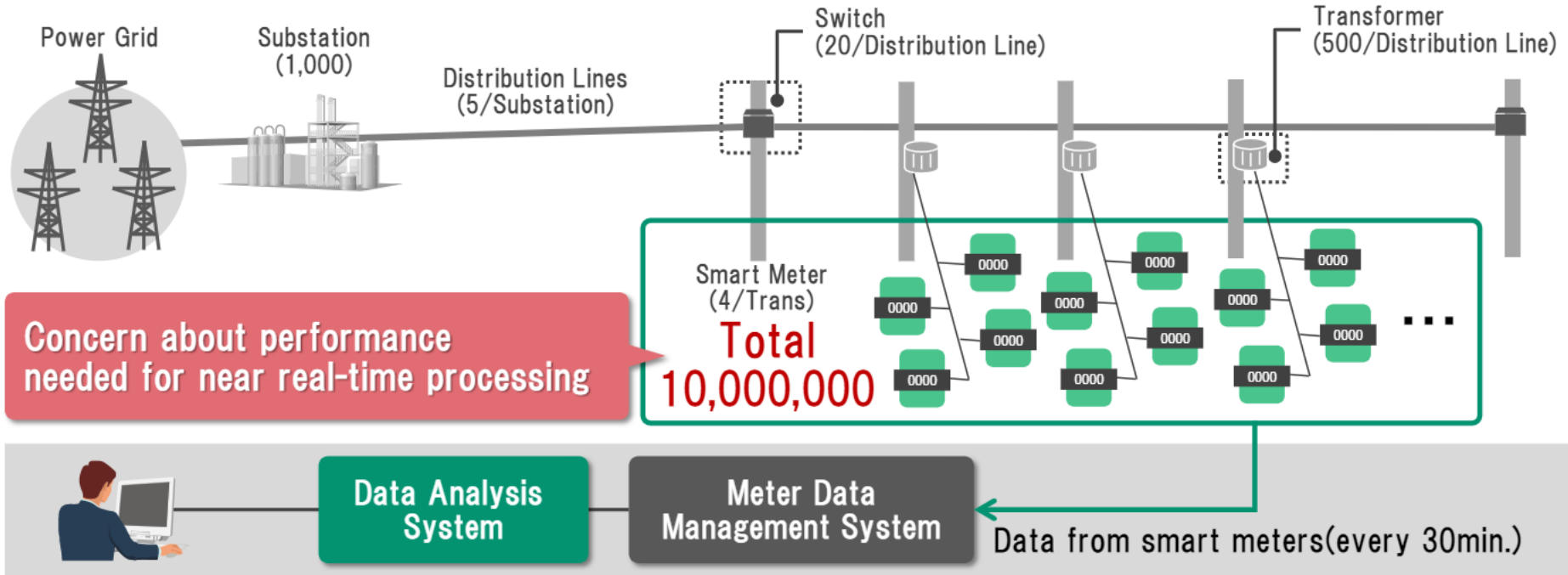
- ◆ **Unstable power supply**
  - Decreasing nuclear plant as a stable power supplier
  - Increasing renewable energy supply
- ◆ **Needs for high level Demand Response**
  - Rates by time zone (current demand response)
  - Many and small renewable energy suppliers
  - Near real-time demand response for each distribution system

**Planning team needs: Obtain near real-time load status of each equipment**



# 1-6 Leveraging big data for electric utilities

## ◆ Meet the needs of electric utilities



Concern about performance needed for near real-time processing

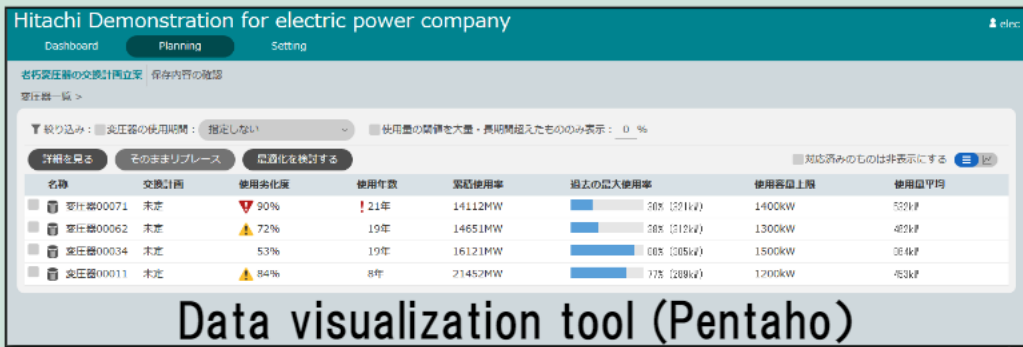
Analyze the data from smart meters to grasp the load status of equipment

## Data Analysis System



Planner  
(Equipment/Demand  
response)

Target of this session



Data processing platform (Hadoop, Spark)

Raw Data

Processed Data

Meter Data  
Management System

Data from smart meters  
(every 30min.)

---

## 2. Performance evaluation of MapReduce and Spark 1.6 (using Hive and Spark SQL)

### ◆ Needs of electric utilities (recap)

#### Needs

- Analyze the data from smart meters to grasp the load status of equipment
- Near real-time

### ◆ Points of analysis

Needs	Point of analysis
Find the equipment that needs to be replaced	Find the equipment that has heavy workload
Estimate the timing for replacement	Check the trend of load status
Select the proper capacity of new equipment	Extract the peak of the load

### ◆ Point of evaluation for near real-time processing

- Concern about performance for processing 10,000,000 meters
- Data comes from each smart meter every 30min (48/Day, spec of smart meter)

Check if MapReduce and Spark can process the data in 30min.

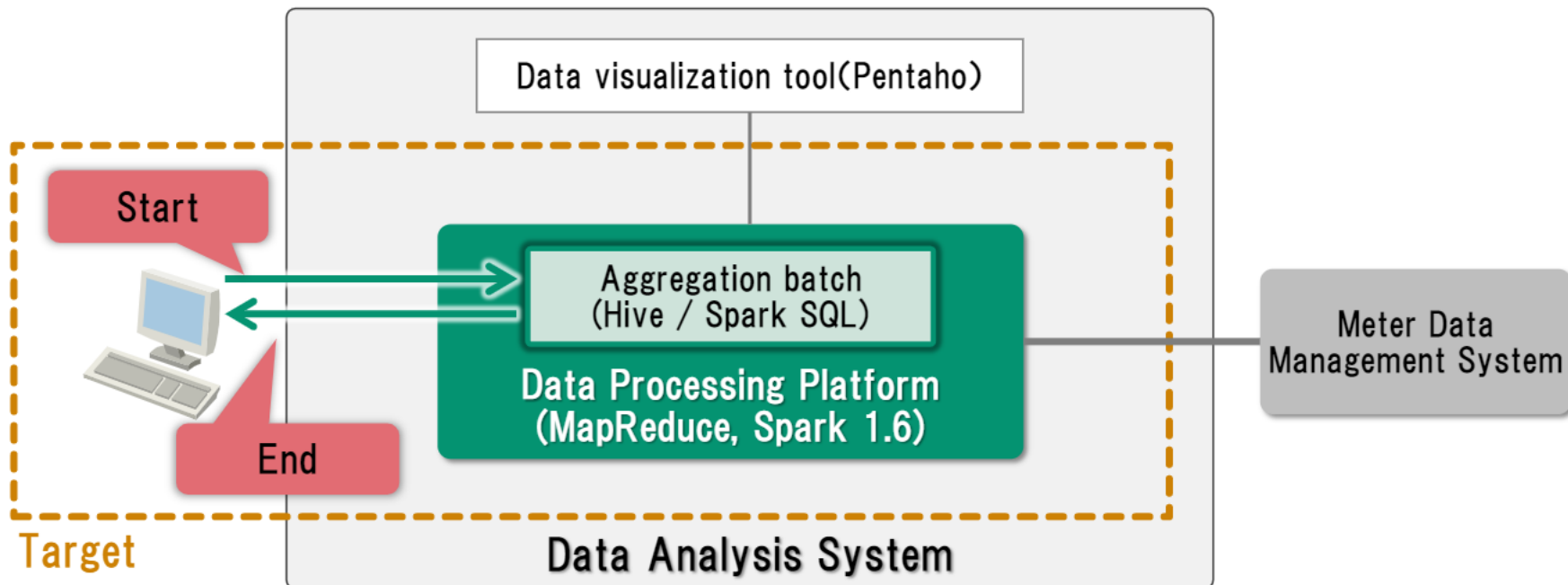
### ◆ Items for evaluation

Point of analysis	Aggregate per	
Find the equipment has heavy workload	Equipment	→ •Distribution system   •Substation   •Switch •Transformer   •Meter
Check the trend of load status	Term	→ •1day   •1month(30days)   •1year(365days)
Extract the peak of the load	Time zone	→ •Specific 30min of each day   •24h

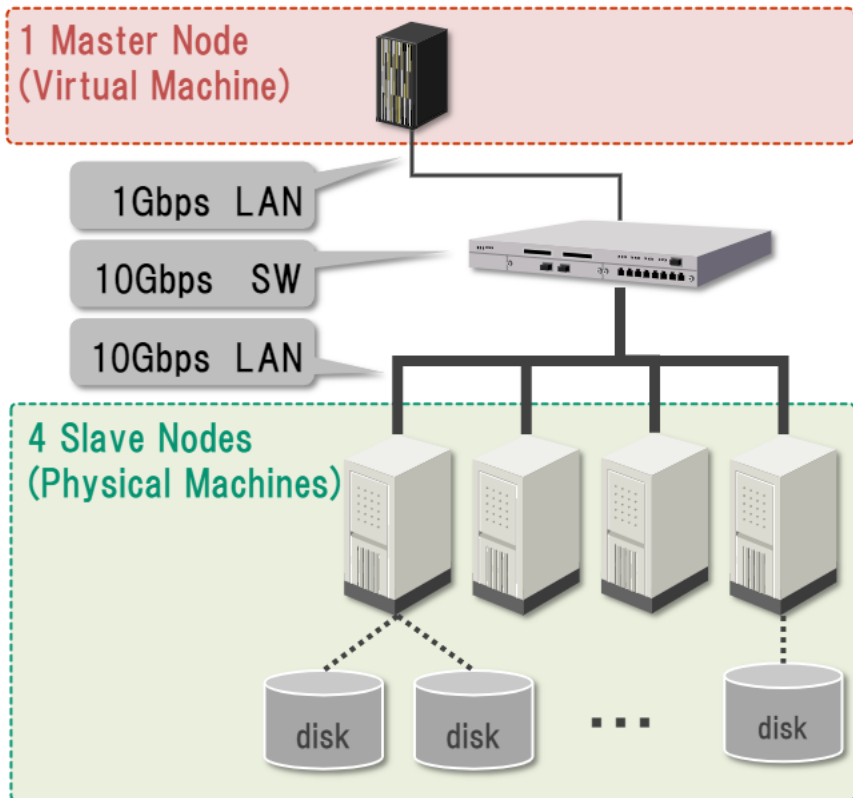
## 2-3 Target of performance evaluation

### ◆ Target of evaluation

Time from start of aggregation batch for meter data to end of the batch



## ◆ System Configuration



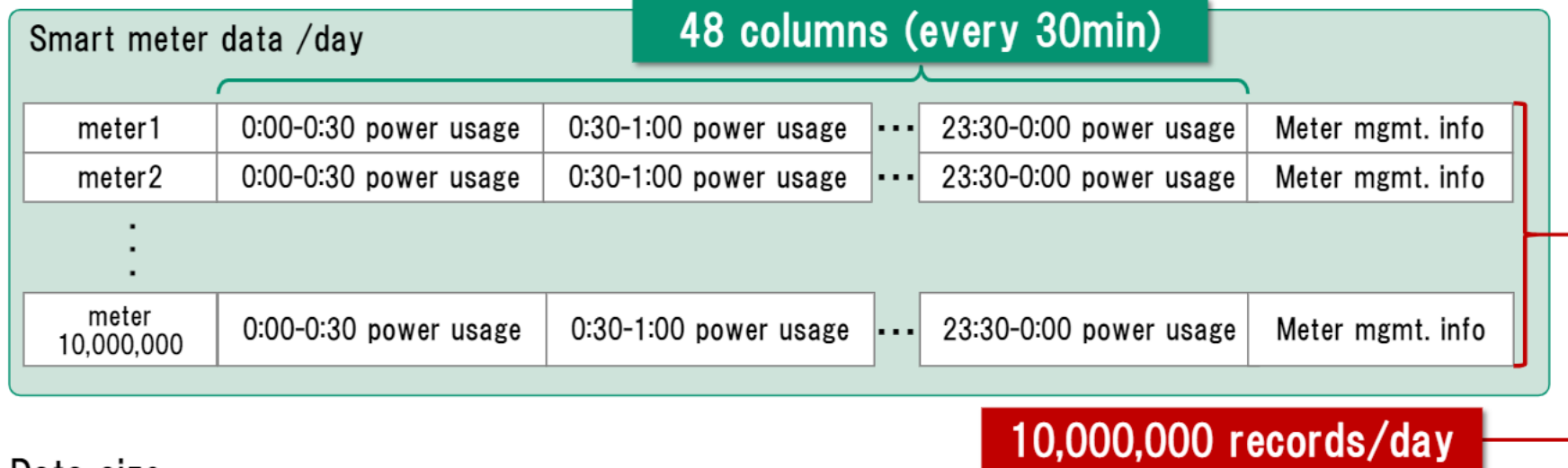
## ◆ Spec

	Master Node
CPU Core	2
Memory	8 GB
Capacity of disk	80 GB
# of disk	1

	Per slave node	total
CPU Core	16	64
Memory	128 GB	512 GB
Capacity of disk	900 GB	-
# of disk	6	24
Total capacity of disks	5.4 TB (5,400 GB)	21.6 TB (21,600 GB)



## ◆ Smart meter data



### Data size

Term	# of records	Size (CSV)	Size (ORCFile)
365days (1year)	3,650 million	2.475 TB	1.325 TB
30days (1month)	300 million	0.205 TB	0.158 TB
1day	10 million	0.007 TB	0.005 TB

## 2-6 Contents of performance evaluation (recap)

### ◆ Point of evaluation

Check if MapReduce and Spark can process the data in 30min.

### ◆ Target of evaluation

Time from start of aggregation batch for meter data to end of the batch

### ◆ Items for evaluation

Point of analysis	Aggregate per
Find the equipment has heavy workload	Equipment
Check the trend of load status	Term
Extract the peak of the load	Time zone



▪ Distribution system   ▪ Substation   ▪ Switch  
▪ Transformer   ▪ Meter



▪ 1day   ▪ 1month(30days)   ▪ 1year(365days)



▪ Specific 30min of each day   ▪ 24h

## 2-6 Contents of performance evaluation (recap)

### ◆ Point of evaluation

Check if MapReduce and Spark can process the data in 30min.

### ◆ Target of evaluation

Time from start of aggregation batch for metadata

+ File type  
▪ Text (CSV)  
▪ ORCFile (Column-based)

### ◆ Items for evaluation

Point of analysis	Aggregate per
Find the equipment has heavy workload	Equipment
Check the trend of load status	Term
Extract the peak of the load	Time zone

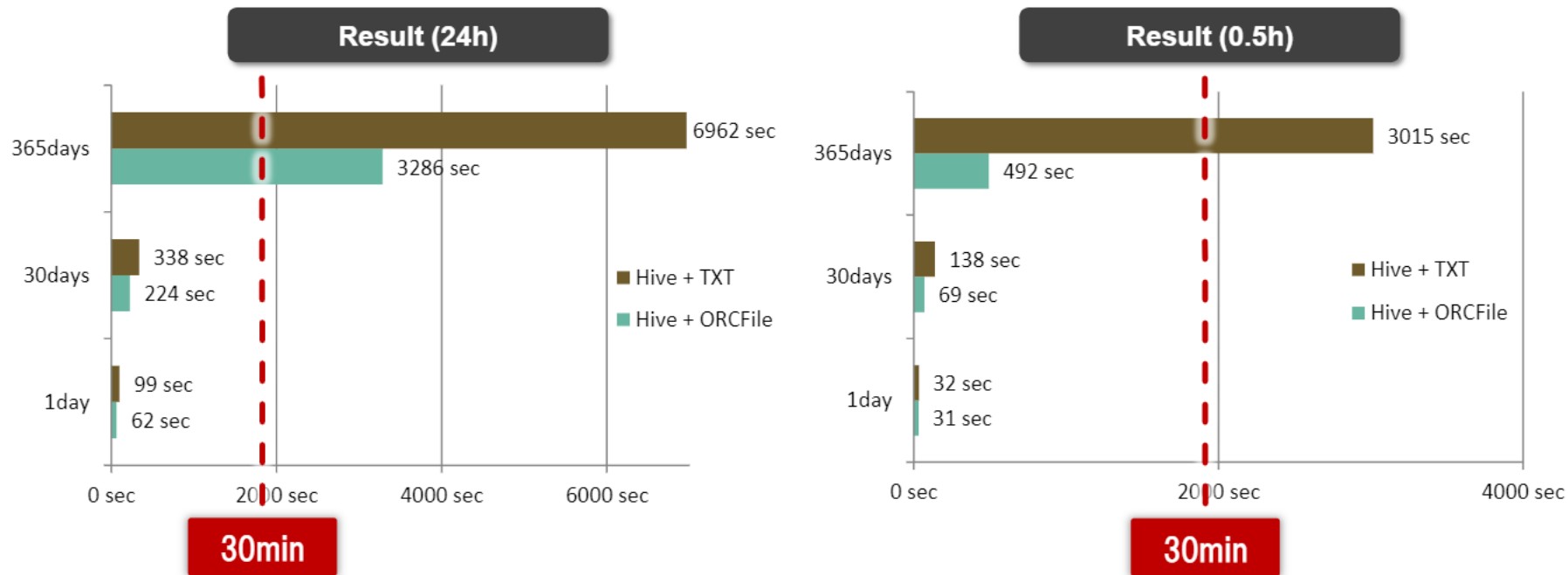
▪ Distribution system    ▪ Substation    ▪ Switch  
▪ Transformer    ▪ Meter

▪ 1day    ▪ 1month(30days)    ▪ 1year(365days)

▪ Specific 30min of each day    ▪ 24h

## 2-7 Comparison of txt with ORCFile (MapReduce)

### ◆ Aggregate meter data of entire Distribution System

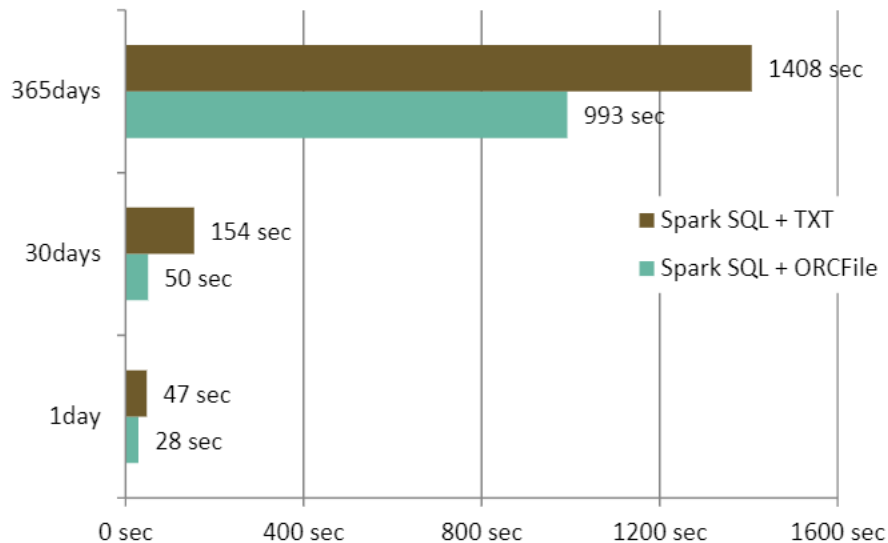


- Couldn't finish processing in 30min
- Performance improvement by ORCFile

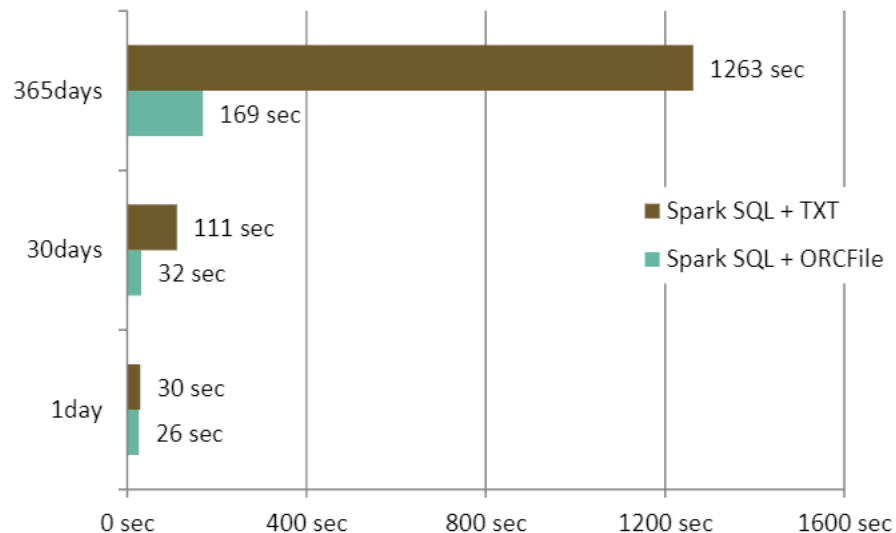
## 2-8 Comparison of txt with ORCFile (Spark 1.6)

### ◆ Aggregate meter data of entire Distribution System

Result (24h)



Result (0.5h)



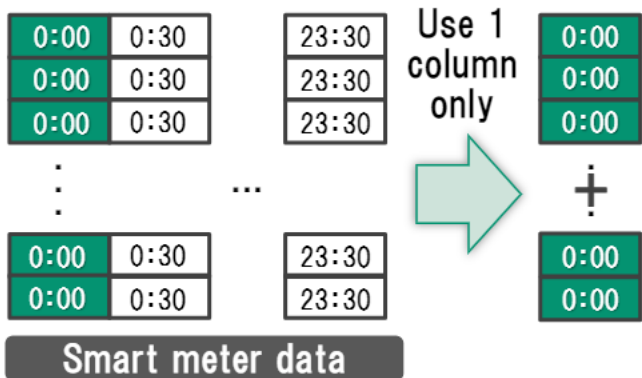
- Could finish processing in 30min (1,800s)
- Performance improvement by ORCFile

## 2-9 Review of the results

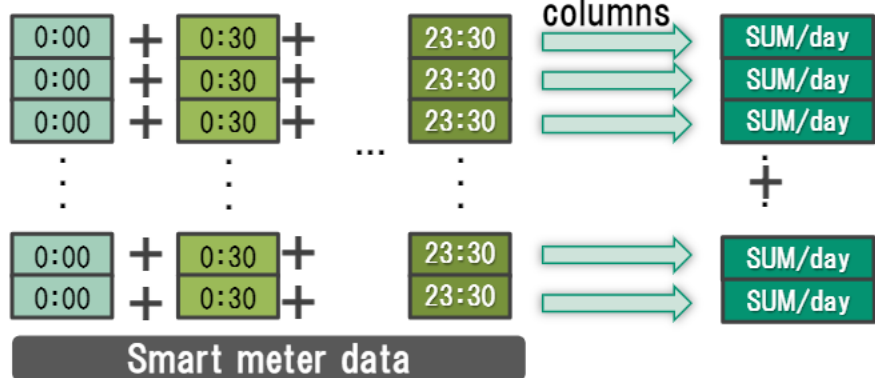
### ◆ Why the processing was fast with ORCFile



### ◆ Processing 0.5h data

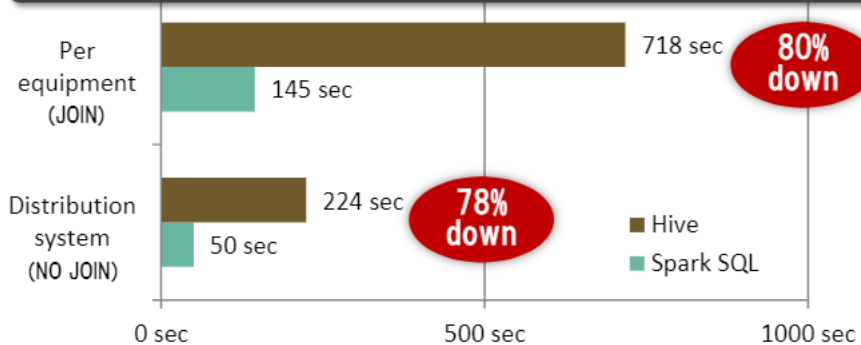


### ◆ Processing 24h data

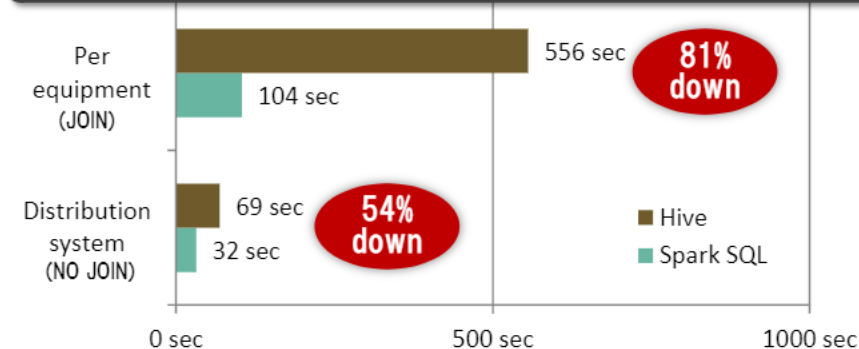


# 2-10 Comparison of MapReduce with Spark 1.6 (ORCFile)

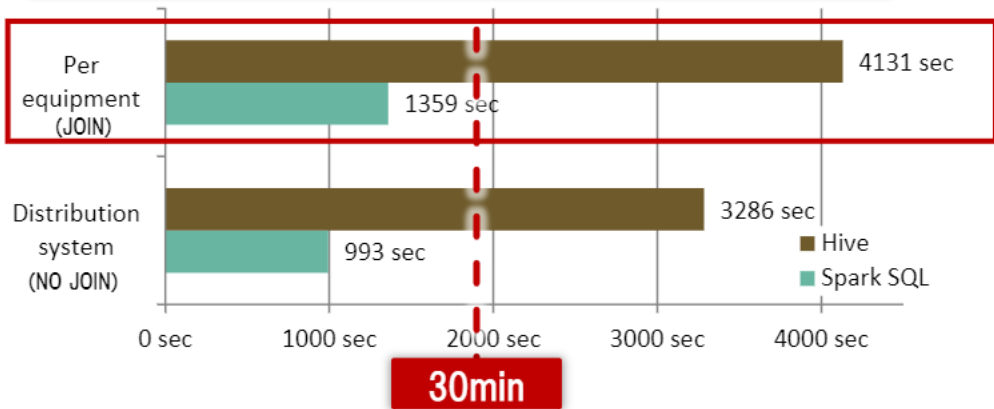
## Result (30days / 24h)



## Result (30days / 0.5h)



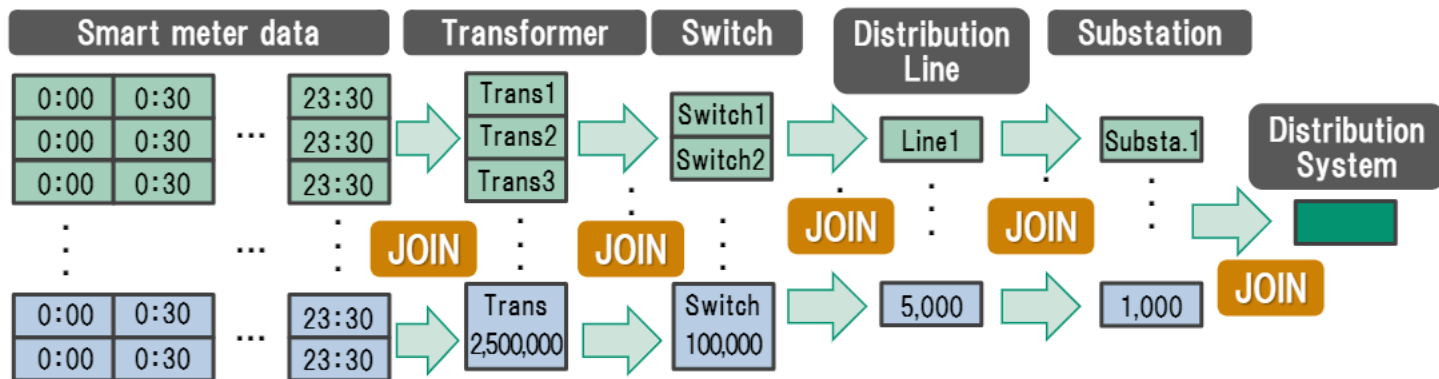
## Result (365days / 24h)



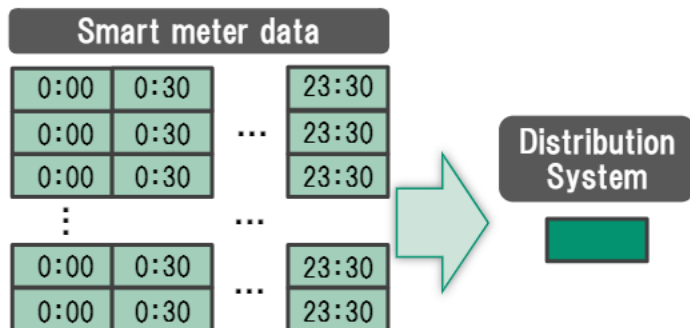
- Could finish processing the data from 10,000,000 meter in 30min using Spark 1.6!  
- Spark's good performance with "per equipment" processing

- ◆ Why the processing per equipment was more effective than the processing for entire distribution system when using spark?

- ◆ Per equipment



- ◆ For entire distribution system



- Less disk I/O than MapReduce
- Smaller data (including re-distributing data) than total memory of cluster



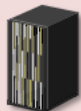
---

## 3. Additional evaluation with Spark 2.1

# 3-1 Evaluation environment

## ◆ System Configuration for additional evaluation

1 Master Node  
(Virtual Machine)



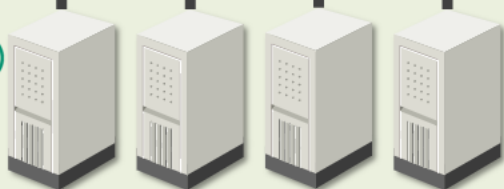
1Gbps LAN

10Gbps SW

10Gbps LAN



4 Slave Nodes  
(Physical Machines)



...



## ◆ Spec

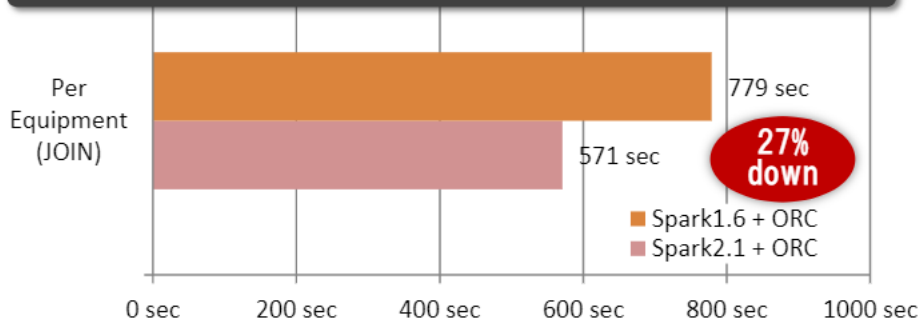
	Master Node
CPU Core	2
Memory	8 GB
Capacity of disk	80 GB
# of disk	1

	Per slave node	total
CPU Core	16	64
Memory	128 GB	512 GB
Capacity of disk	900 GB	-
# of disk	6	24
Total capacity of disks	5.4 TB (5,400 GB)	21.6 TB (21,600 GB)

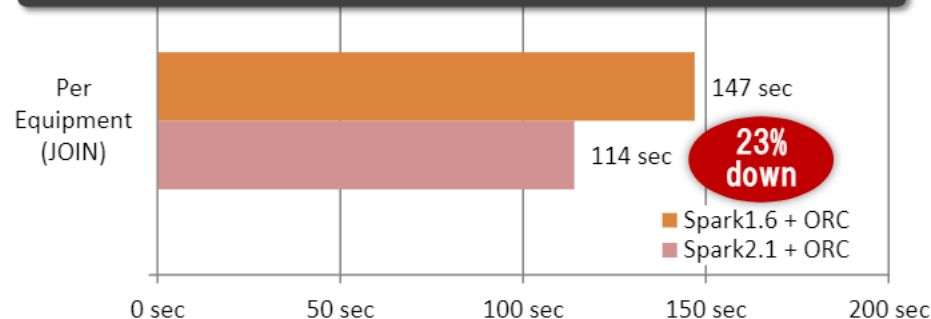
# 3-2 Comparison of Spark 2.1 with Spark 1.6 (ORCFile)

Time from start of aggregation batch for meter data to end of the batch

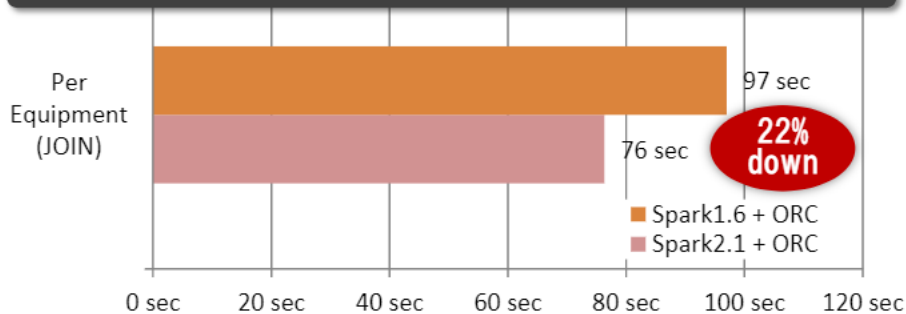
Result (365days / 0.5h)



Result (30days / 0.5h)



Result (1day / 0.5h)

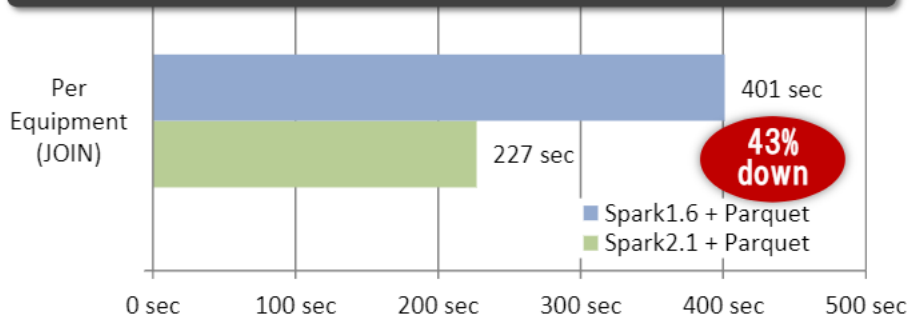


- Performance improvement 22-27% (including disk I/O)
- More effective with large data

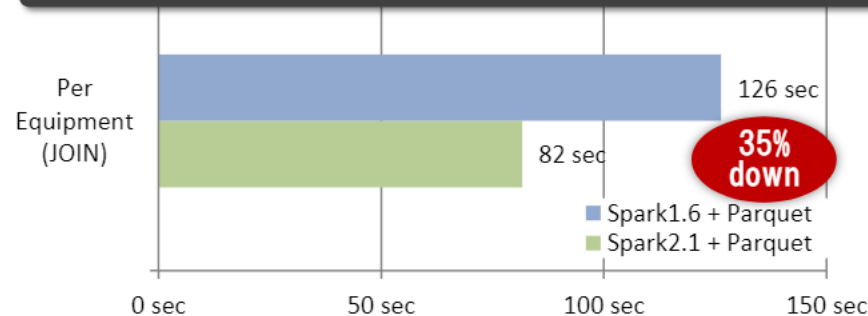
# 3-3 Comparison of Spark 2.1 with Spark 1.6 (Parquet)

Time from start of aggregation batch for meter data to end of the batch

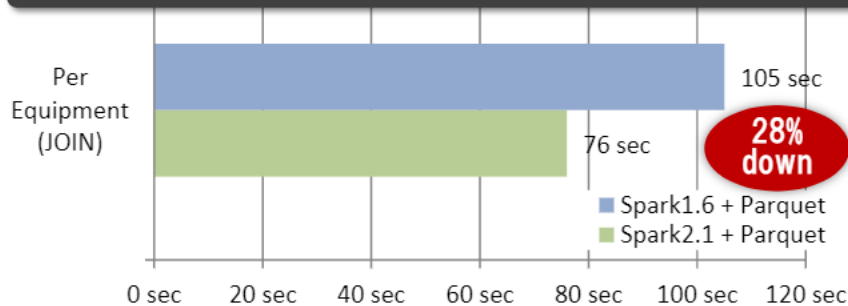
**Result (365days / 0.5h)**



**Result (30days / 0.5h)**



**Result (1day / 0.5h)**

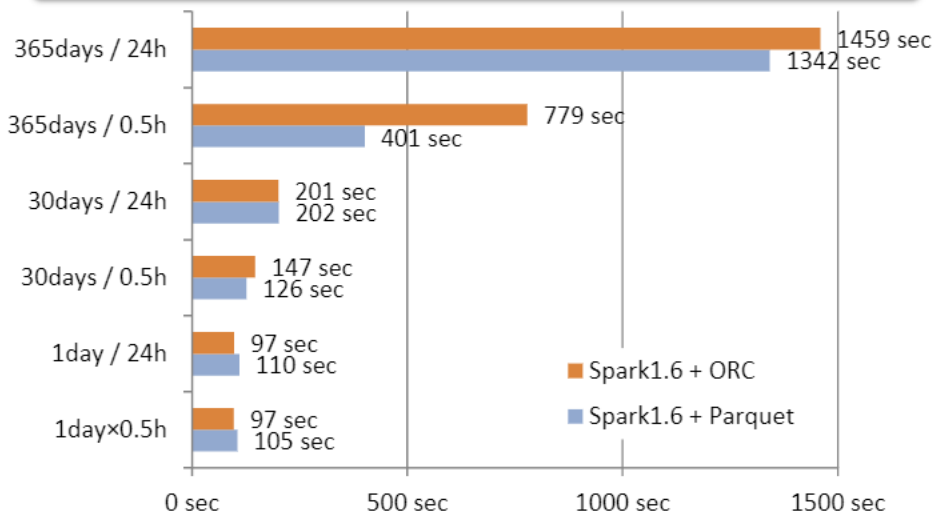


- Performance improvement 28-43% (including disk I/O)
- More effective with large data
- Better improvement than ORCFile

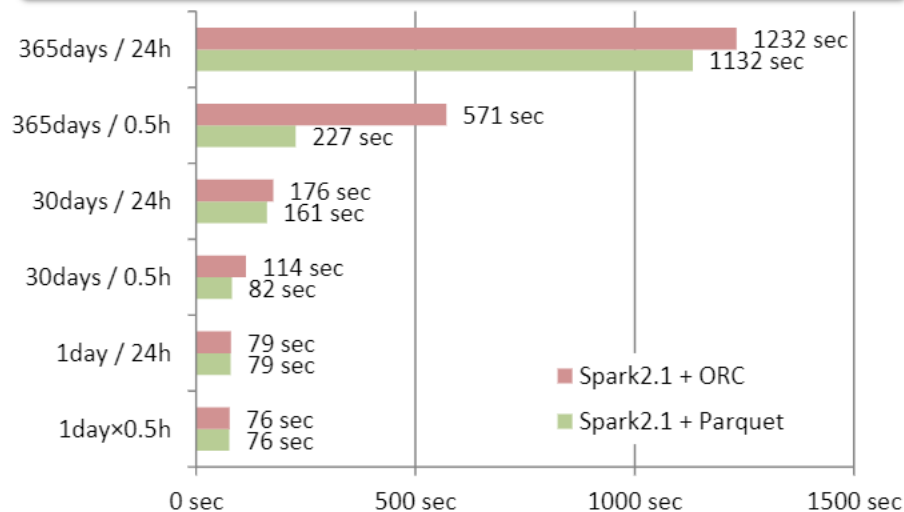
# 3-4 Comparison of ORCFile with Parquet (Spark 1.6/2.1)

Time from start of aggregation batch for meter data to end of the batch

**Result (Spark 1.6 / Per equipment (JOIN))**



**Result (Spark 2.1 / Per equipment (JOIN))**



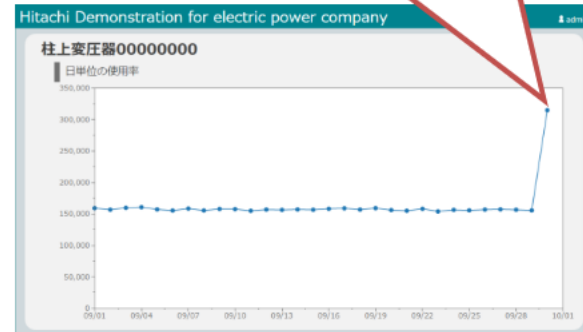
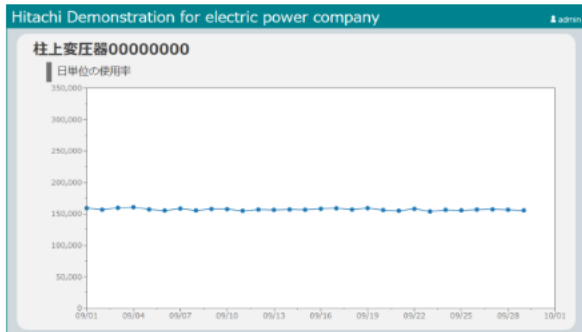
Term	Size (Parquet)	Size (ORCFile)
365days (1year)	1.363 TB	1.328 TB
30days (1month)	0.112 TB	0.109 TB
1day	3.7 GB	3.6 GB

• Basically, performance of Parquet is better than ORCFile  
 • Performance of Parquet with small data is worse than ORCFile in some cases.



## Demo

## Data visualization tool (Pentaho)



① Show 29 days data

② Execute aggregation batch (Spark SQL)

③ Show 30 days data

Aggregated Data (29 days)



Raw Data (1day)

Aggregated Data (30 days)

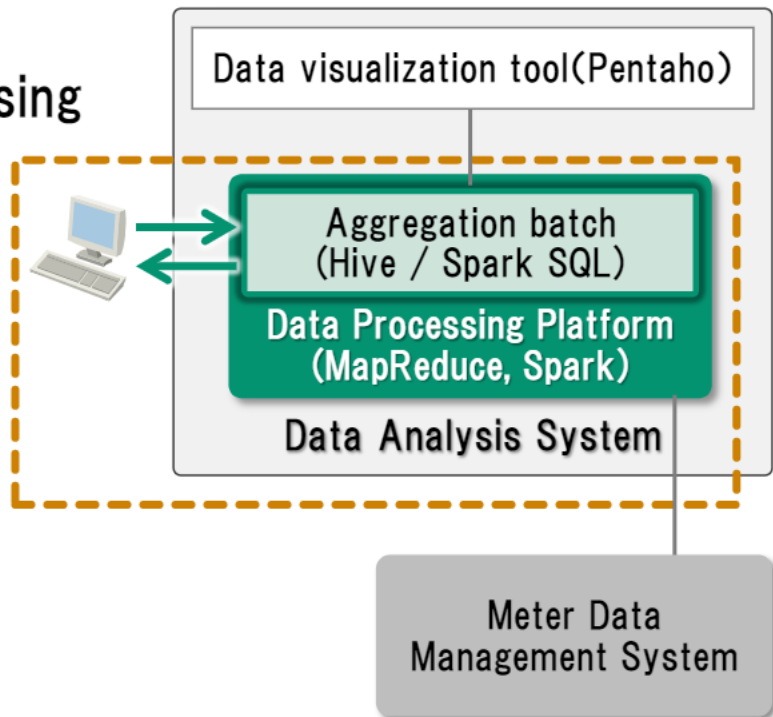
---

## 4. Summary



## 4 Summary

- ◆ Leveraging data from 10,000,000 smart meters for electric utilities in Japan
  - Built data analysis system
  - Concern about performance
- ◆ Evaluate the performance of batch processing
  - Spark could process the data from 10,000,000 meters in 30min (4 slave nodes)
- ◆ Evaluate the performance of Spark 2.1
  - Performance improvement 22-27% (compared to 1.6, ORCFile)
  - Performance improvement 28-43% (compared to 1.6, Parquet)



**END**

---

**Leveraging smart meter data for electric utilities:**  
Comparison of Spark SQL with Hive

5/16/2017

Hitachi, Ltd. OSS Solution Center

**Yusuke Furuyama**  
**Shogo Kinoshita**

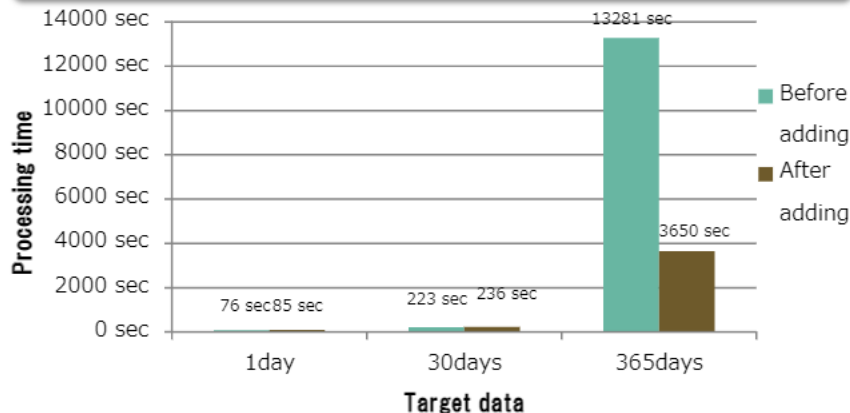
- ◆ Hadoop, Spark and Hive are trademarks or registered trademarks of Apache Software Foundation in the United States and other countries.
- ◆ Other brand names and product names used in this material are trademarks, registered trademarks or trade names of their respective holders.

**HITACHI**  
Inspire the Next 

- ◆ Attempted to aggregate raw (48 columns) meter data per equipment
  - Extremely slow (Spark 2.0) or Job failed (Spark 1.6)
  - Processing: Iteration of JOIN and GROUP BY+SUM
  - Huge data to be shuffled (spilled out from page cache)

Heavy load on a local disk (OS disk) by shuffle

Adding disks for shuffle (Spark 2.0.0)



- ◆ Add HDFS disks as disks for shuffle
  - Performance Improved (365days)
  - Performance degraded (1day/30days)

- Data for Spark (including temporary data) should be smaller than memory.  
- Had better to process as a trial to estimate