# DATASTAX

## Data Velocity *and*
## Continuous Availability

Michael Shaler
Senior Director, Business Development

# What is Big Data's payoff?

**Big Data companies have outperformed their respective markets and have created competitive advantage**

Percent, 10-year CAGR (1999 – 2009)

Legend:
- Big data leader
- Other competitors

## Revenue

| | Big data leader | Other competitors |
|---|---|---|
| Grocers | 12 | 6 |
| Online retailers | 24 | -1 |
| Big box retailers | 9 | 5 |
| Casinos | 11 | 5 |
| Credit cards | 14 | 9 |
| Insurance | 9 | 8 |

## EBITDA

| | Big data leader | Other competitors |
|---|---|---|
| Grocers | 11 | 3 |
| Online retailers | 22 | -15 |
| Big box retailers | 10 | 2 |
| Casinos | 12 | 1 |
| Credit cards | 9 | -1 |
| Insurance | 14 | 5 |

SOURCE: Bloomberg and Datastream; annual reports; McKinsey analysis
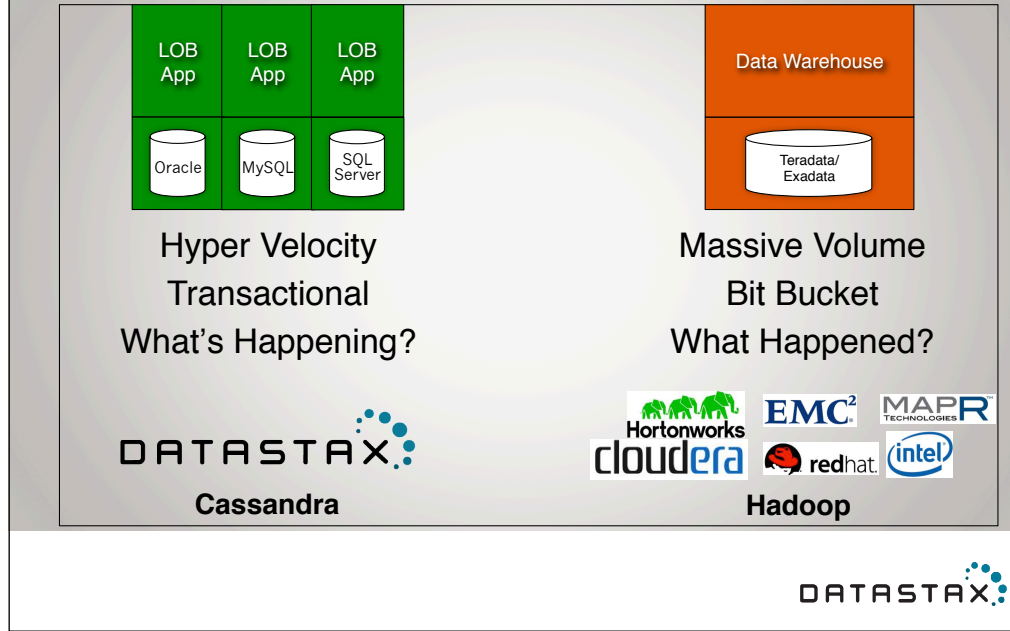
DATASTAX

# Documented Use Cases

# Real Growth In Production

# Our Solution

- DataStax Enterprise powers the big data apps that transform business.
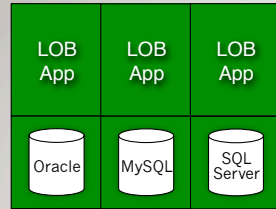
- Velocity at Scale
- Continuous Availability



DATASTAX

The Evolving Data Center

LOB App | LOB App | LOB App

Oracle | MySQL | SQL Server

Hyper Velocity
Transactional
What's Happening?

DATASTAX

Cassandra

Data Warehouse

Teradata/Exadata

Massive Volume
Bit Bucket
What Happened?

Hortonworks | EMC² | MAPR TECHNOLOGIES
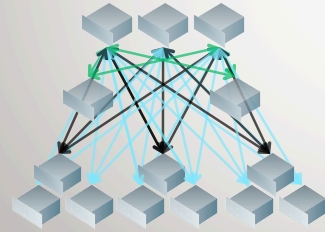cloudera | redhat. | (intel)

Hadoop

DATASTAX

Put the Cloudera logo on the bottom of the their stack and us on ours.  Put the technology names under it.
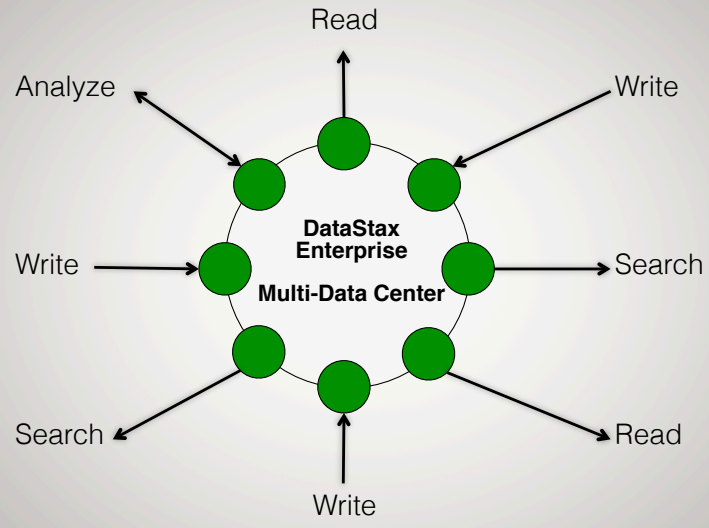
Follow it up with the case studies.
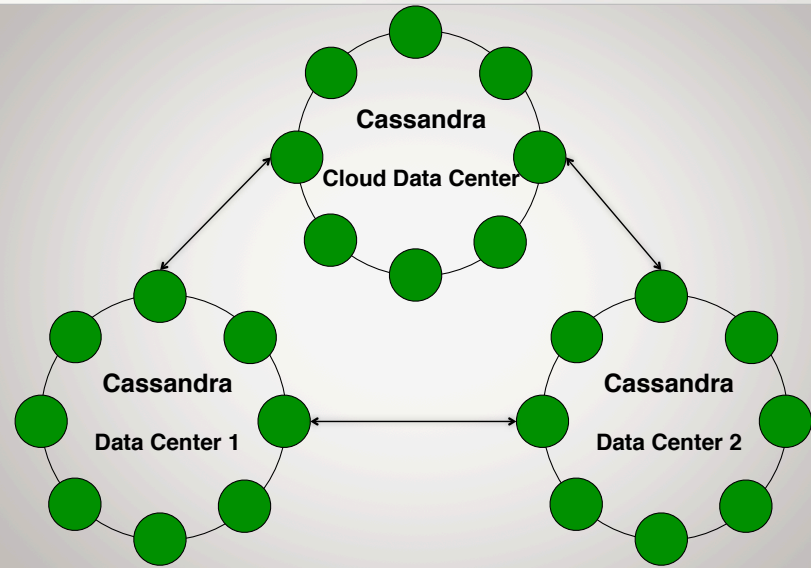
# We Tried… Oh, How We Tried.



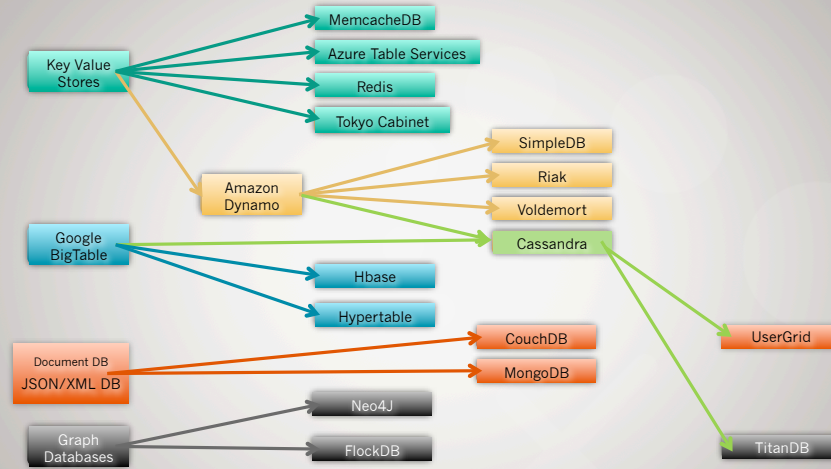LOB App | LOB App | LOB App

Oracle | MySQL | SQL Server

RDBMS Sharding

DATASTAX

# LOB Application Demands



Read

Analyze

Write

**DataStax
Enterprise**

Write

**Multi-Data Center**

Search

Search

Read

Write

DATASTAX

# The New Online Architecture

# Open Source Database Pedigrees



* Courtesy of @GuyHarrison

# Hadoop Use Cases

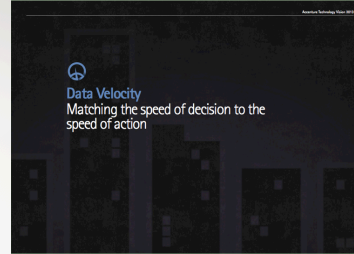Data factory

Variety

Volume

Data reservoir

**Design for Analytics**
Formulate the questions, and design
for the answers

**Relationships at Scale**
Moving beyond transactions
to digital relationships

DATASTAX

# Cassandra Use Cases



Velocity

Complexity

Data Velocity
at Scale

Continuously
available

Data Velocity
Matching the speed of decision to the
speed of action

Relationships at Scale
Moving beyond transactions
to digital relationships

Paul Maritz: "We need to bring consumer-grade back to the enterprise."

DATASTAX

# Common Use Cases

- **Cassandra**
  - Big data OLTP and write intensive systems
  - Time series data management
  - High velocity device data consumption and analysis
  - Healthcare systems input and analysis
  - Media streaming (music, movies, etc.)
  - Online Web retail (shopping carts, user transactions, etc.)
  - Online gaming (real-time messaging, etc.)
  - Real time data analytics
- **Hadoop**
  - Social media input and analysis
  - Web click-stream analysis
  - Buyer event and behavior analytics
  - Fraud detection and analysis
  - Risk analysis and management
  - Supply chain analytics
- **Solr**
  - Web product searches
  - Internal document search (law firms, etc.)
  - Real estate/property searches
  - Social media match ups
  - Web & application log management / analysis

DATASTAX

# Cassandra as Foundation

| Benefit | Feature |
| --- | --- |
| Fully Distributed: no SPOF | Peer-to-peer architecture for continuous availability and operational simplicity |
| Multi-Datacenter | Node-, rack- and DC-aware with tunable consistency |
| Massively Scalable | Multiple customers > 10M writes/second |
| SSD and Cloud optimized | All writes are linear, and all files are immutable |
| Rich Application Data Model | CQL (no joins or 2PC), integration with ODBC/JDBC et al |

DATASTAX

# Continuous Availability Commentary

**mdennis**
@mdennis
Follow

"active/passive", "shared" and "standby" are not phrases found in the description of actual "high availability" systems

Reply   Retweet   Favorite

**Bill de hÒra**
@dehora
Follow

Coming to the conclusion that #cassandra is kind of indestructible. "Robust" doesn't do it justice.

Reply   Retweet   Favorite

**Eric Florenzano**
@ericflo
Follow

"Cassandra ... dealt with the loss of one third of its regional nodes without any loss of data or availability."
techblog.netflix.com/2012/07/lesson... - Nice!

Reply   Retweet   Favorite

**Aaron Turner**
@synfinatic
Follow

took me 10hrs to notice a #cassandra node had a hw failure because everything just kept working. #sweet

Reply   Retweet   Favorite

DATASTAX

The New DR: Simian Army "Dystopia as a Service"

# Time Series Analytics: 70B readings

Smart Grid Proof of Concept: Analyze 2 years of Smart Meter data for 1M households
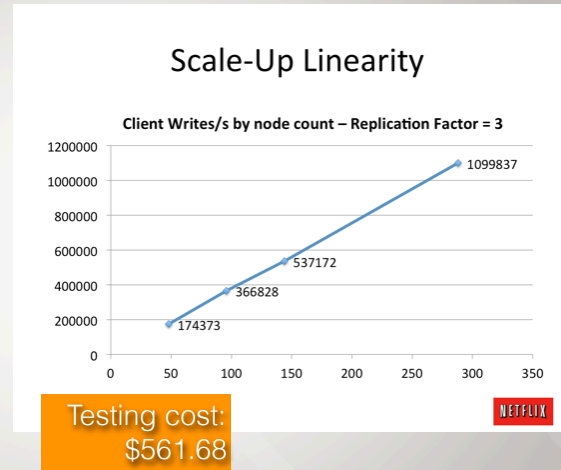Improvements in demand forecasting could yield EBITDA > $100M per GW saved



- $5M CAPEX
- 10 man/months delivery (Deploy, DevOps, Tuning)
- Ongoing OPEX of > $1M

- $450K OPEX
- 2 DevOps running 15 AWS nodes
- Faster performance in 2 weeks
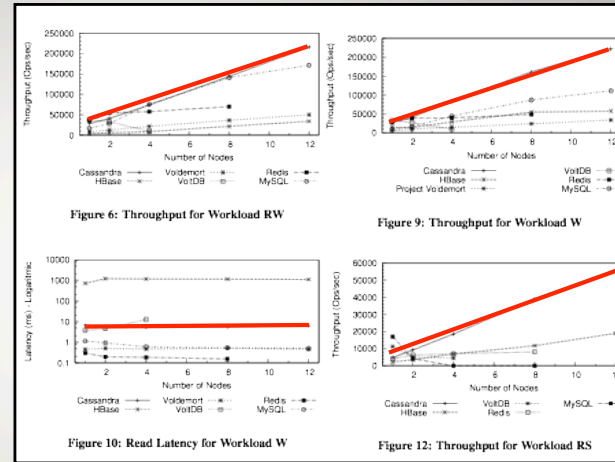- …All in the cloud

DATASTAX

# Linear Scalability on Commodity Hardware

- Yale University: "The Oracle system that is able to achieve 500,000 transactions per second costs a prohibitive $30,000,000!"

### Scale-Up Linearity

**Client Writes/s by node count – Replication Factor = 3**



Testing cost: $561.68

DATASTAX

# Performance: NoSQL Leadership

Cassandra vs. HBase:

- *10x more read throughput*
- *100x faster read latency*
- *8x more write throughput*
- *8x faster scan latency*
- *4x more scan throughput*



Figure 6: Throughput for Workload RW

Figure 9: Throughput for Workload W

Figure 10: Read Latency for Workload W

Figure 12: Throughput for Workload RS

Source: *Solving Big Data Challenges for Enterprise Application Performance Management*
Tillman Rabl, University of Toronto et al VLDB 2012 (August 2012, Istanbul)

DATASTAX

# Performance: NoSQL Leadership

### YCSB Load Process

### YCSB Read-mostly

### YCSB Read-write mix

### YCSB Write-mostly

DATASTAX

# Cassandra Summit SF (June 11-12, 2013)

# Ooyala



**Application/Use Case**

- Enable clients, such as ESPN/Disney, to monetize video streaming properties
- Analytics for all end-user interactions with videos

**Why DataStax?**

- Prior MySQL solution would not scale
- Amount of data coming in too fast and was too large for typical RDBMS
- Online data elasticity for increasing capacity with no downtime a major requirement

*"With a conventional database, we'd have to be really in the trenches, or completely re-architecting how we absorb that data. But because we had Apache Cassandra, we knew we'd have to add a few additional nodes to the cluster, at most, and without having to fundamentally re-architect our solution. It gives us tremendous competitive advantage."*

*- Harry Robertson, Tech Lead*

DATASTAX

# Disney



**Application/Use Case**
•Unified data management for 300+ websites → single platform for internal functional sites
•Real-time, analytic, and search functions

**Why DataStax/Cassandra?**
•DataStax Enterprise chosen because of its ability to support real-time and search
•Needed scalability to serve data needs of hundreds of websites
•Continuous availability for 24x7 uptime
•Search support crucial

*"We use MongoDB, but Cassandra is getting a lot more usage here. We also looked at trying to make HBase work across multiple data centers and couldn't figure out how to do it."*
*- Arun Jacob, Director of Data Services*

DATASTAX

# Walmart



**Application/Use Case**

- Need to maintain global product catalog
- Hundreds of millions of items to manage
- Need to search for products on web site
- Need to analyze product/customer interaction
- Need new product on-line shopping cart

**Why DataStax?**

- Cassandra for fast data input and change (40% of items change per day)
- Hadoop for batch analysis of operations
- Solr for accelerated web site search
- DataStax Enterprise integrates all the above in one package

*"We need a system with low latency, high throughput, and highly available for us to be able to add to our catalog 24 x 7 with staggering updates occurring in that stream of data."*
*-Rajkumar Venkat, Wal-Mart Labs*

DATASTAX

# Expedia



**Application/Use Case**

- Real-time hotel pricing support
- Database requirements exceeded Splunk capacity
- Need to run analytics on log data
- Faster search for Website (500ms SLA = $2.2M revenue)

**Why DataStax?**

- Able to process big data log information much more efficiently than other competitors
- DataStax Enterprise provides built in analytics for inputted log data
- Searches on Expedia.com being implemented for faster search response times with DataStax Enterprise
- Avoid ETL data movement between systems
- Saved $1.5 million over old SQL Server system

DATASTAX

# eBay

**Application/Use Case**

• Variety of use cases supported including write heavy logging and tracking as well as mixed workloads.

• "Social Signal" project, which enables like/own/want features on eBay product pages.

**Why DataStax?**

• Needed database to handle heavy write workloads with continuous availability.

• Multi-data center support essential.

• Easily split out functional area workloads across multiple clusters.

• Required expert Cassandra production support.

DATASTAX

# GNIP



**Application/Use Case**
- Social media provider, ingesting all data from Twitter, Facebook, Tumblr, Wordpress, and others.
- 90% of the Fortune 500 get their social media data from Gnip

**Why DataStax/Cassandra?**
- Sheer data velocity. Twitter alone can generate 20,000 tweets per second
- Massive write load and need for real-time access
- Multi-data center and cloud capable
- Serves system-of-record, compliance, and time series use cases

*"We need a real-time, massively scalable architecture, where no one node is specific, that can easily span multiple data centers and cloud availability zones, and that's Cassandra."*
*- Greg Greenstreet, VP Engineering*

DATASTAX

# Adobe



**Application/Use Case**

•Adobe AudienceManager service, an audience optimization solution that manages customers' digital data assets across the enterprise, helping them to perform web analytics, web content management and online advertising

**Why DataStax?**

•Needed a distributed cache for the ever-changing user profile data Adobe AudienceManager stores for customers

•Looked at HBase, and Membase but chose Cassandra because of linear scale for reads/writes and availability across multiple data centers

*"I can tell you we've been very pleasantly surprised by just how well Cassandra performs ."*
*- Dave Weinstein, Adobe*

**DATASTAX**

# Call to action

DATASTAX

# More Great Resources

Events

- 6/11-12: Cassandra Summit (SF)
- 6/17: Bloomberg Next Big Thing (Half Moon Bay, CA)

- Case Studies and Interviews
  datastax.com/casestudies
- Cassandra Conference Presentations
  and Videos
  datastax.com/events/
  cassandrasummit2012/presentations
- White Papers
  datastax.com/resources/whitepapers
- Webinars
  datastax.com/resources/webinars
- Training
  datastax.com/services/training

DATASTAX

Thank You

We power the big data apps that transform business.

DATASTAX