



# eBay Real-time OLAP Engine Built on Apache Kylin

May 2017

# What is Apache Kylin

kylin 麒麟

--n. (in Chinese art) a mythical animal of composite form

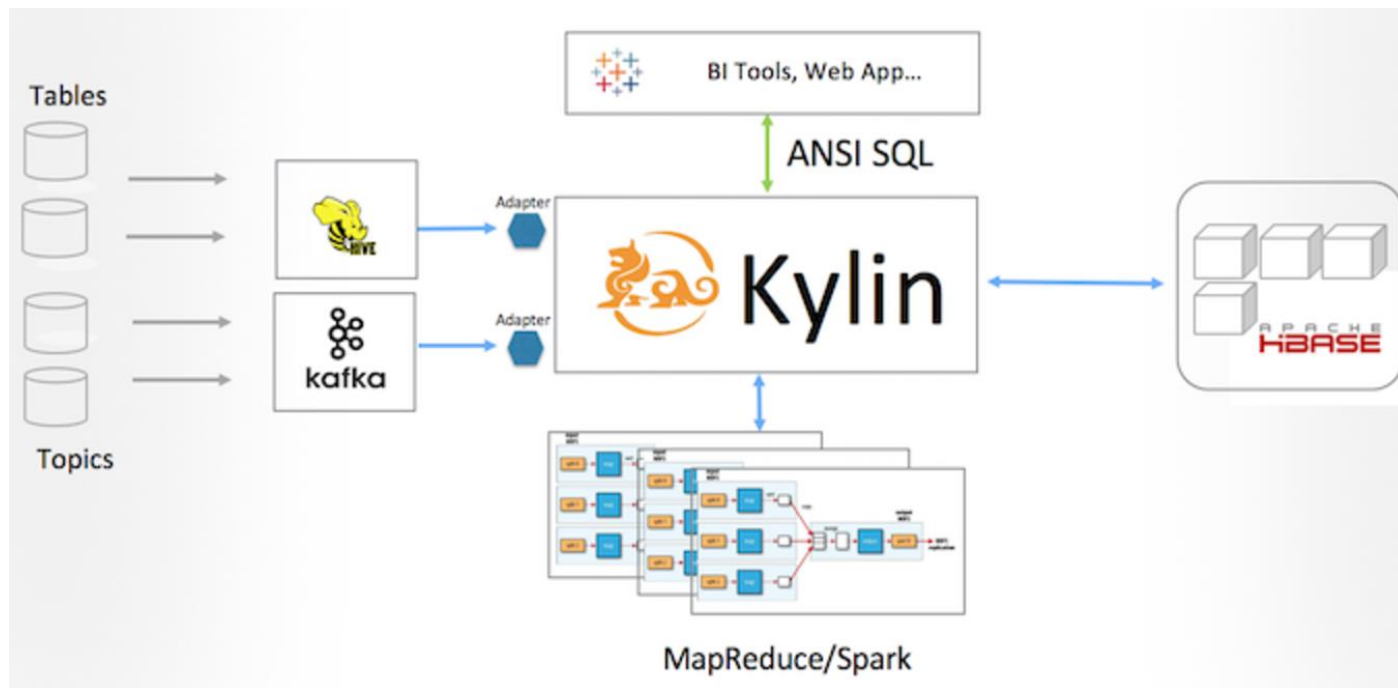


## Extreme OLAP Engine for Big Data

Apache Kylin is an open source Distributed Analytics Engine designed to provide SQL interface and multi-dimensional analysis (OLAP) on Hadoop supporting extremely large datasets, initiated and continuously contributed by eBay Inc.

**Kylin fills an important gap – OLAP on Hadoop**

# NRT Streaming in Apache Kylin 1.6



# Current Design Limitation

- ◆ Minutes level build latency
- ◆ Create many tiny HBase tables
- ◆ Highly depends on Kafka
- ◆ Hard to combine Hive and Kafka together

# Real-time OLAP

Milliseconds Query Latency + Milliseconds Data Preparation Delay

# The Challenge

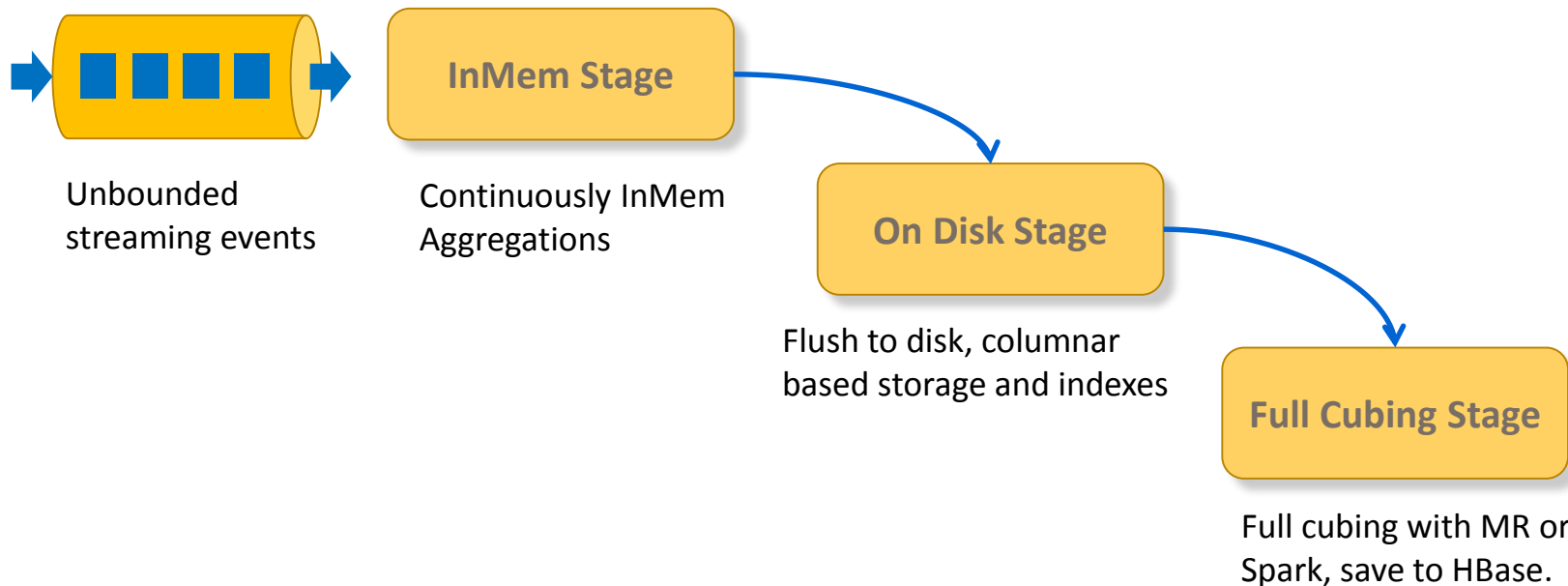
Remember Kylin is about pre-aggregations, it successfully reduced the query latencies by spending more time and resources on data preparation.

But for real-time OLAP, cubing time and data visibility latency is critical. Real-time applications are more sensitive to resource usage also.

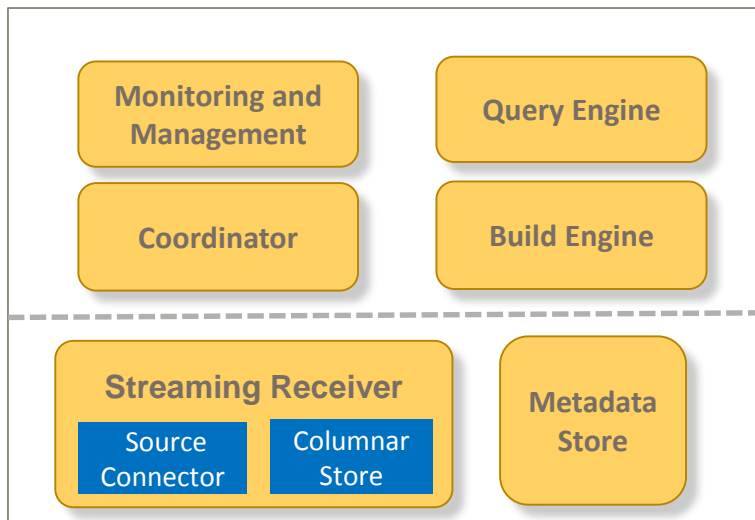
Can we **pre-build the cubes in real-time** in a more **cost effective** manner? This is hard but still doable.

# The New Solution

We divide the unbounded incoming streaming data into 3 stages, the data come into different stages are all queryable.

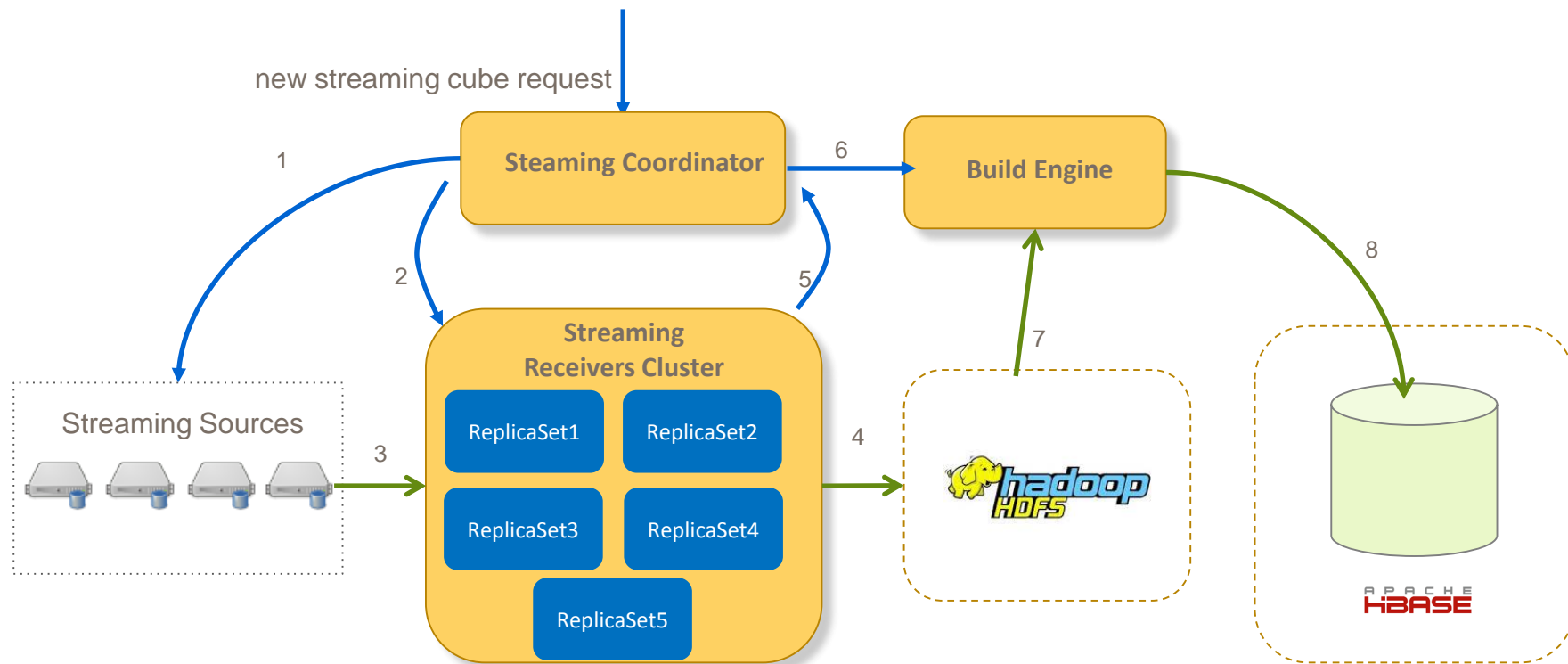


# New Streaming Components





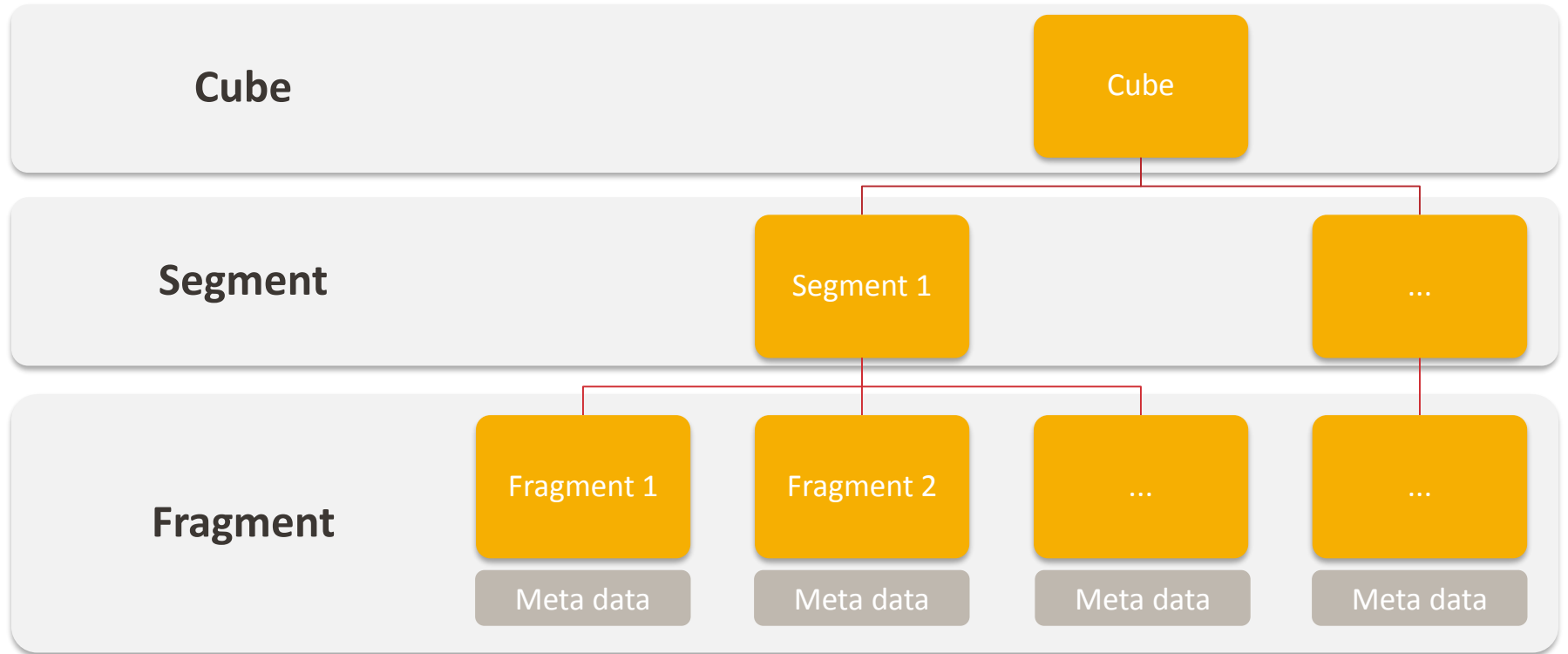
# How Streaming Cube Engine Works



# Core concepts

- ◆ Hierarchy Data Storage
- ◆ On Disk Columnar Store
- ◆ Segment Window
- ◆ Active Segments
- ◆ Immutable Segments
- ◆ Replica Set
- ◆ Check Point

# Hierarchy Data Storage

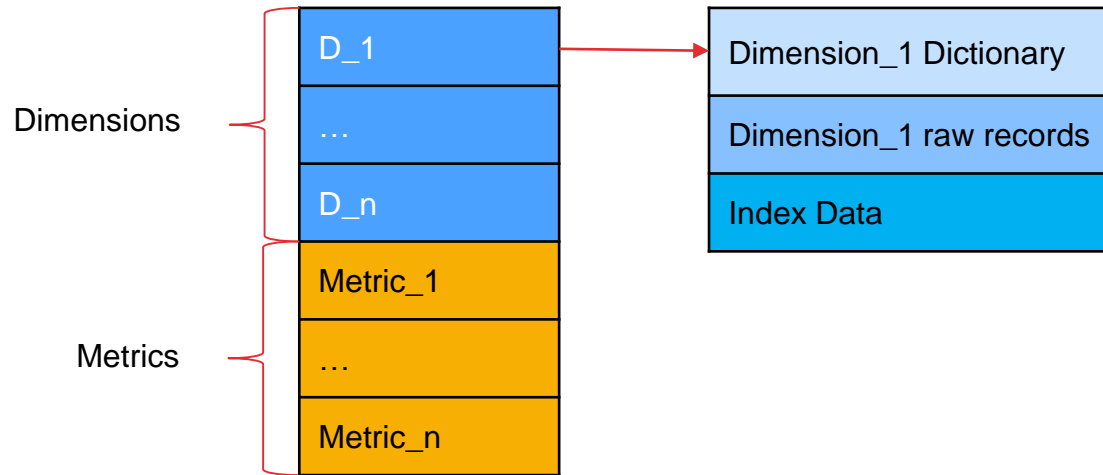


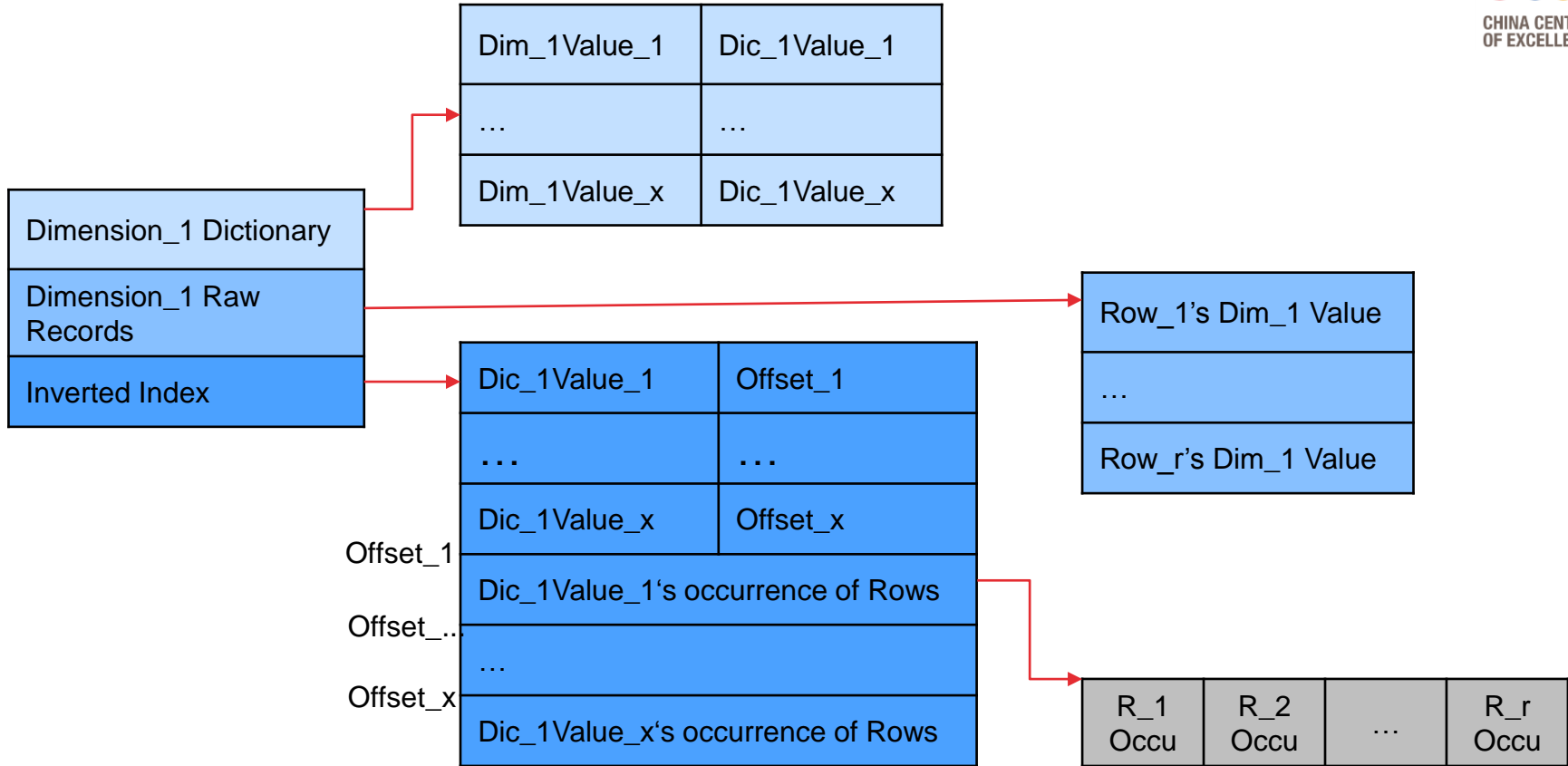
# Column Based Fragment File Design

Dimensions : D1,D2...Dn.

Metrics : M1,M2...Mm.

RawRecords : R1,R2...Rr.



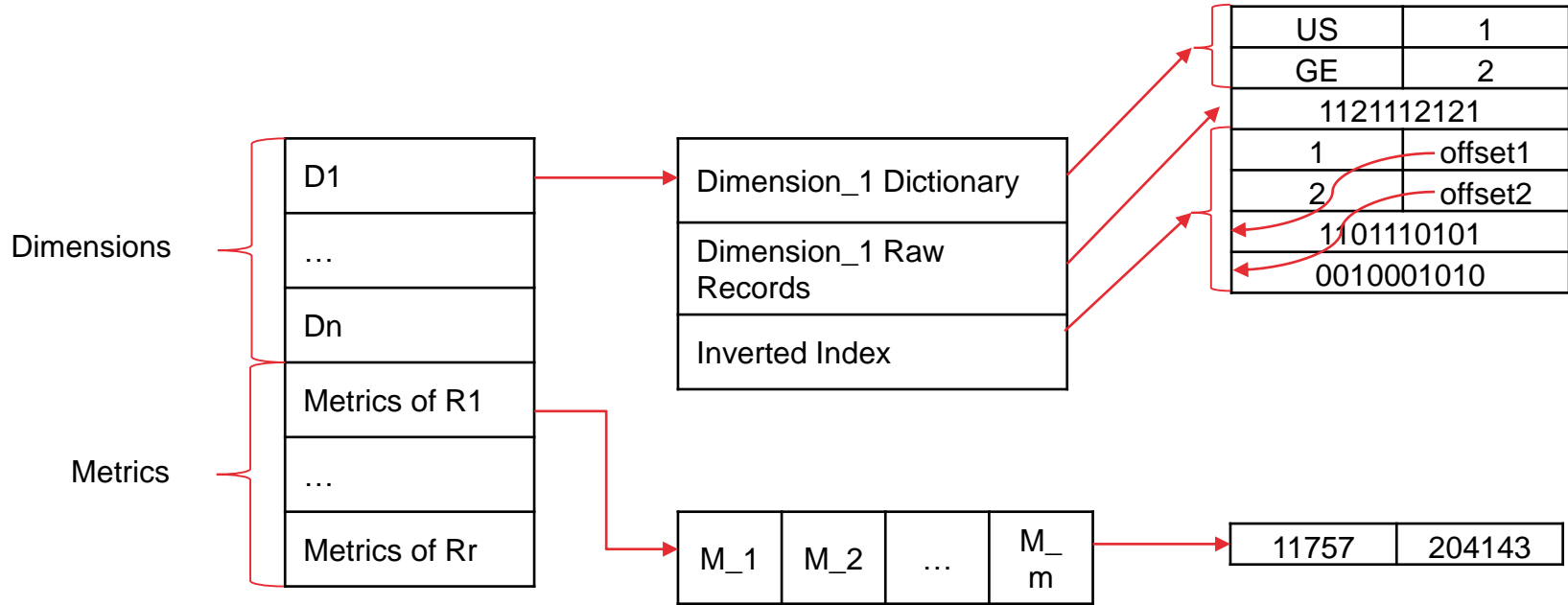


Dimensions : Site, Category, Date

Metrics : SPS, PV.

RawRecords :

Site	Category	Date	SPS	PV
US	Collectibles	8/27/16	11757	204143
US	Collectibles	8/28/16	10949	192350
GE	Fashion	8/28/16	5338	186659
US	Collectibles	8/29/16	9609	186283
US	Fashion	8/27/16	15174	177123
US	Fashion	8/28/16	14612	172552
GE	Home & Garden	8/27/16	7336	171364
US	Fashion	8/29/16	13927	171235
GE	Home & Garden	8/28/16	8615	169846
US	Home & Garden	8/29/16	4560	166506



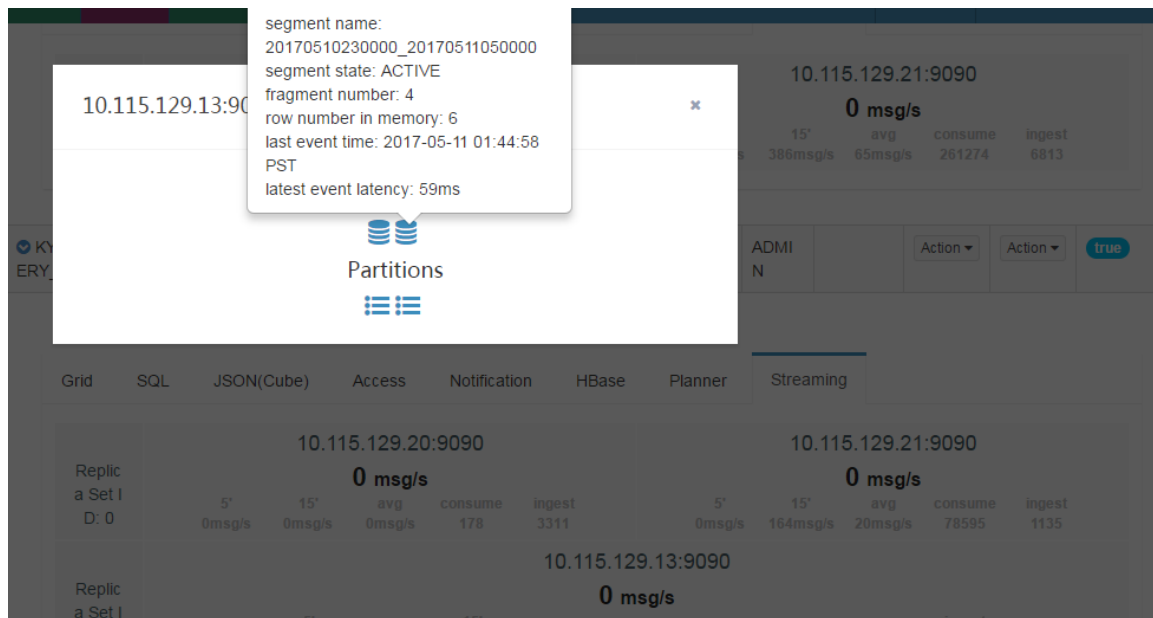
# Segment Window

Segments are divided by event time.

New created segments are first  
In Active state.

After a preconfigured data retention  
period, no further events coming,  
The segment will be changed to  
Immutable state and then write to  
remote HDFS.

Long latency events can't fit into active  
Segments will write to long latency segment.



The screenshot shows a web interface for managing data segments. A modal window titled 'Partitions' is open, displaying the following details for a segment:

- segment name: 20170510230000\_20170511050000
- segment state: ACTIVE
- fragment number: 4
- row number in memory: 6
- last event time: 2017-05-11 01:44:58 PST
- latest event latency: 59ms

The background interface shows a table of segments with columns for IP address, state, and performance metrics. The table is partially obscured by the modal window.

IP Address	State	5' avg	15' avg	consume	ingest
10.115.129.13:9090	ACTIVE	0msg/s	0msg/s	0msg/s	0msg/s
10.115.129.20:9090	ACTIVE	0msg/s	0msg/s	178	3311
10.115.129.21:9090	ACTIVE	0msg/s	164msg/s	20msg/s	78595
10.115.129.13:9090	ACTIVE	0msg/s	0msg/s	0msg/s	0msg/s



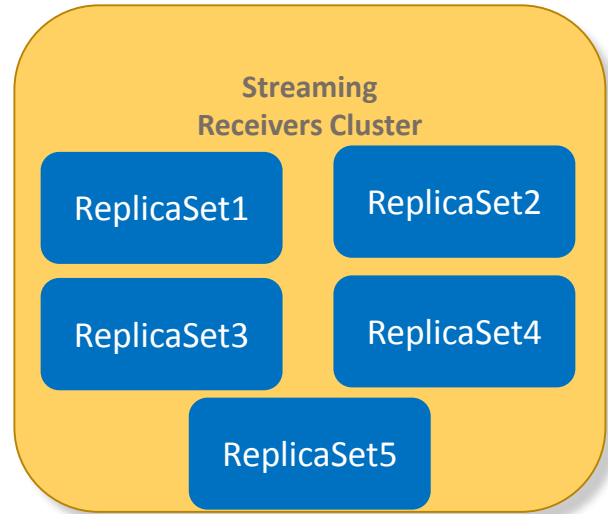
# Replica Set

The concept of Replica Set is similar to many other distributed systems like Kafka, Mongo, Kubernetes, etc.

Replica Set ensures the replications for HA purpose.

Streaming receiver instances in the same replica set have the exact same local state.

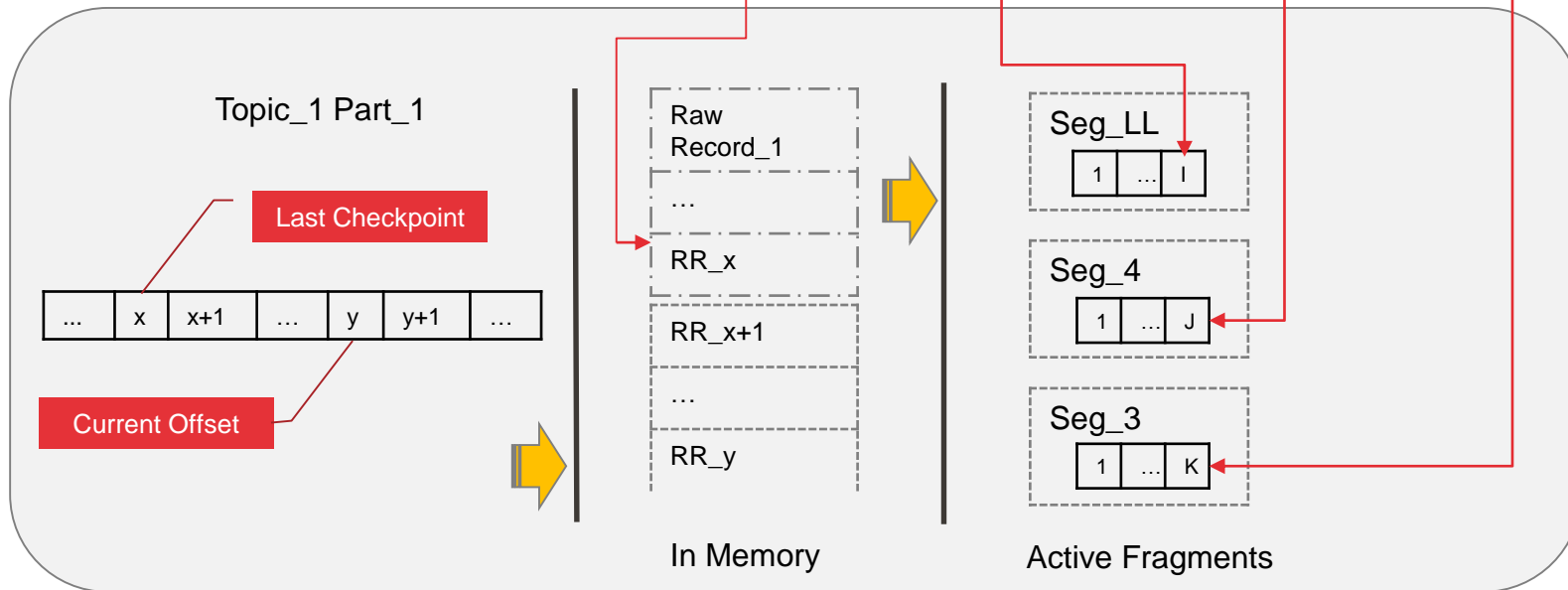
Streaming receive instances are pre-allocated to join the ReplicaSet groups.



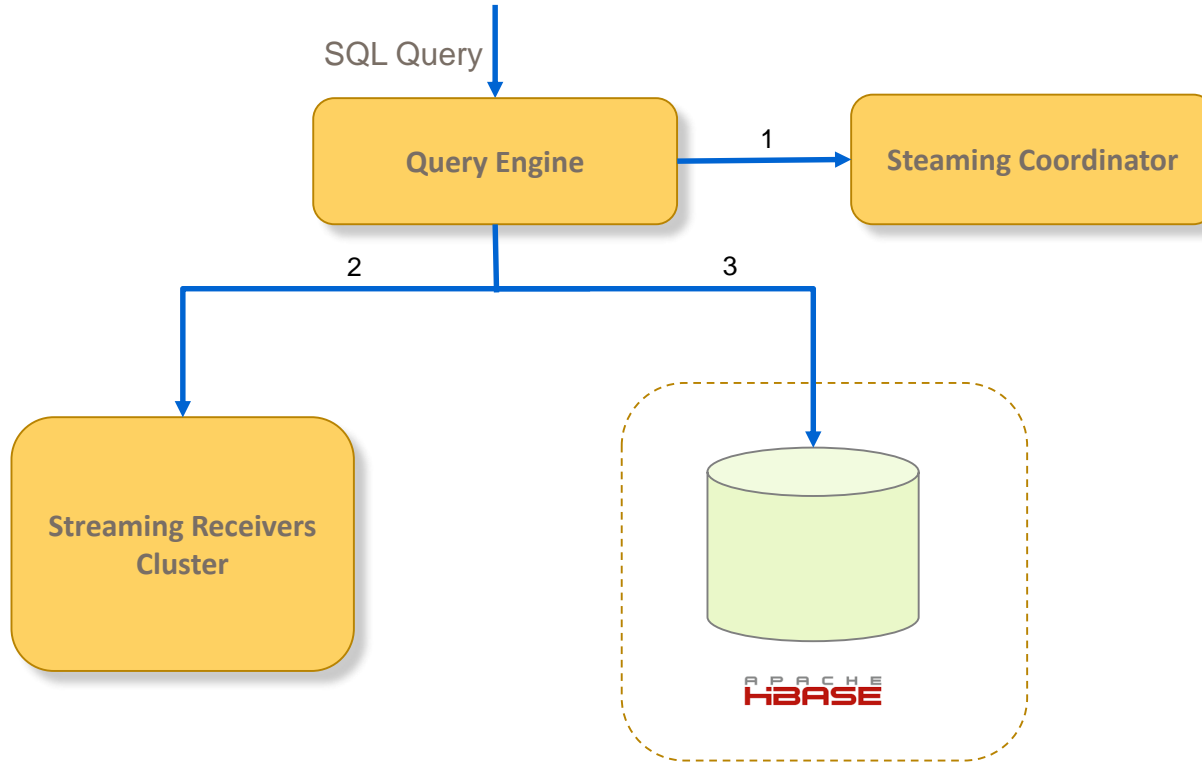
# Check Point

Kafka

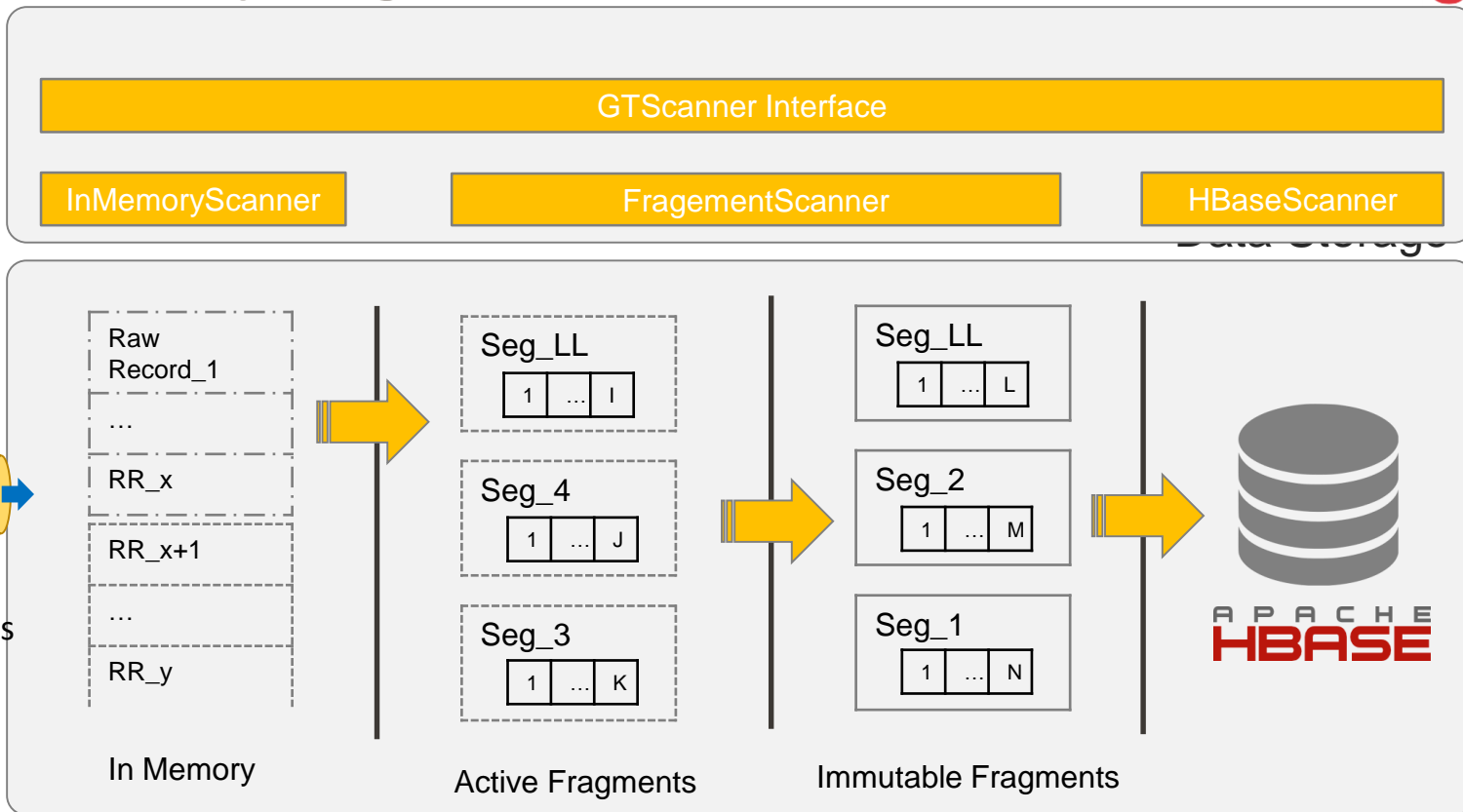
Date Time	Topic Name	Partition #	Offset	SeqID of Seg_LL	SeqID of Active Segments	
2016/10/01 12:00:00	Topic_1	1	x	l	Seg_4, J	Seg_3, K



# Extended Query Engine



# Extended Query Engine



----- Open to Write      ——— Close to Process      - - - - - Processed

# Cubing Job

Coordinator submit MR jobs incrementally to build the full cube when all the streaming receiver instances convert their segment to Immutable for that window.

The MapReduce jobs will do all the dirty work like merge all fragment files together, build dictionaries, build cubes, generate Hfiles and bulkload to Hbase.

Incrementally merge jobs can be configured to merge hourly segments to daily and then to weekly segments to avoid too many Hbase tables

# Cubing Job

KYLIN_RHEOS_METRICS_QUERY_QA2 - 20170509170000_20170509230000 - STREAM - PDT 2017-05-10 23:23:43	KYLIN_RHEOS_METRICS_QUERY_QA2	100%	2017-05-10 23:30:21 PST	6.37 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_QA2 - 20170509110000_20170509170000 - STREAM - PDT 2017-05-10 23:17:43	KYLIN_RHEOS_METRICS_QUERY_QA2	100%	2017-05-10 23:23:37 PST	5.87 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_RPC_QA2 - 20170510110000_20170510170000 - STREAM - PDT 2017-05-10 23:11:43	KYLIN_RHEOS_METRICS_QUERY_RPC_QA2	100%	2017-05-10 23:20:39 PST	8.73 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_CUBE_QA2 - 20170510110000_20170510170000 - STREAM - PDT 2017-05-10 23:13:43	KYLIN_RHEOS_METRICS_QUERY_CUBE_QA2	100%	2017-05-10 23:20:06 PST	6.12 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_QA2 - 20170508110000_20170509110000 - MERGE - PDT 2017-05-10 23:12:58	KYLIN_RHEOS_METRICS_QUERY_QA2	100%	2017-05-10 23:17:18 PST	4.30 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_CUBE_QA2 - 20170509110000_20170510110000 - MERGE - PDT 2017-05-10 23:09:10	KYLIN_RHEOS_METRICS_QUERY_CUBE_QA2	100%	2017-05-10 23:13:40 PST	4.48 min s	Action ▾	
KYLIN_RHEOS_METRICS_QUERY_RPC_QA2 -	KYLIN_RHEOS_METRICS_QUERY_RPC_QA2	100%	2017-05-10	4.17 min	Action ▾	

- 2017-05-10 23:49:45 PST
  - #1 Step Name: Build Dimension Dictionaries For Steaming Job
  - Data Size: 2.23 KB
  - Duration: 1.04 mins
  -
- 2017-05-10 23:50:47 PST
  - #2 Step Name: Save Cube Dictionaries
  - Duration: 0.01 mins
  -
- 2017-05-10 23:50:48 PST
  - #3 Step Name: Build Base Cuboid Data For Streaming Job
  - Data Size: 4.22 KB
  - Duration: 1.08 mins

# Use Case

TOTAL CUBE  
COUNT

11

[More Details](#)

QUERY  
COUNT  
**84,023**

[More Details](#)

AVG QUERY  
LATENCY

**0.15s**

[More Details](#)

JOB  
COUNT

**223**

[More Details](#)

AVG BUILD TIME  
PER MB

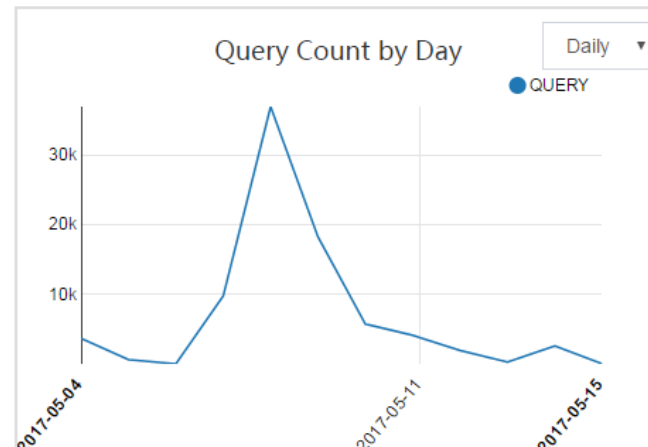
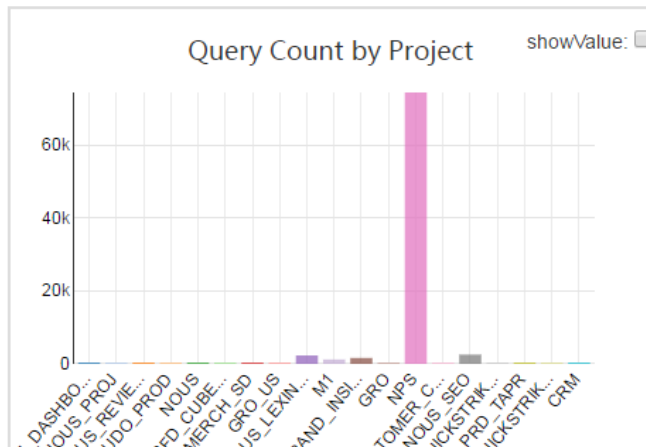
**7.00<sub>sec</sub>**

[More Details](#)

AVG CUBE  
EXPANSION

4.33

[More Details](#)



Thanks!