



Big Data: Beyond the Hype

Why Big Data Matters to You

WHITE PAPER

By DataStax Corporation

March 2012

Contents

Introduction.....	3
Big Data and You	5
Big Data Is More Prevalent Than You Think.....	5
Big Data Formats	6
Competitive Advantages Gained Through Big Data	7
Now What?	10
DataStax Enterprise: The Best Solution for Managing Big Data	11
Powered by Apache Cassandra	11
DataStax Enterprise – Certified Cassandra for Production Applications	13
Big Data Analytics and Enterprise Search	13
Hadoop Analytics.....	13
Search With Solr.....	14
A Complete Big Data Platform	14
Visual Database Management.....	16
Enterprise Production Support and Services	18
Conclusion.....	18
About DataStax	18

Introduction

“Big data” is a big buzz phrase in the IT and business world right now – and there are a dizzying array of opinions on just what these two simple words really mean.

Technology vendors in the legacy database or data warehouse spaces say “big data” simply refers to a traditional data warehousing scenario involving data volumes in either the single or multi-terabyte range. Others disagree: They say “big data” isn’t limited to traditional data warehouse situations, but includes real-time or operational data stores used as the primary data foundation for online applications that power key external or internal business systems.

It used to be that these transactional/real-time databases were typically “pruned” so they could be manageable from a data volume standpoint. Their most recent or “hot” data stayed in the database, and older information was archived to a data warehouse via extract-transform-load (ETL) routines.

But big data has changed dramatically. The evolution of the Web has redefined:

- The speed at which information flows into these primary online systems
- The number of customers a company must deal with
- The acceptable interval between the time that data first enters a system, and its transformation into information that can be analyzed to make key business decisions
- The kind of data that needs to be handled and tracked

Some analysts such as Gartner and others have attempted to categorize these changes by describing big data as:

1. Velocity – how fast the data is coming in
2. Variety – all types are now being captured (structured, semi-structured, unstructured)
3. Volume – potential of terabytes to petabytes of data
4. Complexity – involves everything from moving operational data into big data platforms and the difficulty in managing the data across multiple sites and geographies

Because of these changes, new definitions for big data have been proposed, with a focus on technologies to handle such data. Analyst firms such as IDC say legacy RDBMSs designed to run moderately sized data volumes on single machines do not offer sufficiently powerful engines for the big data scenarios with which modern businesses are wrestling.

Here's how IDC defines "big data":

*Big data technologies describe a new generation of technologies and architectures, designed to economically extract value from very large volumes of a wide variety of data, by enabling high-velocity capture, discovery, and/or analysis.*¹

This definition incorporates all types of data (e.g., real-time, analytic) managed by next-generation systems that must scale to handle constantly increasing user workloads and data volume.

David Kellogg, meanwhile, simply defines big data as being "too big to be reasonably handled by current/traditional technologies."² Consulting and research firm McKinsey & Company agrees with Kellogg's concept of big data and defines it as "datasets whose size is beyond the ability of typical database software tools to capture, store, manage, and analyze."³

Finally, O'Reilly defines big data the following way: "Big data is data that exceeds the processing capacity of conventional database systems. The data is too big, moves too fast, or doesn't fit the strictures of your database architectures. To gain value from this data, you must choose an alternative way to process it."⁴

This paper examines the growing prevalence of big data across nearly every industry; explains why being good at using and understanding big data is critical for firms that want to compete in their chosen market; and details how businesses can use DataStax Enterprise – a solution specifically designed to manage big data easily and effectively – to exploit the benefits derived from handling big data smartly.

¹ *Extracting Value from Chaos*, IDC, June 2011: <http://idcdocserv.com/1142>.

² "Big data' has jumped the shark," DBMS2, September 11, 2011:

<http://www.dbms2.com/2011/09/11/big-data-has-jumped-the-shark/>.

³ *Big Data: The next frontier for innovation, competition, and productivity*, McKinsey Global Institute, May 2011, p. 11:

http://www.mckinsey.com/Insights/MGI/Research/Technology_and_Innovation/Big_data_The_next_frontier_for_innovation.

⁴ "What Is Big Data?," O'Reilly Radar, January 11, 2012, <http://radar.oreilly.com/2012/01/what-is-big-data.html>.

Big Data and You

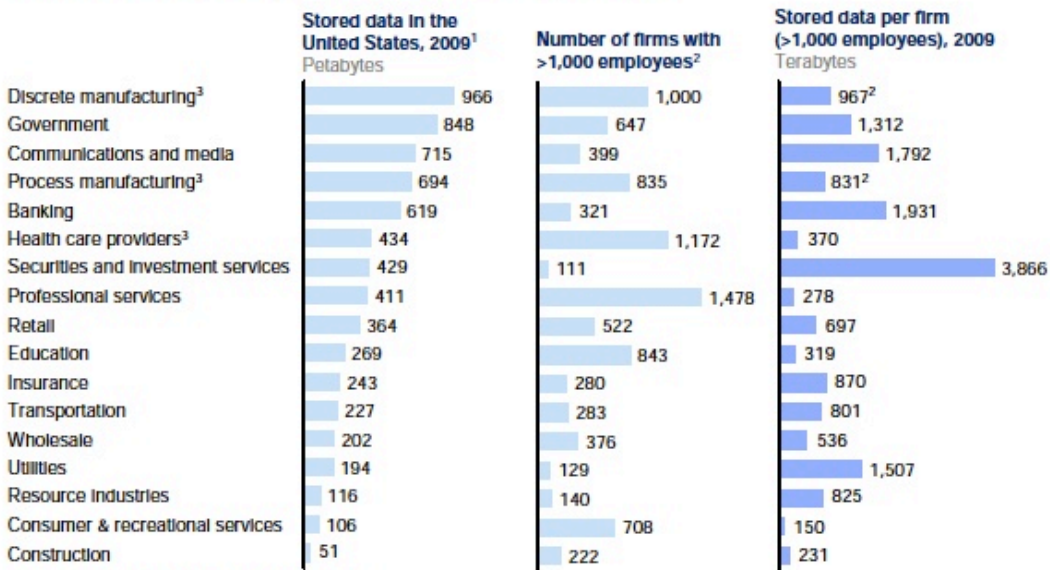
How much should you care about effectively managing big data? A lot – in fact you’re likely already dealing with big data without even knowing it.

Big Data Is More Prevalent Than You Think

Many businesses believe big data is something only companies like Facebook and Google deal with. However, a 2011 McKinsey Global Institute study says otherwise.

For instance, the McKinsey report found that the average investment firm with fewer than 1,000 employees has 3.8 petabytes of data stored, experiences a data growth rate of 40 percent per year, and stores structured, semi-structured, and unstructured data.⁵ Overall, McKinsey found that 15 out of 17 industry sectors in the United States have more data stored per company than the U.S. Library of Congress (which had 235 terabytes of information at the time of McKinsey’s study)⁶ and that companies in all sectors have at least 100 terabytes stored⁷, as Figure 1 shows:

Companies in all sectors have at least 100 terabytes of stored data in the United States; many have more than 1 petabyte



1 Storage data by sector derived from IDC.
 2 Firm data split into sectors, when needed, using employment
 3 The particularly large number of firms in manufacturing and health care provider sectors make the available storage per company much smaller.
 SOURCE: IDC; US Bureau of Labor Statistics; McKinsey Global Institute analysis

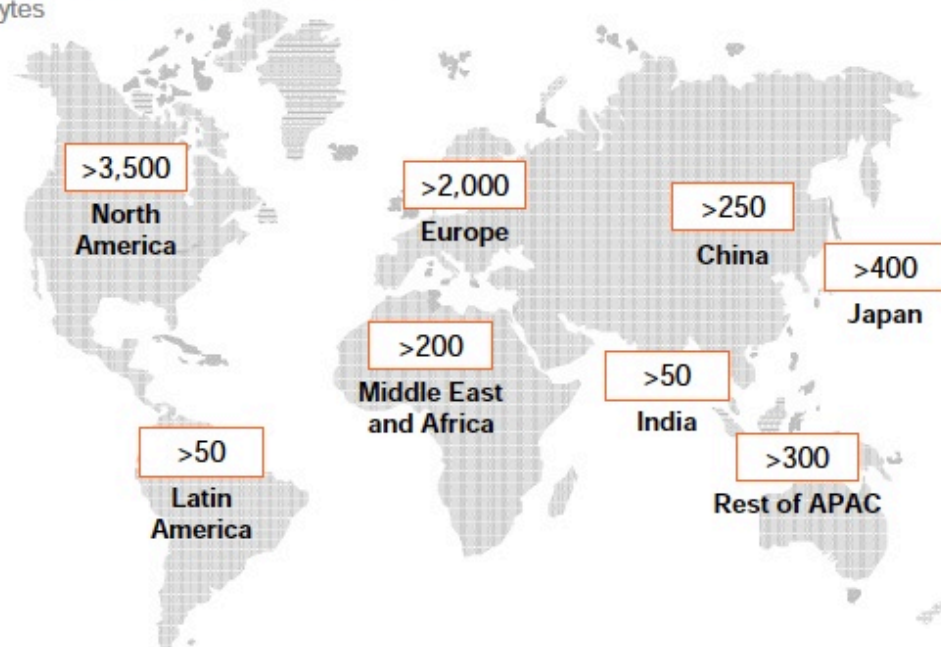
⁵ McKinsey, p. 19.
⁶ McKinsey, p. vi.
⁷ McKinsey, p. 19.

Figure 1: Stored data by industry sector

These data volumes are not confined to enterprise data warehouses that assist only internal decision-makers, but instead exist in the real-time database systems that serve external-facing customers. And that data continues to grow and expand as the underlying business becomes more successful. In 2010, the United States alone stored more than 3,500 petabytes of new information,⁸ as shown in Figure 2:

Amount of new data stored varies across geography

New data stored¹ by geography, 2010
Petabytes



¹ New data stored defined as the amount of available storage used in a given year; see appendix for more on the definition and assumptions.

SOURCE: IDC storage reports; McKinsey Global Institute analysis

Figure 2: Data stores by geography

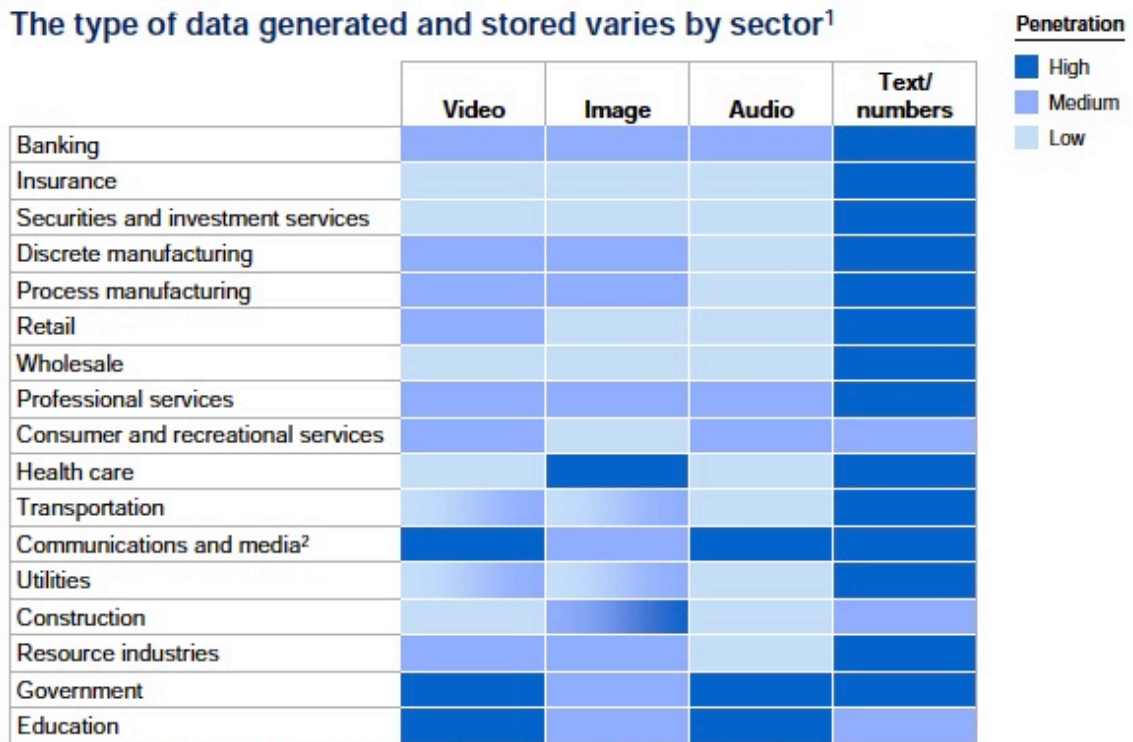
Big Data Formats

Part of the need for new technologies for big data (versus older, legacy RDBMSs) has to do with the format of the data coming in from online applications. A more dynamic, flexible database schema format is needed to handle the structured, semi-structured, and unstructured data that comprises today's big data⁹ (see Figure 3):

⁸ McKinsey, p. 103.

⁹ McKinsey, p. 20.

The type of data generated and stored varies by sector¹



¹ We compiled this heat map using units of data (in files or minutes of video) rather than bytes.

² Video and audio are high in some subsectors.

SOURCE: McKinsey Global Institute analysis

Figure 3: Data stores by sector

Competitive Advantages Gained Through Big Data

Few people today will argue with the fact that a company's data is its most strategic impersonal asset. If you don't use it as a competitive weapon against your market competitors, it's guaranteed you will be at a disadvantage. In a presentation given at the Strata New York conference in September 2011, McKinsey & Company showed the eye-opening, 10-year category growth rate differences (see Figure 4, below) between businesses that smartly use their big data and those that do not.

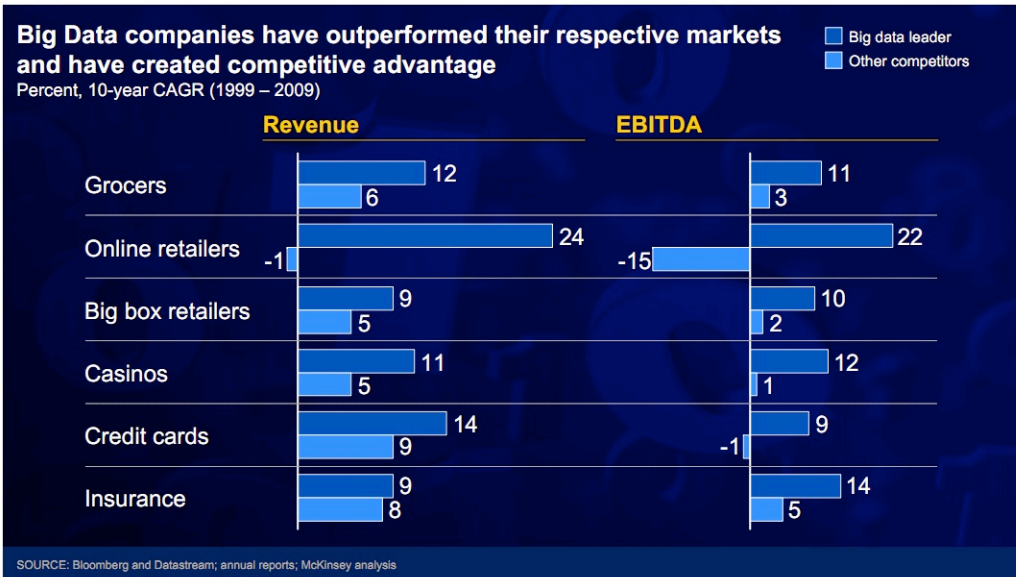


Figure 4: Big data companies have a very real competitive advantage

Clearly, the biggest differences exist between online retailers who don't use big data well, and those that do. And today, online retailing is a business every company is in, whether they're 100 percent web-based or not. Forrester Research predicts that by 2013, over half of all U.S. sales will be online in nature¹⁰ (see Figure 5):

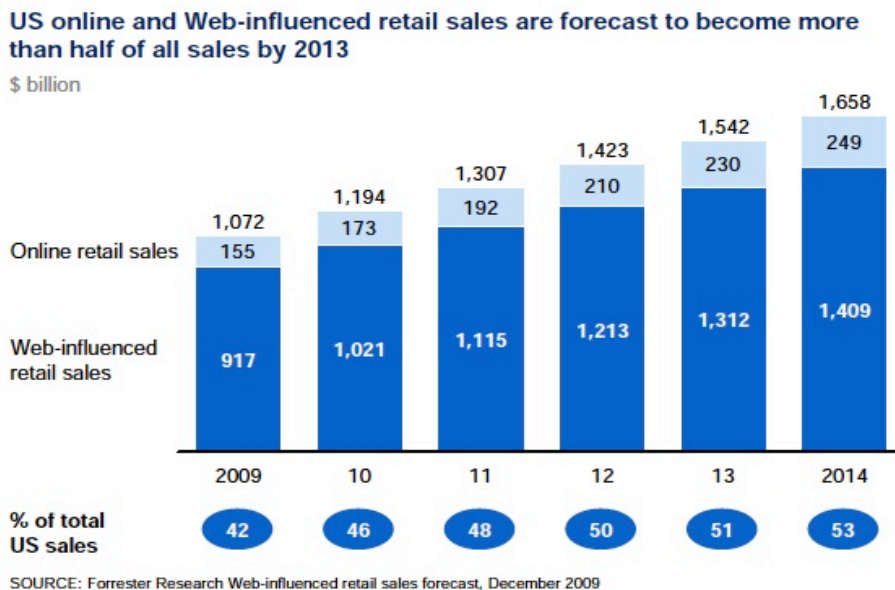


Figure 5: Online retail sales growth

¹⁰ McKinsey, p. 66.

The recognition of these market realities has led smart businesses to concentrate on effectively utilizing their big data. This is reflected in a strong jobs growth trend, as illustrated in the chart below from Indeed.com:

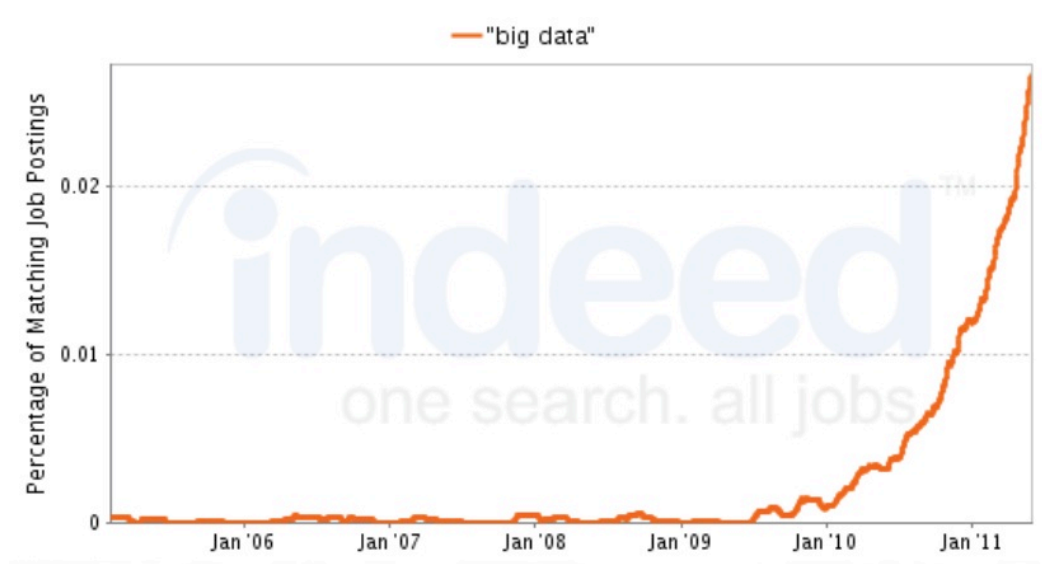


Figure 6: Big data job postings

Once data professionals are hired and put to work, there are many different data-driven projects they can be assigned to. As to the types of processes that can benefit from big data efforts, McKinsey found many different activities across all core internal corporate functions can provide value to a modern business through the use of big data¹¹ (see in Figure 7):

¹¹ McKinsey, p. 67.

Big data retail levers can be grouped by function

Function	Big data lever
Marketing	<ul style="list-style-type: none">▪ Cross-selling▪ Location based marketing▪ In-store behavior analysis▪ Customer micro-segmentation▪ Sentiment analysis▪ Enhancing the multichannel consumer experience
Merchandising	<ul style="list-style-type: none">▪ Assortment optimization▪ Pricing optimization▪ Placement and design optimization
Operations	<ul style="list-style-type: none">▪ Performance transparency▪ Labor inputs optimization
Supply chain	<ul style="list-style-type: none">▪ Inventory management▪ Distribution and logistics optimization▪ Informing supplier negotiations
New business models	<ul style="list-style-type: none">▪ Price comparison services▪ Web-based markets

SOURCE: McKinsey Global Institute analysis

Figure 7: Big data retail levers

Now What?

The facts about big data really do speak for themselves. The nonstop growth of data, new data formats that must be managed, and the competitive advantages that come from managing big data well all underscore why big data should matter to you.

But what should you do about it? If you're an IT professional, you know how difficult it can be to find a solution capable of handling a task like big data management that combines the following benefits:

- Scalability
- Performance
- Ease of use
- Low total cost of ownership (TCO)

Some database offerings may have two of these four features, but it's rare to find one with all four – and that delivers on them well. The good news is there *is* a solution available that confidently provides checkmarks for the four criteria above.

DataStax Enterprise: The Best Solution for Managing Big Data

DataStax is the leading provider of modern enterprise database software products and services based on Apache Cassandra™. It supports businesses that need a progressive data management system that can serve as a primary system of record/real-time datastore for critical production applications, and also deliver built-in analytic capabilities for analyzing that data once it's in Cassandra.

DataStax Enterprise is tailor-made to manage big data effectively. The solution inherits all of Cassandra's powerful feature set for servicing modern real-time applications, and merges in a fault-tolerant, analytics platform that provides Hadoop MapReduce, Hive, and Pig support for business intelligence systems. It also includes enterprise search capabilities via Apache Solr, which is the most popular software in use today where search is concerned.

A key differentiator of DataStax Enterprise over other big data providers is that real-time, analytic, and search workloads are smartly separated across a distributed DataStax Enterprise database cluster, so that no competition for underlying compute resources or data occurs.

Powered by Apache Cassandra

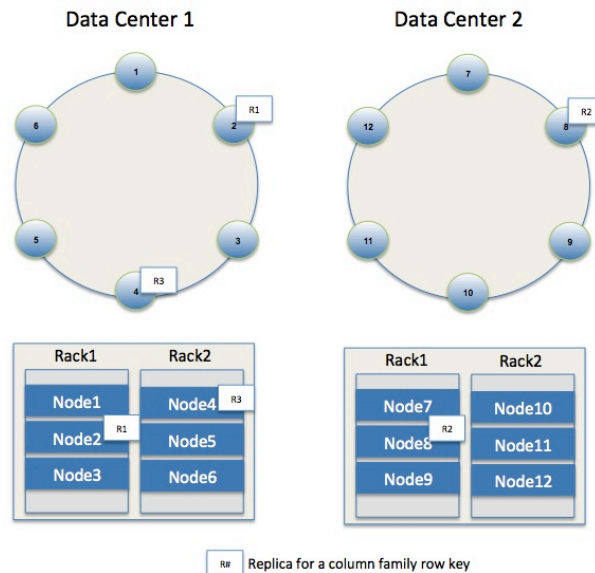
The foundation that enables DataStax Enterprise to tackle big data is Apache Cassandra. Cassandra enjoys an industry reputation for being the only NoSQL database solution able to truly handle big data requirements. It's a highly scalable and high-performance distributed database management system that can handle real-time big data applications that drive key systems for modern and successful businesses.

Key technical differentiators of Cassandra versus its legacy RDBMS predecessors, as well as other NoSQL offerings, include:

- A built-for-scale architecture that can handle petabytes of information and thousands of concurrent users/operations per second as easily as it can manage much smaller amounts of data and user traffic
- Peer-to-peer design that offers no single point of failure for any database process or function
- Online capacity additions that deliver linear performance gains for both read and write operations
- Location independence capabilities that equate to a true network-independent method of storing and accessing data; data can be read and written to anywhere
- Tunable data consistency that allows Cassandra to offer the data durability and protection like an RDBMS, but with the flexible choice of relaxing data consistency when application use cases allow

- Flexible/dynamic schema design that accommodates all formats of big data applications, including structured, semi-structured, and unstructured data; data is represented in Cassandra via column families that are dynamic in nature and accommodate all modifications online
- Simplified replication that provides data redundancy and is capable of being multi-data center and cloud in nature
- Data compression that reduces the footprint of raw big data by over 80 percent in some use cases
- A SQL-like language (CQL) that lessens the learning curve for developers and administrators coming from the RDBMS world
- Support for key developer languages (e.g., Java) and operating systems
- No requirement for any special equipment; Cassandra runs on commodity hardware

Cassandra is built with the assumption that failures can and will occur in a big data infrastructure. Therefore, data redundancy to protect against hardware failure and other data loss scenarios is built into and managed transparently by Cassandra. Furthermore, this capability can be configured so that big data applications can use a single large database that is distributed across multiple, geographically dispersed data centers, between different physical racks in a data center, and between public cloud providers and on-premise managed data centers.



These and other capabilities make Cassandra and DataStax Enterprise the smart choice for modern businesses whose big data management needs have outgrown their traditional RDBMS software.

DataStax Enterprise – Certified Cassandra for Production Applications

Cassandra is a top open source project for the Apache foundation and enjoys strong community support and developer involvement. New community releases and patches are produced very quickly, with the understanding that community builds are not put through any enterprise-styled quality assurance process, and often contain a mixture of enhancements plus bug fixes.

By contrast, DataStax Enterprise only contains selected Cassandra releases that are chosen by the expert staff and committers at DataStax. Each chosen release is then put through a rigorous certification process designed by DataStax engineers and QA staff to ensure it is stable and ready for enterprise production systems. Any found issues are immediately fixed and applied to the DataStax Enterprise server.

In addition, DataStax also provides enterprises with predictable, certified quarterly service pack updates as well as other software benefits such as emergency hot fixes (for production outages) and bug escalation privileges for customers that give priority to their issues over community submitted bugs.

Big Data Analytics and Enterprise Search

A primary benefit that DataStax Enterprise provides to enterprises needing smart big data management capabilities is its ability to service real-time, analytic, and enterprise search data operations in the same database cluster without any of the loads impacting the other. The key to making this possible is the underlying architecture of Cassandra.

Hadoop Analytics

Built into DataStax Enterprise is an enhanced Hadoop distribution that utilizes Cassandra for many of its core services. DataStax Enterprise provides integrated Hadoop MapReduce, Hive, Pig and job/task tracking capabilities, replacing Hadoop's HDFS storage layer with Cassandra (CassandraFS). The end product is a single integrated solution that provides increased reliability, simpler deployment, and lower TCO than a traditional Hadoop solution. DataStax Enterprise also is fully compatible with existing HDFS, Hadoop, and Hive tools and utilities.

Another benefit of using Hadoop in DataStax Enterprise is that it eliminates the complexity and single points of failure of the typical Hadoop HDFS layer. From an operational standpoint, there is no need to set up a Hadoop name node, secondary name node, Zookeeper, and so on.

Instead, DataStax Enterprise provides a single layer in which every node is a peer of the others and automatically knows its position in the cluster. On startup, all DataStax Enterprise nodes automatically start a Hadoop task tracker, and one of the nodes is elected to be the job tracker. If the job tracker node fails, the job tracker is automatically restarted on a different node. DataStax Enterprise utilizes full data locality awareness for Hadoop task assignment.

Search With Solr

DataStax Enterprise includes strong enterprise search support via Lucene and Apache Solr. Coming from the Apache Lucene project, Solr is the most popular open source enterprise search platform in use today.

Solr's primary features include robust full-text search, hit highlighting, faceted search, rich document (e.g., PDF, Microsoft Word) handling, and geospatial search.

By integrating Solr into the DataStax Enterprise big data platform, DataStax extends Solr's native capabilities and provides:

- An easily scalable search platform
- No single point of failure
- No write bottleneck as in community Solr
- Automatic data sharding
- Multi-data center capabilities
- Easy, ad-hoc index rebuilds
- The ability to query search data with Cassandra's CQL

In essence, in the same way that DataStax Enterprise takes Hadoop and delivers a fault-tolerant, no single point of failure, and dynamically scalable Hadoop/analytics system, it automatically does the same thing for Solr and enterprise search operations. Using Cassandra as the underlying foundation, DataStax Enterprise allows search data to be written to any participating search node in a DataStax Enterprise cluster. New search nodes can be added online to increase both fault tolerance and performance, with gains being near linear in nature.

Those currently using Solr will be right at home in DataStax Enterprise, as it is 100 percent Solr compatible, and all Solr utilities, APIs, and so on, are included.

A Complete Big Data Platform

A key benefit of DataStax Enterprise is the tight feedback loop it has between real-time applications and the analytics and search operations that naturally follow. Traditionally, users would be forced to move data between systems via complex ETL processes, or perform both functions on the same system with the risk of one impacting the other. In big data environments, this process can be time-consuming and burdensome.

With DataStax Enterprise, real-time, analytic, and search big data operations take place in the same distributed system, but users have the ability to dedicate certain nodes solely for analytics or search so their workloads don't slow down real-time processing. Users simply define one or more replica groups, and configure the role of each – one or more Cassandra, Hadoop or HDFS (i.e., HDFS without job/task tracker), and search/Solr nodes. Writes are instantly replicated between all nodes.

With DataStax Enterprise, users truly have the best of all worlds for big data management. They have all the power of Cassandra serving their highest-volume and high-velocity real-time applications, the power of Hadoop, Hive, and Pig working directly against the same data for analytics, and Solr for enterprise search in the same distributed database. The result is smart workload isolation for big data application, which is much simpler to manage and more reliable than any of the alternatives.



Figure 8: DataStax Enterprise – real-time, analytic, and search capabilities in one integrated big data platform

Visual Database Management

DataStax Enterprise includes a visual, browser-based management solution named OpsCenter Enterprise. OpsCenter Enterprise allows a developer or administrator to manage and monitor the health of big data clusters from a centralized web console.

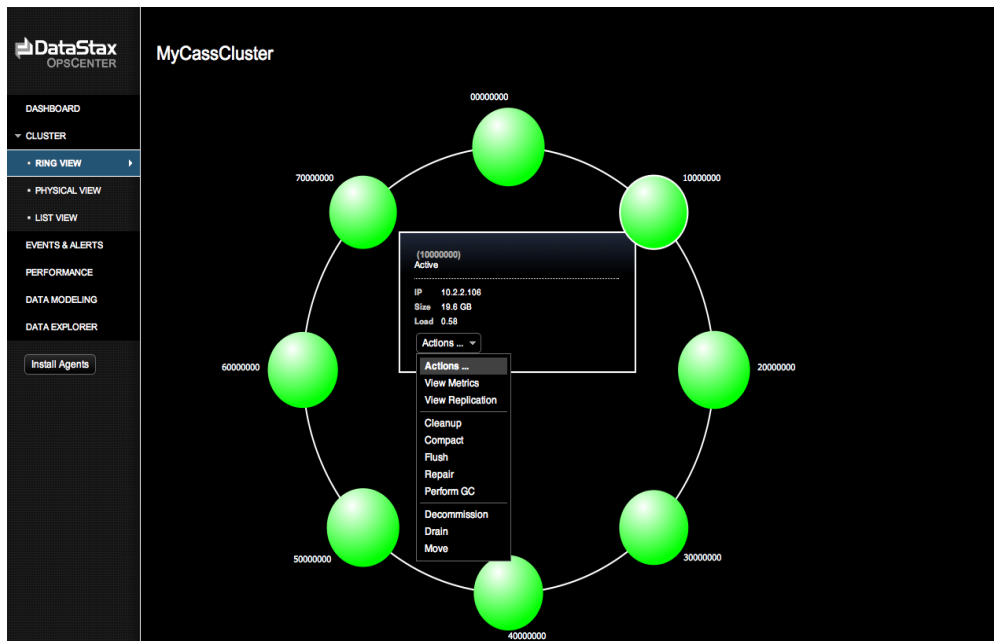


Figure 9: OpsCenter Enterprise database cluster ring view

OpsCenter Enterprise uses an agent-based architecture to monitor and carry out tasks on each node in a DataStax Enterprise cluster. Through a graphical and intuitive point-and-click interface, a user can understand the state of a cluster, which nodes are up and down, and what type of performance users are experiencing. Key events are reported into a centralized dashboard displayed along with other vital statistics.

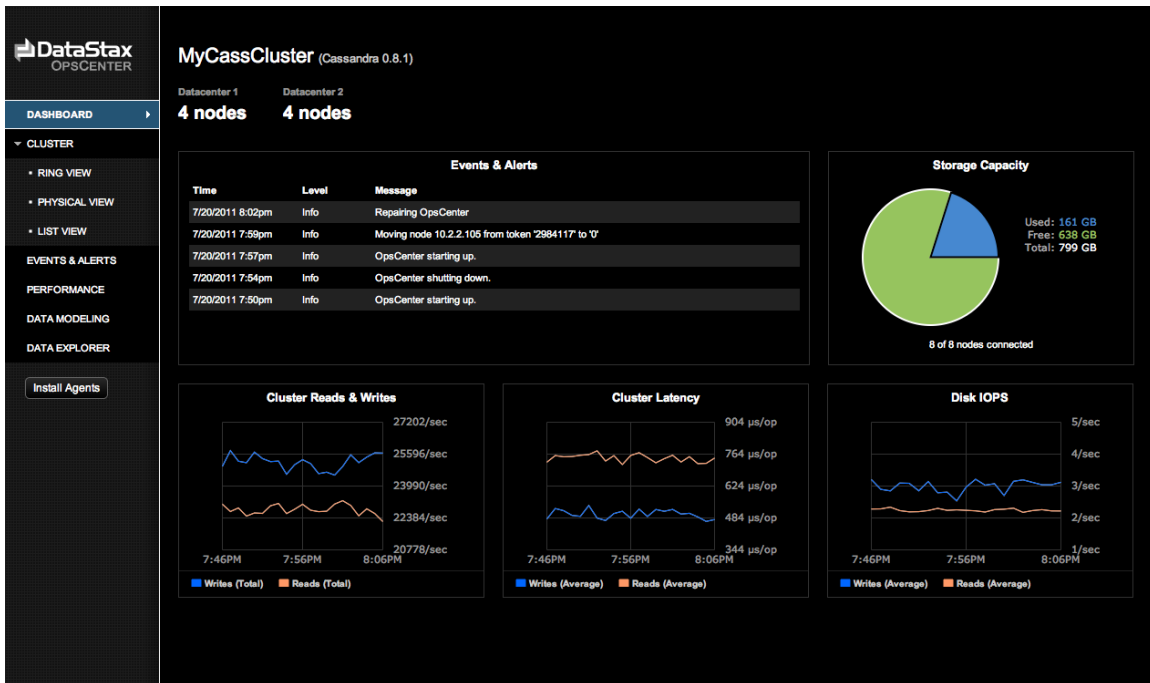


Figure 10: OpsCenter dashboard

Analytic operations also can be monitored and controlled from within OpsCenter Enterprise:

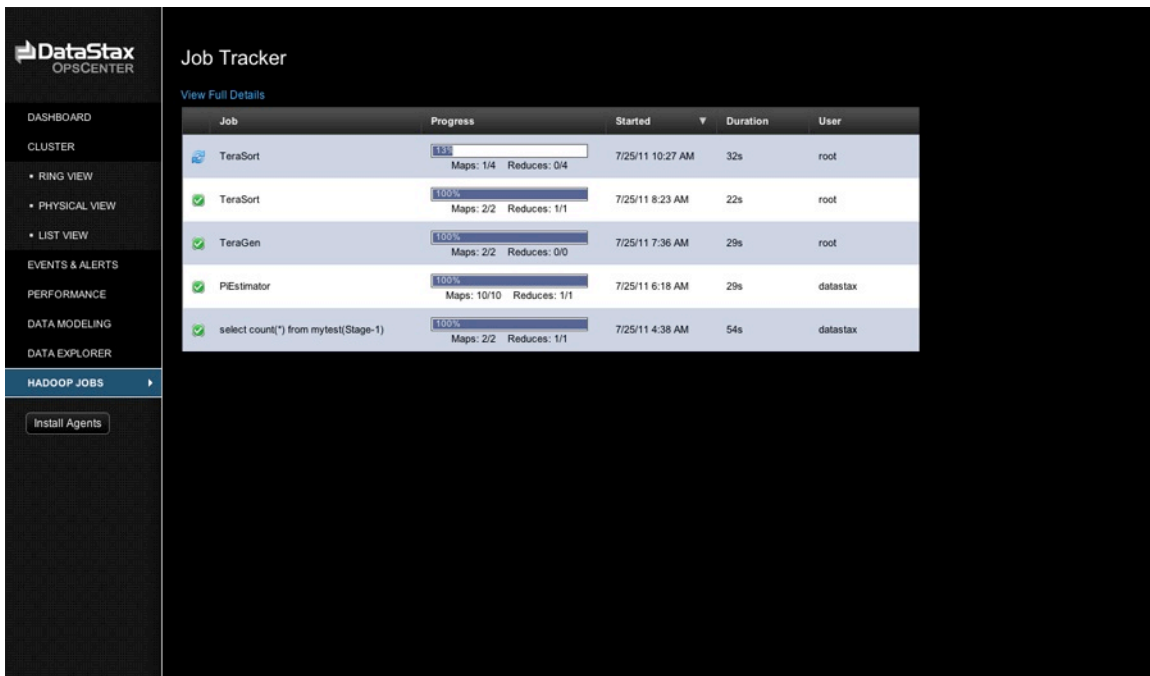


Figure 11: OpsCenter analytic operations monitoring

Enterprise Production Support and Services

Big data situations often require fast access to skilled expertise. DataStax Enterprise includes experienced production support and consultative services from Cassandra experts. You can choose the right production support package for your business needs, including rapid response service level agreements, and consultative help.

Additionally, DataStax offers professional big data training on Cassandra and Hadoop, with classes offered in many major cities as well as on-site for corporations that need many staff members trained at once.

Conclusion

Big data isn't just hype – and it's much more than a buzz phrase. Today, companies across industries are finding they not only need to manage increasingly large data volumes in their real-time systems, but also analyze that information so they can make the right decisions – fast – to compete effectively in the market.

Modern businesses looking for a solution to handle their big data easily and effectively should consider DataStax Enterprise. Its scale-out architecture comfortably scales into the petabytes data range, while offering high performance for reads and writes no matter the data volume size.

DataStax Enterprise also differs from competitors by providing a single integrated database platform that smartly manages real-time, analytic, and enterprise search data. Additionally, DataStax Enterprise does all of this at a fraction of the cost charged by traditional RDBMS vendors.

To find out more about DataStax Enterprise and download the software, please visit www.datastax.com or email info@datastax.com.

About DataStax

DataStax offers products and services based on the popular open-source database, Apache Cassandra™ that solve today's most challenging big data problems. DataStax Enterprise combines the performance of Cassandra with analytics powered by Apache Hadoop and enterprise search with Apache Solr, creating a smartly integrated, big data platform. With DataStax Enterprise, real-time, analytic, and search workloads never conflict, giving you maximum performance with the added benefit of only managing a single database.

The company has over 140 customers, including leaders such as Netflix, Disney, Cisco, Rackspace and Constant Contact, and spans verticals including web, financial services, telecommunications, logistics and government. DataStax is backed by industry leading investors, including Lightspeed Venture Partners and Crosslink Capital and is based in San Mateo, CA.

For more information, visit www.datastax.com.