# Spark and Cassandra

Martin Van Ryswyk, EVP Engineering

# Partnership

DATASTAX **+** databricks

Jonathan Ellis
Chairman, Apache Cassandra

Matei Zaharia
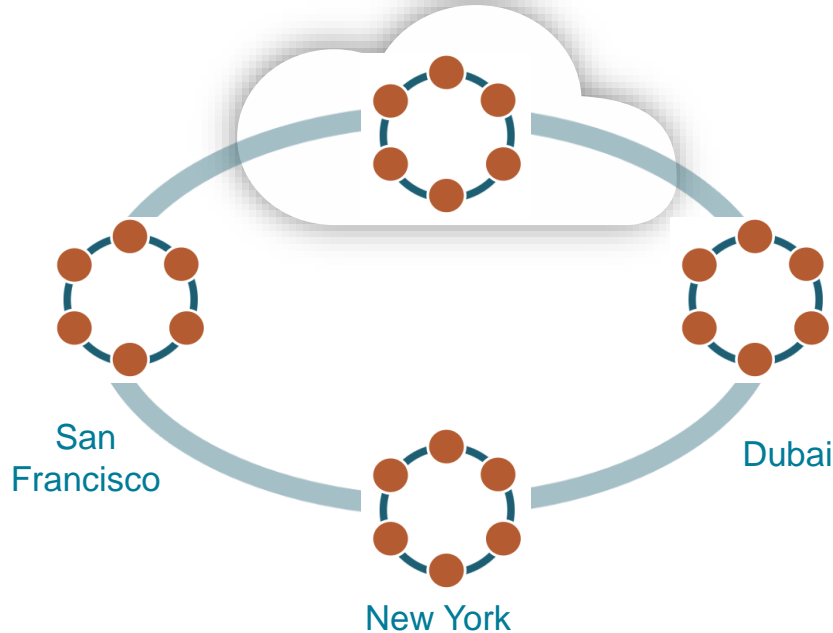Creator, Apache Spark

Announced in May 2014

Based on clear feedback from both communities

Committed to an integration between Spark and Cassandra

# What is Apache Cassandra?

# Apache Cassandra

San Francisco

New York

Dubai

- Massively scalable database
- Always on
- Fully distributed
- Linear scale performance
- Flexible NoSQL data model
- Operationally simple
- SQL-Like language
- Free tools and drivers

# Netflix Ensures Constant Uptime with DataStax Enterprise

World's leading streaming media provider with digital revenue $1.5BN+

95% of all Netflix data stored in DSE

Introduction of 'Profiles' drove throughput to over 10M transactions per second

Does 1 trillion transactions/day with DSE

Replaced Oracle in six data centers, worldwide, 100% in the cloud

# Delivers 150+ Billion Content Recommendations Per Month

Serves content for largest media brands in the world: Reuters, Wall St Journal, USA Today

Needed a massively scalable data store

High velocity of data with 58,000 links to content per second

Always-on data architecture

Lost a data center during Hurricane Sandy but never went offline

# The Weather Channel

*While I have years of experience using Cassandra, my team was mostly new to it;* **CQL made their transition essentially painless**. *But where Cassandra really shines is in* **speed and operational simplicity**, *and I would say those two points were critical."*

ROBBIE STRICKLAND *Software Dev Manager*

# Faster Feedback Loops



Transactional

Analytical

# Announcing cassandra-driver-spark

# cassandra-driver-spark v1.0

- Developed by DataStax with support and review by Databricks

- Free with Apache 2.0 license from DataStax

  https://github.com/datastax/cassandra-driver-spark

- Question to Apache Spark User List

  - user@spark.apache.org

- Offering driver code to Apache Spark community

# cassandra-driver-spark v1.0

- Exposes Cassandra tables as RDD

- Map table rows to CassandraRow objects

- Data type conversions between Cassandra and Scala

- Save RDDs back to Cassandra by implicit saveToCassandra call

- Filter rows on server via CQL WHERE clause (CassandraRDD#where method)

- Select subset of columns (CassandraRDD#select method)

- Optimizations for Cassandra vnodes

# Example

```scala
import com.datastax.driver.spark._

val sparkMasterHost = "127.0.0.1"
val cassandraHost = "127.0.0.1"
val keyspace = "test"
val table = "kv"

// Tell Spark the address of one Cassandra node:
val conf = new SparkConf(true).set("cassandra.connection.host", cassandraHost)

// Connect to the Spark cluster:
val sc = new SparkContext("spark://" + sparkMasterHost + ":7077", "demo-program", conf)

// Read table test.kv and print its contents:
val rdd = sc.cassandraTable("test", "kv").select("key", "value")
rdd.toArray().foreach(println)

// Write two rows to the test.kv table:
val col = sc.parallelize(Seq((1, "value 1"), (2, "value 2")))
col.saveToCassandra("test", "kv", Seq("key", "value"))

sc.stop()
```
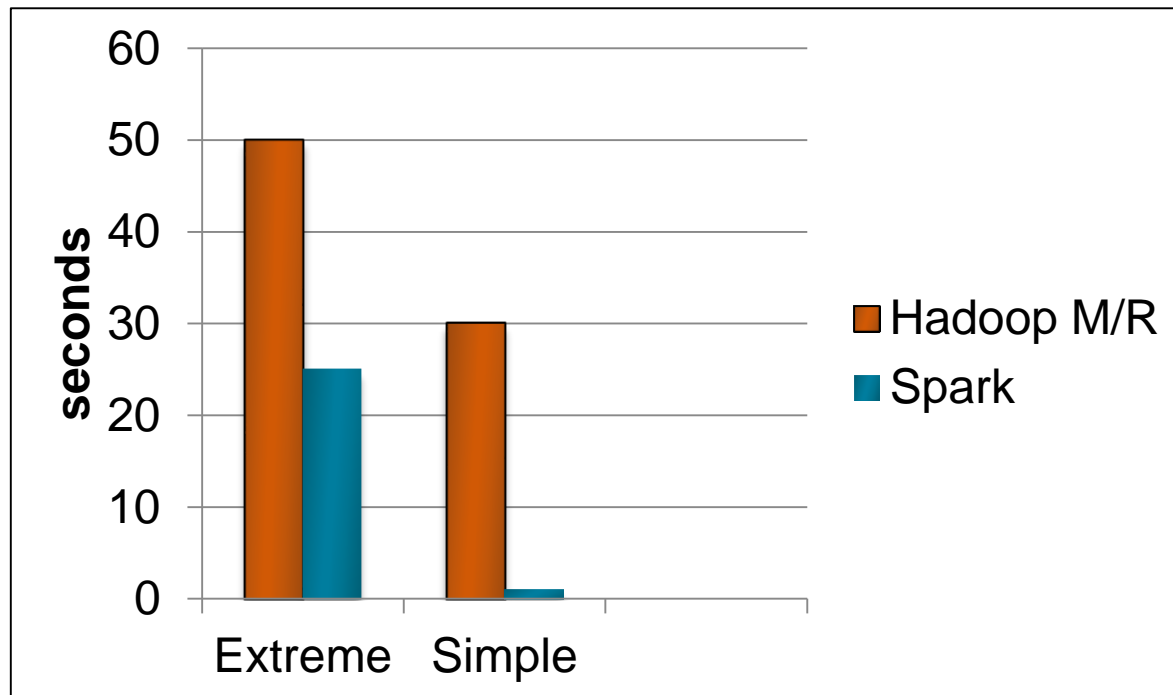
# Performance

2-30x faster than highly optimized "Hadoop on C*" implementation in DataStax Enterprise
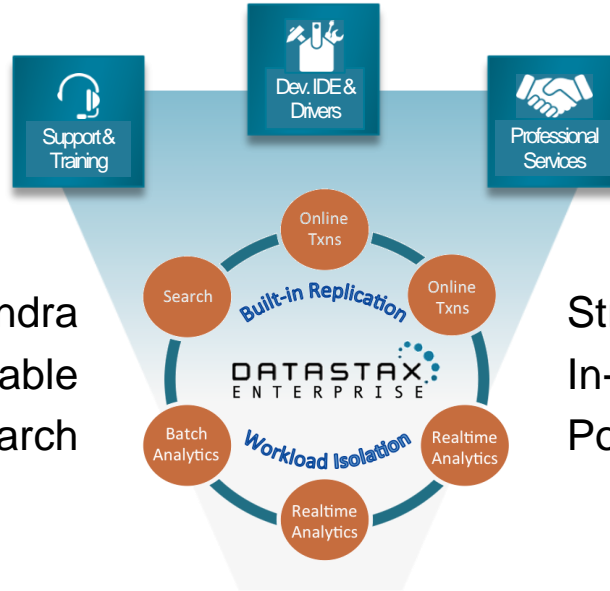
Extreme: data not in memory and fetched from multiple nodes
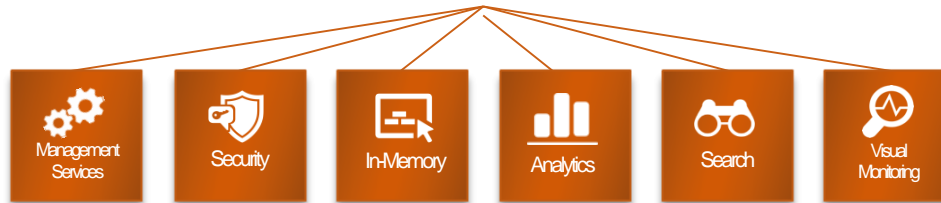Simple:　data set fits in memory

# What is DataStax Enterprise?

# Delivering Apache Cassandra to the Enterprise



Certified Production Cassandra
Multi-Workload/Use Case Capable
Integrated OLTP, Analytics, Search

Strong Data Protection
In-Memory OLTP/Analytics
Point-and-Click/Automated Mgmt

# Announcing



INTRODUCING

**DATASTAX ENTERPRISE 4.5**

The World's Fastest, Most Scalable
Distributed Database Technology

NEW

Certified
Spark
Distribution

# cassandra-driver-spark Roadmap

- Integrate into Apache Spark

- Update for Spark 1.0

- Support for Additional Spark components

  - Spark SQL

  - Spark Streaming

  - GraphX

- Listen to both communities…

# Next Steps

**Today:**

    Attend Tupshin and Al's talk at 1pm

    Visit DataStax booth – meet experts

**Tonight:**

| | |
|---|---|
| Download the integration | github.com/datastax/cassandra-driver-spark |
| Learn about Cassandra | planetcassandra.org |
| Download DataStax Enterprise | www.datastax.com/downloads |

**September:**

**CASSANDRA**SUMMIT**2014**
September 10 - 11 | #CassandraSummit
SAN FRANCISCO

**FREE ADMISSION.**
*SERIOUSLY.*

datastax.com/cassandrasummit14