

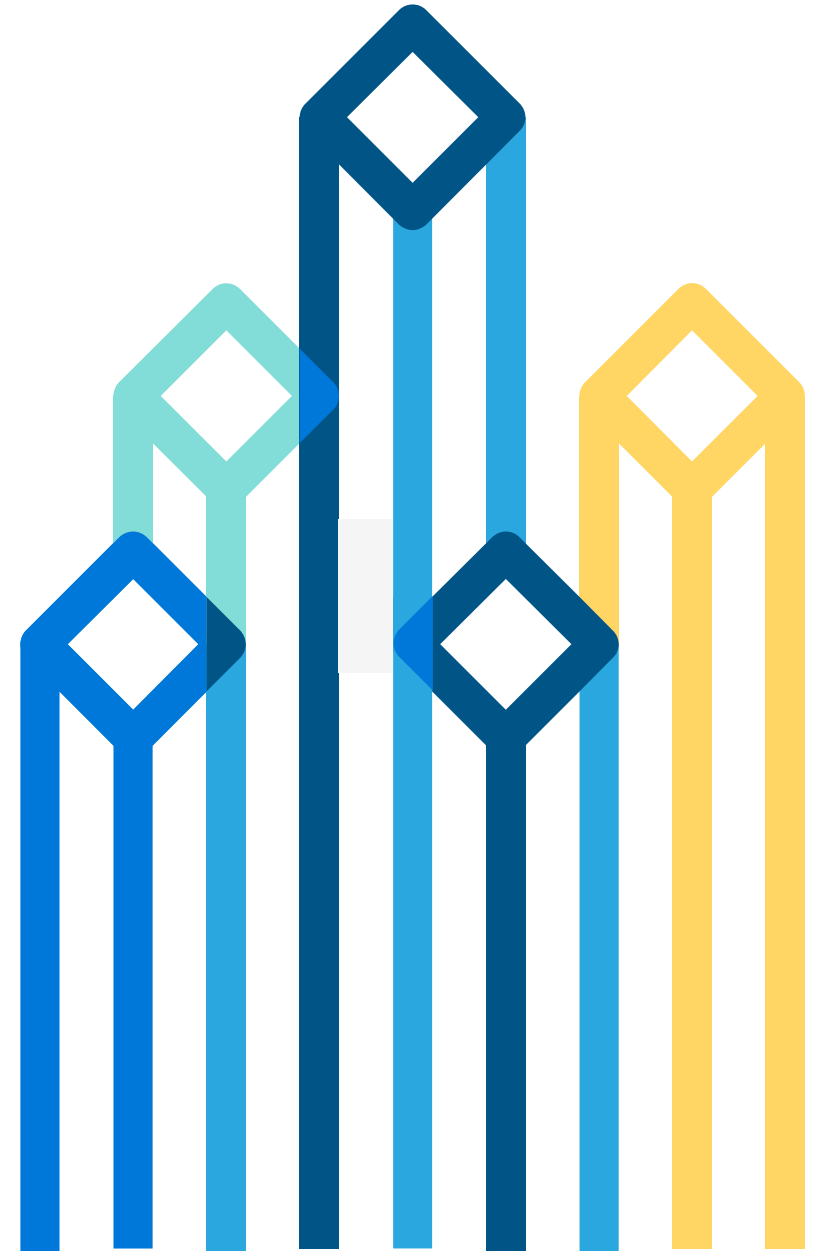


Modern Historian Data Analytics

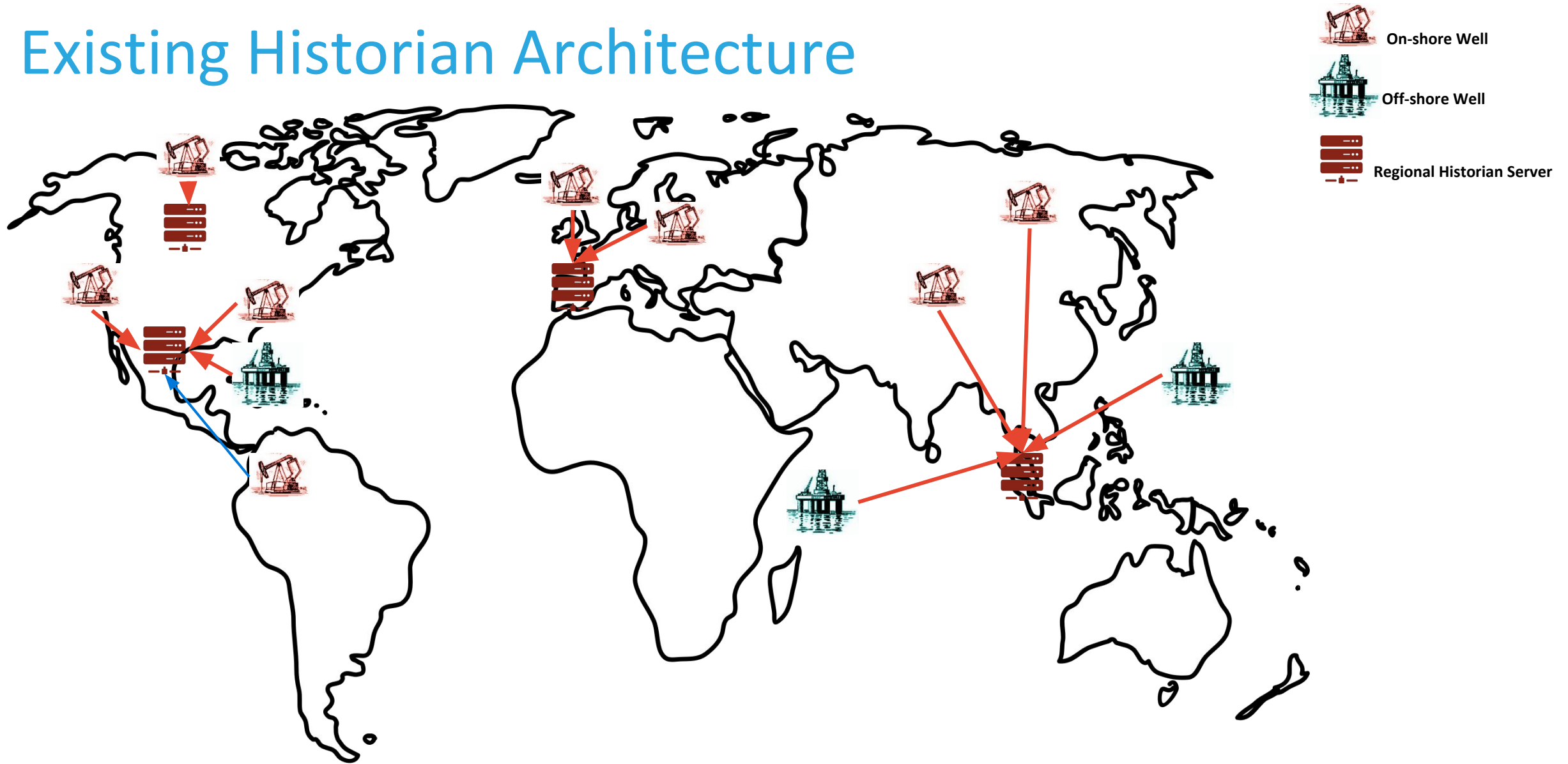
Samir Gupta, Sales Engineer

Sheridan McDonald, Account Executive

January, 2017



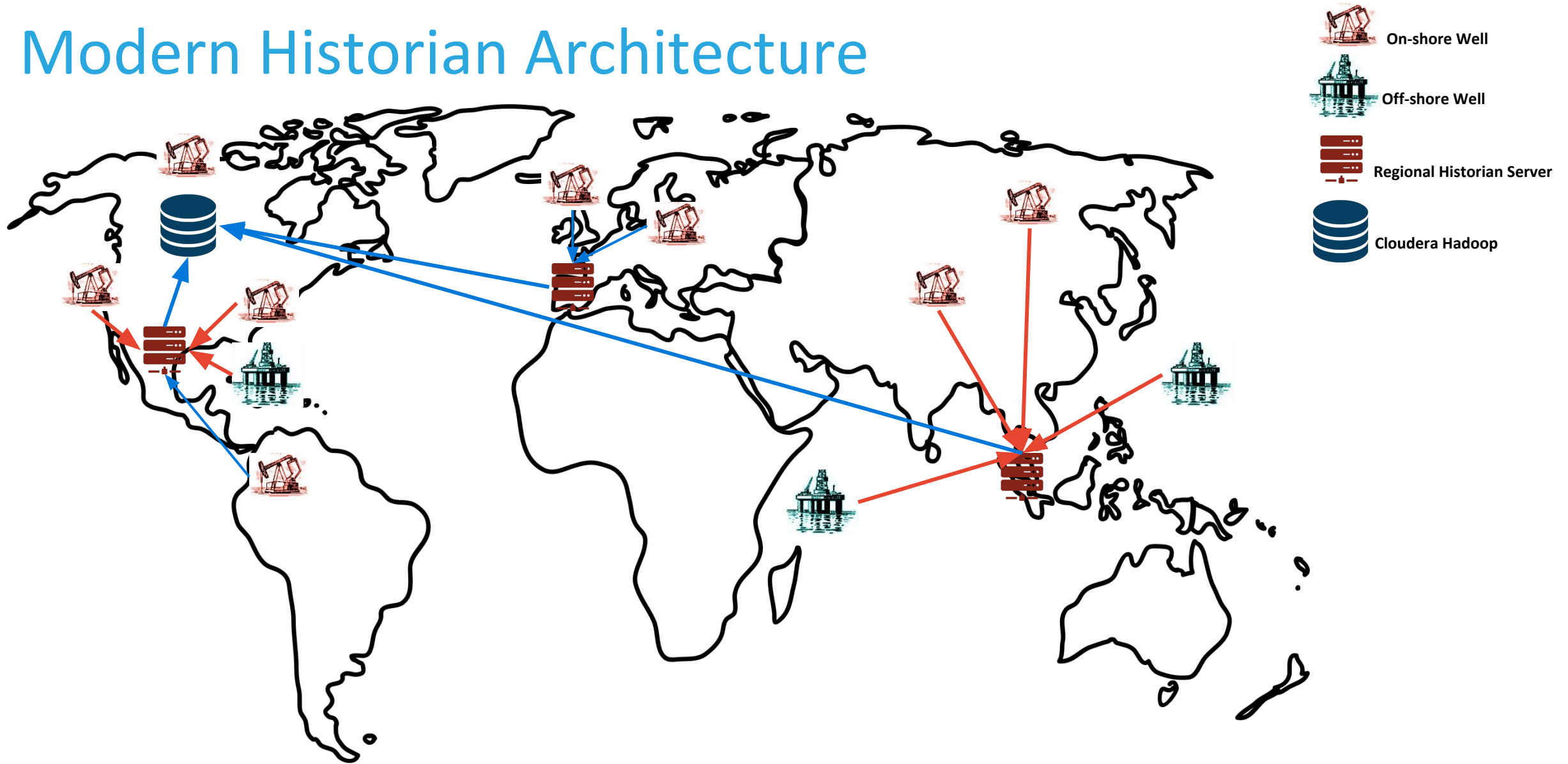
Existing Historian Architecture



Challenges with Existing Architecture

- Expensive
- Proprietary technology
 - Limited analytics capabilities
 - No inherent relationships between data and tags
 - Difficult to get data out in a useful format - ie. by individual tag ID
 - Ancient historian tooling to analyze the data
- No real-time access to the data
 - Raw data needs processing to be easily visualized
- Not well integrated with the data management ecosystem
 - Need to analyze in Excel, which has a 1M record limit (<.001% of data)

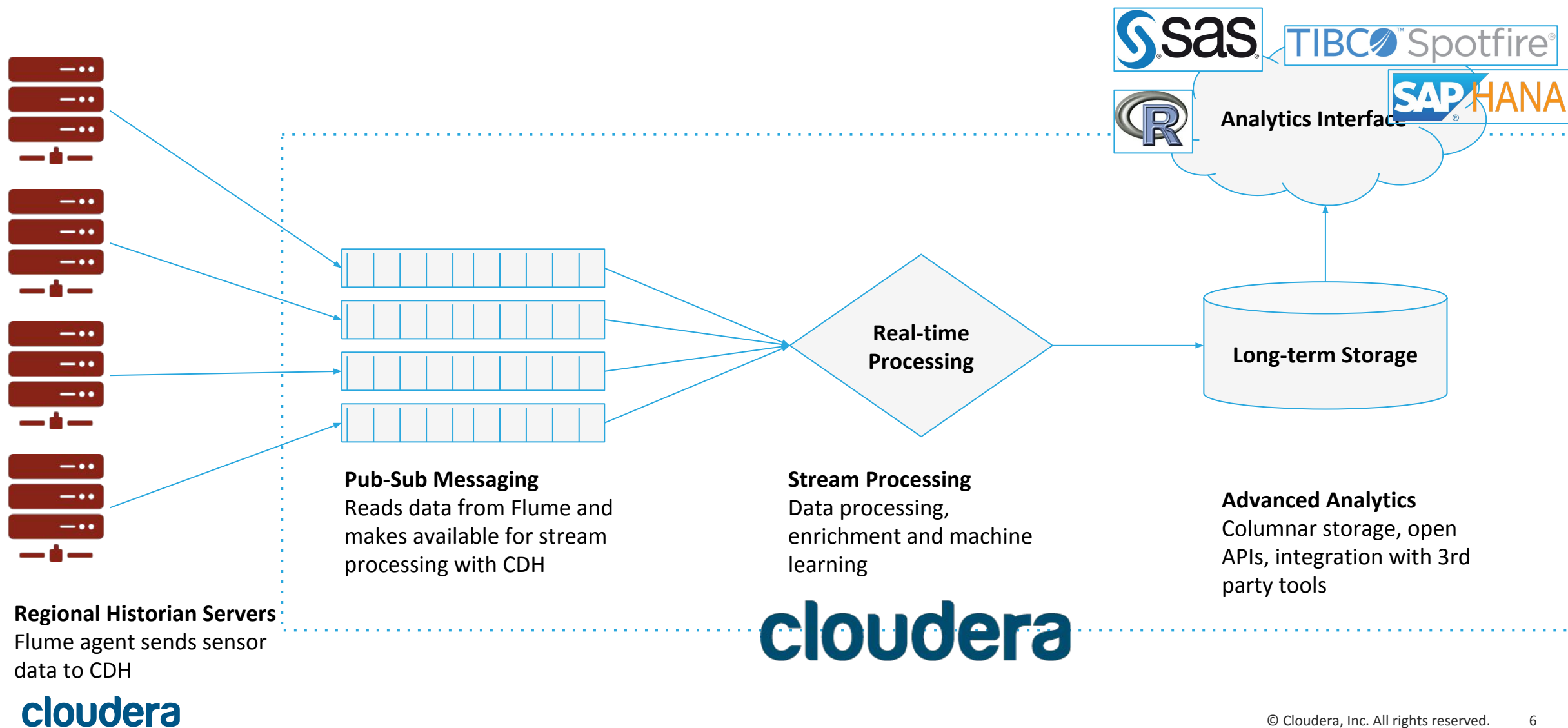
Modern Historian Architecture



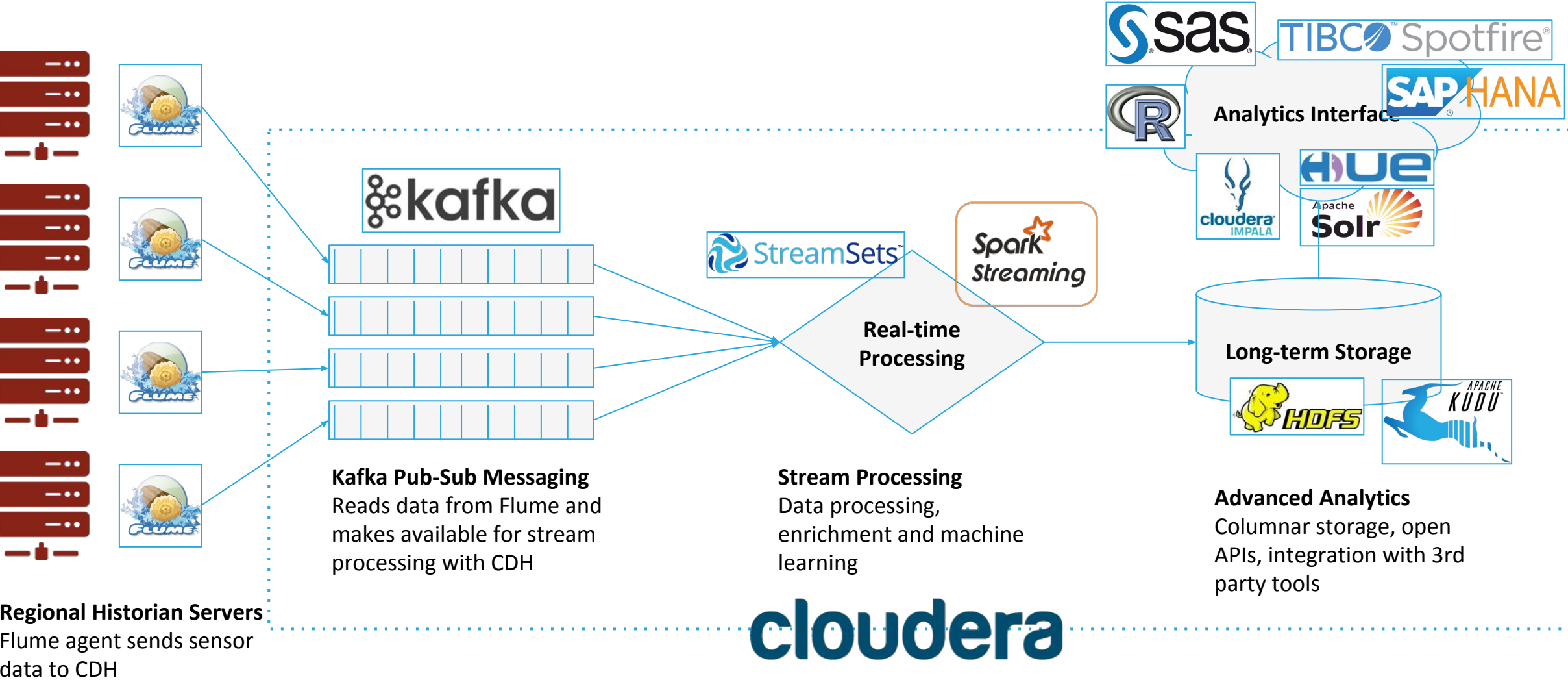
Benefits with Modern Architecture

- Cost Effective
- Open technology
 - Advanced analytics capabilities
 - Easily create relationships between raw data and tag information
 - Integrate with any analytic tooling
- Real-time access to the data
 - Real-time processing of data
- Analyze all sensor data, instead of a small subset
 - 100B+ rows vs. 1M

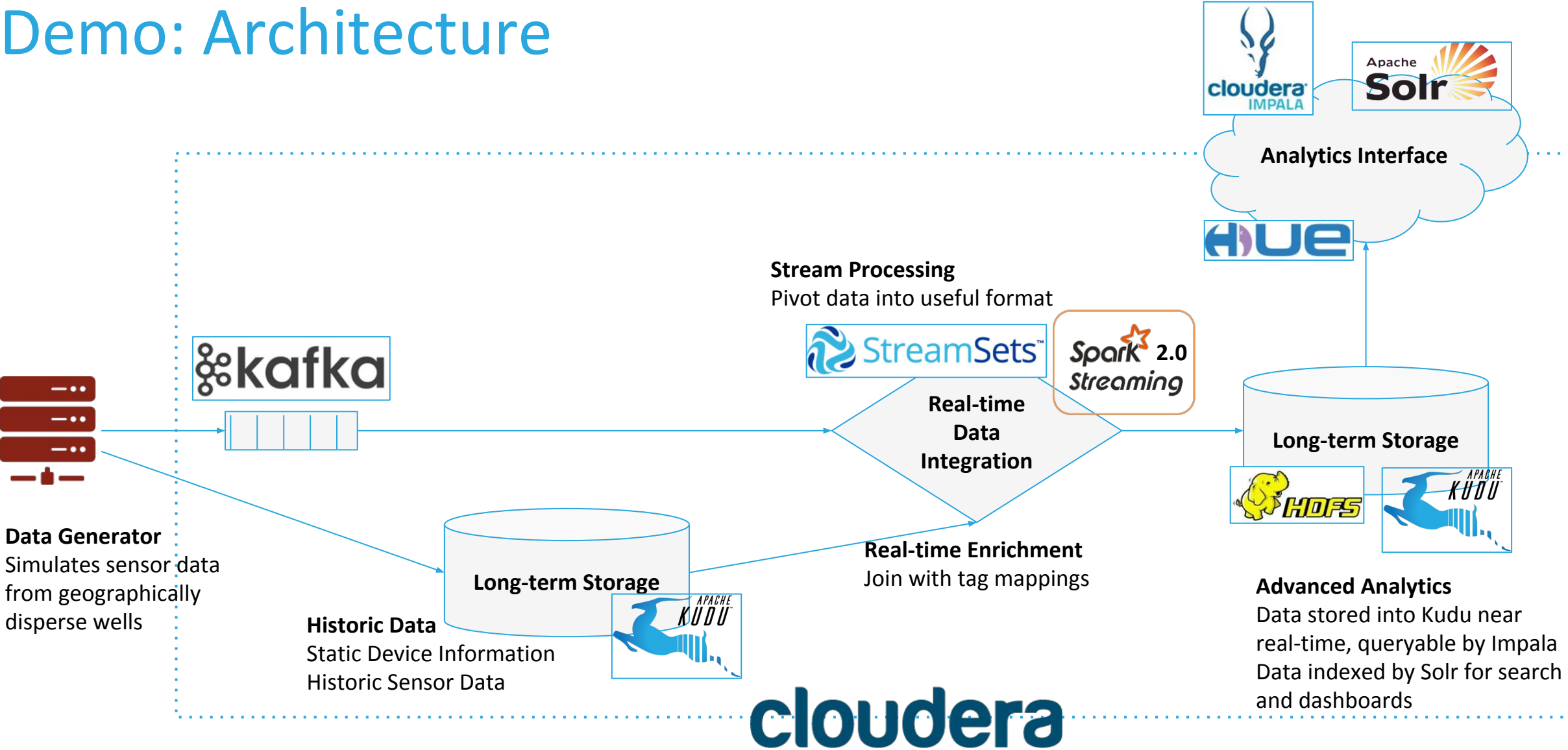
Ingesting PHD Data into Cloudera



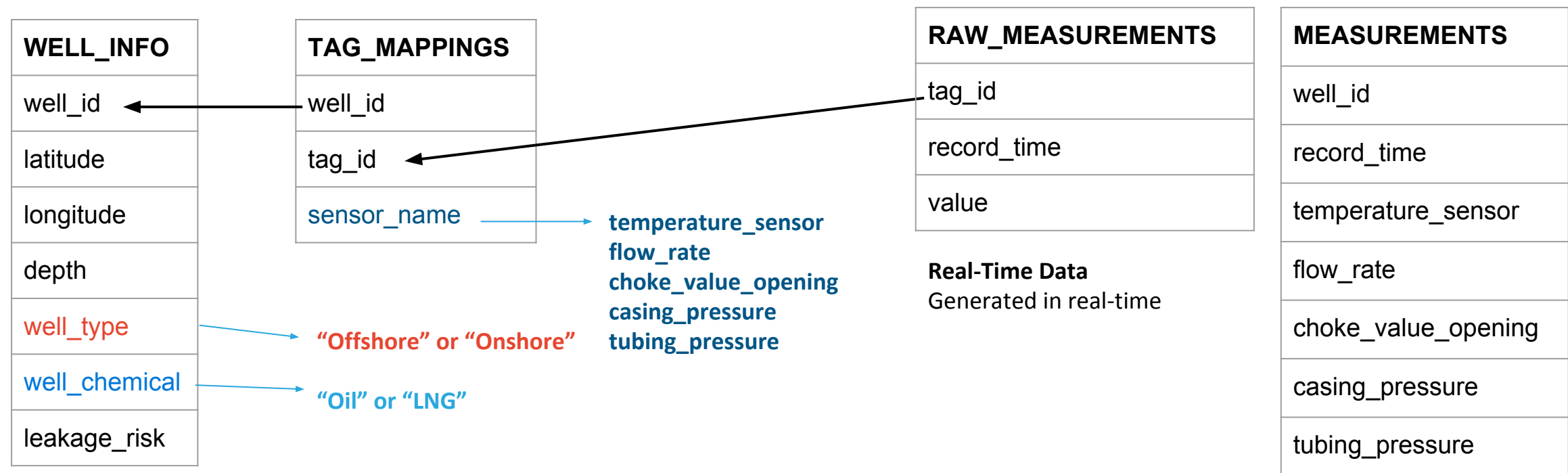
Ingesting Historian Data into Cloudera



Demo: Architecture



Demo: Data Model



Static Data

Generated during initial setup, configurable
Stored in Impala using KUDU for fast lookup

Real-Time Data

Enriched in real-time with Spark
Evaluator in Streamsets
Stored in near-real-time into KUDU

Demo: Pivot in Spark

tag_id	record_time	value
1	9:00:00	50
2	9:00:00	12
3	9:00:00	80
1	9:00:30	48
2	9:00:30	14

tag_id	record_time	value	well	sensor_name
1	9:00:00	50	1	temp_sensor
2	9:00:00	12	1	flow_rate
3	9:00:00	80	1	casing_pressu re
1	9:00:30	48	1	temp_sensor
2	9:00:30	14	1	flow_rate
3	NULL	NULL	NULL	pressure

well_id	record_time	temp	flow	pressure
1	9:00:00	50	12	80
1	9:00:30	48	15	NULL

tag_id	well_id	sensor_name
1	1	temp
2	1	flow
3	1	pressure

RAW_MEASUREMENTS
OUTER JOIN
TAG_MAPPINGS
ON tag_id

PIVOT ON sensor_name

UPSERT INTO KUDU

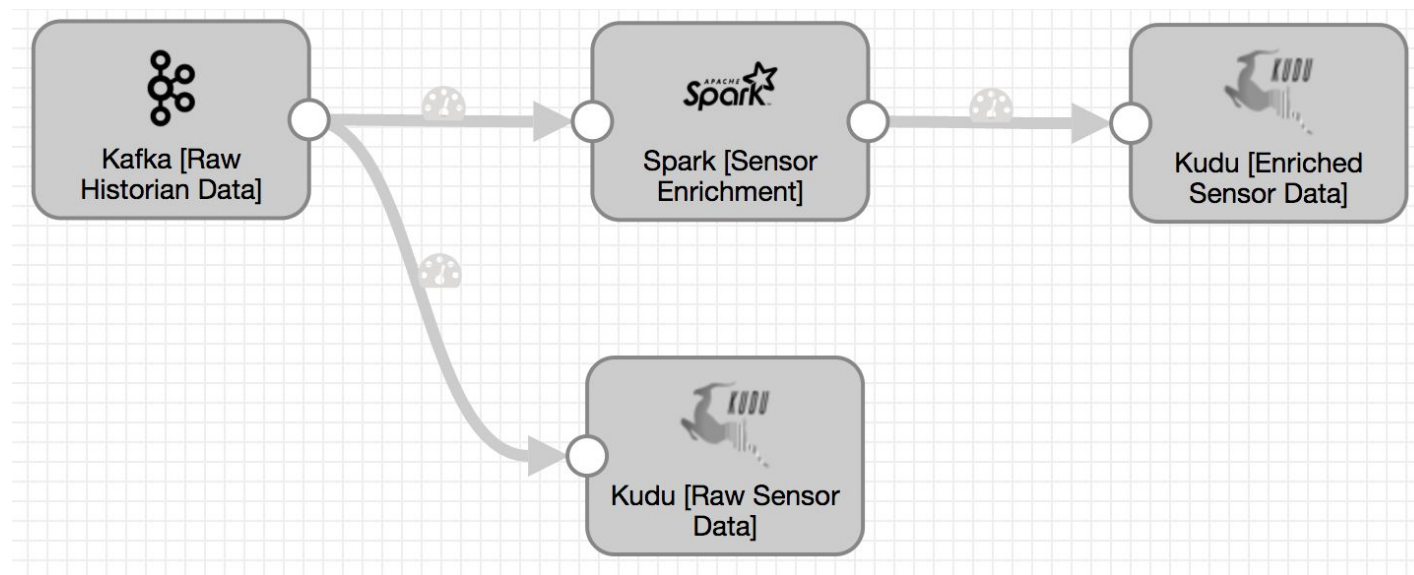


cloudera



well_id	record_time	pressure
1	9:00:00	50
2	9:00:00	80
1	9:00:30	52
2	9:00:30	NULL at W1, 81 at W2
1	9:01:00	49
2	9:01:00	76
1	9:01:30	55
2	9:01:30	86
1	9:02:00	50
2	9:02:00	NULL

Demo: Streamsets Data Collector



Kudu Table

Sends enriched sensor data to Impala table to be queried and used in BI

Kafka Source

Reads raw sensor data from Kafka cluster

Custom Spark Evaluator

Real-time lookup to Kudu to enrich sensor data
Pivot and aggregate before sending to Kudu

Kudu Table

Raw sensor data also stored in Kudu for post-processing analysis/debug

Streamsets running in “Cluster Mode” to read directly from a Kafka cluster. The entire pipeline will run as a Spark streaming application inside CDH.



cloudera

Thank you

Samir Gupta, Systems Engineer
sgupta@cloudera.com

Sheridan McDonald
smcdonald@cloudera.com