



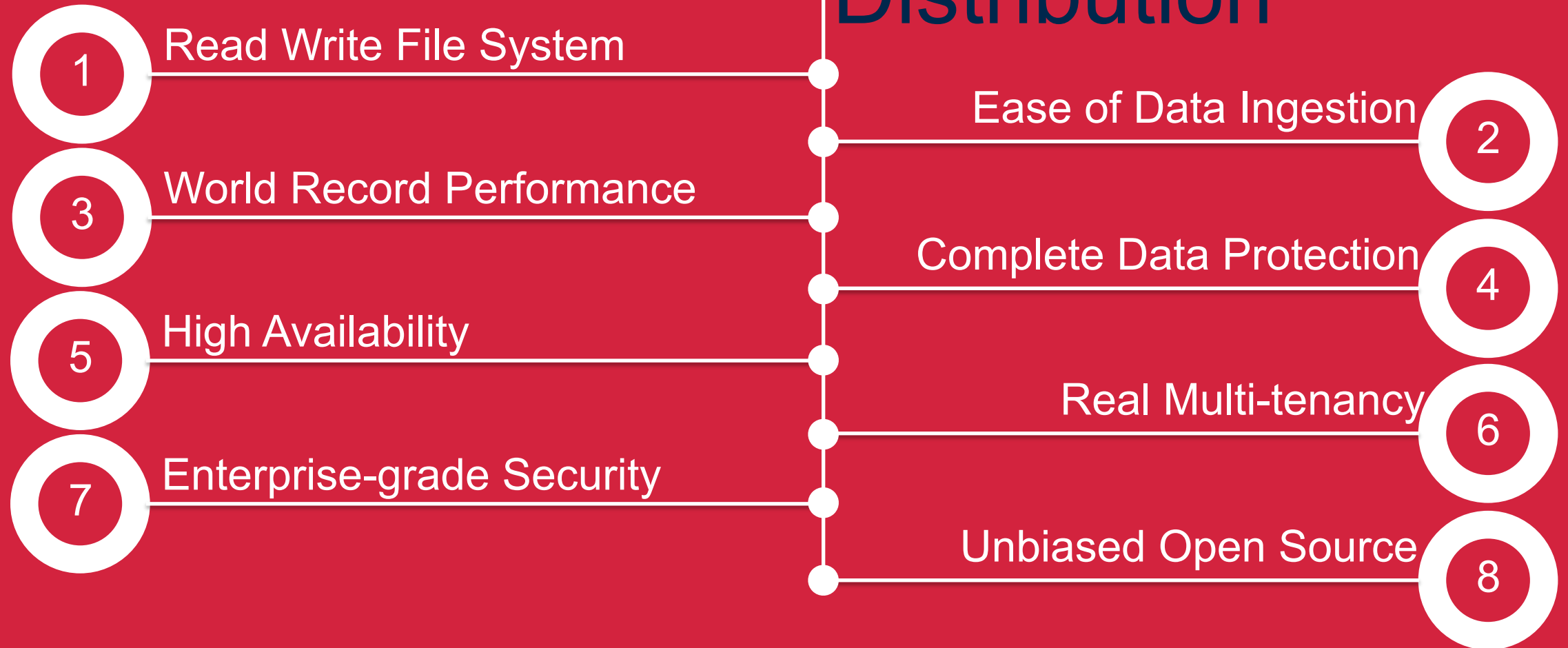
Time-Series Databases and Machine Learning

Jimmy Bates

November 2017



Top-Ranked **Hadoop** Distribution

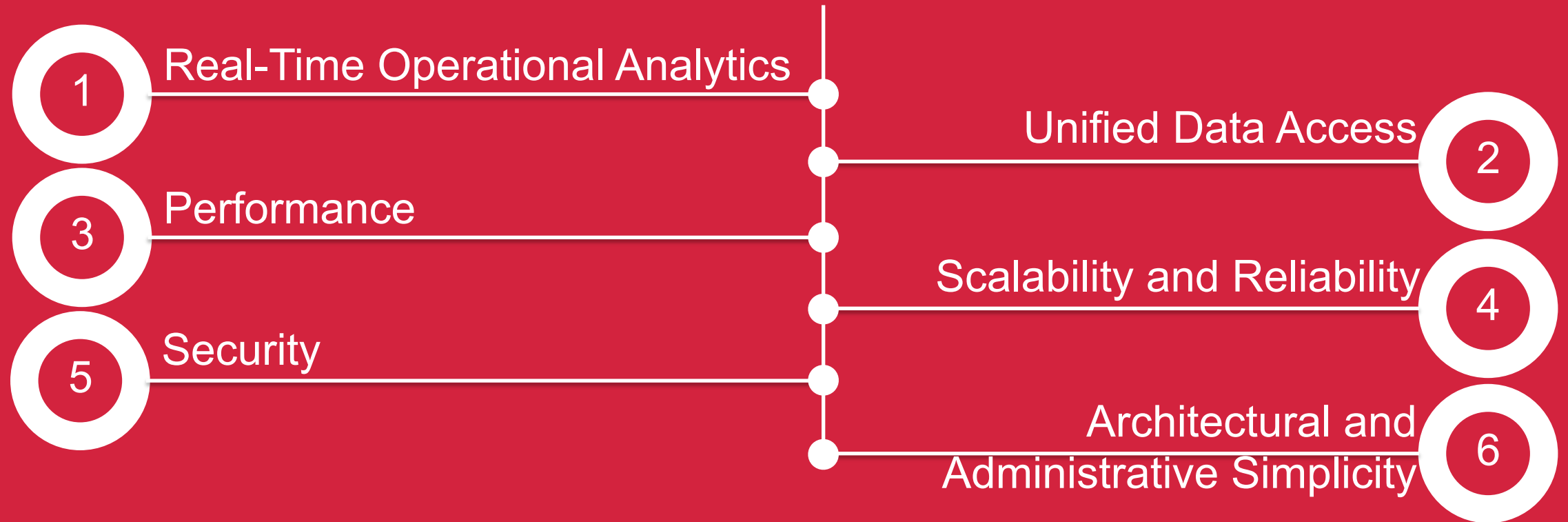




Top-Ranked **Hadoop**
Distribution

Top-Ranked **NoSQL**

Top-Ranked **NoSQL**





Top-Ranked **Hadoop**
Distribution

Top-Ranked **NoSQL**

Top-Ranked **SQL-on-Hadoop**
Solution

Top-Ranked **SQL-on-Hadoop** Solution

1

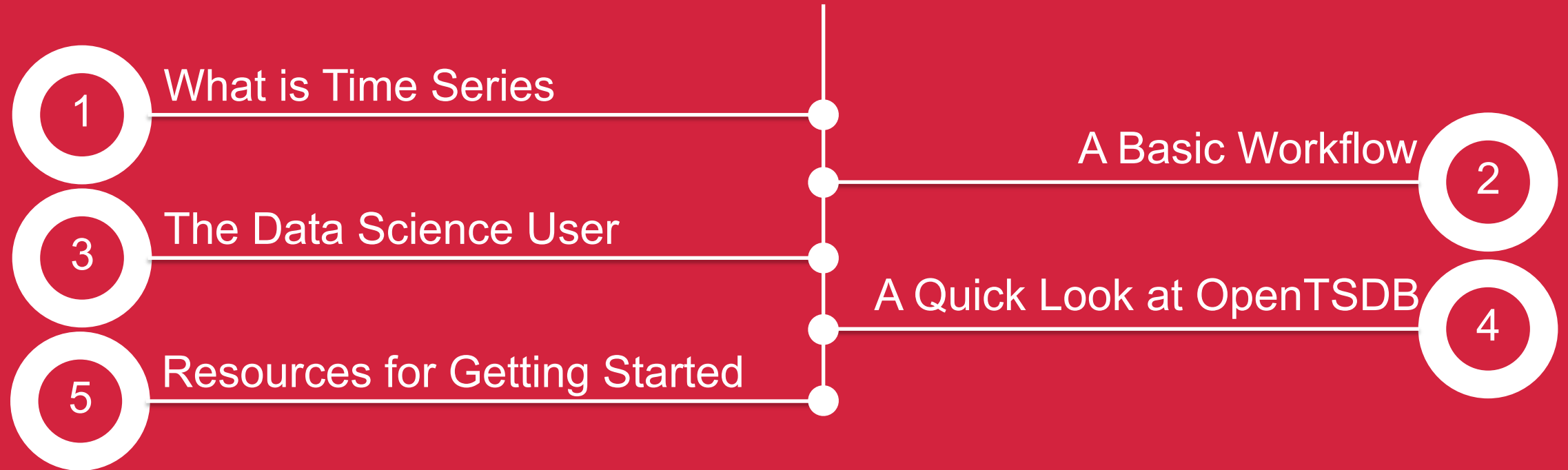
The Most SQL Options on
Hadoop



Including our work with
Apache Drill

2

Agenda



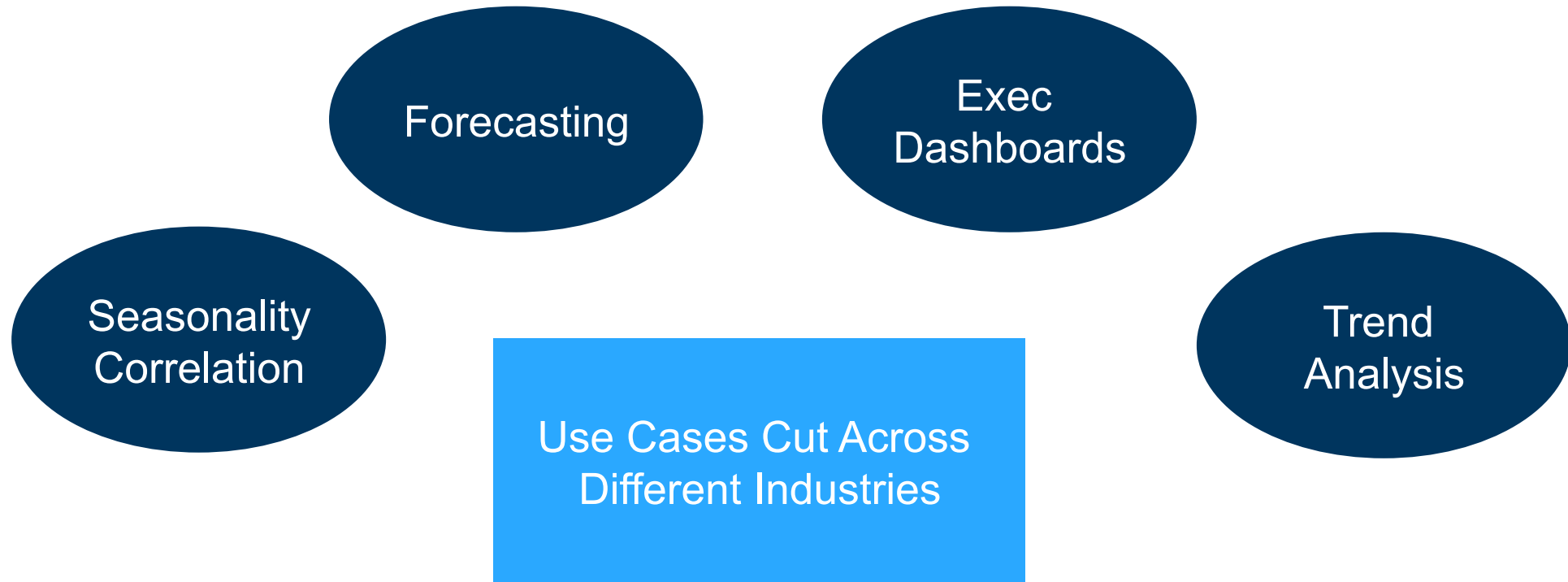
Time Series Data

What is Time Series Data?

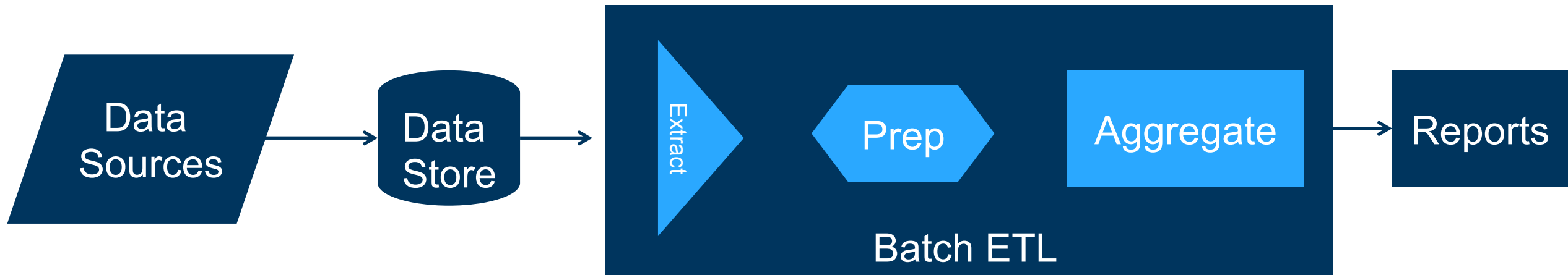
- Any set of data-points that have time associated with it.
- Examples of time series data include
 - weather,
 - web clickstream,
 - product sale,
 - stock trade,
 - machine logs,
 - fleets,
 - sensors,
 - devices,
 - network traffic,
 - system login



Value of Time Series Analytics



Typical Time Series Analytics Flow

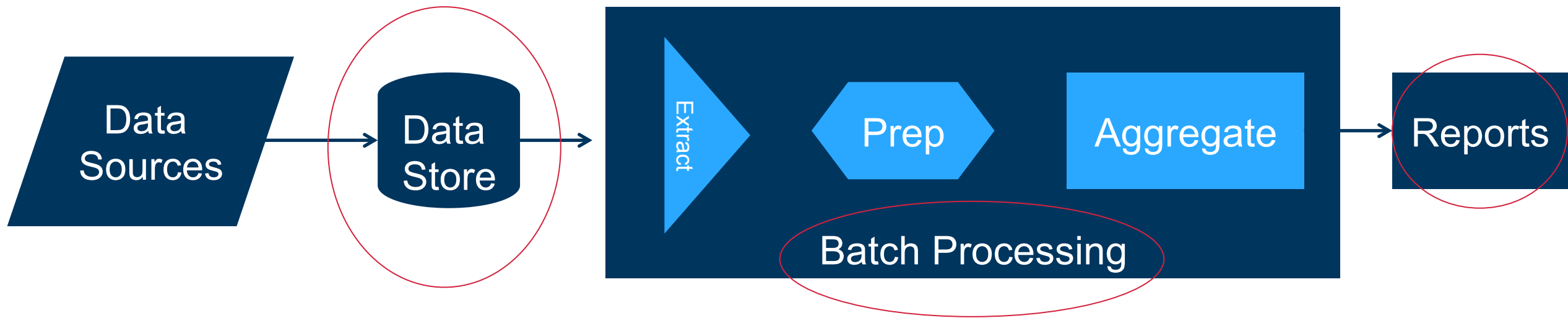


Environment: IT teams/Application Development groups/Analytics groups

Persona: Enterprise Architect, NoSQL Developer, Database developer, Analyst

Task: Conduct analyses or provide reports to management

Typical Time Series Analytics Flow - Challenges



- Larger Volumes of data from newer sources
- Expensive to store

- Takes a lot of time to process information and do simple aggregations

- Not Real-time

Challenges with Status Quo

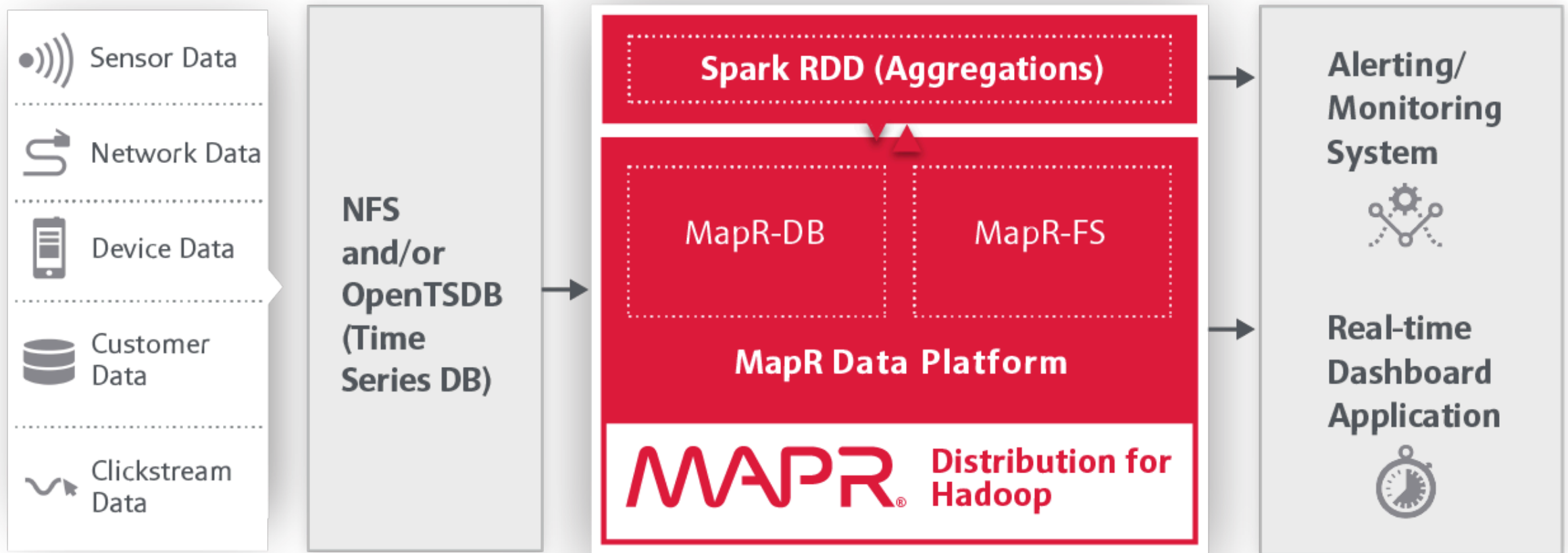
- Existing systems aren't well suited to store high volumes of time-series data
- Building aggregations and statistical computations from time-series data is currently done in batch

Need:

- Cost effective and reliable way to store and analyze large amounts of time series data from various sources.
- Ability to deploy real-time dash-boarding and monitoring capabilities on aggregated data



Time Series Solution Example



MapR Data Exploration Advantages for IoT



Self-Service Data Exploration

Single SQL Interface for Structured
and Semi-Structured Data

Data Agility with Less IT Required

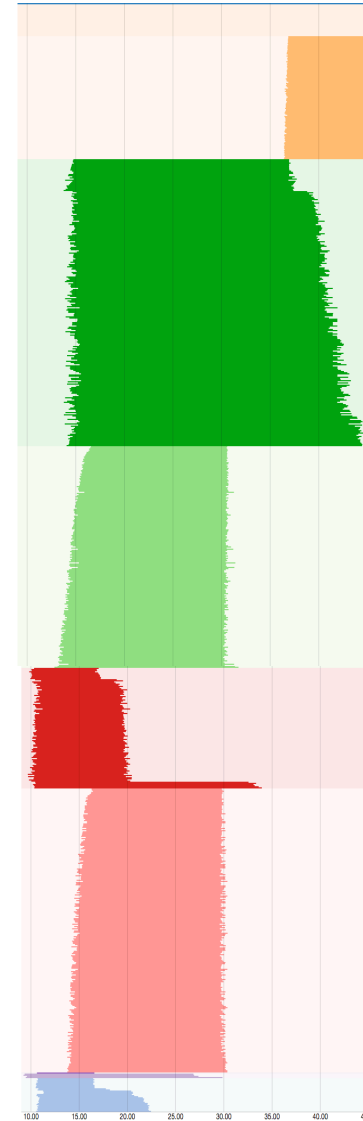
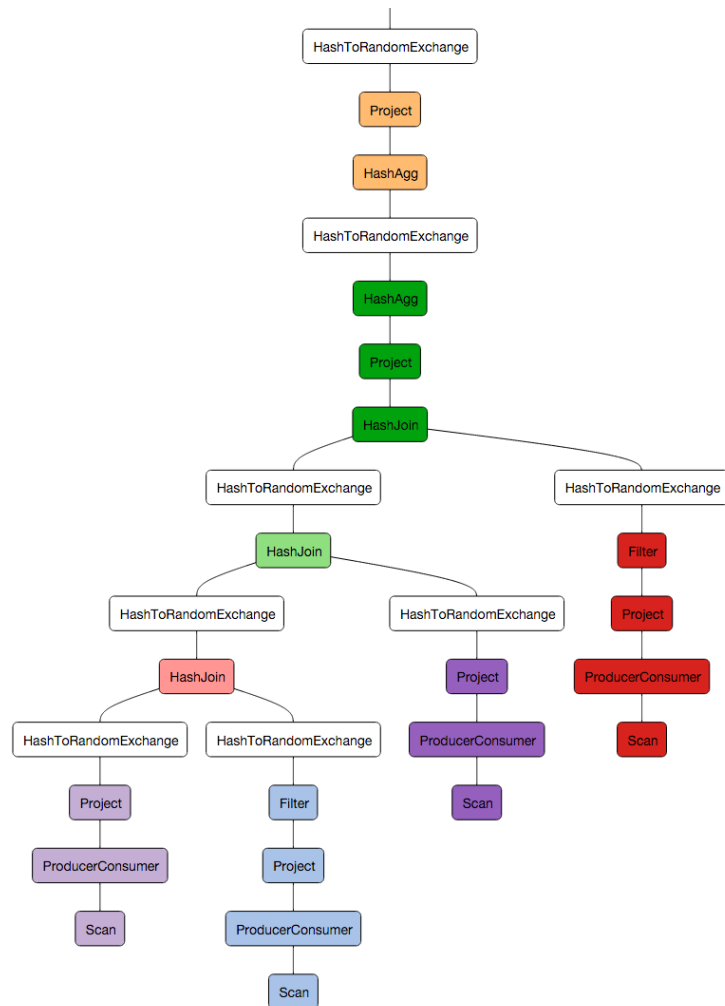


Squaring the Circle

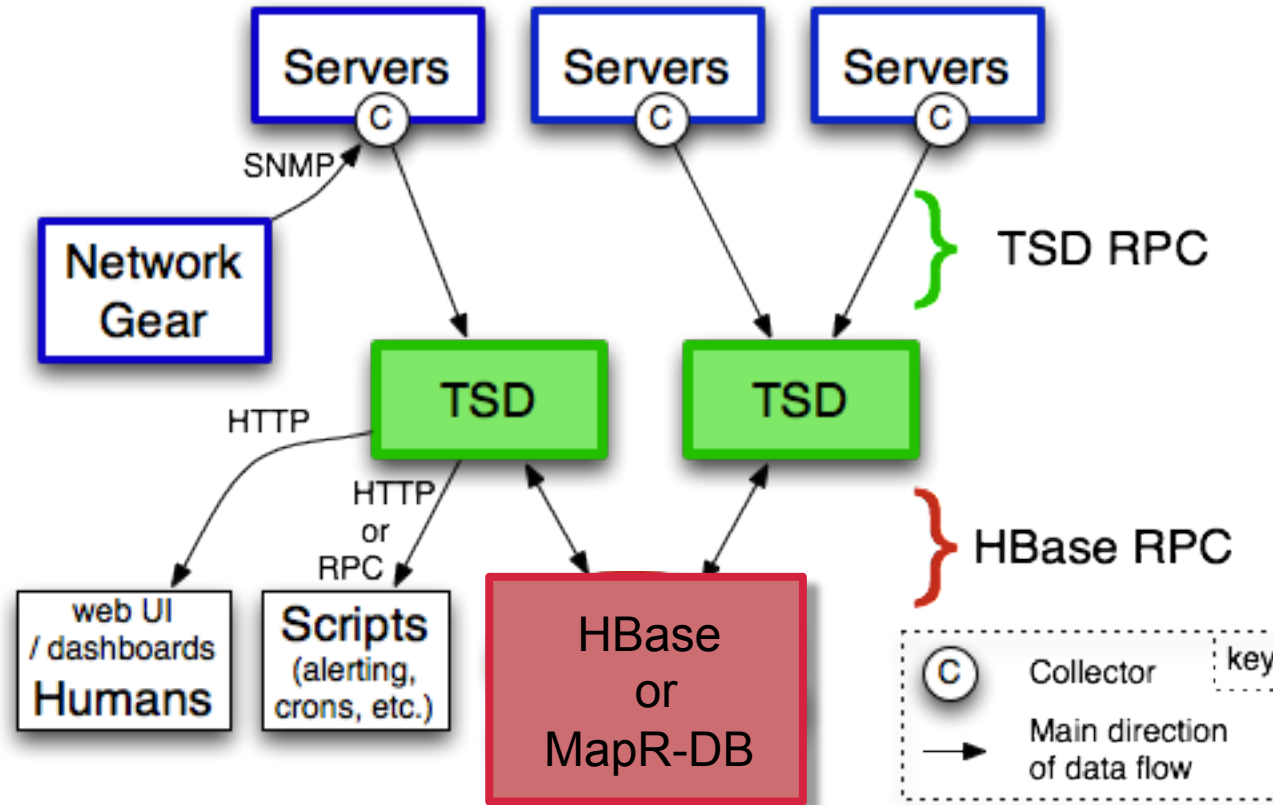
- Enter Apache Drill
- Drill is SQL *compliant*
 - Uses standard syntax and semantics
- Drill extends SQL
 - First class treatment of objects, lists
 - Full support for destructuring, flattening
 - Full power of relational model can be applied to complex data



Drill Provides Scalable and Extended SQL

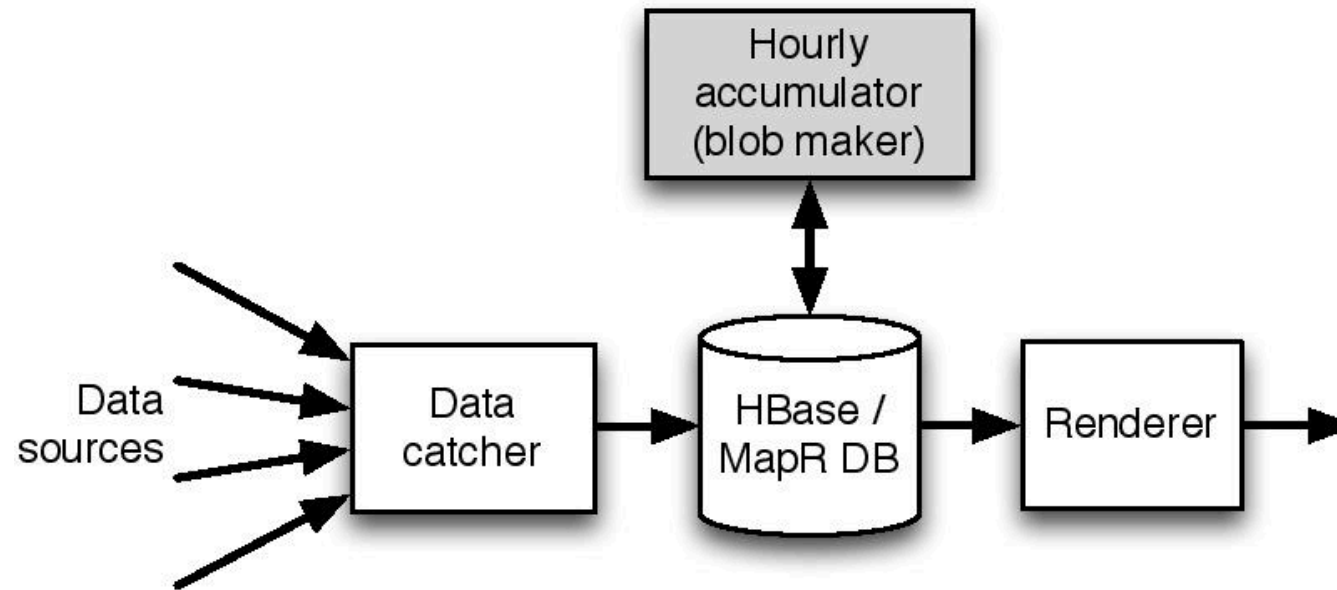


Introduction to Open TSDB



Speeding up OpenTSDB

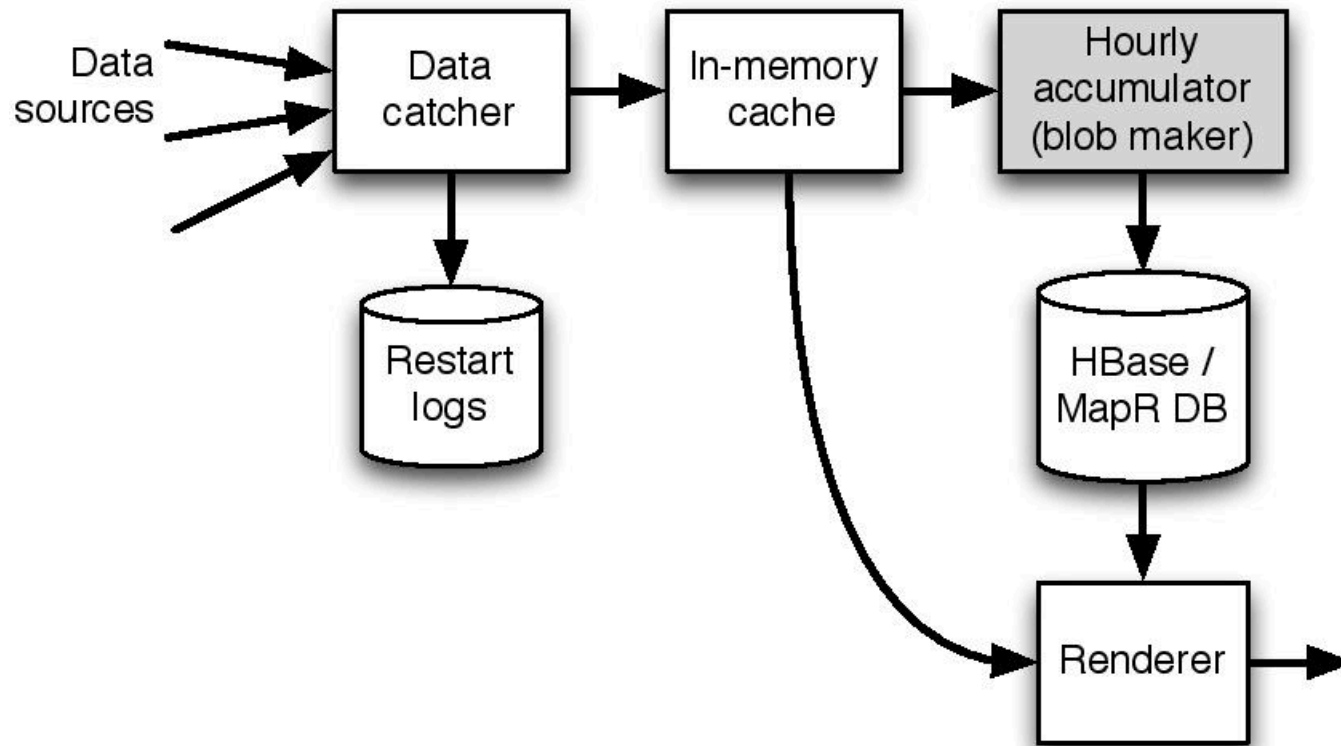
Why can't it be faster ?



20,000 data points per second per node in the cluster



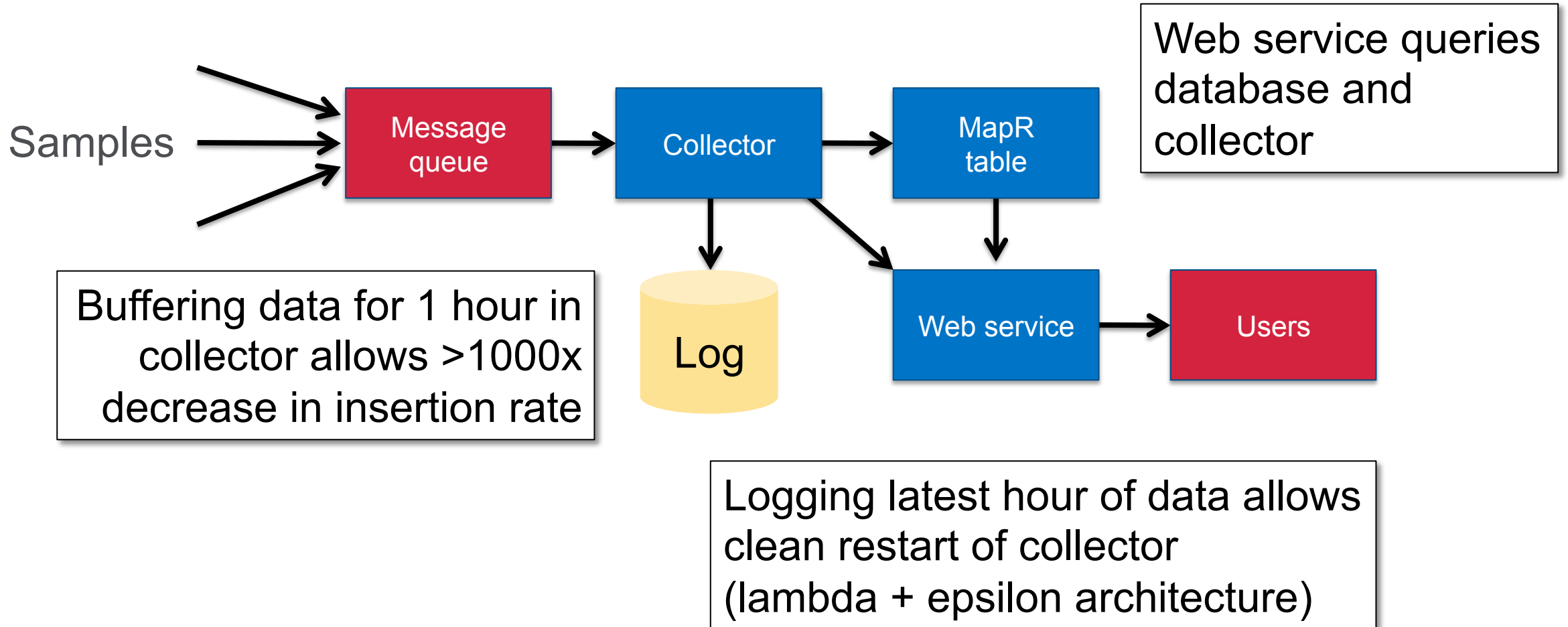
Speeding up OpenTSDB: open source MapR extensions



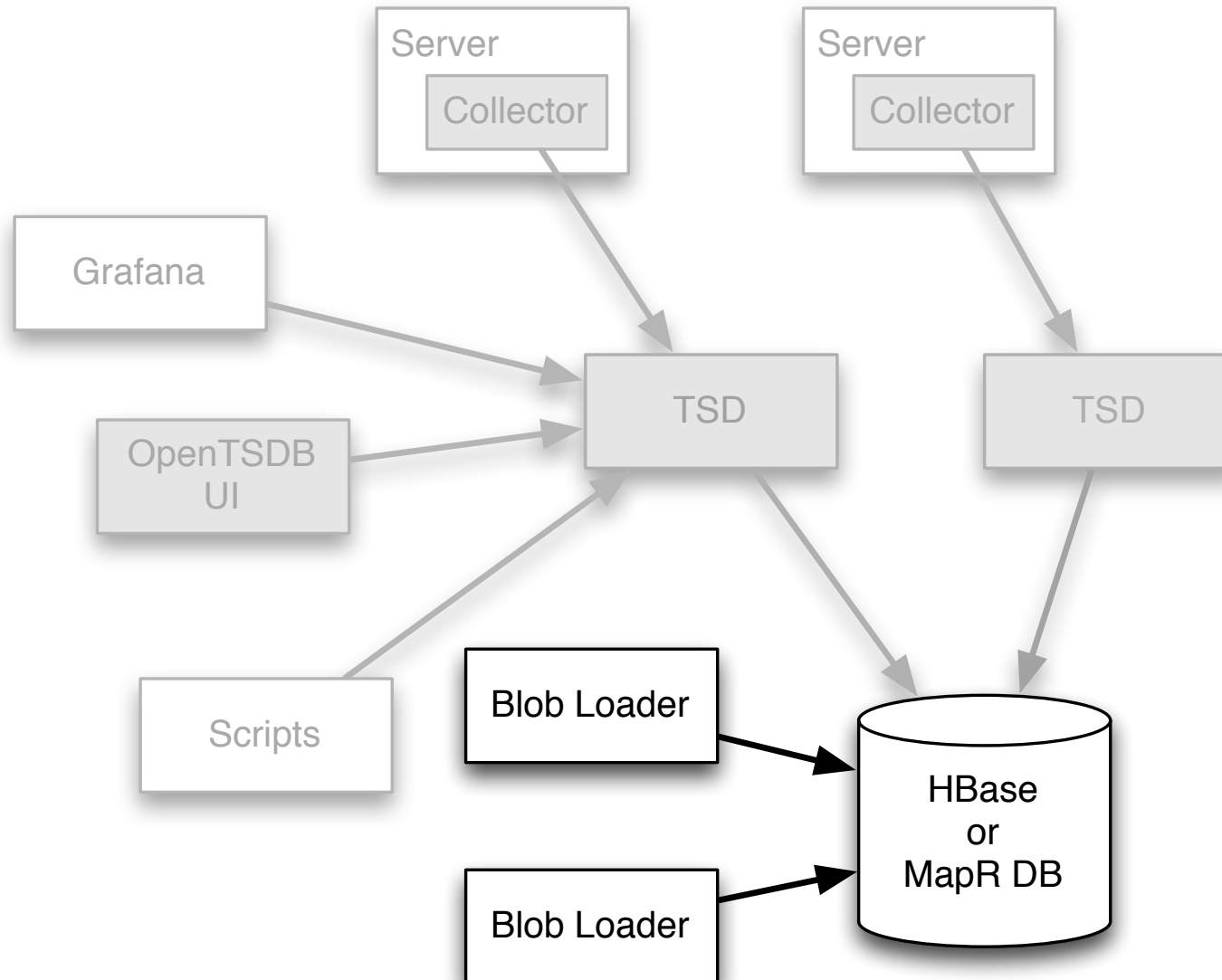
Available on Github: <https://github.com/mapr-demos/opentsdb>



Speeding up OpenTSDB: open source MapR extensions



Direct Blob Loading for Testing



A first example: Time-series data



Column names as data

- When column names are not pre-defined, they can convey information
- Examples
 - Time offsets within a window for time series
 - Top-level domains for web crawlers
 - Vendor id's for customer purchase profiles
- Predefined schema is impossible for this idiom



Relational Model for Time-series

The diagram illustrates a relational model for time-series data. It features a table with three columns: 'Time series ID', 'Sample time', and 'value'. The first two columns are grouped by a bracket labeled 'Row key'. The 'value' column is labeled with the word 'value' in italics. Below the table, there are two brackets: one under the first two columns labeled 'Row key' and another under the 'value' column labeled 'Data values'. Arrows point from the labels 'Time series ID' and 'Sample time' to their respective columns in the table header.

Time series ID	Sample time	<i>value</i>
101	15:51:03	1.16
101	15:52:07	0.04
101	15:52:11	0.08

Table Design: Point-by-Point

Time series ID

Time-window start time

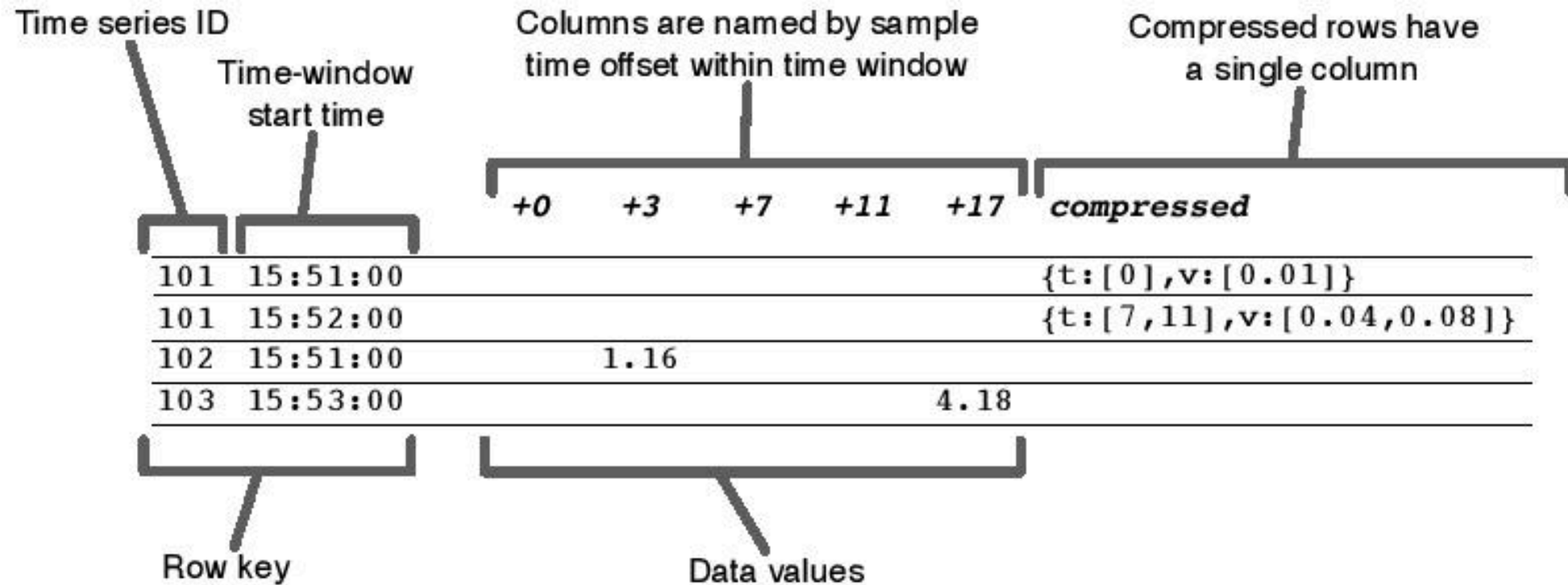
Columns are named by sample time offset within time window

		+0	+3	+7	+11	+17
101	15:51:00	0.01				
101	15:52:00			0.04	0.08	
102	15:51:00		1.16			
103	15:53:00					4.18

Row key

Data values

Table Design: Hybrid Point-by-Point + Sub-table



After close of window, data in row is restated as column-oriented tabular value in different column family.

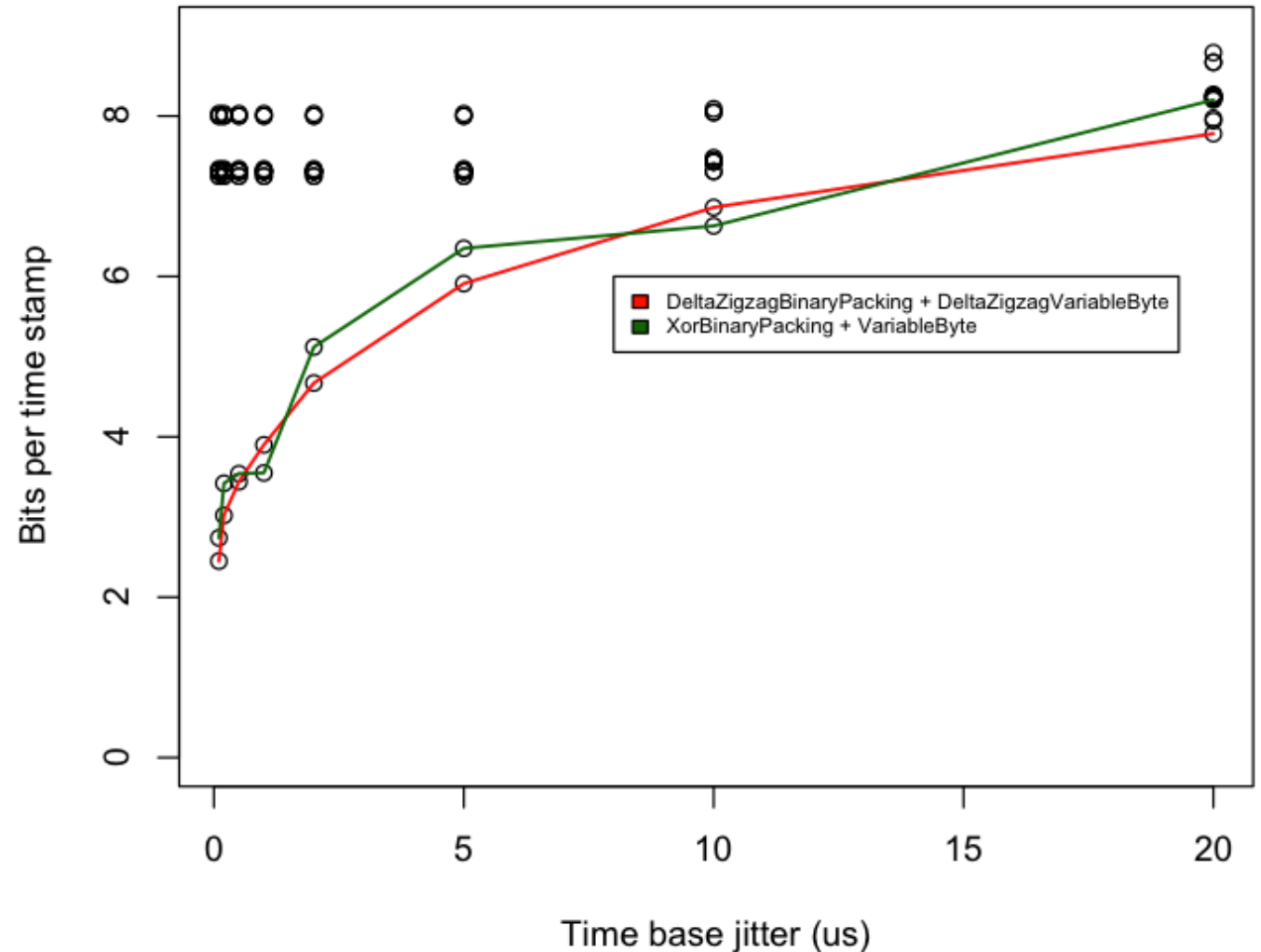
Compression Results

Samples are
64b time, 16 bit sample

Sample time at 10kHz

Sample time jitter makes it
important to keep original
time-stamp

How much overhead to
retain time-stamp?



Insertion Speeds

- Inserting pre-bundled data allows humongous data rates
- >100 M s/s on 4 nodes
- >200 M s/s on 8 nodes
- Linear scaling if enough data sources



How fast is OpenTSDB ?

- OpenTSDB can scale to writing millions of data points per 'second' on commodity servers with regular spinning hard drives
- What if we needed 100-1000x faster writing speed ?



Problem: How Do Load Test Data at Large Scale?

- Testing at large scale is a more realistic measure of performance than on a small sample
- But with high velocity data, how do you set up test?
 - For a sample equivalent to long term data, you need ingest rates of 100 to 1000x faster than normal production
 - OR-
 - You must wait years to load up test data.

What's the solution?



The need for rapid data loading

- Let's say you have 1M samples / second
- Suppose you want to *test* your system
- Perhaps with a year of data
- And you want to load that data in \ll 1 year
- 100x real-time = 100M samples / second





MAPR®

Free on-demand Hadoop training
leading to certification

Start becoming an expert now
mapr.com/training

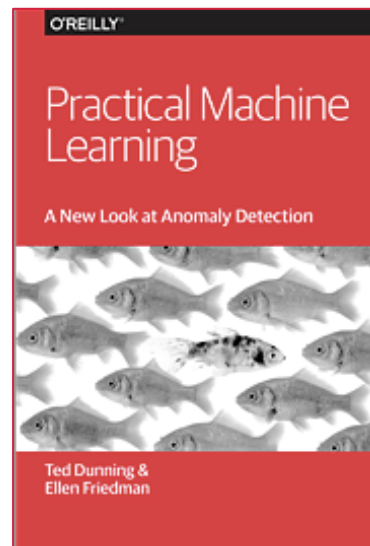


Short Books by Ted Dunning & Ellen Friedman

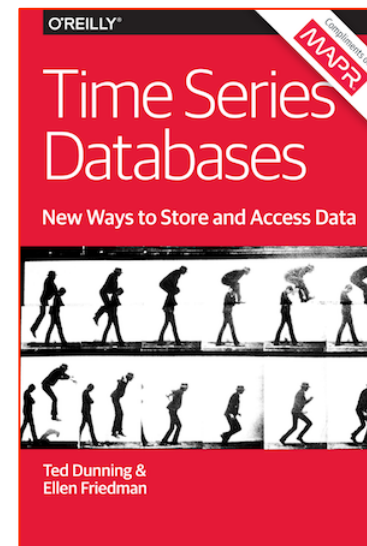
- Published by O'Reilly in 2014 and 2015
- For sale from Amazon or O'Reilly
- Free e-books currently available courtesy of MapR



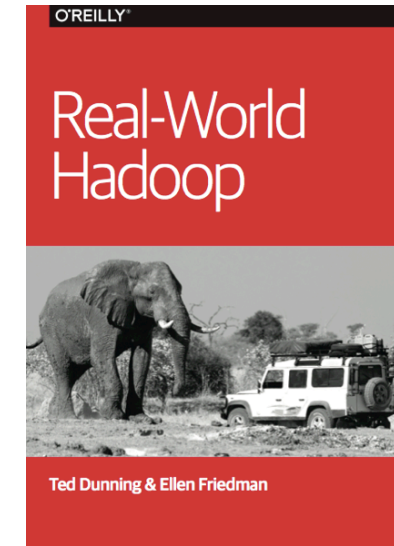
<http://bit.ly/recommendation-ebook>



<http://bit.ly/ebook-anomaly>



<http://bit.ly/mapr-tsdb-ebook>



<http://bit.ly/ebook-real-world-hadoop>



Thank You

@mapr



maprtech

mapr-technologies



MapRTechnologies

jbates@mapr.com



maprtech

