

BIG DATA

The
Definitive
Guide to

TABLE OF CONTENTS

WHAT IS BIG DATA?	3
A Small History of Big Data	4
Structured vs Unstructured vs Semi-structured Data	5
THE ROLE OF BIG DATA IN TODAY'S ENTERPRISE	7
Important Considerations for Big Data	8
BIG DATA PLATFORMS AND APPLICATIONS	9
Data Storage and Management Tools	9
Apache Hadoop	9
NoSQL Databases	10
Hadoop or NoSQL?	11
Other Big Data Applications	11
Splunk on Nutanix	12
The Internet of Things	12
INFRASTRUCTURE CONSIDERATIONS FOR BIG DATA	13
Physical or Virtual?	13
Performance of Server Virtualization	14
Which Hypervisor Should You Choose?	14
Nutanix AHV	14
Other Infrastructure Considerations	14
Which Hardware Should You Choose?	15
Build Your Own	15
Converged Infrastructure	15
Hyperconverged Infrastructure	16
A Better Way to Address the Infrastructure Needs of Big Data	16
BIG DATA AND THE NUTANIX ENTERPRISE CLOUD	17
A Hyperconverged Architecture Based on Big Data Principles	18
Nutanix Enterprise Cloud Architecture	18
Distributed Storage Fabric	18
Nutanix Availability, Data Protection, and Disaster Recovery	20
AHV	21
Managing Your Big Data Environment	21
Nutanix Prism	22
One-Click Management	22
Full REST APIs	22
Top Ten Advantages of a Nutanix Enterprise Cloud for Big Data	22
Online Retailer Deploys Nutanix to Accelerate Elasticsearch	23
Deploying Hadoop on Nutanix	24
Deploying Splunk on Nutanix	25
About Nutanix	26

If you ask people in IT what big data is, chances are you'll get a variety of answers. The answer may depend on an individual's job function and the type and size of their company, but there's also a lot of uncertainty about what constitutes big data, what the benefits are, and how to take advantage of it.

The whole idea of big data is a relatively new addition to the IT lexicon. Mainstream dictionaries have only recognized the term in the last few years. These dictionary definitions provide a useful starting point for thinking about big data.

BIG DATA:

Data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges.

Oxford English Dictionary, 2013

BIG DATA:

An accumulation of data that is too large and complex for processing by traditional database management tools.

Merriam-Webster, 2014

Both mention size, which is certainly a part of the big data equation. However, the absolute size of the dataset may not be the most important factor to consider. The rate at which data is generated or acquired is an important corollary to size. For many companies, the rate of ingest, the speed with which data must be processed, and the ability to extract value from the data are currently bigger challenges than absolute dataset size.

Data complexity is another factor to consider. Many people think of big data primarily as consisting of "unstructured" data from the web and social sources, machine logs, or Internet of Things (IoT) sensor and tracking output. This type of data may not be amenable to conventional database methods. Big data projects frequently include data of multiple types from both traditional and nontraditional sources. Data scientists may incorporate data from relational databases, such as OLTP systems, along with unstructured data sources. The situation becomes complex quickly.

As important as these elements are, for many of us in IT, the central challenge of big data comes down to how best to tackle the infrastructure needs of big data projects—i.e., the "significant logistical challenges" recognized by the OED and the inability to process big data using "traditional database management tools" pointed out by Merriam Webster.

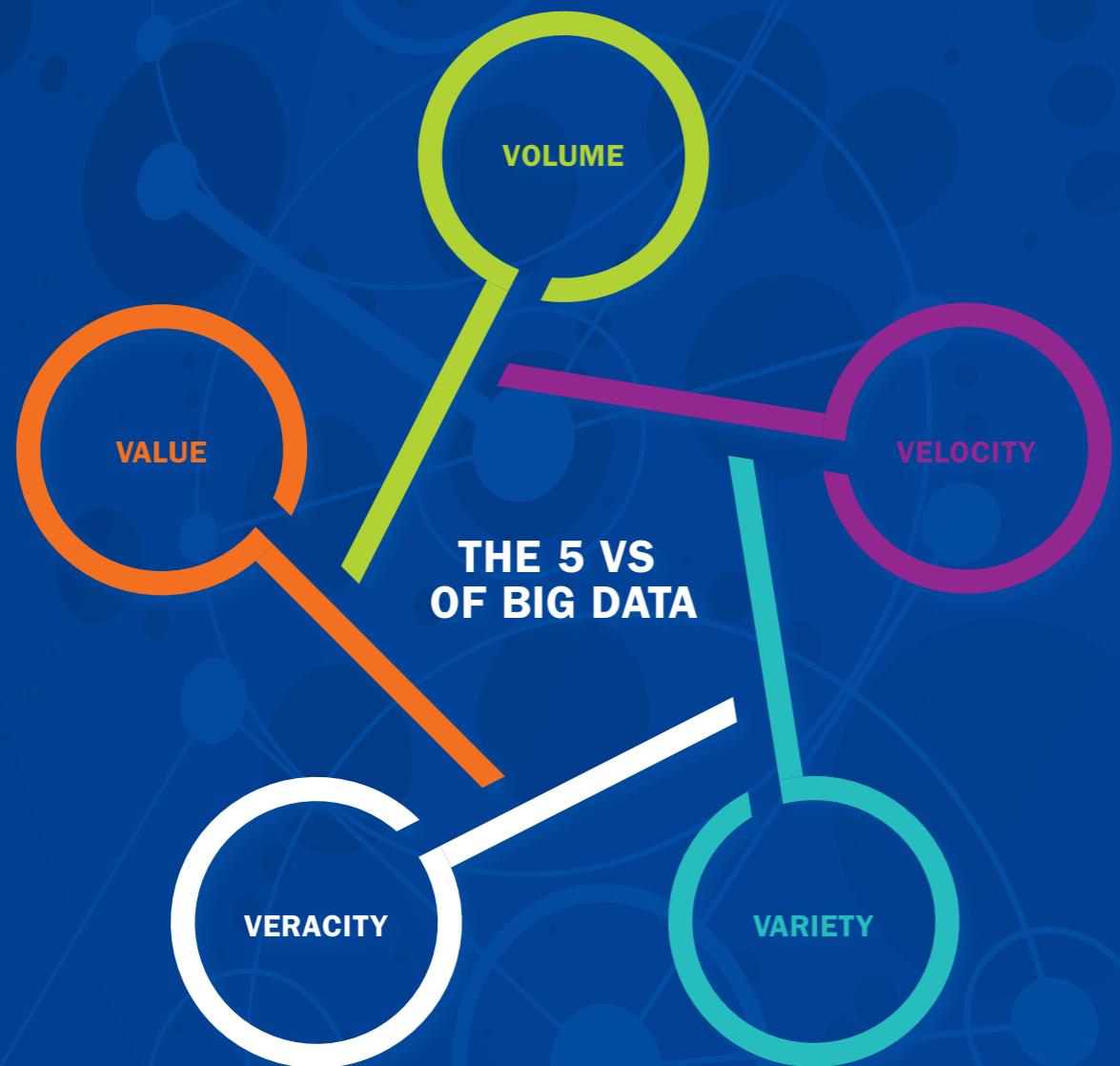


These aspects of big data are often characterized as the “**3 Vs**” – coined in a META Group white paper (now a part of Gartner Group) by Doug Delaney in 2001:

- **VOLUME** of data
- **VELOCITY** at which data must be processed
- **VARIETY** of data types

Increasingly, the 3 Vs have been extended to include:

- **VERACITY** of the data or how reliable and useful it is
- **VALUE** of the data or how value is extracted from it



¹ Oxford English Dictionary; definition added in 2013.
² Merriam-Webster Dictionary; definition added in 2014.

When you’re getting started with a big data project you need infrastructure that allows you to start small for proofs of concept and testing different scenarios. Infrastructure should also be flexible so you can switch workloads as needed. You The Path to Hyperconverged Infrastructure for the Enterprise might try many different big data workloads before settling on the best options for your business.

Finally, you want infrastructure that will scale with the growth of your datasets and ingest rates. Big data projects can rapidly expand as more workloads are added. If the infrastructure stumbles, then the project may lose crucial momentum.

Infrastructure for big data translates to more—and possibly different—hardware choices as well as new software tools that have to be selected and supported. A large part of this book is devoted to architecting infrastructure for big data projects.

A SMALL HISTORY OF BIG DATA

By the end of the 1990s, it was pretty clear that the traditional approach to business intelligence (BI) was no longer getting the job done. Datasets were limited, reporting was too slow, and results were often backward-looking. The “digitalization” of customer interactions and the rise of social media were creating a wealth of new and potentially valuable information that was difficult to incorporate into existing systems.

Around the same time, Google and other web leaders were pioneering new methods for managing the huge quantities of data created by their operations. In 2004, Google described a process called MapReduce that stores and distributes data across multiple servers, using a parallel processing model to process large amounts of data quickly. (MapReduce is in fact integral to the distributed architecture used by the Nutanix enterprise cloud platform.) Apache Hadoop is a well-known implementation of MapReduce.

This was the birth of the modern big data movement. An ecosystem of tools and applications has grown up around Hadoop, and a wide variety of new types of databases and analytics tools is emerging to help make sense of the ever-growing mountain of data. Rather than running on specialized computing hardware, most of these tools are designed to run on scale-out clusters of commodity hardware.

As the practice of big data has grown, it has been further divided into three major areas:

- Analytics that aggregate and analyze data from multiple sources, usually in batch mode
- Real-time or streaming analytics that operate on data as it is received
- Search tools for finding desired information or objects such as documents, videos, etc.

Each of these comes with its own unique requirements.

In discussions of big data, you'll often hear references to structured, unstructured, and even semi-structured data.

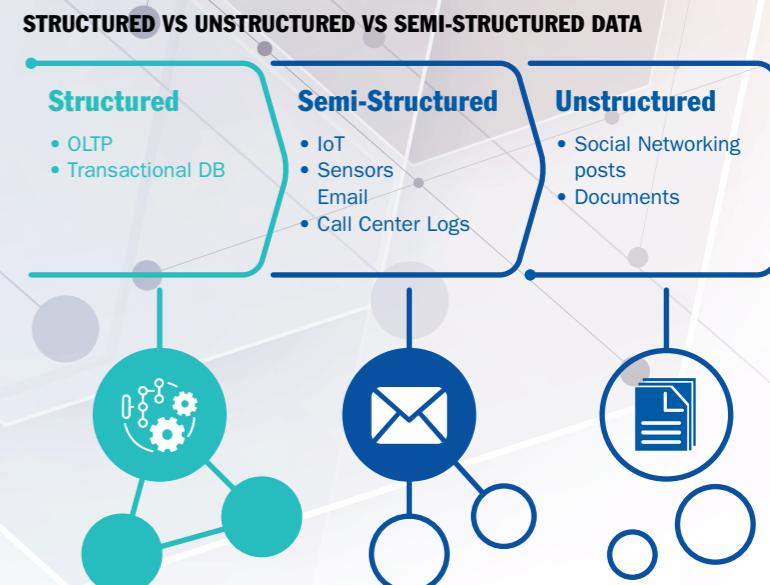
Structured data is the type of well-organized data you find in a typical relational database with regular, well-defined and typed fields such as customer name, account number, phone number, product purchased, date of purchase, etc. Structured data is managed most frequently using relational databases and analyzed with traditional business intelligence tools.

Unstructured data lacks the organization of structured data. Facebook posts, Tweets, and blog posts that make up so much interesting social content are examples of unstructured data. Since this type of data is not easily managed using a conventional relational database, extracting usable information from it requires a different set of tools. For example, big data tools are often used to extract customer sentiment from social media data.

Semi-structured data lies somewhere in the middle—it has a degree of structure, but lacks the well-defined schema of structured data. Not everyone agrees that this is a separate category, preferring instead to lump it in with unstructured data.

An email message is one example of semi-structured data. It includes well-defined data fields in the header such as sender, recipient, and so on, while the actual body of the message is unstructured. If you wanted to find out who is emailing whom and when (information contained in the header), a relational database might be a good choice. But if you're more interested in the message content, big data tools, such as natural language processing, will be a better fit.

Beyond recognizing the difference between types of data, perhaps the most important thing to know is that unstructured data is growing at a much more rapid rate than structured data. About 90% of corporate data is unstructured, and the gap between the amount of structured and unstructured data is likely to widen further. This is why big data tools that can help make sense of unstructured data are so important.



BIG DATA BANDWAGON



If your company hasn't jumped on the big data bandwagon yet, it's likely that you're trying to plan a path forward. Almost 70% of Fortune 1000 firms rate big data as important to their businesses; over 60% already have at least one big data project in place.³ All types of organizations are putting big data to work to improve competitiveness, enhance the customer experience, gain new insights, and achieve better outcomes:

THE ROLE OF BIG DATA IN TODAY'S ENTERPRISE



- **Retailers** use big data to enable advancements such as real-time pricing, product placement, and micro-targeted advertising. In a well-publicized—and controversial—example of this, Target is able to determine when a customer is pregnant based on buying patterns alone.
- **Law enforcement** uses big data to predict where crimes are most likely to occur. The Chicago police department is using data to identify those who are likely to be involved in shootings and reaching out to them proactively to reduce crime.
- **Manufacturers** use big data to decrease time to market, improve supply planning, stock depots, increase product quality, and better forecast product demand.
- **Farmers** use big data to enable “precision agriculture,” maximizing farm yields while minimizing inputs of water, fertilizers, and pesticides. Drones are now often used to gather aerial data.
- **Utilities** gather data from millions of smart meters installed at customer sites to improve energy forecasting and respond more rapidly to changes in supply and demand. Moving from monthly or annual readings to daily or hourly inputs gives utilities a wealth of new data to base decisions on.
- **IT teams** gather data from thousands of devices and applications and use software such as Splunk to identify problems, improve security, and streamline operations.
- **Security teams** at companies of all sizes and across all industries use big data to analyze IT environments, looking for data breaches and other suspicious activity.

Makers of more traditional business intelligence tools are working to enhance existing tools—and adding new tools—to accommodate big data. At the same time, whole new classes of software, such as data discovery, are emerging to make it possible to explore datasets in real time. Rather than forcing you to define what questions are of interest up front, these tools allow you to discover correlations in data you weren't necessarily looking for.

³ New Vantage Partners 2016 Executive Survey

⁴ Inside American Express' Big Data Journey, Forbes, April 2016.

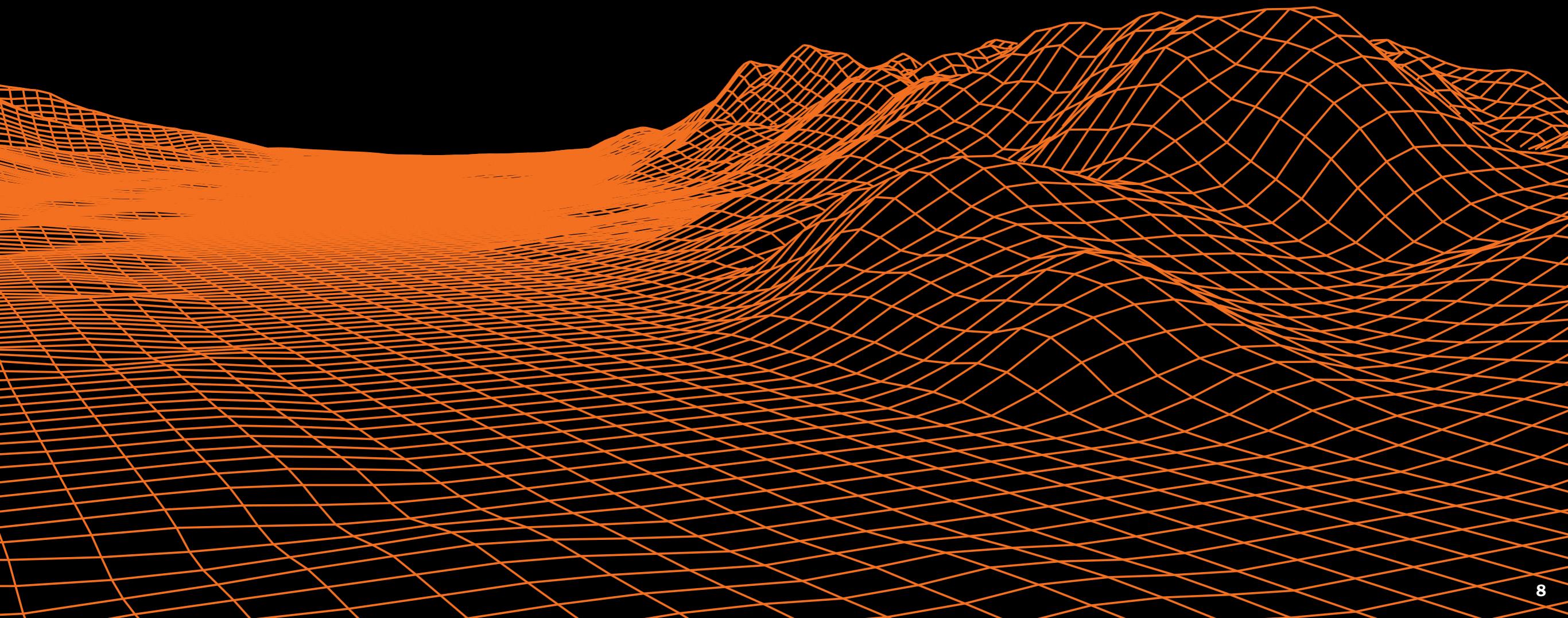
Some big data and analytics companies are focused on refining solutions for particular industries, so make sure you understand the state of the art for your industry.

IMPORTANT CONSIDERATIONS FOR BIG DATA

The big data landscape is changing rapidly, demanding maximum flexibility from anyone embarking on a big data project. But there are some important lessons you can learn from companies that have gone before you. For example, American Express⁴ is already several years into its big data transformation and has been public about sharing what it has learned:

- **Your organization must be prepared to adapt.** New technologies may require changes to both organizational structure and culture.
- **Consider the need for new talent (or help).** You may need to recruit new talent with the necessary big data and analytics skills. In lieu of that—given the shortage of these skills in the job market right now—you may want to choose vendors that are vested in your success and/or engage the right consulting help.

- **Make iterative improvements.** Don't expect to transform your operations and organization all at once. Identify the low hanging fruit, make changes based on what you learn, and build on your successes to do it again.
- **Embrace trial and error.** Foster an environment where experimentation is possible. You may learn as much or more from failures as successes.
- **Focus on speed and agility.** Faced with both a dynamic business environment and a rapidly evolving big data landscape, you need to be able to extract information quickly, and adapt your big data and analytics environment to accommodate business changes and new tools.



BIG DATA PLATFORMS AND APPLICATIONS

If you're new to the big data field, there are several things you will notice almost immediately. One thing that is decidedly different is the amount of open source software in wide use. Many of the tools you hear the most about, such as Hadoop, are open source. Open source is a good way to foster rapid innovation in a rapidly evolving field. As a corollary to that you'll notice that there are a vast number of big data tools from which to choose.

In addition to strictly open source offerings, well-known systems vendors such as IBM, HP, Oracle, and Dell have products for the big data market. Many of these combine open source tools and proprietary solutions. Major cloud providers including Amazon and Google provide big data services in the cloud—often using versions of open source tools.

There are hundreds of big data tools and services. This section, though by no means comprehensive, will acquaint you with some of the most prominent ones.

DATA STORAGE AND MANAGEMENT TOOLS

This book assumes you are already familiar with traditional relational databases such as Oracle and Microsoft SQL Server. As already noted, these databases and related tools (and/or the data they contain) may become a part of your big data efforts.

This section focuses on the relatively new tools available for storing, organizing, and managing the large volumes of data—especially unstructured data—in a big data environment:

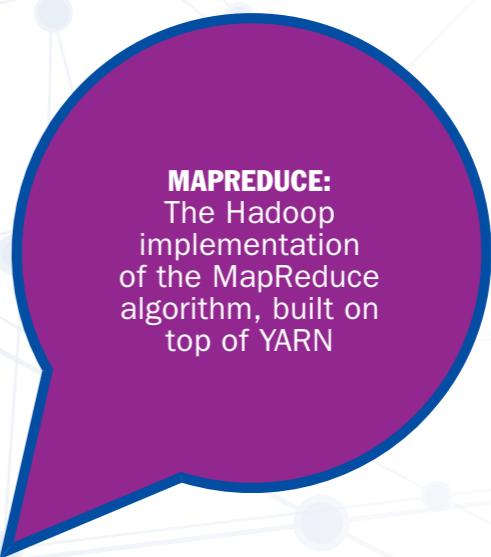
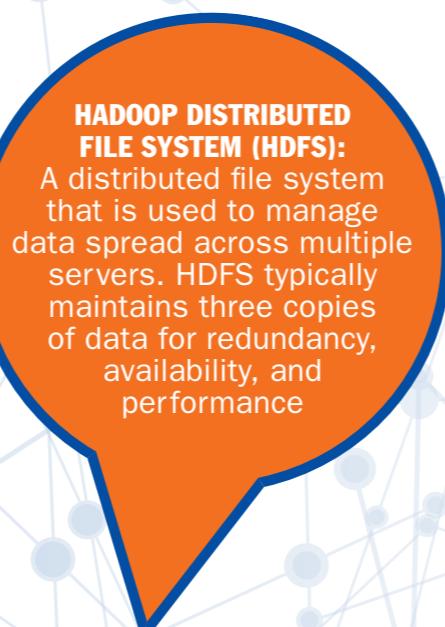
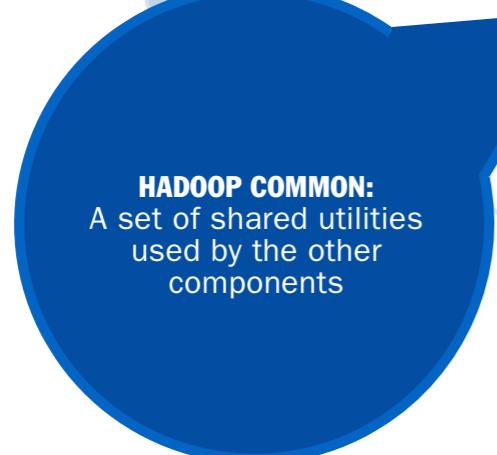
- Apache Hadoop
- NoSQL databases

APACHE HADOOP

As with other open source software, a variety of vendors offer Hadoop distributions. The best known of these are Cloudera, HortonWorks, and MapR. These vendors make Hadoop easier to consume, provide support, and may offer other features and services on top of the base Hadoop software.



Apache Hadoop is not a database. It's a distributed platform for processing large datasets in parallel. Hadoop runs on clusters of machines and can scale up to thousands of nodes.
It consists of FOUR PRIMARY COMPONENTS:



NOSQL DATABASES

New types of databases—collectively dubbed NoSQL databases—have emerged over the last 10 years. The name “NoSQL” started out to mean “non-SQL” but today is widely interpreted as “Not Only SQL” since some of the NoSQL offerings now support structured query language or SQL as used in the relational database world.

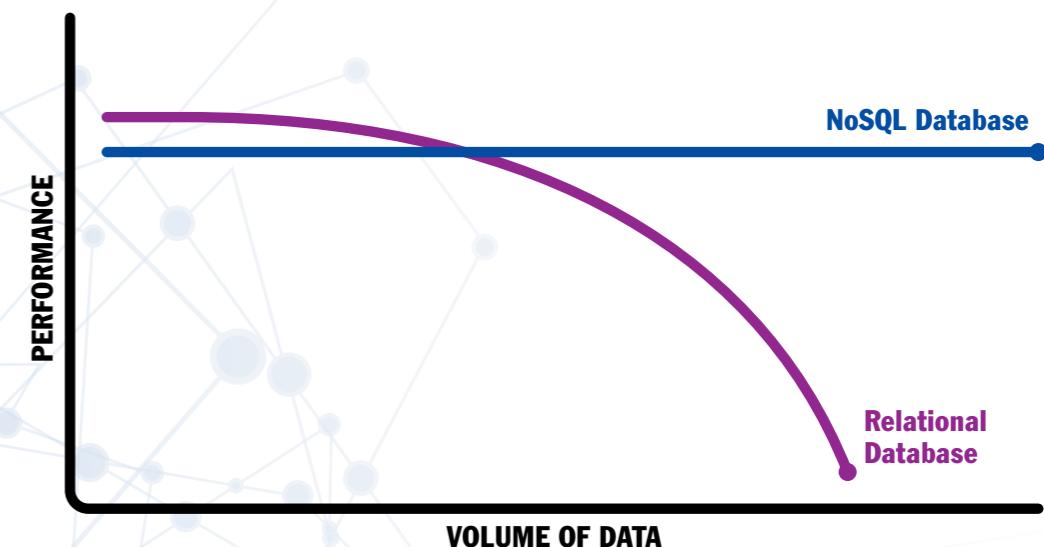
NoSQL databases are perfect for workloads that need a predictable response time no matter how much data they contain. Many companies use NoSQL databases as the backend for mobile applications with custom-coded front ends.

While the particulars of NoSQL databases differ widely, they share a number of characteristics that make them well suited for big data, especially in comparison to traditional relational databases. These databases are built to handle massive amounts of data by scaling out across many servers and are highly available by design. In many cases, they are designed for “eventual consistency” (prioritizing availability over the strict consistency of relational databases) and maximum speed. Where relational databases are designed for structured data, NoSQL databases are much better suited for the unstructured and semi-structured data types prevalent in the big data world.

In addition to the core components, an ecosystem of open source projects related to Hadoop has been created. Infrastructure-related offerings include:

- **Ambari:** A web-based tool for provisioning, managing, and monitoring Apache Hadoop clusters
- **ZooKeeper:** A coordination service for distributed applications
- **Cassandra:** A NoSQL database (see following section)
- **HBase:** A NoSQL database (see following section)

Nutanix has incorporated and customized several of these components, including Cassandra, Zookeeper and MapReduce.



HADOOP OR NOSQL?

Should you choose Hadoop or NoSQL for your big data project? The two have a number of attributes in common. Both use a scale-out architecture, provide high availability, and are good at handling large amounts of data and a variety of data formats. And, as you may have noticed, the two aren't mutually exclusive. Some NoSQL databases can be incorporated into the Hadoop framework, and organizations increasingly deploy both to serve different big data use cases.

Hadoop is most often deployed for processing and analyzing large amounts of data in parallel in a batch mode of operation. You submit a job and wait for your result to be computed which could take minutes or hours.

NoSQL is better suited for real-time, interactive access to data including realtime analytics, fraud detection, intrusion detection, and other streaming applications.



OTHER BIG DATA APPLICATIONS

There are a variety of big data applications and application frameworks designed to support activities such as search, fraud detection, intrusion detection, and more. Applications are now driven by a need for continuous innovation. Businesses are using these tools to facilitate frequent—weekly or even daily—software adjustments, making them a key part of the ongoing digital transformation of business. Many of these may be complementary to existing data warehouse and business intelligence applications.

Popular open source applications include:

For Hadoop:

- **Hive.** A data warehouse infrastructure that provides data summarization, query, and analysis. Originally developed by Facebook, Hive is now widely used by other companies including Netflix
- **Mahout.** A scalable machine learning and data mining library that incorporates a broad and growing library of algorithms for a wide range of uses such as clustering, classification, and collaborative filtering, as well as math and statistics
- **Spark.** A compute engine with a programming model that supports a wide range of applications including streaming. In addition to Hadoop, Spark runs on Mesos, standalone, or in the cloud. Data sources can include HDFS, Cassandra, HBase, and S3
- **Drill.** An SQL query engine that supports a wide range of data sources in addition to HDFS, including HBase, MongoDB, MapR-DB, HDFS, MapR-FS, Amazon S3, Azure Blob Storage, Google Cloud Storage, Swift, NAS, and local files. A query can join data from multiple sources

Other:

- **Lucene Core.** Search software that provides Java-based indexing and search, spellchecking, hit highlighting, and advanced analysis/tokenization. PyLucene is a port of Lucene core in Python.
- **Elastic.** A search server based on Lucene, designed to take data from any source and search, analyze, and visualize it in real time. Currently the most popular enterprise search engine
- **Solr.** A high performance search server built using Lucene Core. Second most popular enterprise search engine
- **Splunk.** A software system that captures, indexes, and correlates machine generated data. It is used widely in IT for a variety of purposes including improving IT operations, compliance, security, and auditing. Because it operates on any kind of machine data, it also finds uses in other industries such as oil and gas, manufacturing, and to support the Internet of Things



SPLUNK ENTERPRISE



RUNNING SPLUNK ON NUTANIX

Splunk Enterprise is the leading software platform for unleashing the power of machine data gathered from IT infrastructure and equipment of all types. As the amount of machine data increases without bound, the job of the Splunk architect has become more challenging. Reducing complexity, improving data security, and eliminating bottlenecks are top priorities. Traditional IT infrastructure is ill-suited to address the needs of growing Splunk installations.

The Nutanix enterprise cloud platform takes the complexity out of managing infrastructure for Splunk, allowing Splunk experts to spend more time extracting insight from data. Nutanix allows Splunk to take full advantage of server virtualization without the limitations of other solutions. By ensuring data is accessed locally by all Splunk indexers, Nutanix eliminates the “I/O Blender” effect that can plague conventional infrastructure.

THE INTERNET OF THINGS

The buzz around applications for the Internet of Things (IoT) seems poised to surpass the general buzz in the industry around big data. Though IoT is just getting started, it is already creating big expectations in the enterprise, particularly in retail and manufacturing organizations. Because of the sheer amount of data that may need to be processed, big data methods are becoming a requirement.

In terms of data and analytics, IoT projects require a big data back end that is no different than other big data projects, so the information in this book is directly applicable. When everything is wired, there's no end to the data that can be tracked to create strategic advantage or simply to streamline processes and mundane tasks. Here are a few examples of IoT projects:

-  **Airplane manufacturers** instrument the latest planes with sensors to provide a large volume of data for every flight. This data is being used to improve servicing and detect problems before they occur.
 -  **Delivery services** use IoT to gather data on vehicles in order to improve fuel efficiency and reduce costs.
 -  Some **entertainment companies** and theme parks use intelligent wristbands that make it easier for customers to check in, order food, and make purchases. These wristbands allow customer activities to be tracked and monitored to detect trends and patterns.
 -  **Livestock producers** use IoT to keep better track of animals.
 -  **Farmers** can track exactly where a particular crop was grown, packed, shipped, and processed, allowing a suspect crop to be more easily pinpointed in the case of contamination.



INFRASTRUCTURE CONSIDERATIONS FOR BIG DATA

When thinking about the infrastructure needed to support a big data project, it's important to keep in mind the considerations mentioned at the end of chapter 2:

- Be prepared to adapt**
- Consider the need for new talent**
- Make iterative improvements**
- Embrace trial and error**
- Focus on speed and agility**

The process is likely to be an iterative one with some trial and error involved, and you need to consider how agile the infrastructure you choose will be. As you add or change processes and tools, you want to be able to adapt and scale with a minimum of effort and without the need to start from scratch. Server virtualization can simplify the process dramatically versus deploying big data software directly on bare metal.

PHYSICAL OR VIRTUAL?

Most of the infrastructure initially deployed for big data utilized bare metal deployments on physical servers with local storage. In part this was because many of the tools originated at companies such as Google, Amazon, and Facebook that had standardized on this type of infrastructure, and in part it was because bare metal was thought necessary to achieve maximum performance. Ten years ago, the state of the art in virtualization had not yet reached the high levels we have now.

Today it makes sense to deploy big data projects in virtual server environments whenever possible to take advantage of the well-known advantages of virtualization:

- Rapid virtual machine (VM) provisioning
- Simplified infrastructure management
- Improved availability
- Greater flexibility
- Increased efficiency and decreased cost



WHICH HYPERVISOR SHOULD YOU CHOOSE?

A wide variety of hypervisors are available. The three options are supported with Nutanix systems: VMware vSphere, Microsoft Hyper-V, and Nutanix AHV. Because most open source big data software is designed to run on Linux, support for Linux as a guest operating system is an important consideration.

- **VMware vSphere** remains the most widely deployed hypervisor and in all likelihood you have staff with VMware expertise. It offers broad support for various versions of Linux. A drawback of VMware is the additional cost of licensing the software.

- **Microsoft Hyper-V** added support for Linux as of Windows Server 2012 R2. If you're already running Windows Server and/or Hyper-V, it may be worth considering. Otherwise, the Linux support may be limited compared to other options.

- **Nutanix AHV** is a next-generation hypervisor that integrates closely with the Nutanix hyperconverged architecture and is included at no additional charge with a Nutanix purchase. Based on the Linux KVM hypervisor, AHV offers broad Linux support.

Virtualization gives you the ability to deploy and redeploy resources rapidly, which fits well with the needs of a big data project. If a given infrastructure configuration turns out to be non-optimal, you can adjust it easily, without making physical changes, and you can more easily try different node counts and storage configurations to find the best fit. Some Nutanix customers run VDI during the day and run Hadoop on the same infrastructure at night.

PERFORMANCE OF SERVER VIRTUALIZATION

In the past, it was generally assumed that bare metal servers would always deliver better performance than virtualized servers. However, hypervisor designers have worked hard to minimize overhead and today a virtual machine delivers performance very close to that of a similarly configured physical server. In some cases, virtualized environments have even been shown to deliver better performance than physical servers.⁵

Of particular interest, VMware found that a 32-node Linux cluster running Hadoop delivered performance close to bare metal when a single VM per node was configured. Running two VMs per node matched bare metal performance, and running four VMs per node exceeded bare metal performance. They concluded that multiple VMs were able to more efficiently use memory and disk resources. The paper, written in 2013, also concluded that this advantage would likely become greater as servers continued to get more powerful.⁶

Your mileage may vary, but it's probably safe to conclude that the advantages of virtualization outweigh any performance impact—and by careful sizing and configuration, you may achieve better performance than a bare metal deployment.

OTHER INFRASTRUCTURE CONSIDERATIONS

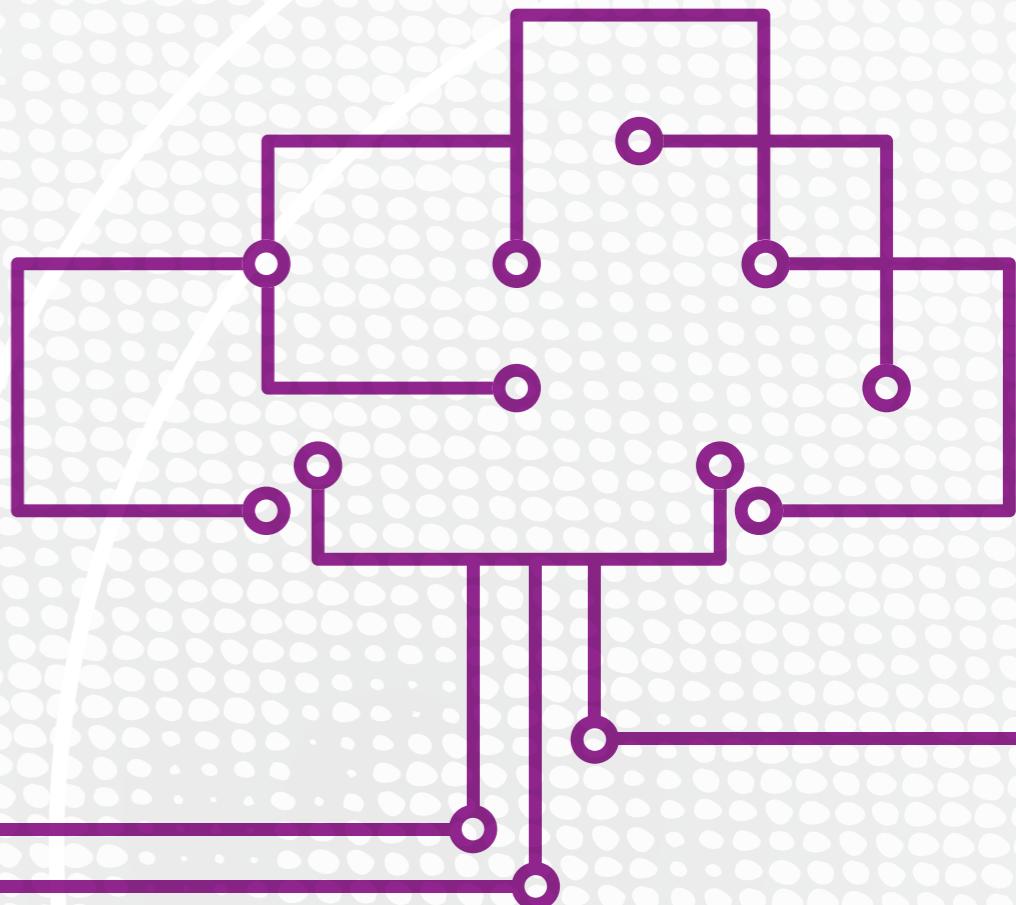
Some organizations deploying infrastructure for big data—whether on physical or virtualized hardware—utilize conventional IT infrastructure designs with servers and storage arrays separated by storage networks. This is seen as a way to reduce storage overhead and increase infrastructure flexibility. However, this yields mixed results for a number of reasons:

- **Lack of data locality.** Hadoop and other big data frameworks are designed around the idea that the data being operated on will be local to the server performing the operation. Introducing networked storage creates delays and bottlenecks.

- **I/O blender effect.** Especially in virtual environments, the stream of I/O requests coming from demanding big data jobs results in a random mix of I/O that can defeat storage read-ahead algorithms, making it difficult for a central storage system to operate efficiently. This can result in further I/O delays.

- **Mixed I/O needs.** Some big data jobs require high streaming performance while others create random I/O. Most storage arrays don't handle both types equally, especially when simultaneous.

As a result, many of the expected benefits of conventional infrastructure fail to emerge. As illustrated in Figure 1, when you're configuring multiple targets per server, the virtual environment can become extremely complicated. The ability to rearchitect storage quickly in response to changing needs may only be marginally better than local storage, if at all.



⁵ Yes, virtualization is faster (sometimes) than native hardware, ZDNet, 2013.

⁶ Virtualized Hadoop Performance with VMware vSphere® 5.1, VMware, 2013.

WHICH HARDWARE SHOULD YOU CHOOSE?

A final decision you need to consider when initiating a big data project is what hardware to choose.

Essentially, you have three choices:

- Build your own (BYO)
- Converged infrastructure (CI)
- Hyperconverged infrastructure (HCI)

BUILD YOUR OWN

The BYO infrastructure alternative is just what the name implies: your team makes all hardware selections, which could include servers, storage, and networks. You have to choose between:

- **Servers with internal storage.** A bare-metal big data deployment means you must manage hardware configuration tasks manually as described above.
- **Conventional infrastructure.** Choosing separate servers and storage can result in a significant increase in upfront planning and research as your team must evaluate each product separately and also assess how products will work together. You can deploy conventional infrastructure for either bare metal or virtualized infrastructure, keeping in mind the potential limitations as described earlier in this chapter.

In addition to all the decisions to be made, your team is also responsible for identifying and managing compatible firmware levels across many, many devices.

CONVERGED INFRASTRUCTURE

There are a variety of converged infrastructure solutions on the market that prepackage servers, networking, and storage to make conventional IT infrastructure easier to consume. Usually, these solutions are deployed in conjunction with virtualization. If you've settled on conventional architecture, a CI solution can speed up deployment, but it won't eliminate the limitations of traditional infrastructure for big data projects.

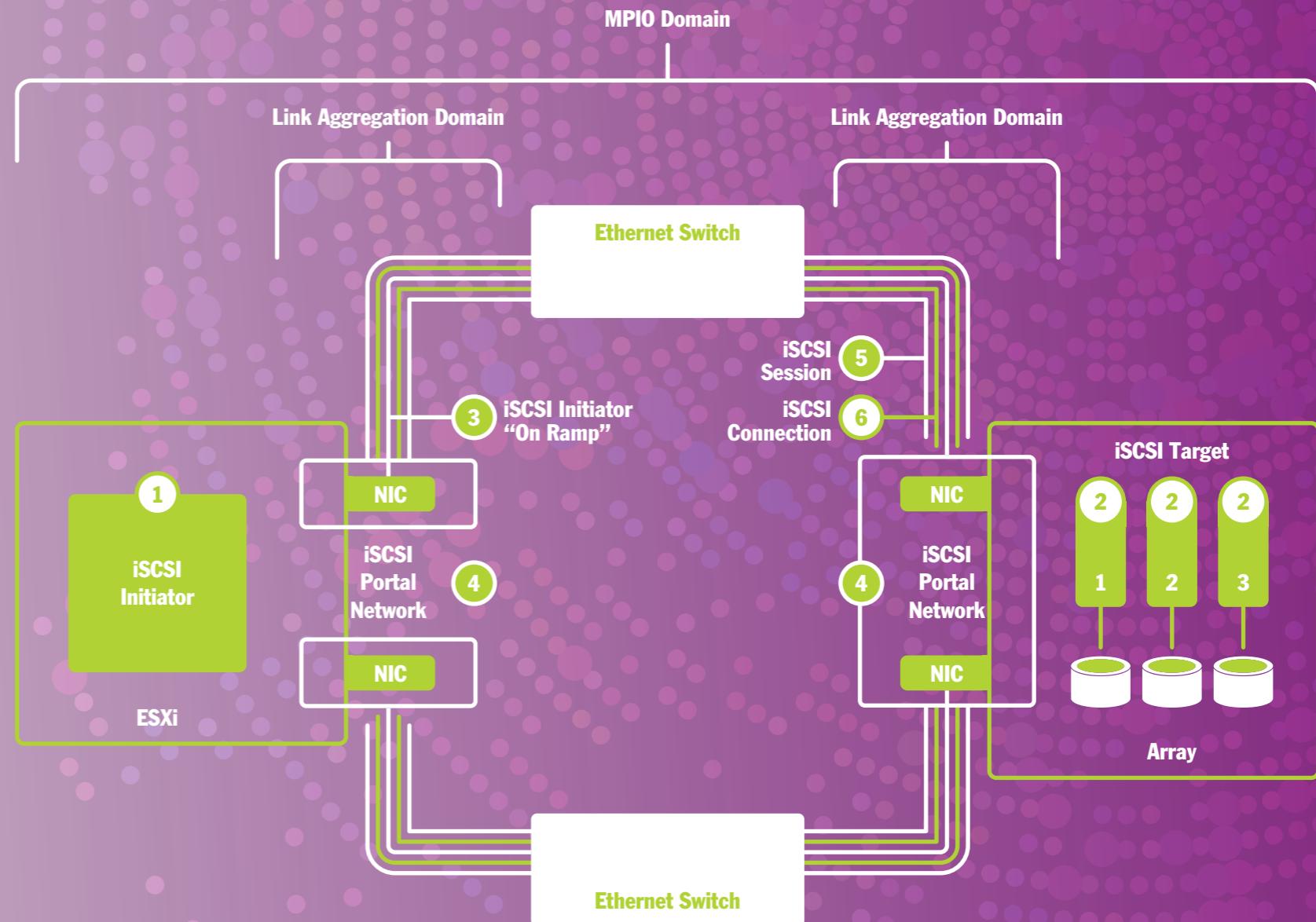


Figure 1. Complexity of conventional infrastructure for big data

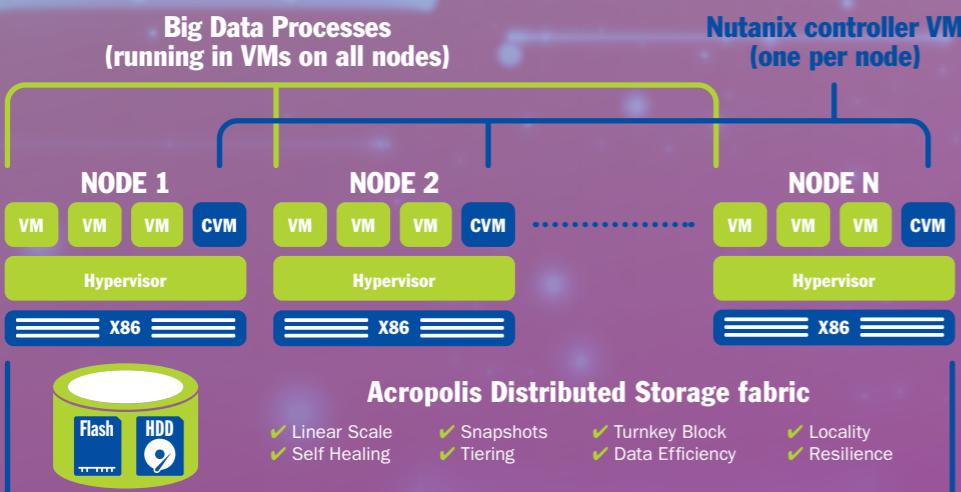
HYPERCONVERGED INFRASTRUCTURE

The final infrastructure option converges servers, storage, and virtualization in easy-to-deploy building blocks. HCI offers the potential to combine the virtues of bare metal deployment with the flexibility of virtualization for ease of management, scaling, and availability. HCI solutions typically use a hybrid approach to storage combining flash solid-state drives (SSDs) and hard disk drives (HDDs) in each node.

There are a few things to keep in mind when considering HCI:

- **Data locality.** As you've seen, big data software is built to expect local data. An HCI solution that provides data locality automatically satisfies this requirement, improving performance while simplifying data management tasks. However, typical HCI solutions don't provide data locality.
- **Data reduction and erasure coding.** Big data solutions such as Hadoop usually create three copies of all data for redundancy and performance, but this can consume a lot of extra storage. A correctly architected HCI solution includes data reduction technologies, such as compression, that can reduce storage consumption without affecting performance. Erasure coding can further reduce storage consumption for cold data without affecting redundancy or performance.
- **High availability.** A virtue of most big data software frameworks is that they are resilient to failure. However, you still need underlying infrastructure that is as resilient as possible to ensure that operations aren't disrupted. Features such as non-disruptive upgrades and self-healing can improve availability versus bare metal deployments. Failed virtual machines can be automatically restarted to minimize the impact to big data jobs.

Nutanix believes that a hyperconverged infrastructure approach can significantly reduce time to value and lower total cost of ownership (TCO), putting you on a faster path to big data success.



A BETTER WAY TO ADDRESS THE INFRASTRUCTURE NEEDS OF BIG DATA

Wouldn't it be great if a big data infrastructure solution could:

- Simplify and accelerate infrastructure provisioning for big data projects?
- Combine the flexibility and ease-of-management of virtualization with the performance of bare metal?
- Make on-premises infrastructure as easy to consume as public cloud services?
- Reduce the cost of purchasing, installing, and managing infrastructure?
- Increase availability and eliminate planned downtime?
- Allow you to quickly redeploy existing infrastructure for big data tasks during off hours such as weekends and evenings?

The unique, hyperconverged design of a Nutanix enterprise cloud makes all this possible and more.



BIG DATA AND

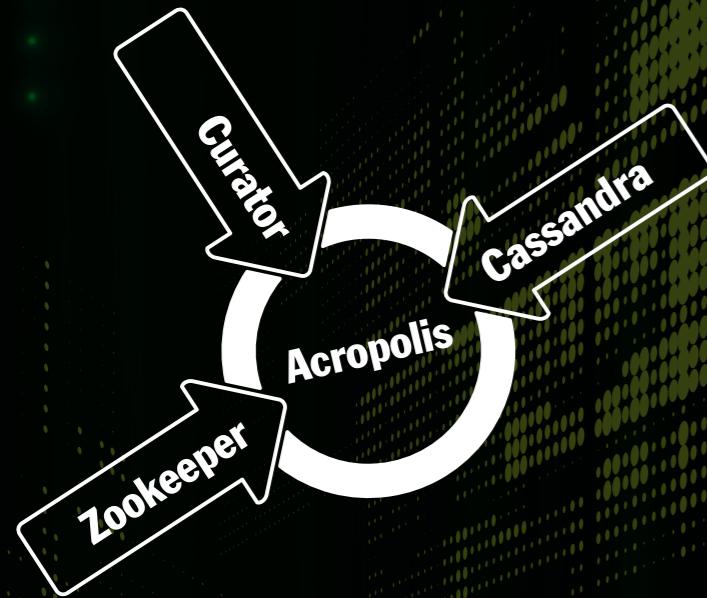
THE NUTANIX ENTERPRISE CLOUD PLATFORM

The Nutanix enterprise cloud platform is perfect for hosting big data applications of all types. In today's dynamic business environment, IT plays a larger role in businesses of all sizes and all industries. Nutanix was founded with the goal of simplifying IT infrastructure and reducing routine management, elevating IT teams to focus attention on critical applications and processes.

By combining industry standard hardware with innovative software—and focusing tirelessly on customer needs and customer service—Nutanix has created a new, simpler IT infrastructure designed to propel IT into the next era of computing. Nutanix enterprise cloud solutions offer the agility of the public cloud without sacrificing security, predictable costs, or high levels of service.

Because the Nutanix distributed architecture incorporates many of the same design principles as big data applications, including scale-out and data locality, it is a perfect fit for big data projects, enabling you to get big data clusters up and running more quickly at lower cost—and to flexibly scale and adapt your infrastructure as your needs change and evolve.

Nutanix delivers the flexibility and simplicity of server virtualization and shared storage without sacrificing the performance of bare-metal. Like bare-metal, Nutanix hardware building blocks combine compute and storage. Unlike bare metal, Nutanix simplifies the process of deploying and managing a large number of nodes. In addition, Nutanix uses a security development lifecycle to design security into Enterprise Cloud from the ground up.



A HYPERCONVERGED ARCHITECTURE BASED ON BIG DATA PRINCIPLES

One of the things that makes Nutanix such a good fit for big data and analytics applications is that it's built on many of the same design principles.

For instance, Nutanix utilizes a scale-out, shared nothing design that incorporates data locality for performance and replication factor for availability. Because Hadoop and other big data software frameworks use the same principles, they run very well on Nutanix infrastructure and benefit immediately from Nutanix architectural features.

To support scaling, Nutanix also incorporates many of the well-known algorithms from the big data open source movement. (In keeping with the open source tradition, these have been customized for our needs):

- **MapReduce** is used throughout the Nutanix operating software, distributing tasks uniformly across the cluster
- **Zookeeper** is used to keep track of Nutanix cluster configuration information
- **Cassandra** is used as a metadata store

As a result, Nutanix services scale-out in lockstep with your big data project.

THE NUTANIX ENTERPRISE CLOUD PLATFORM

Hyperconverged infrastructure is the foundation of the Nutanix enterprise cloud platform. Scale-out clusters of high-performance servers (nodes) contain processors, memory, and local storage. Clusters can scale from three nodes up to a very large number.

In addition to standard virtual machines, each Nutanix node runs a special controller VM (CVM) that acts as a storage controller for the drives on that node, providing data services for VMs running locally, as well as other nodes as needed.

In hybrid clusters, each node contains a combination of flash SSDs and hard disk drives (HDDs) for performance and capacity. All-Flash clusters are also available to support the most stringent real-time and streaming data projects.

Each node in a Nutanix cluster runs a standard hypervisor and runs virtual machines just like any other virtualized server. Nutanix provides support for VMware vSphere and Microsoft Hyper-V. Our own Nutanix AHV (Acropolis Hypervisor) provides a next-generation, native hypervisor that integrates closely with the Nutanix architecture and is included with your Nutanix purchase.

DISTRIBUTED STORAGE FABRIC

The Acropolis Distributed Storage Fabric (DSF) virtualizes local storage from all nodes in a cluster into a unified pool. DSF uses the local SSDs and HDDs from all nodes to store virtual machines and data.

DSF provides a number of features that are extremely valuable in big data environments:

Data locality. The data used by each virtual machine is preferentially kept on local storage on the node where the virtual machine is running.

Auto-tiering. Hot data is stored preferentially in flash, while cold data resides in the HDD tier. For optimal performance, data is automatically moved between tiers based on access patterns.

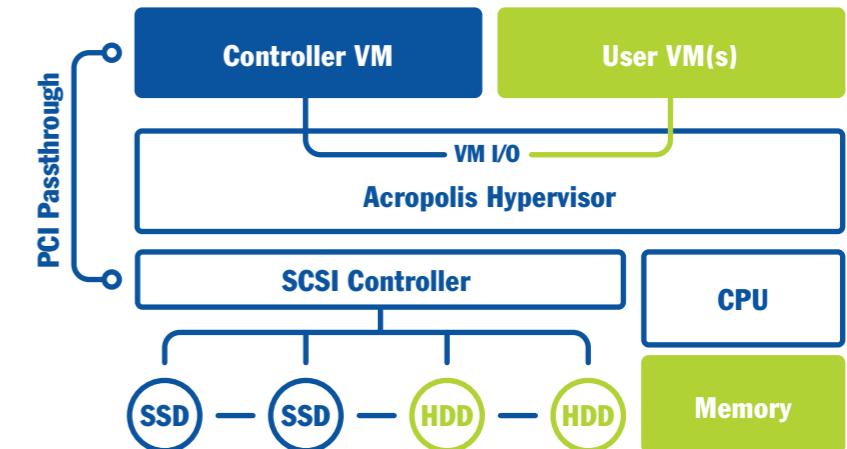


Figure 2. Nutanix node architecture

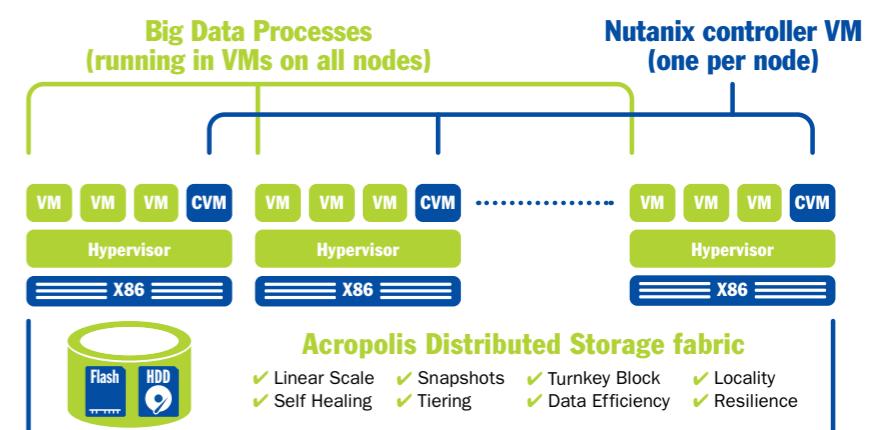


Figure 3. Acropolis Distributed Storage Fabric (DSF)

Auto-tuning. Nutanix delivers excellent performance for both random and sequential I/O— without the need for constant performance tuning, even when multiple and/or different workloads are running simultaneously.

Auto-leveling. Advanced algorithms ensure that data is balanced across nodes. Data from any virtual machine can utilize storage on other nodes when needed, eliminating the possibility of running out of storage on any single node and simplifying data management. If a particular node has more hot data than it has local SSD, available SSD capacity on other nodes is used.

Auto-archiving. A Nutanix cluster can include capacity-only nodes to increase available storage. When these nodes are deployed, cold data preferentially finds its way to them. In effect, this provides auto-archiving of cold data. Whenever data becomes active again, it moves back to the node where it is needed without administrator intervention. Because these capacity nodes don't run big data applications, there is typically no licensing associated with them. For instance, neither Cloudera nor Hortonworks charges licensing for these nodes.

Replication factor. A Nutanix enterprise cloud relies on a replication factor (RF) for data protection and availability. Either two or three copies of all data are maintained on different nodes within a cluster. Nutanix's patented erasure coding algorithm, ECX, can be used to reduce the storage overhead resulting from RF. EC-X operates on write-cold data, making it an excellent choice for write-once-read-many (WORM) workloads. Using EC can increase usable space in the cluster by up to 70%.

Self-healing. Replication factor also enables a Nutanix cluster to be self-healing. When a disk or node fails, full data redundancy is quickly and automatically restored to protect against additional failures. VMs are restarted on other nodes as necessary. A larger Nutanix cluster can withstand the failure of an entire four-node enclosure (referred to as a block). By architecting a big data application such as Hadoop to span multiple clusters, even a full rack failure can be sustained without interrupting operations.

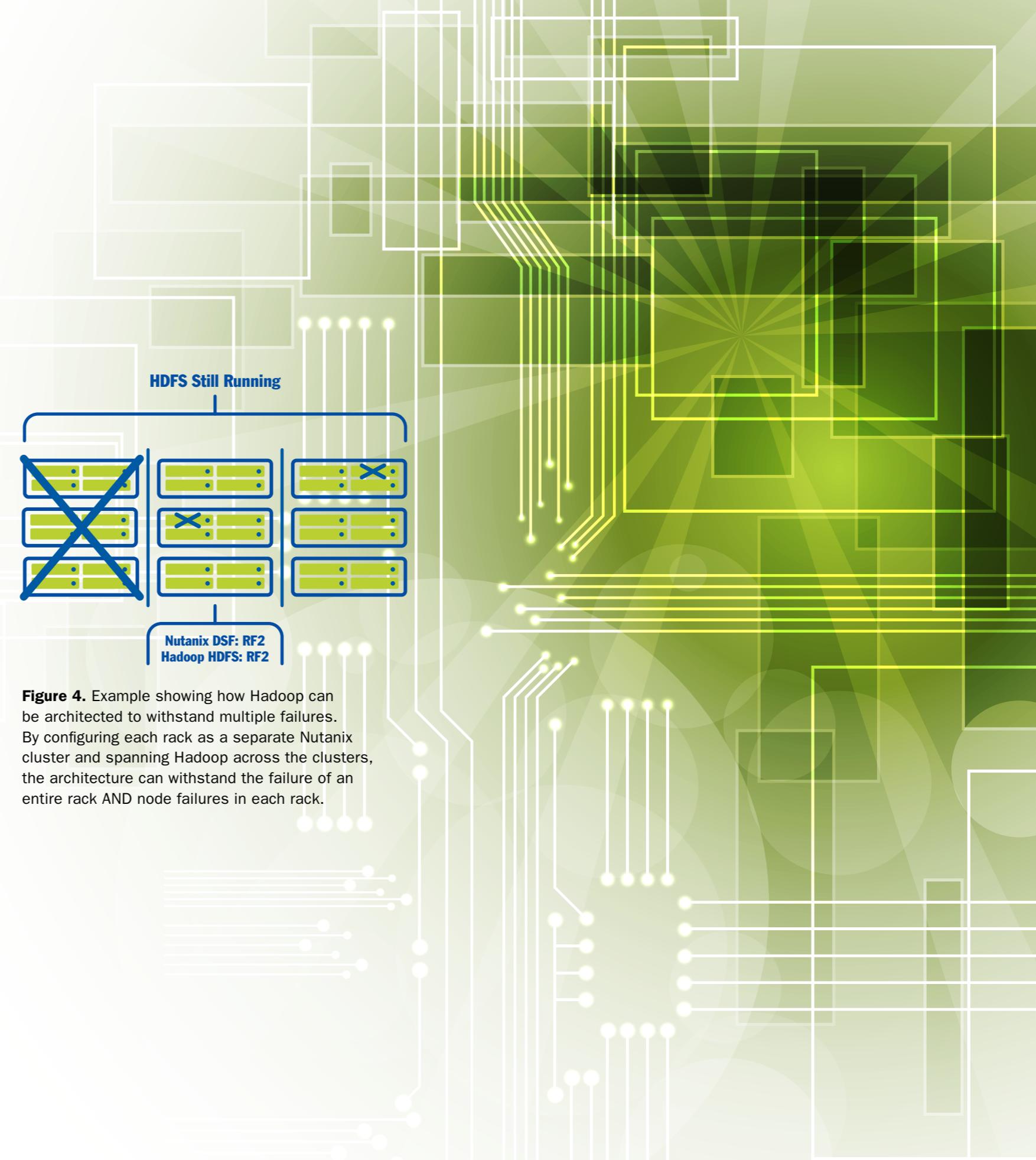


Figure 4. Example showing how Hadoop can be architected to withstand multiple failures. By configuring each rack as a separate Nutanix cluster and spanning Hadoop across the clusters, the architecture can withstand the failure of an entire rack AND node failures in each rack.

NUTANIX AVAILABILITY, DATA PROTECTION, AND DISASTER RECOVERY

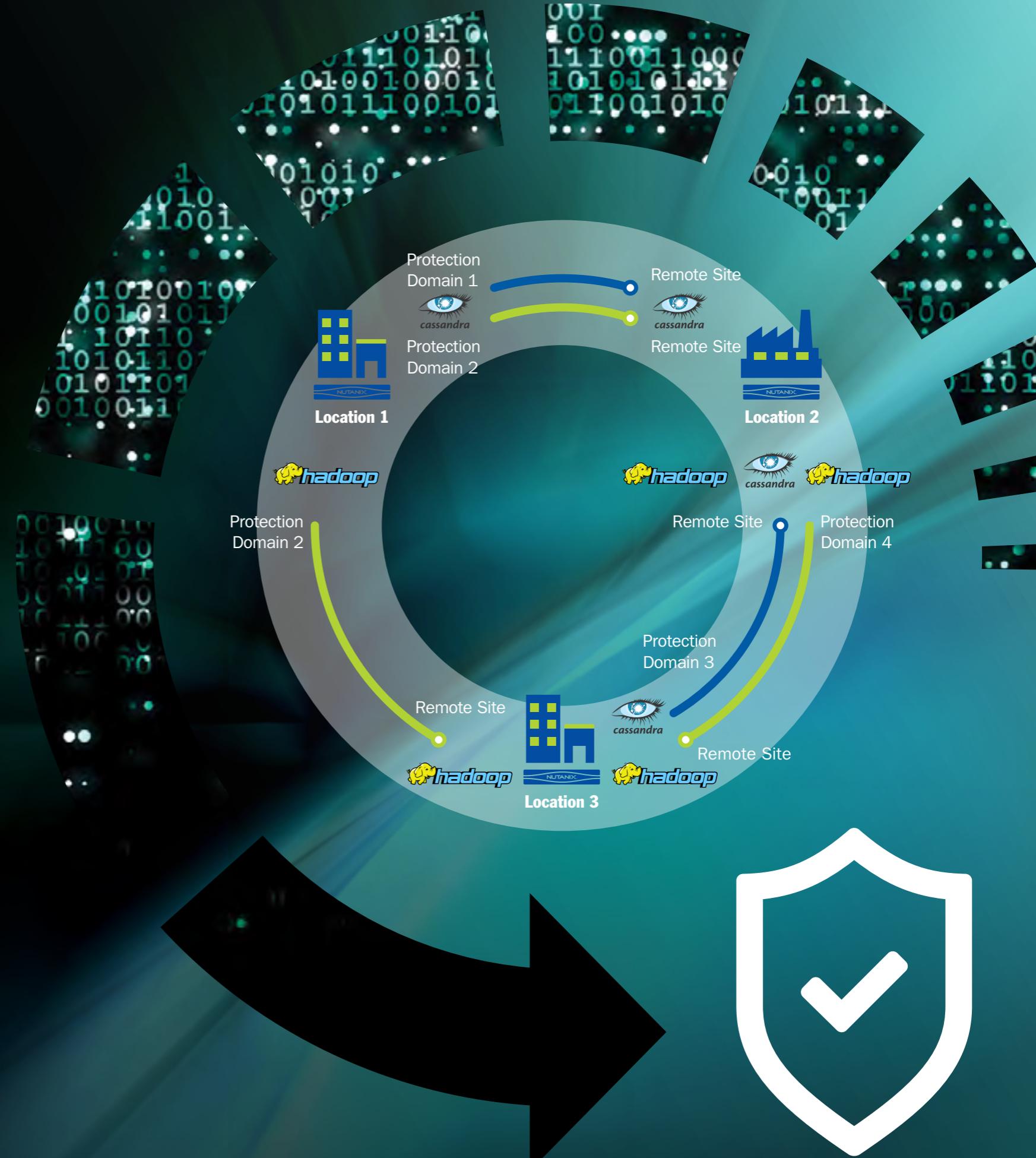
Rather than rely on RAID as conventional IT architectures do, the Distributed Storage Fabric maintains two or three copies of all data for resilience. If one copy becomes un-available for any reason, the alternative copy is accessed. With self-healing, full resiliency is restored automatically in the background. Because the work required for self-healing is distributed across the cluster, there is no long period of “degraded mode” operation that impacts application performance and creates added risk.

For big data applications, the best practice is to configure the application to create only as many copies as you need for performance and use the Nutanix replication factor to provide resiliency. This ensures that data copies are optimally located on separate hardware for resilience.

Nutanix offers a range of integrated data protection options to ensure continued access to important big data applications. Backups can be performed locally, to a remote site, or to a cloud service provider. Nutanix Prism simplifies backup management by giving you centralized control of backup policies for multiple clusters and sites.

For disaster recovery, choose from asynchronous or synchronous replication. By compressing and deduplicating data before it is sent over the wire, these technologies can reduce storage footprint and network bandwidth significantly, depending on the dataset.

For critical workloads requiring zero RPO and fast recovery, Nutanix offers Metro Availability, which uses synchronous replication to ensure continuous data availability across different sites during planned maintenance or disasters. Metro Availability greatly simplifies disaster recovery and eliminates the need for secondary solutions.





AHV

If you're going to virtualize your big data application on Nutanix, AHV is the preferred choice because of its low cost and native integration with Nutanix Prism. Nutanix is developing reference architectures specifically for deployment of big data applications on AHV.

Traditional hypervisors were designed for a world of monolithic non-VM-aware storage arrays and switch fabrics; they were built to accommodate thousands of combinations of servers, NICs, and drivers. They require multi-pathing policies and complex designs to mitigate issues such as storage congestion and application resource contention while still accommodating high availability and scalability. Acceptable performance often requires segregating workloads.

Acropolis Hypervisor (AHV) was built from the ground up to provide a much simpler and more scalable hypervisor and associated management platform by leveraging the software intelligence of the hyperconverged architecture. AHV changes the core building block of the virtualized big data application from the hypervisor to the application. It liberates virtualization from the domain of specialists, making it easier to deploy for big data applications. AHV is perfect for big data applications operated by non-traditional IT teams such as security teams. Nutanix provides an all-in-one compliance system without the need to share storage, compute, virtualization space, or licenses with other teams.

AHV is based on the proven Linux KVM hypervisor to ensure support for all popular workloads and is hardened to meet the most stringent enterprise security requirements. It is fully supported by Nutanix, which means that you get full infrastructure and virtualization support from a single vendor with no hidden costs.

MANAGING YOUR BIG DATA ENVIRONMENT

Among the biggest challenges with infrastructure for big data is a painful and disjointed management experience. Servers, storage systems, and storage networks come with their own management tools, and you have to be an expert on both the technology and the tools to use them effectively. If you deploy your big data environment on bare metal with local storage, you're faced with a variety of management tasks on each server: software upgrades and changes, firmware upgrades on drives, and so on. This makes it very difficult to keep everything in sync and can turn any troubleshooting into a lengthy and painful process.



A Nutanix enterprise cloud incorporates management as part of a complete solution. The Nutanix Prism management platform delivers consumer-grade simplicity for infrastructure management and makes it easy to keep infrastructure up and running. With full integration of AHV management, infrastructure changes such as adding or removing VMs can be accomplished quickly.

Powered by advanced data analytics and heuristics, Prism streamlines common IT workflows, providing a single interface for managing servers, storage, data protection, and virtualization. Prism makes configuring, monitoring, and managing Nutanix solutions remarkably simple. One-click management reduces the administrative burden and the potential for operator error while eliminating the need for planned downtime.

Prism Central allows you to manage multiple Nutanix clusters in dispersed locations from a single pane of glass. Prism Central is an ideal tool to help manage big data operations that span multiple datacenter locations.

ONE-CLICK MANAGEMENT

One-click software upgrades. A consistent pain point for any IT environment is keeping system software and firmware up to date. IT administrators spend countless evening and weekend hours on upgrade tasks. Prism takes the pain and disruption out of upgrades, allowing them to be executed during normal business hours. Intelligent software does all the heavy lifting, eliminating the need for detailed upfront planning.

Nutanix operating software and hypervisor software on each node is updated using a rolling methodology that eliminates disruption to running jobs.

One-click remediation. In the event of alerts or failures, Prism suggests remediation actions that you can initiate to correct problems quickly. With one-click remediation, the mean time to repair and restore services is greatly reduced, significantly improving availability.

FULL REST APIs

Any task that can be performed via Prism can also be performed using REST APIs or a library of PowerShell cmdlets. As a result, you can incorporate Nutanix management as part of the existing management processes for your big data application, or programmers can automate tasks such as scaling up and scaling down in response to demand as part of big data applications.

ONLINE RETAILER DEPLOYS NUTANIX TO ACCELERATE ELASTIC

A top 50 Internet Retailer was looking for a new infrastructure to support their Elastic deployment which managed their inventory, consisting of hundreds of thousands of items.

Challenges. Existing infrastructure was not delivering the predictable responses needed from database searches. The IT team struggled to identify the compute, storage, and network constraints that were hindering performance. An IT project was initiated to replatform existing product and search engines and add new product and search capabilities.

Why Nutanix? Although the initial capital outlay was slightly higher, the capabilities Nutanix provided for automation, visibility, and ease of management offset the difference. In the end, velocity and density were the most important criteria in the decision to go with Nutanix.

Benefits:

- **Faster inventory updates.** Once a night, the company builds the indexes for 1.6 million different products using Elastic. The data is indexed and optimized for consumption by front-end applications. The combination of Elastic and Nutanix delivers near real-time updates.
- **Higher conversion rates through faster search.** Query performance has improved from seconds to milliseconds —two orders of magnitude faster than the previous infrastructure. This results in more click-throughs, leading to an increase in sales conversion rate.
- **Simplified management.** With Nutanix hyperconverged infrastructure, they can quickly scale up or down without worrying about capacity management. The overhead is minimal with the highly intuitive Nutanix Prism interface.
- **Excellent support.** Within 30 minutes of the initial call, the Nutanix team isolated the problem and helped the company move workloads to other nodes so that a module could be replaced with no disruption.

DEPLOYING HADOOP ON NUTANIX

Architecting and setting up infrastructure to support Hadoop can be a time-consuming and complicated chore. Nutanix takes the pain out of Hadoop deployment, accelerating time to value while delivering exceptional performance. Virtualizing Hadoop as part of a Nutanix enterprise cloud offers many advantages:

- **Manage Hadoop like any other app.** One-click management simplifies infrastructure management and upgrades.
- **Avoid the hypervisor tax.** Nutanix includes AHV at no additional cost.
- **Improve availability and security.** VM-HA and automated Security Technical Implementation Guides (STIGs) keep data available and secure.
- **Take advantage of data locality.** Nutanix is the only HCI vendor that provides data locality.
- **Increase hardware utilization.** Nutanix delivers higher CPU utilization and greater flexibility.
- **Change Hadoop economics.** Downtime and underutilized hardware can jeopardize big data projects. Virtualizing Hadoop changes the economics.
- **Increase data efficiency.** Enable compression at the VM or file level.
- **Benefit from auto-leveling and auto-archive.** Data is spread evenly across the cluster. Cold data moves automatically from compute nodes to cold storage nodes.

COMPLETE REFERENCE ARCHITECTURE

Working with HortonWorks, Nutanix has created a full reference architecture for Hadoop on Nutanix. This reference architecture not only further streamlines Hadoop deployment, it has undergone extensive performance testing.

The figure shows Nutanix performance relative to a traditional deployment with 24 1TB SATA drives for the popular TeraSort benchmark. Each 2U appliance added to a cluster adds a predictable 500MB/s of additional throughput in this scenario.

SPLUNK ENTERPRISE



Logfiles Configs Messages Traps Alerts Metrics Scripts Changes Tickets

A Nutanix enterprise cloud lets you operate and scale Splunk on dedicated hardware or in conjunction with other hosted services. For existing sources and platforms, machine data can be sent to the Splunk platform on Nutanix over the network. Modular scale-out enables you to start with a modest deployment and grow in granular increments, eliminating the need for large up-front infrastructure purchases that may take months or years to grow into, ensuring a faster time-to-value for Splunk.

Virtualizing Splunk Enterprise as part of a Nutanix enterprise cloud offers many advantages:

- **Start small and scale incrementally.** Match performance and capacity with demand to minimize upfront costs. Linear scaling yields predictable results.
- **Increase data efficiency.** Enable compression at the VM or file level.
- **Reduce infrastructure footprint.** Nutanix can increase density up to 4x relative to traditional infrastructure.
- **Improve lifecycle management.** Nutanix continuously monitors data access patterns and places data in the most appropriate location, complementing the Splunk life cycle.
- **Improve availability and security.** VM-HA and automated Security Technical Implementation Guides (STIGs) keep data available and secure.

COMPLETE REFERENCE ARCHITECTURE

Nutanix has created a full reference architecture for Splunk on Nutanix. This reference architecture not only streamlines Splunk deployment, it features extensive performance testing.

- **Ingest terabytes of data per day.** A compact 4-node, 2U cluster provides sequential throughput of 3 GB/s or more.
- **Process millions of events per second.** A 4-node cluster can process 500,000 events per second.

SUCCESSFUL START

The journey to Big Data and cloud technologies has truly changed datacenter models. IT leaders now have the ability to focus on strategic ways of solving business problems and to spend less time worrying about the challenges of managing infrastructure. It's not that organizations who continue with legacy architectures will fail, rather it's about moving beyond traditional outcomes and creating exponential business growth. The enterprise cloud platform provides IT leaders with leverage to move past error-prone manual day-to-day operations and enable pay-as-you-grow economics and fractional consumption models.

If you have any questions, contact us at info@nutanix.com, follow us on [Twitter @nutanix](#), or send us a request at www.nutanix.com/demo to set up your own customized briefing and demonstration to see how validated and certified solutions from Nutanix can help your organization make the most of its Big Data technologies. You can also stay engaged with Nutanix experts and customers on the Nutanix Next online community (next.nutanix.com).



Nutanix makes infrastructure invisible, elevating IT to focus on the applications and services that power their business. The Nutanix enterprise cloud platform leverages web-scale engineering and consumer-grade design to natively converge compute, virtualization and storage into a resilient, software-defined solution with rich machine intelligence. The result is predictable performance, cloud-like infrastructure consumption, robust security, and seamless application mobility for a broad range of enterprise applications.