

SPARK AT OOYALA

EVAN CHAN
AUG 29, 2013



Who is this guy?

- Staff Engineer, Compute and Data Services, Ooyala
- Building multiple web-scale real-time systems on top of C*, Kafka, Storm, etc.
- Scala/Akka guy
- Very excited by open source, big data projects
- @evanfchan



Agenda

- Ooyala and Big Data
- What problem are we trying to solve?
- Spark and Shark
- Our Spark/Cassandra Architecture
- Spark Job Server



OOYALA AND BIG DATA



OOYALA
Powering **personalized** video
experiences across all
screens.



COMPANY OVERVIEW

Founded in 2007

Commercially launch in 2009

230+ employees in Silicon Valley, LA, NYC,
London, Paris, Tokyo, Sydney & Guadalajara

Global footprint, 200M unique users,
110+ countries, and more than 6,000 websites

Over 1 billion videos played per month
and 2 billion analytic events per day

25% of U.S. online viewers watch video
powered by Ooyala



CONFIDENTIAL—DO NOT DISTRIBUTE

TRUSTED VIDEO PARTNER

CUSTOMERS



STRATEGIC PARTNERS



CONFIDENTIAL—DO NOT DISTRIBUTE

We have a large Big Data stack

- > 250GB of fresh logs every day
- Total of 28TB of data managed over ~200 Cassandra nodes
- Traditional stack: Hadoop, Ruby, Cassandra, Ruby...
- Real-time stack: Kafka, Storm, Scala, Cassandra
- New stack: Kafka, Akka, Cassandra, Spark, Scala/Go



Becoming a big Spark user...

- Started investing in Spark beginning of 2013
- 2 teams of developers doing stuff with Spark
- Actively contributing to Spark developer community
- Deploying Spark to a large (>100 node) production cluster
- Spark community very active, huge amount of interest



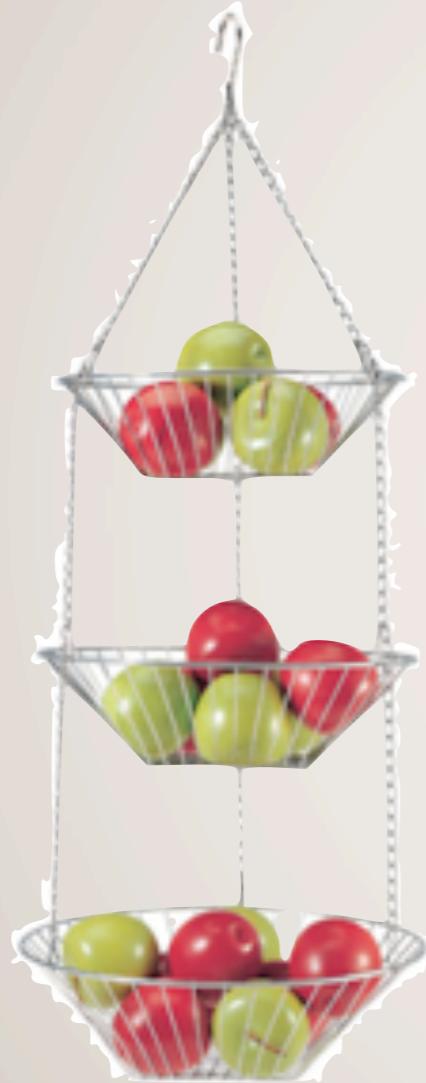
WHAT PROBLEM ARE WE TRYING TO SOLVE?



From mountains of raw data...



To nuggets of truth...



- **Quickly**
- **Painlessly**
- **At scale?**



TODAY: PRECOMPUTED AGGREGATES

- Video metrics computed along several high cardinality dimensions
- Very fast lookups, but inflexible, and hard to change
- Most computed aggregates are never read
- What if we need more dynamic queries?
 - Top content for mobile users in France
 - Engagement curves for users who watched recommendations
 - Data mining, trends, machine learning



THE STATIC - DYNAMIC CONTINUUM



100% Precomputation

100% Dynamic

- Super fast lookups
 - Inflexible,
wasteful
 - Best for 80% most
common queries
- Always compute results
from raw data
 - Flexible but slow



WHERE WE WANT TO BE



Partly dynamic

- Pre-aggregate most common queries
- Flexible, fast dynamic queries
- Easily generate many materialized views



INDUSTRY TRENDS

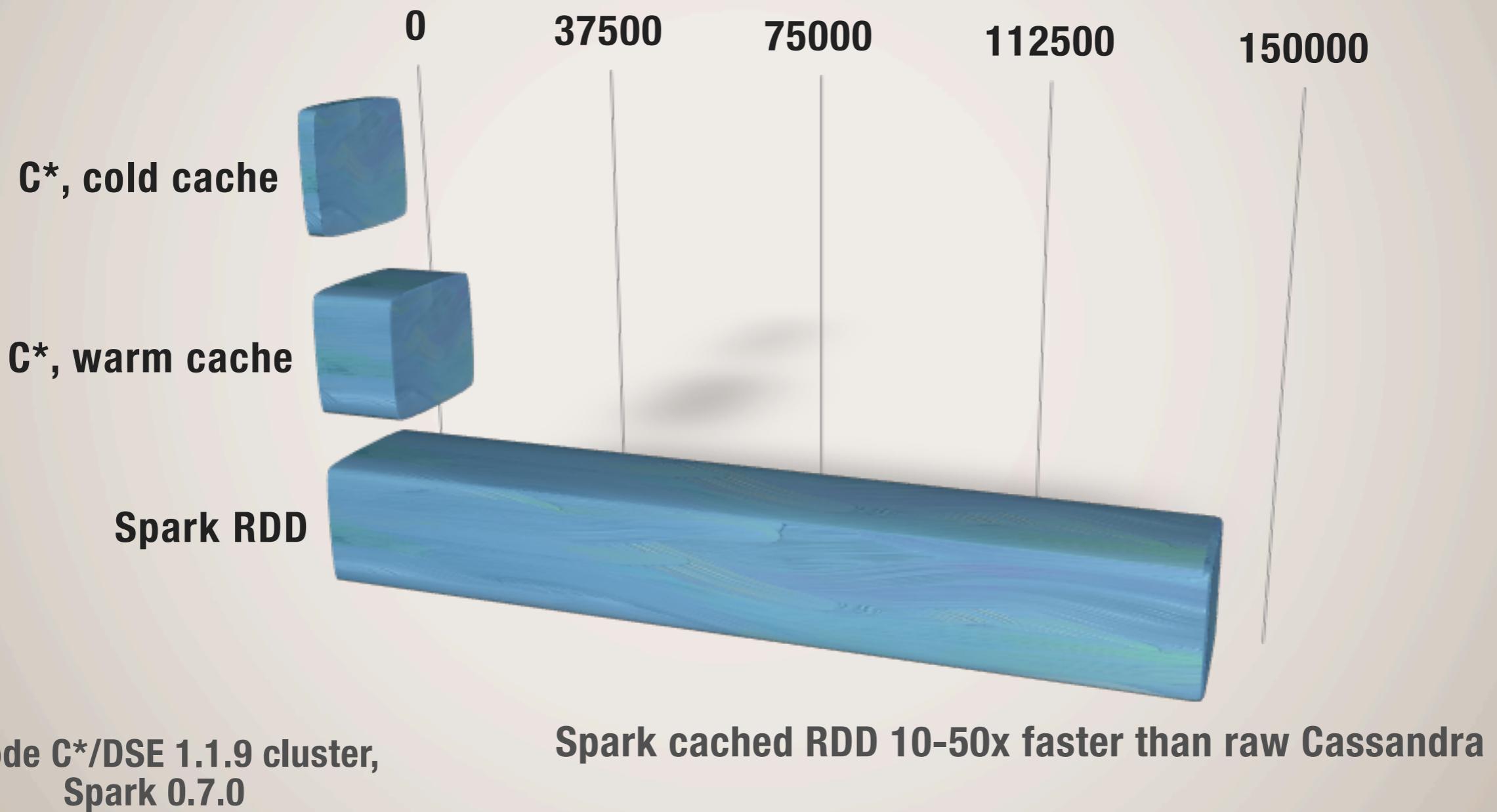
- Fast execution frameworks
 - Impala
- In-memory databases
 - VoltDB, Druid
- Streaming and real-time
- Higher-level, productive data frameworks



WHY SPARK?



THROUGHPUT: MEMORY IS KING



DEVELOPERS LOVE IT

- “I wrote my first aggregation job in 30 minutes”
- High level “distributed collections” API
- No Hadoop cruft
- Full power of Scala, Java, Python
- Interactive REPL shell



SPARK VS HADOOP WORD COUNT

```
file = spark.textFile("hdfs://...")  
  
file.flatMap(line => line.split(" "))  
  .map(word => (word, 1))  
  .reduceByKey(_ + _)
```

```
1 package org.myorg;  
2  
3 import java.io.IOException;  
4 import java.util.*;  
5  
6 import org.apache.hadoop.fs.Path;  
7 import org.apache.hadoop.conf.*;  
8 import org.apache.hadoop.io.*;  
9 import org.apache.hadoop.mapreduce.*;  
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
11 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;  
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;  
13 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;  
14  
15 public class WordCount {  
16  
17     public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {  
18         private final static IntWritable one = new IntWritable(1);  
19         private Text word = new Text();  
20  
21         public void map(LongWritable key, Text value, Context context) throws  
IOException, InterruptedException {  
22             String line = value.toString();  
23             StringTokenizer tokenizer = new StringTokenizer(line);  
24             while (tokenizer.hasMoreTokens()) {  
25                 word.set(tokenizer.nextToken());  
26                 context.write(word, one);  
27             }  
28         }  
29     }  
30  
31     public static class Reduce extends Reducer<Text, IntWritable, Text,  
IntWritable> {  
32  
33         public void reduce(Text key, Iterable<IntWritable> values, Context context)  
throws IOException, InterruptedException {  
34             int sum = 0;  
35             for (IntWritable val : values) {  
36                 sum += val.get();  
37             }  
38             context.write(key, new IntWritable(sum));  
39         }  
40     }  
41  
42     public static void main(String[] args) throws Exception {  
43         Configuration conf = new Configuration();  
44  
45         Job job = new Job(conf, "wordcount");  
46  
47         job.setOutputKeyClass(Text.class);  
48         job.setOutputValueClass(IntWritable.class);  
49  
50         job.setMapperClass(Map.class);  
51         job.setReducerClass(Reduce.class);  
52  
53         job.setInputFormatClass(TextInputFormat.class);  
54         job.setOutputFormatClass(TextOutputFormat.class);  
55  
56         FileInputFormat.addInputPath(job, new Path(args[0]));  
57         FileOutputFormat.setOutputPath(job, new Path(args[1]));  
58  
59         job.waitForCompletion(true);  
60     }  
61 }  
62  
63 }
```



ONE PLATFORM TO RULE THEM ALL



Bagel -
Pregel on
Spark

SHARK
HIVE on Spark

Spark Streaming -
discretized stream
processing

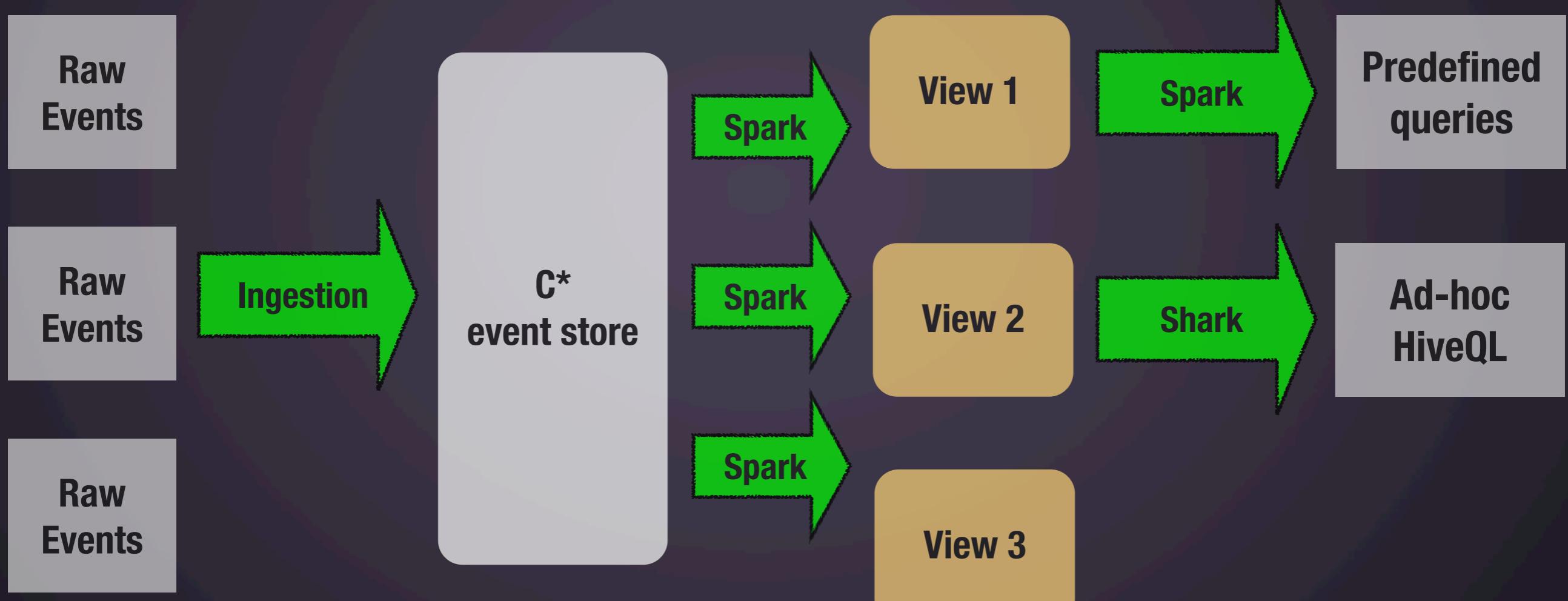
- Fewer platforms == lower TCO
- Much higher code sharing/reuse
- Spark/Shark/Streaming can replace Hadoop, Storm, and Impala
- Integration with Mesos, YARN helps



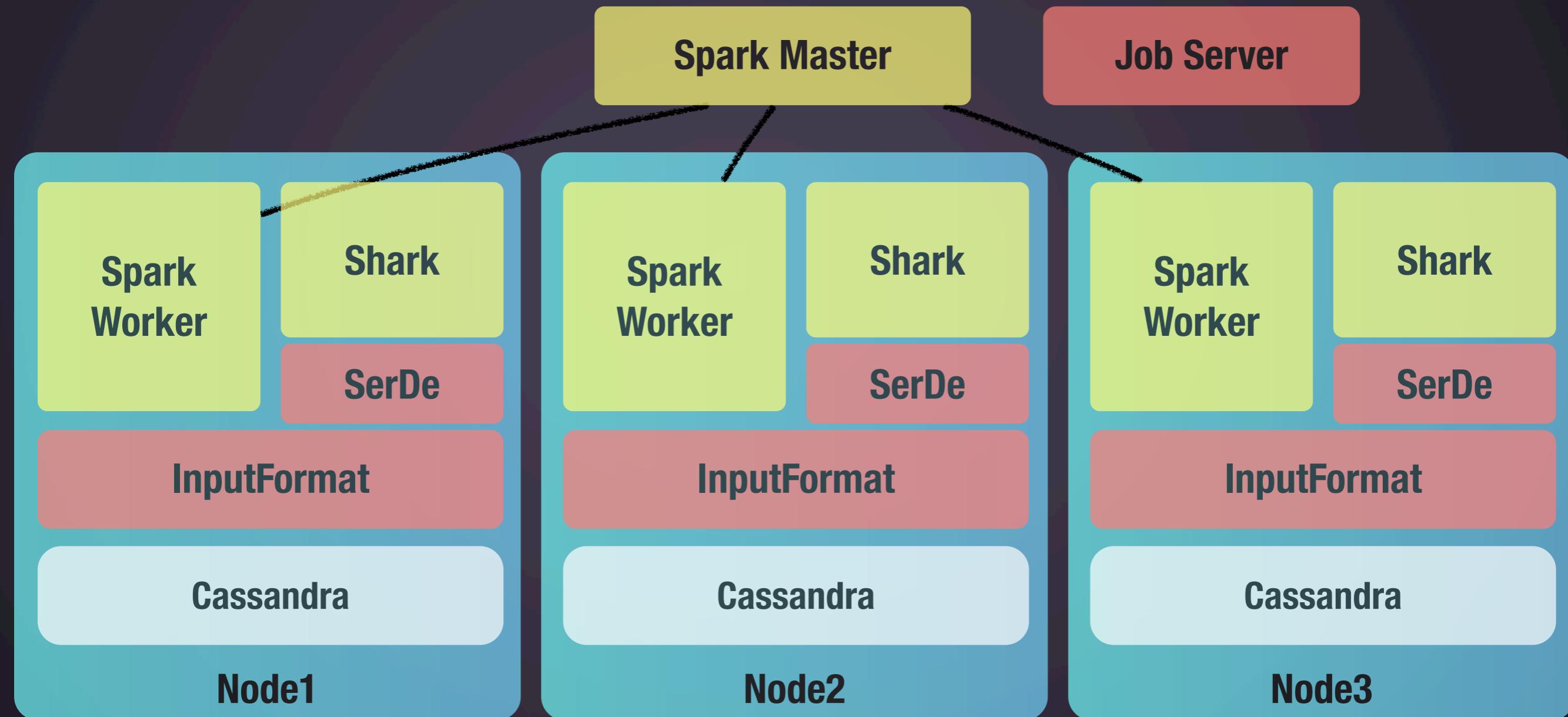
OUR SPARK ARCHITECTURE



From raw events to fast queries



Our Spark/Shark/Cassandra Stack



CASSANDRA SCHEMA

Event CF

	t0	t1	t2	t3	t4
2013-04-05 T00:00Z#id1	{event0: : a0}	{event1: a1}	{event2: a2}	{event3: a3}	{event4: a4}

EventAttr CF

	ipaddr: 10.20.30.40:t1	videoid:45678:t1	providerId:500:t0
2013-04-05 T00:00Z#id1			



INPUTFORMAT VS RDD

- You can easily use InputFormats in Spark using newAPIHadoopRDD().
- Writing a custom RDD could have saved us lots of time.

InputFormat	RDD
Supports Hadoop, HIVE, Spark, Shark	Spark / Shark only
Have to implement multiple classes - InputFormat, RecordReader, Writeable, etc. Clunky API.	One class - simple API.
Two APIs, and often need to implement both (HIVE needs older...)	Just one API.



UNPACKING RAW EVENTS

	t0	t1	UserID	Video	Type
2013-04-0 5T00:00Z#i	{video: 10, type: 5}	{video: 11, type: 11}	id1	10	5
2013-04-05 T00:00Z#id	{video: 20, type: 5}	{video: 25, type: 9}			



UNPACKING RAW EVENTS

	t0	t1
2013-04-05 T00:00Z#id	{video: 10, +vne·5}	{video: 11, +vne·11}
2013-04-05 T00:00Z#id	{video: 20, +vne·5}	{video: 25, type: 9}

UserID	Video	Type
id1	10	5
id1	11	1



UNPACKING RAW EVENTS

	t0	t1
2013-04-05 T00:00Z#id	{video: 10, type: 5}	{video: 11, type: 11}
2013-04-0 5T00:00Z#i	{video: 20, type: 5}	{video: 25, type: 9}



UserID	Video	Type
id1	10	5
id1	11	1
id2	20	5



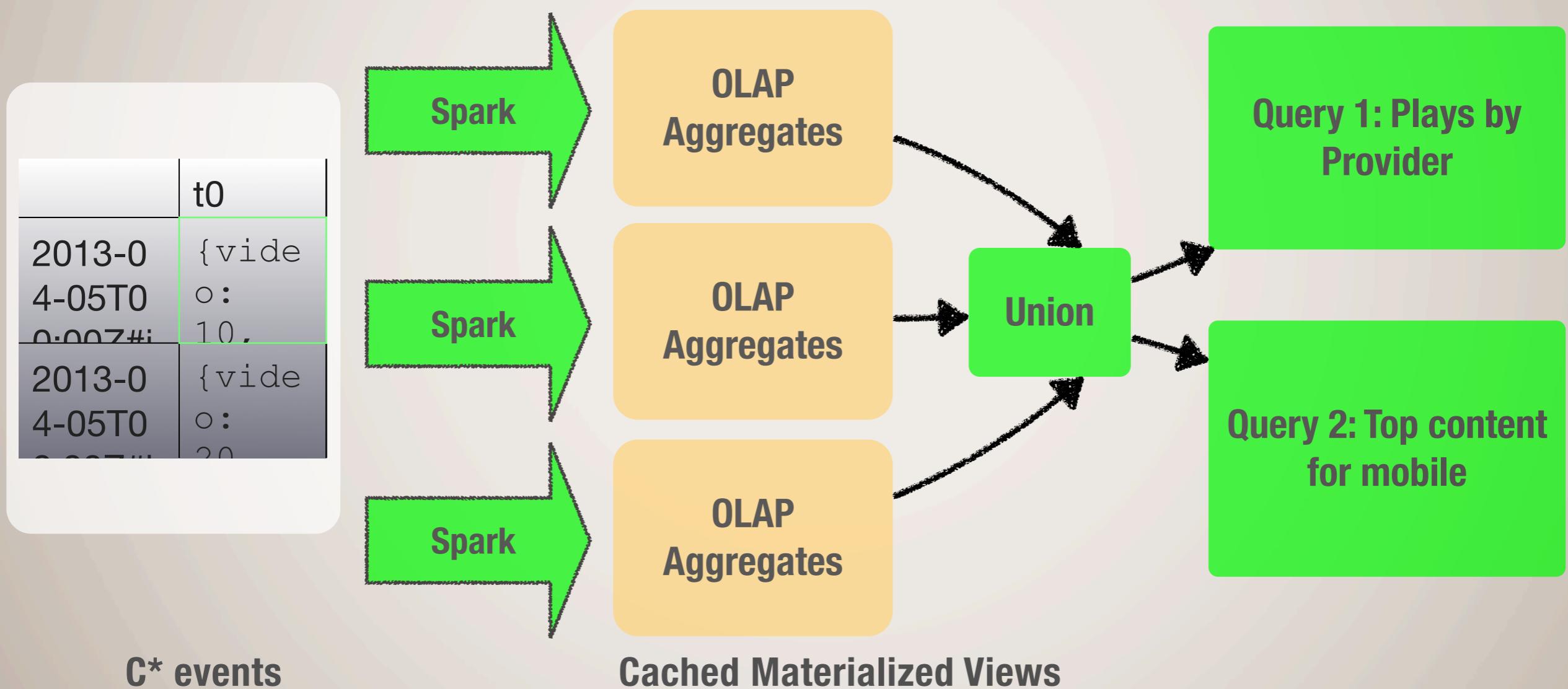
UNPACKING RAW EVENTS

	t0	t1
2013-04-05 T00:00Z#id	{video: 10, +vne·5}	{video: 11, type: 11}
2013-04-05 T00:00Z#id	{video: 20, +vne·5}	{video: 25, +vne·9}

UserID	Video	Type
id1	10	5
id1	11	1
id2	20	5
id2	25	9



EXAMPLE: OLAP PROCESSING



PERFORMANCE #'S

Spark: C* -> OLAP aggregates
cold cache, 1.4 million events 130 seconds

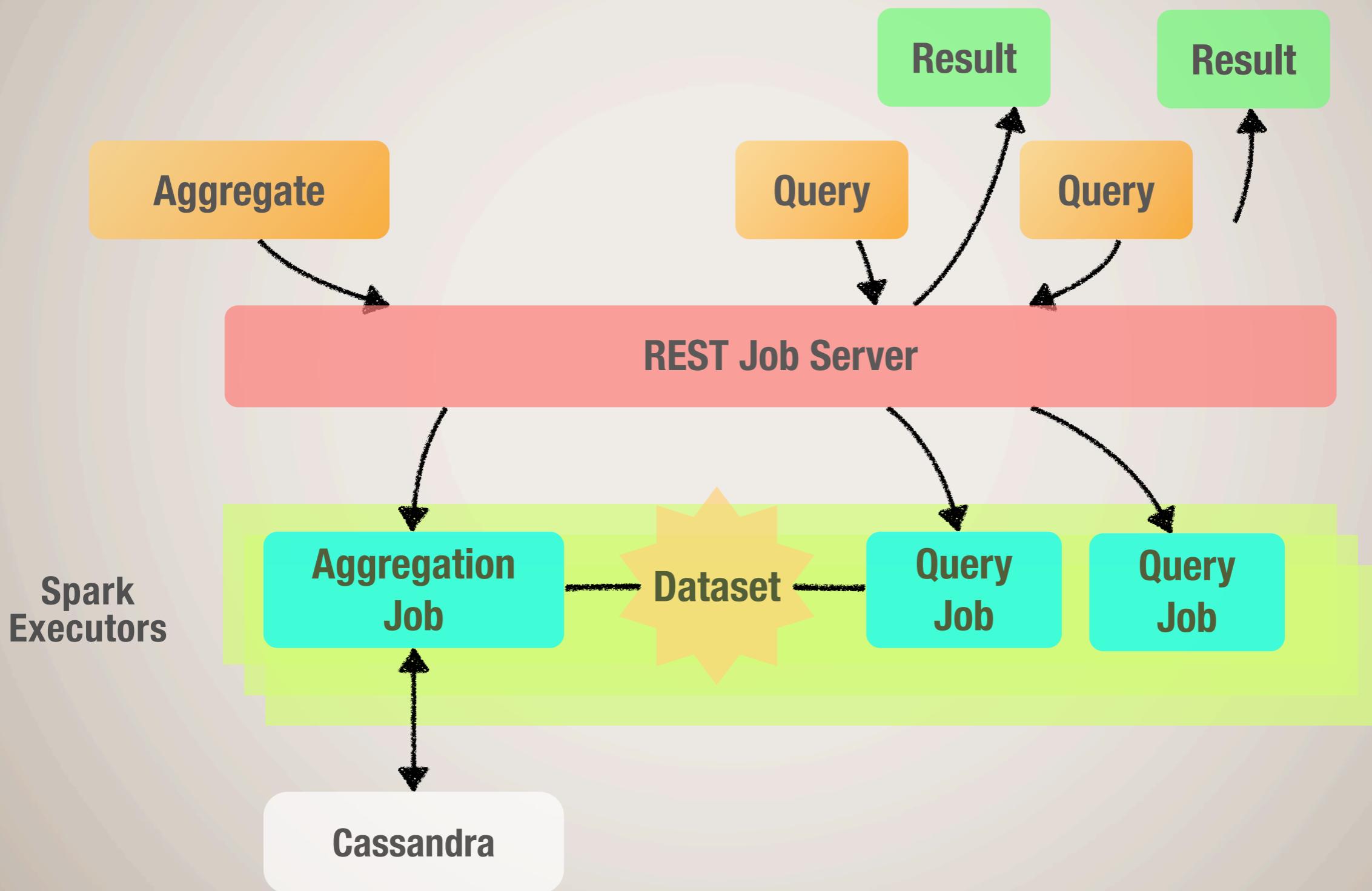
C* -> OLAP aggregates
warmed cache 20-30 seconds

OLAP aggregate query via
Spark
(56k records) 60 ms

6-node C*/DSE 1.1.9 cluster,
Spark 0.7.0



OLAP WORKFLOW



FAULT TOLERANCE

- Cached dataset lives in Java Heap only – what if process dies?
- Spark lineage – automatic recomputation from source, but this is expensive!
 - Can also replicate cached dataset to survive single node failures
- Persist materialized views back to C*, then load into cache -- now recovery path is much faster



SPARK JOB SERVER



JOB SERVER OVERVIEW

- Spark as a Service – Job, Jar, and Context management
- Run ad-hoc Spark jobs
- Great support for sharing cached RDDs across jobs and low-latency jobs
- Works with Standalone Spark as well as Mesos
- Jars and job history is persisted via pluggable API
- Async and sync API, JSON job results
- Contributing back to Spark community in the near future



EXAMPLE JOB SERVER JOB

```
/**  
 * A super-simple Spark job example that implements the SparkJob trait and  
 * can be submitted to the job server.  
 */  
  
object WordCountExample extends SparkJob {  
    override def validate(sc: SparkContext, config: Config): SparkJobValidation = {  
        Try(config.getString("input.string"))  
            .map(x => SparkJobValid)  
            .getOrElse(SparkJobInvalid("No input.string"))  
    }  
  
    override def runJob(sc: SparkContext, config: Config): Any = {  
        val dd = sc.parallelize(config.getString("input.string").split(" ")).toSeq  
        dd.map((_, 1)).reduceByKey(_ + _).collect().toMap  
    }  
}
```



SUBMITTING AND RUNNING A JOB

```
◆ curl --data-binary @../target/mydemo.jar localhost:8090/jars/demo  
OK[11:32 PM] ~  
  
◆ curl -d "input.string = A lazy dog jumped mean dog" 'localhost:8090/jobs?  
appName=demo&classPath=WordCountExample&sync=true'  
{  
  "status": "OK",  
  "RESULT": {  
    "lazy": 1,  
    "jumped": 1,  
    "A": 1,  
    "mean": 1,  
    "dog": 2  
  }  
}
```



THANK YOU

And YES, We're HIRING!!

ooyala.com/careers