

**NIST Special Publication 1500-1**

---

**DRAFT NIST Big Data Interoperability  
Framework:  
Volume 1, Definitions**

---

NIST Big Data Public Working Group  
Definitions and Taxonomies Subgroup

Draft Version 1  
April 6, 2015

<http://dx.doi.org/10.6028/NIST.SP.1500-1>

**NIST**  
National Institute of  
Standards and Technology  
U.S. Department of Commerce

NIST Special Publication 1500-1  
Information Technology Laboratory

**DRAFT NIST Big Data Interoperability  
Framework:  
Volume 1, Definitions  
Draft Version 1**

NIST Big Data Public Working Group (NBD-PWG)  
Definitions and Taxonomies Subgroup  
National Institute of Standards and Technology  
Gaithersburg, MD 20899

April 2015



U. S. Department of Commerce  
*Penny Pritzker, Secretary*

National Institute of Standards and Technology  
*Dr. Willie E. May, Under Secretary of Commerce for Standards and Technology and Director*

**National Institute of Standards and Technology NIST Special Publication 1500-1**  
33 pages (April 6, 2015)

Certain commercial entities, equipment, or materials may be identified in this document in order to describe an experimental procedure or concept adequately. Such identification is not intended to imply recommendation or endorsement by NIST, nor is it intended to imply that the entities, materials, or equipment are necessarily the best available for the purpose.

There may be references in this publication to other publications currently under development by NIST in accordance with its assigned statutory responsibilities. The information in this publication, including concepts and methodologies, may be used by Federal agencies even before the completion of such companion publications. Thus, until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain operative. For planning and transition purposes, Federal agencies may wish to closely follow the development of these new publications by NIST.

Organizations are encouraged to review all draft publications during public comment periods and provide feedback to NIST. All NIST Information Technology Laboratory publications, other than the ones noted above, are available at <http://www.nist.gov/publication-portal.cfm>.

**Public comment period: April 6, 2015 through May 21, 2015**

**Comments on this publication may be submitted to Wo Chang**

National Institute of Standards and Technology  
Attn: Wo Chang, Information Technology Laboratory  
100 Bureau Drive (Mail Stop 8900) Gaithersburg, MD 20899-8930  
Email: [SP1500comments@nist.gov](mailto:SP1500comments@nist.gov)

## Reports on Computer Systems Technology

The Information Technology Laboratory (ITL) at NIST promotes the U.S. economy and public welfare by providing technical leadership for the Nation's measurement and standards infrastructure. ITL develops tests, test methods, reference data, proof of concept implementations, and technical analyses to advance the development and productive use of information technology. ITL's responsibilities include the development of management, administrative, technical, and physical standards and guidelines for the cost-effective security and privacy of other than national security-related information in Federal information systems. This document reports on ITL's research, guidance, and outreach efforts in Information Technology and its collaborative activities with industry, government, and academic organizations.

### Abstract

Big Data is a term used to describe the new deluge of data in our networked, digitized, sensor-laden, information-driven world. While great opportunities exist with Big Data, it can overwhelm traditional technical approaches and its growth is outpacing scientific and technological advances in data analytics. To advance progress in Big Data, the NIST Big Data Public Working Group (NBD-PWG) is working to develop consensus on important, fundamental questions related to Big Data. The results are reported in the *NIST Big Data Interoperability Framework* series of volumes. This volume, Volume 1, contains a definition of Big Data and related terms necessary to lay the groundwork for discussions surrounding Big Data.

### Keywords

Big Data, Data Science, Reference Architecture, System Orchestrator, Data Provider, Big Data Application Provider, Big Data Framework Provider, Data Consumer, Security and Privacy Fabric, Management Fabric, Big Data taxonomy, use cases, Big Data characteristics

## Acknowledgements

This document reflects the contributions and discussions by the membership of the NBD-PWG, co-chaired by Wo Chang of the NIST ITL, Robert Marcus of ET-Strategies, and Chaitanya Baru, University of California San Diego Supercomputer Center.

The document contains input from members of the NBD-PWG Definitions and Taxonomies Subgroup, led by Nancy Grady (SAIC), Natasha Balac (SDSC), and Eugene Luster (R2AD).

NIST SP1500-1, Version 1 has been collaboratively authored by the NBD-PWG. As of the date of this publication, there are over six hundred NBD-PWG participants from industry, academia, and government. Federal agency participants include the National Archives and Records Administration (NARA), National Aeronautics and Space Administration (NASA), National Science Foundation (NSF), and the U.S. Departments of Agriculture, Commerce, Defense, Energy, Health and Human Services, Homeland Security, Transportation, Treasury, and Veterans Affairs.

NIST would like to acknowledge the specific contributions to this volume by the following NBD-PWG members:

Deborah Blackstock  
*MITRE Corporation*

David Boyd  
*L3 Data Tactics*

Pw Carey  
*Compliance Partners, LLC*

Wo Chang  
*NIST*

Yuri Demchenko  
*University of Amsterdam*

Frank Farance  
*Consultant*

Geoffrey Fox  
*University of Indiana*

Ian Gorton  
*CMU*

Nancy Grady  
*SAIC*

Karen Guertler  
*Consultant*

Keith Hare  
*JCC Consulting, Inc.*

Christine Hawkinson  
*U.S. Bureau of Land  
Management*

Thomas Huang  
*NASA*

Philippe Journeau  
*ResearXis*

Pavithra Kenjige  
*PK Technologies*

Orit Levin  
*Microsoft*

Eugene Luster  
*U.S. Defense Information  
Systems Agency/R2AD LLC*

Ashok Malhotra  
*Oracle*

Bill Mandrick  
*L3 Data Tactics*

Robert Marcus  
*ET-Strategies*

Lisa Martinez  
*Consultant*

Gary Mazzaferro  
*AlloyCloud, Inc.*

William Miller  
*MaCT USA*

Sanjay Mishra  
*Verizon*

Bob Natale  
*Mitre*

Rod Peterson  
*U.S. Department of Veterans  
Affairs*

Ann Racuya-Robbins  
*World Knowledge Bank*

Russell Reinsch  
*Calibrum*

John Rogers  
*HP*

Arnab Roy  
*Fujitsu*

Mark Underwood  
*Krypton Brothers LLC*

William Vorhies  
*Predictive Modeling LLC*

Tim Zimmerman  
*Consultant*

Alicia Zuniga-Alvarado  
*Consultant*

The editors for this document were Nancy Grady and Wo Chang.

## Notice to Readers

NIST is seeking feedback on the proposed working draft of the *NIST Big Data Interoperability Framework: Volume 1, Definitions*. Once public comments are received, compiled, and addressed by the NBD-PWG, and reviewed and approved by NIST internal editorial board, Version 1 of this volume will be published as final. Three versions are planned for this volume, with Versions 2 and 3 building on the first. Further explanation of the three planned versions and the information contained therein is included in Section 1.5 of this document.

Please be as specific as possible in any comments or edits to the text. Specific edits include, but are not limited to, changes in the current text, additional text further explaining a topic or explaining a new topic, additional references, or comments about the text, topics, or document organization. These specific edits can be recorded using one of the two following methods.

1. **TRACK CHANGES**: make edits to and comments on the text directly into this Word document using track changes
2. **COMMENT TEMPLATE**: capture specific edits using the Comment Template : ([http://bigdatawg.nist.gov/uploadfiles/SP1500-1-to-7\\_comment\\_template.docx](http://bigdatawg.nist.gov/uploadfiles/SP1500-1-to-7_comment_template.docx)), which includes space for Section number, page number, comment, and text edits

Submit the edited file from either method 1 or 2 to [SP1500comments@nist.gov](mailto:SP1500comments@nist.gov) with the volume number in the subject line (e.g., Edits for Volume 1.)

Please contact Wo Chang ([wchang@nist.gov](mailto:wchang@nist.gov)) with any questions about the feedback submission process.

Big Data professionals continue to be welcome to join the NBD-PWG to help craft the work contained in the volumes of the NIST Big Data Interoperability Framework. Additional information about the NBD-PWG can be found at <http://bigdatawg.nist.gov>.

# Table of Contents

---

<b>EXECUTIVE SUMMARY</b> .....	<b>VII</b>
<b>1 INTRODUCTION</b> .....	<b>1</b>
1.1 BACKGROUND .....	1
1.2 SCOPE AND OBJECTIVES OF THE DEFINITIONS AND TAXONOMIES SUBGROUP .....	2
1.3 REPORT PRODUCTION .....	2
1.4 REPORT STRUCTURE .....	3
1.5 FUTURE WORK ON THIS VOLUME .....	3
<b>2 BIG DATA AND DATA SCIENCE DEFINITIONS</b> .....	<b>4</b>
2.1 BIG DATA DEFINITIONS .....	4
2.2 DATA SCIENCE DEFINITIONS .....	7
2.3 OTHER BIG DATA DEFINITIONS .....	9
<b>3 BIG DATA FEATURES</b> .....	<b>12</b>
3.1 DATA ELEMENTS AND METADATA .....	12
3.2 DATA RECORDS AND NON-RELATIONAL MODELS .....	12
3.3 DATASET CHARACTERISTICS AND STORAGE .....	13
3.4 DATA IN MOTION .....	15
3.5 DATA SCIENCE LIFECYCLE MODEL FOR BIG DATA .....	16
3.6 BIG DATA ANALYTICS .....	16
3.7 BIG DATA METRICS AND BENCHMARKS .....	17
3.8 BIG DATA SECURITY AND PRIVACY .....	17
3.9 DATA GOVERNANCE .....	18
<b>4 BIG DATA ENGINEERING PATTERNS (FUNDAMENTAL CONCEPTS)</b> .....	<b>19</b>
<b>APPENDIX A: INDEX OF TERMS</b> .....	<b>A-1</b>
<b>APPENDIX B: TERMS AND DEFINITIONS</b> .....	<b>B-1</b>
<b>APPENDIX C: ACRONYMS</b> .....	<b>C-1</b>
<b>APPENDIX D: REFERENCES</b> .....	<b>D-1</b>
 <b>FIGURE</b>	
FIGURE 1: SKILLS NEEDED IN DATA SCIENCE .....	8
 <b>TABLE</b>	
TABLE 1: SAMPLING OF CONCEPTS ATTRIBUTED TO BIG DATA .....	10

# 1 Executive Summary

---

2 The NIST Big Data Public Working Group (NBD-PWG) Definitions and Taxonomy Subgroup prepared  
3 this *NIST Big Data Interoperability Framework: Volume 1, Definitions* to address fundamental concepts  
4 needed to understand the new paradigm for data applications, collectively known as Big Data, and the  
5 analytic processes collectively known as data science. While Big Data has been defined in a myriad of  
6 ways, the shift to a Big Data paradigm occurs when the scale of the data leads to the need for a cluster of  
7 computing and storage resources to provide cost-effective data management. Data science combines  
8 various technologies, techniques, and theories from various fields, mostly related to computer science and  
9 statistics, to obtain actionable knowledge from data. This report seeks to clarify the underlying concepts  
10 of Big Data and data science to enhance communication among Big Data producers and consumers. By  
11 defining concepts related to Big Data and data science, a common terminology can be used among Big  
12 Data practitioners.

13 The *NIST Big Data Interoperability Framework* consists of seven volumes, each of which addresses a  
14 specific key topic, resulting from the work of the NBD-PWG. The seven volumes are as follows:

- 15 • Volume 1, Definitions
- 16 • Volume 2, Taxonomies
- 17 • Volume 3, Use Cases and General Requirements
- 18 • Volume 4, Security and Privacy
- 19 • Volume 5, Architectures White Paper Survey
- 20 • Volume 6, Reference Architecture
- 21 • Volume 7, Standards Roadmap

22 The *NIST Big Data Interoperability Framework* will be released in three versions, which correspond to  
23 the three stages of the NBD-PWG work. The three stages aim to achieve the following:

- 24 Stage 1: Identify the high-level Big Data reference architecture key components, which are  
25 technology, infrastructure, and vendor agnostic
- 26 Stage 2: Define general interfaces between the NIST Big Data Reference Architecture (NBDRA)  
27 components
- 28 Stage 3: Validate the NBDRA by building Big Data general applications through the general interfaces

29 Potential areas of future work for the Subgroup during stage 2 are highlighted in Section 1.5 of this  
30 volume. The current effort documented in this volume reflects concepts developed within the rapidly  
31 evolving field of Big Data.

32



# 33 1 INTRODUCTION

---

## 34 1.1 BACKGROUND

35 There is broad agreement among commercial, academic, and government leaders about the remarkable  
 36 potential of Big Data to spark innovation, fuel commerce, and drive progress. Big Data is the common  
 37 term used to describe the deluge of data in today's networked, digitized, sensor-laden, and information-  
 38 driven world. The availability of vast data resources carries the potential to answer questions previously  
 39 out of reach, including the following:

- 40 • How can a potential pandemic reliably be detected early enough to intervene?
- 41 • Can new materials with advanced properties be predicted before these materials have ever been  
 42 synthesized?
- 43 • How can the current advantage of the attacker over the defender in guarding against cyber-  
 44 security threats be reversed?

45 There is also broad agreement on the ability of Big Data to overwhelm traditional approaches. The growth  
 46 rates for data volumes, speeds, and complexity are outpacing scientific and technological advances in data  
 47 analytics, management, transport, and data user spheres.

48 Despite widespread agreement on the inherent opportunities and current limitations of Big Data, a lack of  
 49 consensus on some important, fundamental questions continues to confuse potential users and stymie  
 50 progress. These questions include the following:

- 51 • What attributes define Big Data solutions?
- 52 • How is Big Data different from traditional data environments and related applications?
- 53 • What are the essential characteristics of Big Data environments?
- 54 • How do these environments integrate with currently deployed architectures?
- 55 • What are the central scientific, technological, and standardization challenges that need to be  
 56 addressed to accelerate the deployment of robust Big Data solutions?

57 Within this context, on March 29, 2012, the White House announced the Big Data Research and  
 58 Development Initiative.<sup>1</sup> The initiative's goals include helping to accelerate the pace of discovery in  
 59 science and engineering, strengthening national security, and transforming teaching and learning by  
 60 improving the ability to extract knowledge and insights from large and complex collections of digital  
 61 data.

62 Six federal departments and their agencies announced more than \$200 million in commitments spread  
 63 across more than 80 projects, which aim to significantly improve the tools and techniques needed to  
 64 access, organize, and draw conclusions from huge volumes of digital data. The initiative also challenged  
 65 industry, research universities, and nonprofits to join with the federal government to make the most of the  
 66 opportunities created by Big Data.

67 Motivated by the White House initiative and public suggestions, the National Institute of Standards and  
 68 Technology (NIST) has accepted the challenge to stimulate collaboration among industry professionals to  
 69 further the secure and effective adoption of Big Data. As one result of NIST's Cloud and Big Data Forum  
 70 held on January 15–17, 2013, there was strong encouragement for NIST to create a public working group  
 71 for the development of a Big Data Standards Roadmap. Forum participants noted that this roadmap  
 72 should define and prioritize Big Data requirements, including interoperability, portability, reusability,  
 73 extensibility, data usage, analytics, and technology infrastructure. In doing so, the roadmap would  
 74 accelerate the adoption of the most secure and effective Big Data techniques and technology.

75 On June 19, 2013, the NIST Big Data Public Working Group (NBD-PWG) was launched with extensive  
76 participation by industry, academia, and government from across the nation. The scope of the NBD-PWG  
77 involves forming a community of interests from all sectors—including industry, academia, and  
78 government—with the goal of developing consensus on definitions, taxonomies, secure reference  
79 architectures, security and privacy, and—from these—a standards roadmap. Such a consensus would  
80 create a vendor-neutral, technology- and infrastructure-independent framework that would enable Big  
81 Data stakeholders to identify and use the best analytics tools for their processing and visualization  
82 requirements on the most suitable computing platform and cluster, while also allowing value-added from  
83 Big Data service providers.

84 The *NIST Big Data Interoperability Framework* consists of seven volumes, each of which addresses a  
85 specific key topic, resulting from the work of the NBD-PWG. The seven volumes are as follows:

- 86 • Volume 1, Definitions
- 87 • Volume 2, Taxonomies
- 88 • Volume 3, Use Cases and General Requirements
- 89 • Volume 4, Security and Privacy
- 90 • Volume 5, Architectures White Paper Survey
- 91 • Volume 6, Reference Architecture
- 92 • Volume 7, Standards Roadmap

93 The *NIST Big Data Interoperability Framework* will be released in three versions, which correspond to  
94 the three stages of the NBD-PWG work. The three stages aim to achieve the following:

- 95 Stage 1: Identify the high-level Big Data reference architecture key components, which are  
96 technology, infrastructure, and vendor agnostic
- 97 Stage 2: Define general interfaces between the NIST Big Data Reference Architecture (NBDRA)  
98 components
- 99 Stage 3: Validate the NBDRA by building Big Data general applications through the general interfaces

100 Potential areas of future work for the Subgroup during stage 2 are highlighted in Section 1.5 of this  
101 volume. The current effort documented in this volume reflects concepts developed within the rapidly  
102 evolving field of Big Data.

## 103 **1.2 SCOPE AND OBJECTIVES OF THE DEFINITIONS AND TAXONOMIES SUBGROUP**

104 This volume was prepared by the NBD-PWG Definitions and Taxonomy Subgroup, which focused on  
105 identifying Big Data concepts and defining related terms in areas such as data science, reference  
106 architecture, and patterns.

107 The aim of this volume is to provide a common vocabulary for those involved with Big Data. For  
108 managers, the terms in this volume will distinguish the concepts needed to understand this changing field.  
109 For procurement officers, this document will provide the framework for discussing organizational needs,  
110 and distinguishing among offered approaches. For marketers, this document will provide the means to  
111 promote solutions and innovations. For the technical community, this volume will provide a common  
112 language to better differentiate the specific offerings.

## 113 **1.3 REPORT PRODUCTION**

114 *Big Data* and *data science* are being used as buzzwords and are composites of many concepts. To better  
115 identify those terms, the NBD-PWG Definitions and Taxonomy Subgroup first addressed the individual  
116 concepts needed in this disruptive field. Then, the two over-arching buzzwords—Big Data and data  
117 science—and the concepts they encompass were clarified.

118 To keep the topic of data and data systems manageable, the Subgroup attempted to limit discussions to  
 119 differences affected by the existence of Big Data. Expansive topics such as data type or analytics  
 120 taxonomies and metadata were only explored to the extent that there were issues or effects specific to Big  
 121 Data. However, the Subgroup did include the concepts involved in other topics that are needed to  
 122 understand the new Big Data methodologies.

123 Terms were developed independent of a specific tool or implementation, to avoid highlighting specific  
 124 implementations, and to stay general enough for the inevitable changes in the field.

125 The Subgroup is aware that some fields, such as legal, use specific language that may differ from the  
 126 definitions provided herein. The current version reflects the breadth of knowledge of the Subgroup  
 127 members. During the comment period, the broader community is requested to address any domain  
 128 conflicts caused by the terminology used in this volume.

## 129 **1.4 REPORT STRUCTURE**

130 This volume seeks to clarify the meanings of the broad terms Big Data and data science, which are  
 131 discussed at length in Section 2. The more elemental concepts and terms that provide additional insights  
 132 are discussed in Section 3. Section 4 explores several concepts that are more detailed. This first version of  
 133 *NIST Big Data Interoperability Framework: Volume 1, Definitions* describes some of the fundamental  
 134 concepts that will be important to determine categories or functional capabilities that represent  
 135 architecture choices.

136 Tightly coupled information can be found in the other volumes of the *NIST Big Data Interoperability*  
 137 *Framework. Volume 2, Taxonomies* provides a description of the more detailed components of the NIST  
 138 Big Data Reference Architecture (NBDRA) presented in *Volume 6, Reference Architecture*. Security and  
 139 privacy related concepts are described in detail in *Volume 4, Security and Privacy*. To understand how  
 140 these systems are architected to meet users' needs, the reader is referred to *Volume 3, Use Cases and*  
 141 *General Requirements. Volume 7, Standards Roadmap* recaps the framework established in Volumes 1  
 142 through 6 and discusses NBDRA related standards. Comparing related sections in these volumes will  
 143 provide a more comprehensive understanding of the consensus of the NBD-PWG.

## 144 **1.5 FUTURE WORK ON THIS VOLUME**

145 This volume represents the beginning stage of the NBD-PWG's effort to provide order and clarity to an  
 146 emerging and rapidly changing field. Big Data encompasses a large range of data types, fields of study,  
 147 technologies, and techniques. Distilling from the varied viewpoints a consistent, core set of definitions to  
 148 frame the discussion has been challenging. However, through discussion of the varied viewpoints a  
 149 greater understanding of the Big Data paradigm will emerge. As the field matures, this document will also  
 150 need to mature to accommodate innovations in the field. To ensure the concepts are accurate, future  
 151 NBD-PWG tasks will consist of the following:

- 152 • Defining the different patterns of communications between Big Data resources to better clarify  
 153 the different approaches being taken
- 154 • Updating Volume 1 taking into account the efforts of other working groups such as International  
 155 Organization for Standardization (ISO) Joint Technical Committee 1 (JTC 1) and the Transaction  
 156 Processing Performance Council.
- 157 • Improve the discussions of governance and data ownership
- 158 • Develop the Management section
- 159 • Develop the Security and Privacy section
- 160 • Add a discussion of the value of data

161

## 162 2 BIG DATA AND DATA SCIENCE DEFINITIONS

---

163 The rate of growth of data generated and stored has been increasing exponentially. In a 1965 paper<sup>2</sup>,  
 164 Gordon Moore estimated that the density of transistors on an integrated circuit board was doubling every  
 165 two years. Known as “Moore’s Law”, this rate of growth has been applied to all aspects of computing,  
 166 from clock speeds to memory. The growth rates of data volumes are considered faster than Moore’s Law,  
 167 with data volumes more than doubling every eighteen months. This data explosion is creating  
 168 opportunities for new ways of combining and using data to find value, as well as providing significant  
 169 challenges due to the size of the data being managed and analyzed. One significant shift is in the amount  
 170 of unstructured data. Historically, structured data has typically been the focus of most enterprise analytics,  
 171 and has been handled through the use of the relational data model. Recently, the quantity of unstructured  
 172 data, such as micro-texts, web pages, relationship data, images and videos, has exploded and the trend  
 173 indicates an increase in the incorporation of unstructured data to generate value. The central benefit of  
 174 Big Data analytics is the ability to process large amounts and various types of information. Big Data does  
 175 not imply that the current data volumes are simply “bigger” than before, or bigger than current techniques  
 176 can efficiently handle. The need for greater performance or efficiency happens on a continual basis.  
 177 However, Big Data represents a fundamental change in the architecture needed to efficiently handle  
 178 current datasets.

179 In the evolution of data systems, there have been a number of times when the need for efficient, cost  
 180 effective data analysis has forced a change in existing technologies. For example, the move to a relational  
 181 model occurred when methods to reliably handle changes to structured data led to the shift toward a data  
 182 storage paradigm that modeled relational algebra. That was a fundamental shift in data handling. The  
 183 current revolution in technologies referred to as Big Data has arisen because the relational data model can  
 184 no longer efficiently handle all the current needs for analysis of large and often unstructured datasets. It is  
 185 not just that data is bigger than before, as it has been steadily getting larger for decades. The Big Data  
 186 revolution is instead a one-time fundamental shift in architecture, just as the shift to the relational model  
 187 was a one-time shift. As relational databases evolved to greater efficiencies over decades, so too will Big  
 188 Data technologies continue to evolve. Many of the conceptual underpinnings of Big Data have been  
 189 around for years, but the last decade has seen an explosion in their maturation and application to scaled  
 190 data systems.

191 The term Big Data has been used to describe a number of concepts, in part because several distinct  
 192 aspects are consistently interacting with each other. To understand this revolution, the interplay of the  
 193 following four aspects must be considered: the characteristics of the datasets, the analysis of the datasets,  
 194 the performance of the systems that handle the data, and the business considerations of cost effectiveness.

195 In the following sections, the two broad concepts, Big Data and data science, are broken down into  
 196 specific individual terms and concepts.

### 197 2.1 BIG DATA DEFINITIONS

198 Big Data refers to the inability of traditional data architectures to efficiently handle the new datasets.  
 199 Characteristics of Big Data that force new architectures are *volume* (i.e., the size of the dataset) and  
 200 *variety* (i.e., data from multiple repositories, domains, or types), and the data in motion characteristics of  
 201 *velocity* (i.e., rate of flow) and *variability* (i.e., the change in other characteristics). These  
 202 characteristics—volume, variety, velocity, and variability—are known colloquially as the ‘Vs’ of Big  
 203 Data and are further discussed in Section 3. Each of these characteristics influences the overall design of a  
 204 Big Data system, resulting in different data system architectures or different data lifecycle process  
 205 orderings to achieve needed efficiencies.

206 *Big Data* consists of extensive datasets—primarily in the characteristics of volume,  
 207 variety, velocity, and/or variability—that require a scalable architecture for efficient  
 208 storage, manipulation, and analysis.

209 Note that this definition contains the interplay between the characteristics of the data and the need for a  
 210 system architecture that can scale to achieve the needed performance and cost efficiency. There are two  
 211 fundamentally different methods for system scaling, often described metaphorically as “vertical” or  
 212 “horizontal” scaling. **Vertical scaling** implies increasing the system parameters of processing speed,  
 213 storage, and memory for greater performance. This approach is limited by physical capabilities whose  
 214 improvements have been described by Moore’s Law, requiring ever more sophisticated elements (e.g.,  
 215 hardware, software) that add time and expense to the implementation. The alternate method is to use  
 216 **horizontal scaling**, to make use of a cluster of individual (usually commodity) resources integrated to act  
 217 as a single system. It is this horizontal scaling that is at the heart of the Big Data revolution.

218 *The Big Data paradigm* consists of the distribution of data systems across horizontally  
 219 coupled, independent resources to achieve the scalability needed for the efficient  
 220 processing of extensive datasets.

221 This new paradigm leads to a number of conceptual definitions that suggest Big Data exists when the  
 222 scale of the data causes the management of the data to be a significant driver in the design of the system  
 223 architecture. This definition does not explicitly refer to the horizontal scaling in the Big Data paradigm.

224 As stated above, fundamentally, the Big Data paradigm is a shift in data system architectures from  
 225 monolithic systems with vertical scaling (i.e., adding more power, such as faster processors or disks, to  
 226 existing machines) into a parallelized, “horizontally scaled”, system (i.e., adding more machines to the  
 227 available collection) that uses a loosely coupled set of resources in parallel. This type of parallelization  
 228 shift began over 20 years ago in the simulation community, when scientific simulations began using  
 229 massively parallel processing (MPP) systems.

230 *Massively parallel processing* refers to a multitude of individual processors working in  
 231 parallel to execute a particular program.

232 In different combinations of splitting the code and data across independent processors, computational  
 233 scientists were able to greatly extend their simulation capabilities. This, of course, introduced a number of  
 234 complications in such areas as message passing, data movement, latency in the consistency across  
 235 resources, load balancing, and system inefficiencies, while waiting on other resources to complete their  
 236 computational tasks.

237 The Big Data paradigm of today is similar. Data systems need a level of extensibility that matches the  
 238 scaling in the data. To get that level of extensibility, different mechanisms are needed to distribute data  
 239 and data retrieval processes across loosely coupled resources.

240 While the methods to achieve efficient scalability across resources will continually evolve, this paradigm  
 241 shift (in analogy to the prior shift in the simulation community) is a one-time occurrence. Eventually, a  
 242 new paradigm shift will likely occur beyond this distribution of a processing or data system that spans  
 243 multiple resources working in parallel. That future revolution will need to be described with new  
 244 terminology.

245 Big Data focuses on the self-referencing viewpoint that data is big because it requires scalable systems to  
 246 handle it. Conversely, architectures with better scaling have come about because of the need to handle Big  
 247 Data. It is difficult to delineate a size requirement for a dataset to be considered Big Data. Data is usually  
 248 considered “big” if the use of new scalable architectures provides a cost or performance efficiency over  
 249 the traditional vertically scaled architectures (i.e., if similar performance cannot be achieved in a  
 250 traditional, single platform computing resource.) This circular relationship between the characteristics of

251 the data and the performance of data systems leads to different definitions for Big Data if only one aspect  
252 is considered.

253 Some definitions for Big Data focus on the systems innovations required because of the characteristics of  
254 Big Data.

255 *Big Data engineering includes advanced techniques that harness independent resources*  
256 *for building scalable data systems when the characteristics of the datasets require new*  
257 *architectures for efficient storage, manipulation, and analysis.*

258 Once again the definition is coupled, so that Big Data engineering is used when the characteristics of the  
259 data require it. New engineering techniques in the data layer have been driven by the growing prominence  
260 of datasets that cannot be handled efficiently in a traditional relational model. The need for scalable  
261 access in structured data has led to software built on the key-value pair paradigm. The rise in importance  
262 of document analysis has spawned a document-oriented database paradigm, and the increasing  
263 importance of relationship data has led to efficiencies in the use of graph-oriented data storage.

264 The new non-relational model database paradigms are typically referred to as *NoSQL* (Not Only or No  
265 Structured Query Language [SQL]) systems, which are further discussed in Section 3. The problem with  
266 identifying Big Data storage paradigms as NoSQL is, first, that it describes the storage of data with  
267 respect to a set theory-based language for query and retrieval of data, and, second, that there is a growing  
268 capability in the application of the SQL query language against the new non-relational data repositories.  
269 While NoSQL is in such common usage that it will continue to refer to the new data models beyond the  
270 relational model, it is hoped the term itself will be replaced with a more suitable term, since it is unwise to  
271 name a set of new storage paradigms with respect to a query language currently in use against that  
272 storage.

273 *Non-relational models, frequently referred to as NoSQL, refer to logical data models*  
274 *that do not follow relational algebra for the storage and manipulation of data.*

275 Another related engineering technique is the federated database system, which is related to the variety  
276 characteristic of Big Data.

277 *A federated database system is a type of meta-database management system, which*  
278 *transparently maps multiple autonomous database systems into a single federated*  
279 *database.*

280 A federated database is thus a database system comprised of underlying database systems. Big Data  
281 systems can likewise pull a variety of data from many sources, but the underlying repositories do not all  
282 have to conform to the relational model.

283 Note that for systems and analysis processes, the Big Data paradigm shift also causes changes in the  
284 traditional data lifecycle processes. One description of the end-to-end data lifecycle categorizes the  
285 process steps as collection, preparation, analysis, and action. Different Big Data use cases can be  
286 characterized in terms of the dataset characteristics and in terms of the time window for the end-to-end  
287 data lifecycle. Dataset characteristics change the data lifecycle processes in different ways, for example in  
288 the point in the lifecycle at which the data is placed in persistent storage. In a traditional relational model,  
289 the data is stored after preparation (for example, after the extract-transform-load and cleansing processes).  
290 In a high velocity use case, the data is prepared and analyzed for alerting, and only then is the data (or  
291 aggregates of the data) given a persistent storage. In a volume use case, the data is often stored in the raw  
292 state in which it was produced—before being cleansed and organized (sometimes referred to as extract-  
293 load-transform). The consequence of persistence of data in its raw state is that a schema or model for the  
294 data is only applied when the data is retrieved for preparation and analysis. This Big Data concept is  
295 described as schema-on-read.

296            *Schema-on-read* is the application of a data schema through preparation steps such as  
 297            transformations, cleansing, and integration at the time the data is read from the  
 298            database.

299            Another concept of Big Data is often referred to as *moving the processing to the data, not the data to the*  
 300            *processing.*

301            *Computational portability* is the movement of the computation to the location of the data.

302            The implication is that data is too extensive to be queried and moved into another resource for analysis, so  
 303            the analysis program is instead distributed to the data-holding resources, with only the results being  
 304            aggregated on a remote resource. This concept of data locality is actually a critical aspect of parallel data  
 305            architectures. Additional system concepts are the interoperability (ability for tools to work together),  
 306            reusability (ability to apply tools from one domain to another), and extendibility (ability to add or modify  
 307            existing tools for new domains). These system concepts are not specific to Big Data, but their presence in  
 308            Big Data can be understood in the examination of a Big Data reference architecture, which is discussed in  
 309            *NIST Big Data Interoperability Framework: Volume 6, Reference Architecture* of this series.

310            Additional concepts used in reference to the term Big Data refer to changes in analytics, which will be  
 311            discussed in Section 2.2. A number of other terms (particularly terms starting with the letter V) are also  
 312            used, several of which refer to the data science process or its benefit, instead of new Big Data  
 313            characteristics. Some of these additional terms include *veracity* (i.e., accuracy of the data), *value* (i.e.,  
 314            value of the analytics to the organization), *volatility* (i.e., tendency for data structures to change over  
 315            time), and *validity* (i.e., appropriateness of the data for its intended use). While these characteristics and  
 316            others—including quality control, metadata, and data provenance—long pre-dated Big Data, their impact  
 317            is still important in Big Data systems. Several of these terms are discussed with respect to Big Data  
 318            analytics in Section 3.4.

319            Essentially, Big Data refers to the extensibility of data repositories and data processing across resources  
 320            working in parallel, in the same way the compute-intensive simulation community embraced massively  
 321            parallel processing two decades ago. By working out methods for communication among resources, the  
 322            same scaling is now available to data-intensive applications.

## 323            2.2 DATA SCIENCE DEFINITIONS

324            In its purest form, data science is the *fourth paradigm* of science, following theory, experiment, and  
 325            computational science. The fourth paradigm is a term coined by Dr. Jim Gray in 2007. It refers to the  
 326            conduct of data analysis as an empirical science, learning directly from data itself. Data science as a  
 327            paradigm would refer to the formulation of a hypothesis, the collection of the data—new or pre-  
 328            existing—to address the hypothesis, and the analytical confirmation or denial of the hypothesis (or the  
 329            determination that additional information or study is needed.) In many data science projects, the raw data  
 330            is browsed first, which informs a hypothesis, which is then investigated. As in any experimental science,  
 331            the end result could be that the original hypothesis itself needs to be reformulated. The key concept is that  
 332            data science is an empirical science, performing the scientific process directly on the data. Note that the  
 333            hypothesis may be driven by a business need, or can be the restatement of a business need in terms of a  
 334            technical hypothesis.

335            *The data science paradigm* is extraction of actionable knowledge directly from data  
 336            through a process of discovery, hypothesis, and hypothesis testing.

337            Data science can be understood as the activities happening in the processing layer of the system  
 338            architecture, against data stored in the data layer, in order to extract knowledge from the raw data through  
 339            the complete data lifecycle.

340 *The **data lifecycle** is the set of processes that transform raw data into actionable*  
 341 *knowledge.*

342 Traditionally, the term analytics has been used as one of the steps in the data lifecycle of collection,  
 343 preparation, analysis, and action.

344 ***Analytics** is the synthesis of knowledge from information.*

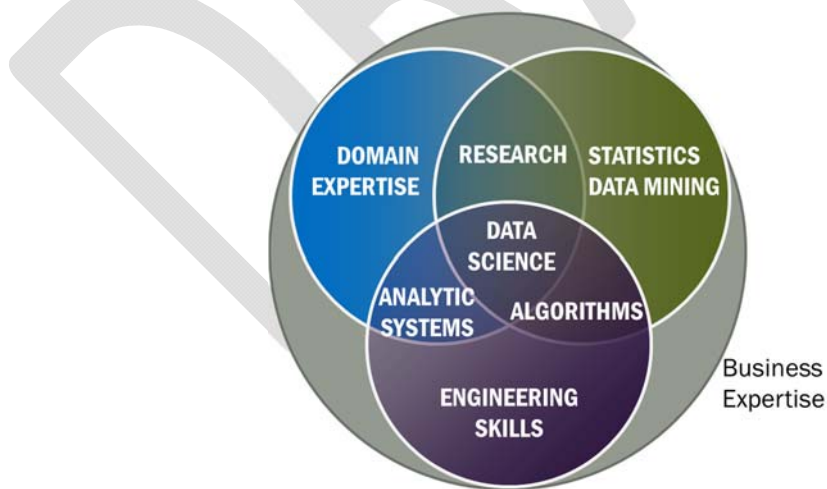
345 With the new Big Data paradigm, analytics are no longer separable from the data model and the  
 346 distribution of that data across parallel resources. When structured data was almost exclusively stored as  
 347 organized information in a relational model, the analytics could be designed for this structure. While the  
 348 working definition of the data science paradigm refers to learning directly from data, in the Big Data  
 349 paradigm this learning must implicitly involve all steps in the data lifecycle, with analytics being only a  
 350 subset.

351 ***Data science** is the empirical synthesis of actionable knowledge from raw data through*  
 352 *the complete data lifecycle process.*

353 Data science across the entire data lifecycle now incorporates principles, techniques, and methods from  
 354 many disciplines and domains, including the analytics domains of mathematics, data mining (specifically  
 355 machine learning and pattern recognition), statistics, operations research, and visualization, along with the  
 356 domains of systems, software, and network engineering. Data scientists and data science teams solve  
 357 complex data problems by employing deep expertise in one or more of these disciplines, in the context of  
 358 business strategy, and under the guidance of domain knowledge. Personal skills in communication,  
 359 presentation, and inquisitiveness are also very important given the complexity of interactions within Big  
 360 Data systems.

361 *A **data scientist** is a practitioner who has sufficient knowledge in the overlapping regimes*  
 362 *of business needs, domain knowledge, analytical skills, and software and systems*  
 363 *engineering to manage the end-to-end data processes through each stage in the data*  
 364 *lifecycle.*

365 While this full collection of skills can be present in a single individual, it is also possible that these skills,  
 366 as shown in Figure 1, are covered in the members of a team.



367  
 368 *Figure 1: Skills Needed in Data Science*

369 Data science is not solely concerned with analytics, but also with the end-to-end experimental lifecycle,  
 370 where the data system is essentially the scientific equipment. The implication is that the data scientist  
 371 must be aware of the sources and provenance of the data, the appropriateness and accuracy of the



372 transformations on the data, the interplay between the transformation algorithms and processes, and the  
 373 data storage mechanisms. This end-to-end overview role ensures that everything is performed correctly to  
 374 meaningfully address the hypothesis. These analytics concepts are discussed further in Section 3.4.

375 Data science is increasingly used to influence business decisions. In Big Data systems, identifying a  
 376 correlation is often sufficient for a business to take action. As a simple example, if it can be determined  
 377 that using the color blue on a website leads to greater sales than using green, then this correlation can be  
 378 used to improve the business. The reason for the preference is not needed—it is enough to determine  
 379 correlation.

380 Several issues are currently being debated within the data science community, two of which are data  
 381 sampling, and the idea that more data is superior to better algorithms

382 Data sampling, a central concept of statistics, involves the selection of a subset of data from the larger  
 383 data population. The subset of data can be used as input for analytical processes, to determine  
 384 methodology to be used for experimental procedures, or to address questions. For example, it is possible  
 385 to calculate the data needed to determine an outcome for an experimental procedure (e.g., during a  
 386 pharmaceutical clinical trial).

387 When the data mining community began, the emphasis was typically on re-purposed data (i.e., data used  
 388 to train models was sampled from a larger dataset that was originally collected for another purpose). The  
 389 often-overlooked critical step was to ensure that the analytics were not prone to over-fitting (i.e., the  
 390 analytical pattern matched the data sample but did not work well to answer questions of the overall data  
 391 population). In the new Big Data paradigm, it is implied that data sampling from the overall data  
 392 population is no longer necessary since the Big Data system can theoretically process all the data without  
 393 loss of performance. However, even if all of the available data is used, it still only represents a population  
 394 subset whose behaviors led them to produce the data, which might not be the true population of interest.  
 395 For example, studying Twitter data to analyze people’s behaviors does not represent all people, as not  
 396 everyone uses Twitter. While less sampling may be used in data science processes, it is important to be  
 397 aware of the implicit sampling when trying to address business questions.

398 The assertion that more data is superior to better algorithms implies that better results can be achieved by  
 399 analyzing larger samples of data rather than refining the algorithms used in the analytics. The heart of this  
 400 debate states that a few bad data elements are less likely to influence the analytical results in a large  
 401 dataset than if errors are present in a small sample of that dataset. If the analytics needs are correlation  
 402 and not causation, then this assertion is easier to justify. Outside the context of large datasets in which  
 403 aggregate trending behavior is all that matters, the data quality rule remains “garbage-in, garbage-out”,  
 404 where you cannot expect accurate results based on inaccurate data.

405 For descriptive purposes, analytics activities can be broken into different categories, including discovery,  
 406 exploratory analysis, correlation analysis, predictive modeling, and machine learning. Again, these  
 407 analytics categories are not specific to Big Data, but some have gained more visibility due to their greater  
 408 application in data science.

409 Data science is tightly linked to Big Data, and refers to the management and execution of the end-to-end  
 410 data processes, including the behaviors of the data system. As such, data science includes all of analytics,  
 411 but analytics does not include all of data science.

## 412 **2.3 OTHER BIG DATA DEFINITIONS**

413 A number of Big Data definitions have been suggested as efforts have been made to understand the extent  
 414 of this new field. Several Big Data concepts, discussed in previous sections, were observed in a sample of  
 415 definitions taken from blog posts<sup>3 4 5 6</sup>. The sample of formal and informal definitions offer a sense of the  
 416 spectrum of concepts applied to the term Big Data. The sample of Big Data concepts and definitions are

417 aligned in Table 1. The NBD-PWG’s definition is closest to the Gartner definition, with additional  
 418 emphasis that the horizontal scaling is the element that provides the cost efficiency. The Big Data  
 419 concepts and definitions in Table 1 are not comprehensive, but rather illustrate the inter-related concepts  
 420 attributed to the catch-all term Big Data.

421

**Table 1: Sampling of Concepts Attributed to Big Data**

Concept	Author	Definition
<b>4Vs (Volume, Variety, Velocity, and Variability) and Engineering</b>	Gartner <sup>7,8</sup>	“Big data is high-volume, high-velocity and high-variety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making.”
<b>Volume</b>	Techtarget <sup>9</sup>	“Although Big data doesn't refer to any specific quantity, the term is often used when speaking about petabytes and exabytes of data.”
	Oxford English Dictionary (OED) <sup>10</sup>	“big data n. Computing (also with capital initials) data of a very large size, typically to the extent that its manipulation and management present significant logistical challenges; (also) the branch of computing involving such data.”
<b>Bigger Data</b>	Annette Greiner <sup>9</sup>	“Big data is data that contains enough observations to demand unusual handling because of its sheer size, though what is unusual changes over time and varies from one discipline to another.”
<b>Not Only Volume</b>	Quentin Hardy <sup>9</sup>	“What’s ‘big’ in big data isn’t necessarily the size of the databases, it’s the big number of data sources we have, as digital sensors and behavior trackers migrate across the world.”
	Chris Neumann <sup>9</sup>	“...our original definition was a system that (1) was capable of storing 10 TB of data or more ... As time went on, diversity of data started to become more prevalent in these systems (particularly the need to mix structured and unstructured data), which led to more widespread adoption of the “3 Vs” (volume, velocity, and variety) as a definition for big data.”
<b>Big Data Engineering</b>	IDC <sup>11</sup> [16]	“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.”
	Hal Varian <sup>9</sup>	“Big data means data that cannot fit easily into a standard relational database.”
	McKinsey <sup>12</sup>	“Big Data refers to a dataset whose size is beyond the ability of typical database software tools to capture, store, manage, and analyze.”
<b>Less Sampling</b>	John Foreman <sup>9</sup>	“Big data is when your business wants to use data to solve a problem, answer a question, produce a product, etc., ...crafting a solution to the problem that leverages the data without simply sampling or tossing out records.”
	Peter Skomoroch <sup>9</sup>	“Big data originally described the practice in the consumer Internet industry of applying algorithms to increasingly large amounts of disparate data to solve problems that had suboptimal solutions with smaller datasets.”
<b>New Data Types</b>	Tom Davenport <sup>13</sup>	“The broad range of new and massive data types that have appeared over the last decade or so.”
	Mark van Rijmenam <sup>9</sup>	“Big data is not all about volume, it is more about combining different data sets and to analyze it in real-time to get insights for your organization. Therefore, the right definition of big data should in fact be: mixed data.”

Concept	Author	Definition
<b>Analytics</b>	Ryan Swanson <sup>9</sup>	“Big data used to mean data that a single machine was unable to handle. Now big data has become a buzzword to mean anything related to data analytics or visualization.”
<b>Data Science</b>	Joel Gurin <sup>9</sup>	“Big data describes datasets that are so large, complex, or rapidly changing that they push the very limits of our analytical capability.”
	Josh Ferguson <sup>9</sup>	“Big data is the broad name given to challenges and opportunities we have as data about every aspect of our lives becomes available. It’s not just about data though; it also includes the people, processes, and analysis that turn data into meaning.”
<b>Value</b>	Harlan Harris <sup>9</sup>	“To me, ‘big data’ is the situation where an organization can (arguably) say that they have access to what they need to reconstruct, understand, and model the part of the world that they care about.”
	Jessica Kirkpatrick <sup>9</sup>	“Big data refers to using complex datasets to drive focus, direction, and decision making within a company or organization.”
	Hilary Mason <sup>9</sup>	“Big data is just the ability to gather information and query it in such a way that we are able to learn things about the world that were previously inaccessible to us.”
	Gregory Piatetsky-Shapiro <sup>9</sup>	“The best definition I saw is, “Data is big when data size becomes part of the problem.” However, this refers to the size only. Now the buzzword “big data” refers to the new data-driven paradigm of business, science and technology, where the huge data size and scope enables better and new services, products, and platforms.”
<b>Cultural Change</b>	Drew Conway <sup>9</sup>	“Big data, which started as a technological innovation in distributed computing, is now a cultural movement by which we continue to discover how humanity interacts with the world—and each other—at large-scale.”
	Daniel Gillick <sup>9</sup>	““Big data’ represents a cultural shift in which more and more decisions are made by algorithms with transparent logic, operating on documented immutable evidence. I think ‘big’ refers more to the pervasive nature of this change than to any particular amount of data.”
	Cathy O’Neil <sup>9</sup>	““Big data’ is more than one thing, but an important aspect is its use as a rhetorical device, something that can be used to deceive or mislead or overhype.”

422

423

## 424 3 BIG DATA FEATURES

---

425 The diversity of Big Data concepts discussed in Section 2 is similarly reflected in the discussion of Big  
 426 Data features in Section 3. Some Big Data terms and concepts are discussed in Section 3 to understand  
 427 new aspects brought about by the Big Data paradigm in the context of existing data architecture and  
 428 analysis context.

### 429 3.1 DATA ELEMENTS AND METADATA

430 Individual data elements have not changed with Big Data and are not discussed in detail in this document.  
 431 For additional information on data types, readers are directed to the ISO standard ISO/IEC 11404:2007  
 432 General Purpose Datatypes<sup>14</sup>, and, as an example, its extension into healthcare information data types in  
 433 ISO 21090:2011 Health Informatics<sup>15</sup>.

434 One important concept to Big Data is metadata, which is often described as “data about data.” Metadata  
 435 describes additional information about the data such as how and when data was collected and how it has  
 436 been processed. Metadata should itself be viewed as data with all the requirements for tracking, change  
 437 management, and security. Many standards are being developed for metadata, for general metadata  
 438 coverage (e.g., ISO/IEC 11179-x<sup>16</sup>) and discipline specific metadata (e.g., ISO 19115-x<sup>17</sup> for geospatial  
 439 data).

440 Metadata that describes the history of a dataset is called its *provenance*, which is discussed in Section 3.6.  
 441 As *open data* (data available to others) and *linked data* (data that is connected to other data) become the  
 442 norm, it is increasingly important to have information about how data was collected, transmitted, and  
 443 processed. Provenance type of metadata guides users to correct data utilization when the data is  
 444 repurposed from its original collection process in an effort to extract additional value.

445 *Semantic metadata*, another type of metadata, refers to the definitional description of a data element to  
 446 assist with proper interpretation. An *ontology* can be conceptualized as a graphic model, representing a  
 447 semantic relationship between entities. Ontologies are semantic models constrained to follow different  
 448 levels of logic models. Ontologies and semantic models predated Big Data and not discussed in depth this  
 449 document. Ontologies can be very general or extremely domain specific in nature. A number of  
 450 mechanisms exist for implementing these unique definitional descriptions, and the reader is referred to the  
 451 World Wide Web Consortium (W3C) efforts on the semantic web<sup>1819</sup> for additional information. Semantic  
 452 data is important in the new Big Data Paradigm since the Semantic Web represents a Big Data attempt to  
 453 provide cross-cutting meanings for terms. Again, semantic metadata is especially important for linked  
 454 data efforts.

455 *Taxonomies* represent in some sense metadata about data element relationships. Taxonomy is a  
 456 hierarchical relationship between entities, where a data element is broken down into smaller component  
 457 parts. While these concepts are important, they predated the Big Data paradigm shift.

### 458 3.2 DATA RECORDS AND NON-RELATIONAL MODELS

459 Data elements are collected into records that describe a particular observation, event, or transaction.  
 460 Previously, most of the data in business systems was *structured* data, where each record was consistently  
 461 structured and could be described efficiently in a *relational model*. Records are conceptualized as the  
 462 rows in a table where data elements are in the cells. Unstructured data types, such as text, image, video,  
 463 and relationship data, have been increasing in both volume and prominence. While modern relational  
 464 databases tend to have support for these types of data elements, their ability to directly analyze, index, and  
 465 process them has tended to be both limited and accessed via non-standard SQL extensions. The need to

466 analyze *unstructured* or *semi-structured* data has been present for many years. However, the Big Data  
 467 paradigm shift has increased the emphasis on the value of unstructured or relationship data, and also on  
 468 different engineering methods that can handle data more efficiently.

469 Again, semantic metadata is, Big Data Engineering refers to the new ways data is stored in records. In  
 470 some cases the records are still in the concept of a table structure. One storage paradigm is a key-value  
 471 structure, with a record consisting of a key and a string of data together in the value. The data is retrieved  
 472 through the key, and the non-relational database software handles accessing the data in the value. This can  
 473 be viewed as a subset/simplification of a relational database table with a single index field and column. A  
 474 variant on this is the document store, where the document has multiple value fields, any of which can be  
 475 used as the index/key. The difference from the relational table model is that the set of documents do not  
 476 all need to have same value fields.

477 Another type of new Big Data record storage is in a graphical model. A graphical model represents the  
 478 relationship between data elements. The data elements are nodes, and the relationship is represented as a  
 479 link between nodes. Graph storage models represent each data element as a series of subject, predicate,  
 480 and object triples. Often, the available types of objects and relationships are described via ontologies as  
 481 discussed above.

482 Another data element relationship concept that is not new in the Big Data paradigm shift is the presence  
 483 of *complexity* between the data elements. There are systems where data elements cannot be analyzed  
 484 outside the context of other data elements. This is evident, for example, in the analytics for the Human  
 485 Genome Project, where it is the relationship between the elements and their position and proximity to  
 486 other elements that matters. The term *complexity* is often attributed to Big Data, but it refers to this inter-  
 487 relationship between data elements or across data records, independent of whether the dataset has the  
 488 characteristics of Big Data

### 489 **3.3 DATASET CHARACTERISTICS AND STORAGE**

490 Data records are grouped into datasets, which can have the Big Data characteristics of volume, velocity,  
 491 variety, and variability. Dataset characteristics can refer to the data itself, or *data at rest*, while  
 492 characteristics of the data that is traversing a network or temporarily residing in computer memory to be  
 493 read or updated is referred to as *data in motion*, which is discussed in Section 3.4.

494 ***Data at Rest:*** Typical characteristics of data at rest that are notably different in the era of Big Data are  
 495 volume and variety. Volume is the characteristic of data at rest that is most associated with Big Data.  
 496 Estimates show that the amount of data in the world doubles every two years.<sup>20</sup> Should this trend  
 497 continue, by 2020 there would be 500 times the amount of data as existed in 2011. The sheer volume of  
 498 the data is colossal. The data volumes have stimulated new ways for scalable storage across a collection  
 499 of horizontally coupled resources, as described in Section 2.1.

500 The second characteristic of data at rest is the increasing need to use a variety of data, meaning the data  
 501 represents a number of data domains and a number of data types. Traditionally, a variety of data was  
 502 handled through transformations or pre-analytics to extract features that would allow integration with  
 503 other data. The wider range of data formats, logical models, timescales, and semantics, which is desirous  
 504 to use in analytics, complicates the integration of the variety of data. For example, data to be integrated  
 505 could be text from social networks, image data, or a raw feed directly from a sensor source. To deal with  
 506 a wider range of data formats, a federated database model was designed as a database across the  
 507 underlying databases. Data to be integrated for analytics could now be of such volume that it cannot be  
 508 moved to integrate, or it may be that some of the data is not under control of the organization creating the  
 509 data system. In either case, the variety of Big Data forces a range of new Big Data engineering solutions  
 510 to efficiently and automatically integrate data that is stored across multiple repositories, in multiple  
 511 formats, and in multiple logical data models.

512 Big Data engineering has spawned data storage models that are more efficient for unstructured data than  
513 the traditional relational model, causing a derivative issue for the mechanisms to integrate this data. New  
514 scalable techniques have arisen to manage and manipulate Big Data not stored in traditional expensive  
515 high-performance “vertically” scaled systems, but rather spread across a number of less expensive  
516 resources. For example, the document store was developed specifically to support the idea of storing and  
517 indexing heterogeneous data in a common repository for analysis. New types of non-relational storage for  
518 data records are discussed below.

519 **Shared-disk File Systems:** These approaches, such as Storage Area Networks (SANs) and Network  
520 Attached Storage (NAS), use a single storage pool, which is accessed from multiple computing resources.  
521 While these technologies solved many aspects of accessing very large datasets from multiple nodes  
522 simultaneously, they suffered from issues related to data locking and updates and, more importantly,  
523 created a performance bottleneck (from every input/output [I/O] operation accessing the common storage  
524 pool) that limited their ability to scale up to meet the needs of many Big Data applications. These  
525 limitations were overcome through the implementation of fully *distributed file systems*.

526 **Distributed File Systems:** In distributed file storage systems, multi-structured (object) datasets are  
527 distributed across the computing nodes of the server cluster(s). The data may be distributed at the  
528 file/dataset level, or more commonly, at the block level, allowing multiple nodes in the cluster to interact  
529 with different parts of a large file/dataset simultaneously. Big Data frameworks are frequently designed to  
530 take advantage of data locality to each node when distributing the processing, which avoids any need to  
531 move the data between nodes. In addition, many distributed file systems also implement file/block level  
532 replication where each file/block is stored multiple times on different machines for both  
533 reliability/recovery (data is not lost if a node in the cluster fails), as well as enhanced data locality. Any  
534 type of data and many sizes of files can be handled without formal extract, transformation, and load  
535 conversions, with some technologies performing markedly better for large file sizes.

536 **Distributed Computing:** The popular framework for distributed computing consists of a storage layer and  
537 processing layer combination that implements a multiple-class, algorithm-programming model. Low cost  
538 servers supporting the distributed file system that stores the data can dramatically lower the storage costs  
539 of computing on a large scale of data (e.g., web indexing). **MapReduce** is the default processing  
540 component in data-distributed computing. Processing results are typically then loaded into an analysis  
541 environment.

542 The use of inexpensive servers is appropriate for slower, batch-speed Big Data applications, but do not  
543 provide good performance for applications requiring low latency processing. The use of basic MapReduce  
544 for processing places limitations on updating or iterative access to the data during computation. Bulk  
545 Synchronous Parallelism systems or newer MapReduce developments can be used when repeated  
546 updating is a requirement. Improvements and “generalizations” of MapReduce have been developed that  
547 provide additional functions lacking in the older technology, including fault tolerance, iteration flexibility,  
548 elimination of middle layer, and ease of query.

549 **Resource Negotiation:** The common distributed computing system has little in the way of built-in data  
550 management capabilities. In response, several technologies have been developed to provide the necessary  
551 support functions, including operations management, workflow integration, security, and governance. Of  
552 special importance to resource management development, are new features for supporting additional  
553 processing models (other than MapReduce) and controls for multi-tenant environments, higher  
554 availability, and lower latency applications.

555 In a typical implementation, the resource manager is the hub for several node managers. The client or user  
556 accesses the resource manager which in turn launches a request to an application master within one or  
557 many node managers. A second client may also launch its own requests, which will be given to other  
558 application masters within the same or other node managers. Tasks are assigned a priority value allocated  
559 based on available CPU and memory, and provided the appropriate processing resource in the node.

560 Data movement is normally handled by transfer and application program interface (API) technologies  
 561 other than the resource manager. In rare cases, peer-to-peer (P2P) communications protocols can also  
 562 propagate or migrate files across networks at scale, meaning that technically these P2P networks are also  
 563 distributed file systems. The largest social networks, arguably some of the most dominant users of Big  
 564 Data, move binary large objects (BLOBs) of over 1 gigabyte (GB) in size internally over large numbers of  
 565 computers via such technologies. The internal use case has been extended to private file synchronization,  
 566 where the technology permits automatic updates to local folders whenever two end users are linked  
 567 through the system.

568 In external use cases, each end of the P2P system contributes bandwidth to the data movement, making  
 569 this currently the fastest way to leverage documents to the largest number of concurrent users. For  
 570 example, NASA (U.S. National Aeronautics and Space Administration) uses this technology to make  
 571 3GB images available to the public. However, any large bundle of data (e.g., video, scientific data) can be  
 572 quickly distributed with lower bandwidth cost.

573 There are additional aspects of Big Data that are changing rapidly and are not fully explored in this  
 574 document, including cluster management and other mechanisms for providing communication among the  
 575 cluster resources holding the data in the non-relational models. Discussion of the use of multiple tiers of  
 576 storage (e.g., in-memory, cache, solid state drive, hard drive, network drive) in the newly emerging  
 577 software defined storage can be found in other industry publications. Software defined storage is the use  
 578 of software to determine the dynamic allocation of tiers of storage to reduce storage costs while  
 579 maintaining the required data retrieval performance.

### 580 3.4 DATA IN MOTION

581 Another important characteristic of Big Data is the time window in which the analysis can take place.  
 582 Data in motion is processed and analyzed in real time, or near-real time, and has to be handled in a very  
 583 different way than data at rest (i.e., persisted data). Data in motion tends to resemble event-processing  
 584 architectures, and focuses on real-time or operational intelligence applications.

585 Typical characteristics of data in motion that are significantly different in the era of Big Data are velocity  
 586 and variability. The velocity is the rate of flow at which the data is created, stored, analyzed, and  
 587 visualized. Big Data velocity means a large quantity of data is being processed in a short amount of time.  
 588 In the Big Data era, data is created and passed on in real time or near real time. Increasing data flow rates  
 589 create new challenges to enable real- or near real-time data usage. Traditionally this concept has been  
 590 described as *streaming data*. While these aspects are new for some industries, other industries (e.g.,  
 591 telecommunications) have processed high volume and short time interval data for years. However, the  
 592 new in-parallel scaling approaches do add new Big Data engineering options for efficiently handling this  
 593 data.

594 The second characteristic for data in motion is variability, which refers to any change in data over time,  
 595 including the flow rate, the format, or the composition. Given that many data processes generate a surge  
 596 in the amount of data arriving in a given amount of time, new techniques are needed to efficiently handle  
 597 this data. The data processing is often tied up with the automatic provisioning of additional virtualized  
 598 resources in a cloud environment. Detailed discussions of the techniques used to process data can be  
 599 found in other industry publications that focus on operational cloud architectures.<sup>21 22</sup> Early Big Data  
 600 systems built by Internet search providers and others were frequently deployed on bare metal to achieve  
 601 the best efficiency at distributing I/O across the clusters and multiple storage devices. While cloud (i.e.,  
 602 virtualized) infrastructures were frequently used to test and prototype Big Data deployments, there are  
 603 recent trends, due to improved efficiency in I/O virtualization infrastructures, of production solutions  
 604 being deployed on cloud or Infrastructure-as-a-Service (IaaS) platforms. A high velocity system with high  
 605 variability may be deployed on a cloud infrastructure, because of the cost and performance efficiency of  
 606 being able to add or remove nodes to handle the peak performance. Being able to release those resources

607 when they are no longer needed provides significant cost savings for operating this type of Big Data  
 608 system. Very large implementations and in some cases cloud providers are now implementing this same  
 609 type of elastic infrastructure on top of their physical hardware. This is especially true for organizations  
 610 that already need extensive infrastructure but simply need to balance resources across application  
 611 workloads that can vary.

### 612 3.5 DATA SCIENCE LIFECYCLE MODEL FOR BIG DATA

613 As was introduced in Section 2.1, the data lifecycle consists of the following four stages:

- 614 1. **Collection:** This stage gathers and stores data in its original form (i.e., raw data.)
- 615 2. **Preparation:** This stage involves the collection of processes that convert raw data into cleansed,  
 616 organized information.
- 617 3. **Analysis:** This stage involves the techniques that produce synthesized knowledge from organized  
 618 information.
- 619 4. **Action:** This stage involves processes that use the synthesized knowledge to generate value for  
 620 the enterprise.

621 In the traditional data warehouse, the data handling process followed the order above (i.e., collection,  
 622 preparation, storage, and analysis.) The relational model was designed in a way that optimized the  
 623 intended analytics. The different Big Data characteristics have influenced changes in the ordering of the  
 624 data handling processes. Examples of these changes are as follows:

- 625 • **Data warehouse:** Persistent storage occurs after data preparation
- 626 • **Big Data volume system:** Data is stored immediately in raw form before preparation; preparation  
 627 occurs on read, and is referred to as ‘schema on read’
- 628 • **Big Data velocity application:** The collection, preparation, and analytics (alerting) occur on the  
 629 fly, and possibly includes some summarization or aggregation prior to storage

630 Just as simulations split the analytical processing across clusters of processors, data processes are  
 631 redesigned to split data transformations across data nodes. Because the data may be too big to move, the  
 632 transformation code may be sent in parallel across the data persistence nodes, rather than the data being  
 633 extracted and brought to the transformation servers.

### 634 3.6 BIG DATA ANALYTICS

635 Analytic processes are often characterized as **discovery** for the initial hypothesis formulation,  
 636 **development** for establishing the analytics process for a specific hypothesis, and **applied** for the  
 637 encapsulation of the analysis into an operational system. While Big Data has touched all three types of  
 638 analytic processes, the majority of the changes is observed in development and applied analytics. New  
 639 Big Data engineering technologies change the types of analytics that are possible, but do not result in  
 640 completely new types of analytics. However, given the retrieval speeds, analysts are able to interact with  
 641 their data in ways that were not previously possible. Traditional statistical analytic techniques downsize,  
 642 sample, or summarize the data before analysis. This was done to make analysis of large datasets  
 643 reasonable on hardware that could not scale to the size of the dataset. Big Data analytics often emphasize  
 644 the value of computation across the entire dataset, which gives analysts better chances to determine  
 645 causation, rather than just correlation. Correlation, though, is still useful when knowing the direction or  
 646 trend of something is enough to take action. Today, most analytics in statistics and data mining focus on  
 647 causation—being able to describe why something is happening. Discovering the cause aids actors in  
 648 changing a trend or outcome. Actors, which in system development can represent individuals,  
 649 organizations, software, or hardware, are discussed in *NIST Big Data Interoperability Framework:  
 650 Volume 2, Taxonomy*. Big Data solutions make it more feasible to implement causation type of complex  
 651 analytics for large, complex, and heterogeneous data.



652 In addition to volume, velocity, variety, and variability, several terms, many beginning with V, have been  
 653 used in connection with Big Data requirements for the system architecture. Some of these terms strongly  
 654 relate to analytics on the data. Veracity and provenance are two such terms and are discussed below.

655 Veracity refers to the completeness and accuracy of the data and relates to the vernacular “garbage-in,  
 656 garbage-out” description for data quality issues in existence for a long time. If the analytics are causal,  
 657 then the quality of every data element is extremely important. If the analytics are correlations or trending  
 658 over massive volume datasets, then individual bad elements could be lost in the overall counts and the  
 659 trend will still be accurate. As mentioned in Section 2.2, many people debate whether “more data is  
 660 superior to better algorithms,” but that is a topic better discussed elsewhere.

661 As discussed in Section 3.1, the provenance, or history of the data, is increasingly an essential factor in  
 662 Big Data analytics, as more and more data is being repurposed for new types of analytics in completely  
 663 different disciplines from which the data was created. As the usage of data persists far beyond the control  
 664 of the data producers, it becomes ever more essential that metadata about the full creation and processing  
 665 history is made available along with the data. In addition, it is vital to know what analytics may have  
 666 produced the data, since there are always confidence ranges, error ranges, and precision/recall limits  
 667 associated with analytic outputs.

668 Another analytics consideration is the speed of interaction between the analytics processes and the person  
 669 or process responsible for delivering the actionable insight. Analytic data processing speed can fall along  
 670 a continuum between batch and streaming oriented processing. Although the processing continuum  
 671 existed prior to the era of Big Data, the desired location on this continuum is a large factor in the choice  
 672 of architectures and component tools to be used. Given the greater query and analytic speeds within Big  
 673 Data due to the scaling across a cluster, there is an increasing emphasis on interactive (i.e., real-time)  
 674 processing. Rapid analytics cycles allow an analyst to do exploratory discovery on the data, browsing  
 675 more of the data space than might otherwise have been possible in any practical time frame. The  
 676 processing continuum is further discussed in *NIST Big Data Interoperability Framework: Volume 6,  
 677 Reference Architecture*.

### 678 **3.7 BIG DATA METRICS AND BENCHMARKS**

679 Initial considerations in the use of Big Data engineering include the determination, for a particular  
 680 situation, of the size threshold after which data should be considered Big Data. Multiple factors must be  
 681 considered in this determination and the outcome is particular to each application. As described in Section  
 682 2.1, Big Data characteristics lead to use of Big Data engineering techniques to allow the data system to  
 683 operate affordably and efficiently. Whether a performance or cost efficiency can be attained for a  
 684 particular application requires a design analysis, which is beyond the scope of this report.

685 There is a significant need for metrics and benchmarking to provide standards for the performance of Big  
 686 Data systems. This topic is being addressed by the Transaction Processing Performance Council TCP-  
 687 xHD Big Data Committee, and available information from their efforts may be included in future versions  
 688 of this report.

### 689 **3.8 BIG DATA SECURITY AND PRIVACY**

690 Security and privacy have also been affected by the emergence of the Big Data paradigm. A detailed  
 691 discussion of the influence of Big Data on security and privacy is included in NIST Big Data  
 692 Interoperability Framework: Volume 4, Security and Privacy. Some of the effects of Big Data  
 693 characteristics on security and privacy summarized below:

- 694 • **Variety:** Retargeting traditional relational database security to non-relational databases has been  
 695 a challenge. An emergent phenomenon introduced by Big Data variety that has gained

696 considerable importance is the ability to infer identity from anonymized datasets by correlating  
 697 with apparently innocuous public databases.

- 698 • **Volume:** The volume of Big Data has necessitated storage in multi-tiered storage media. The  
 699 movement of data between tiers has led to a requirement of systematically analyzing the threat  
 700 models and research and development of novel techniques.
- 701 • **Velocity:** As with non-relational databases, distributed programming frameworks such as Hadoop  
 702 were not developed with security as a primary objective.
- 703 • **Veracity:** Complex challenges have been introduced in protecting data integrity as well as  
 704 maintaining privacy policies as data moves across individual boundaries to groups, communities  
 705 of interest, state, national, and international boundaries.
- 706 • **Volatility:** Security and privacy requirements can shift according to the time dependent nature of  
 707 roles that collected, processed, aggregated, and stored it. Governance can shift as responsible  
 708 organizations merge or even disappear

709 Privacy concerns, and frameworks to address these concerns, predate Big Data. While bounded in  
 710 comparison to Big Data, past solutions considered legal, social, and technical requirements for privacy in  
 711 distributed systems, very large databases, and in HPCC. The addition of variety, volume, velocity,  
 712 veracity, volatility, and value to the mix has amplified these concerns to the level of a national  
 713 conversation, with unanticipated impacts on privacy frameworks.

### 714 3.9 DATA GOVERNANCE

715 Data governance is a fundamental element in the management of data and data systems.

716 *Data governance refers to administering, or formalizing, discipline (e.g., behavior*  
 717 *patterns) around the management of data.*

718 The definition of data governance includes management across the complete data lifecycle, whether the  
 719 data is at rest, in motion, in incomplete stages, or transactions. To maximize its benefit, data governance  
 720 must also consider the issues of privacy and security of individuals of all ages, individuals as companies,  
 721 and companies as companies.

722 Data governance is needed to address important issues in the new global Internet Big Data economy. For  
 723 example, many businesses provide a data hosting platform for data that is generated by the users of the  
 724 system. While governance policies and processes from the point of view of the data hosting company are  
 725 commonplace, the issue of governance and control rights of the data providers is new. Many questions  
 726 remain including the following. Do they still own their data, or is the data owned by the hosting  
 727 company? Do the data producers have the ability to delete their data? Can they control who is allowed to  
 728 see their data?

729 The question of governance resides between the value that one party (e.g., the data hosting company)  
 730 wants to generate versus the rights that the data provider wants to retain to obtain their own value. New  
 731 governance concerns arising from the Big Data Paradigm need greater discussion, and will be discussed  
 732 during the development of the next version of this document.

733

734 **4 BIG DATA ENGINEERING PATTERNS (FUNDAMENTAL**  
735 **CONCEPTS)**

---

736 To define the differences between Big Data technologies, different ‘scenarios’ and ‘patterns’ are needed  
737 to illustrate relationships between Big Data characteristics (Section 2.1) and between the NBDRA  
738 components found in *NIST Big Data Interoperability Framework: Volume 6, Reference Architecture*. The  
739 scenarios would describe the high-level functional processes that can be used to categorize and, therefore,  
740 provide better understanding of the different use cases presented in *NIST Big Data Interoperability*  
741 *Framework: Volume 3, Use Cases and General Requirements*, as well as help to clarify the differences in  
742 specific implementations of components listed in the *NIST Big Data Interoperability Framework: Volume*  
743 *6, Reference Architecture*.

744 The topics surrounding the relaxation of the principles of a relational model in non-relational systems are  
745 very important. These topics are discussed in industry publications on concurrency, and will be addressed  
746 more fully in of future additions to this document.

747

DRAFT

748 **Appendix A: Index of Terms**

---

749  
750  
751

- A**  
analytics, 7
- B**  
Big Data, 5, 6  
Big Data engineering, 5  
Big Data paradigm, 4  
Big Data velocity application, 14  
Big Data volume system, 14
- C**  
complexity, 10  
Computational portability, 6
- D**  
data lifecycle, 6  
data sampling, 8  
data science, 8  
data science paradigm, 6  
data scientist, 7  
data warehouse, 14
- F**  
federated database system, 5  
fourth paradigm, 6
- M**  
massively parallel processing, 4  
metadata, 10
- N**  
non-relational models, 5NoSQL, 5
- O**  
ontologies, 10
- P**  
provenance, 13
- R**  
relational model, 10
- S**  
Schema-on-read, 6  
semantic data, 10  
semi-structured data, 10  
streaming data, 12  
structured data, 10
- T**  
taxonomies, 10
- U**  
unstructured data, 10
- V**  
validity, 8  
value, 8  
variability, 4  
variety, 4  
velocity, 10  
veracity, 8  
volatility, 8  
volume, 4

## 731 **Appendix B: Terms and Definitions**

---

732 *Analytics* is the synthesis of knowledge from information.

733 *Big Data* consists of extensive datasets—primarily in the characteristics of volume, variety, velocity,  
734 and/or variability—that require a scalable architecture for efficient storage, manipulation, and analysis.

735 *Big Data engineering* includes advanced techniques that harness independent resources for building  
736 scalable data systems when the characteristics of the datasets require new architectures for efficient  
737 storage, manipulation, and analysis.

738 The *Big Data paradigm* consists of the distribution of data systems across horizontally coupled,  
739 independent resources to achieve the scalability needed for the efficient processing of extensive  
740 datasets.

741 *Computational portability* is the movement of the computation to the location of the data.

742 *Data governance* refers to the overall management of the availability, usability, integrity, and  
743 security of the data employed in an enterprise.

744 The *data lifecycle* is the set of processes that transforms raw data into actionable knowledge, which  
745 includes data collection, preparation, analytics, visualization, and access.

746 *Data science* is the empirical synthesis of actionable knowledge from raw data through the complete data  
747 lifecycle process.

748 The *data science paradigm* is extraction of actionable knowledge directly from data through a process of  
749 discovery, hypothesis, and hypothesis testing.

750 A *Latency* is a practitioner who has sufficient knowledge in the overlapping regimes of business needs,  
751 domain knowledge, analytical skills, and software and systems engineering to manage the end-to-end  
752 data processes through each stage in the data lifecycle.

753 *Distributed Computing* is a computing system in which components located on networked  
754 computers communicate and coordinate their actions by passing messages.

755 *Distributed File Systems* contain multi-structured (object) datasets that are distributed across the  
756 computing nodes of the server cluster(s).

757 A *federated database system* is a type of meta-database management system, which transparently maps  
758 multiple autonomous database systems into a single federated database.

759 *horizontal scaling* implies the coordination of individual resources (e.g., server) that are integrated to act  
760 in parallel as a single system (i.e., operate as a cluster).

761 *Latency* refers to the delay in processing or in availability.

762 *Massively parallel processing* refers to a multitude of individual processors working in parallel to execute  
763 a particular program.

764 *Non-relational models*, frequently referred to as NoSQL, refer to logical data models that do not follow  
765 relational algebra for the storage and manipulation of data.

766 *Resource Negotiation* consists of built-in data management capabilities that provide the necessary  
767 support functions, such as operations management, workflow integration, security, governance, support  
768 for additional processing models, and controls for multi-tenant environments, providing higher  
769 availability and lower latency applications.

770 **Schema-on-read** is the application of a data schema through preparation steps such as transformations,  
771 cleansing, and integration at the time the data is read from the database.

772 **Shared-disk File Systems**, such as Storage Area Networks (SANs) and Network Attached Storage (NAS),  
773 use a single storage pool, which is accessed from multiple computing resources.

774 **validity** refers to appropriateness of the data for its intended use

775 **value** refers to the inherent wealth, economic and social, embedded in any data set

776 **variability** refers to the change in other data characteristics

777 **variety** refers to data from multiple repositories, domains, or types

778 **velocity** refers to the rate of data flow

779 **veracity** refers to the accuracy of the data

780 **Vertical scaling** implies increasing the system parameters of processing speed, storage, and memory for  
781 greater performance.

782 **volatility** refers to the tendency for data structures to change over time

783 **volume** refers to the size of the dataset

784

## 785 **Appendix C: Acronyms**

---

786	API	application program interface
787	BLOBs	binary large objects
788	GB	gigabyte
789	I/O	input/output
790	ISO	International Organization for Standardization
791	ITL	Information Technology Laboratory
792	JTC 1	Joint Technical Committee 1
793	MPP	massively parallel processing
794	NARA	National Archives and Records Administration
795	NAS	Network Attached Storage
796	NASA	National Aeronautics and Space Administration
797	NBD-PWG	NIST Big Data Public Working Group
798	NBDRA	NIST Big Data Reference Architecture
799	NIST	National Institute of Standards and Technology
800	NSF	National Science Foundation
801	OED	Oxford English Dictionary
802	P2P	peer-to-peer
803	SANs	Storage Area Networks
804	SQL	Structured Query Language
805	NoSQL	Not Only or No Structured Query Language
806	W3C	World Wide Web Consortium
807		

## Appendix D: References

---

### DOCUMENT REFERENCES

- <sup>1</sup> The White House Office of Science and Technology Policy, "Big Data is a Big Deal," *OSTP Blog*, accessed February 21, 2014, <http://www.whitehouse.gov/blog/2012/03/29/big-data-big-deal>.
- <sup>2</sup> Gordon Moore, "Cramming More Components Onto Integrated Circuits," *Electronics*, Volume 38, Number 8 (1965), pages 114-117.
- <sup>3</sup> Jenna Dutcher, "What is Big Data," *Data Science at Berkeley Blog*, September 3, 2014, <http://datascience.berkeley.edu/what-is-big-data/>.
- <sup>4</sup> Emerging Technology From the arXiv (Contributor), "The Big Data Conundrum: How to Define It?," MIT Technology Review, October 3, 2013, <http://www.technologyreview.com/view/519851/the-big-data-conundrum-how-to-define-it/>.
- <sup>5</sup> ISO/IEC JTC 1 Study Group on Big Data (SGBD), "N0095 Final SGBD Report to JTC1," September 3, 2014, [http://jtc1bigdatasg.nist.gov/uploadfiles/N0095\\_Final\\_SGBD\\_Report\\_to\\_JTC1.docx](http://jtc1bigdatasg.nist.gov/uploadfiles/N0095_Final_SGBD_Report_to_JTC1.docx).
- <sup>6</sup> Gil Press (Contributor), "12 Big Data Definitions: What's Yours?," *Forbes.com*, accessed November 17, 2014, <http://www.forbes.com/sites/gilpress/2014/09/03/12-big-data-definitions-whats-yours/>.
- <sup>7</sup> ISO/IEC JTC 1 Study Group on Big Data (SGBD), "N0095 Final SGBD Report to JTC1," September 3, 2014, [http://jtc1bigdatasg.nist.gov/uploadfiles/N0095\\_Final\\_SGBD\\_Report\\_to\\_JTC1.docx](http://jtc1bigdatasg.nist.gov/uploadfiles/N0095_Final_SGBD_Report_to_JTC1.docx).
- <sup>8</sup> Gartner IT Glossary, "Big Data" (definition), *Gartner.com*, accessed November 17, 2014, <http://www.gartner.com/it-glossary/big-data>.
- <sup>9</sup> Jenna Dutcher, "What is Big Data," *Data Science at Berkeley Blog*, September 3, 2014, <http://datascience.berkeley.edu/what-is-big-data/>.
- <sup>10</sup> Oxford English Dictionary, "Big Data" (definition), *OED.com*, accessed November 17, 2014, <http://www.oed.com/view/Entry/18833#eid301162178>.
- <sup>11</sup> John Gantz and David Reinsel, "Extracting Value from Chaos," *IDC iView sponsored by EMC Corp*, accessed November 17, 2014, <http://www.emc.com/collateral/analyst-reports/idc-extracting-value-from-chaos-ar.pdf>.
- <sup>12</sup> James Manyika et al., "Big data: The next frontier for innovation, competition, and productivity," McKinsey Global Institute, May 2011.
- <sup>13</sup> Tom Davenport, "Big Data@Work," Harvard Business Review Press, February 25, 2014.
- <sup>14</sup> ISO/IEC 11404:2007, "Information technology -- General-Purpose Datatypes (GPD)," *International Organization for Standardization*, [http://www.iso.org/iso/home/store/catalogue\\_ics/catalogue\\_detail\\_ics.htm?csnumber=39479](http://www.iso.org/iso/home/store/catalogue_ics/catalogue_detail_ics.htm?csnumber=39479).
- <sup>15</sup> ISO 21090:2011, "Health informatics -- Harmonized data types for information interchange," *International Organization for Standardization*, [http://www.iso.org/iso/home/store/catalogue\\_tc/catalogue\\_detail.htm?csnumber=35646](http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=35646).
- <sup>16</sup> ISO/IEC 11179-2004, Information technology -- "Metadata registries (MDR) -- Part 1: Framework," *International Organization for Standardization*, [http://www.iso.org/iso/home/store/catalogue\\_tc/catalogue\\_detail.htm?csnumber=35343](http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=35343).
- <sup>17</sup> ISO 19115-2014, "Geographic information -- Metadata -- Part 1: Fundamentals," *International Organization for Standardization*, [http://www.iso.org/iso/home/store/catalogue\\_tc/catalogue\\_detail.htm?csnumber=53798](http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=53798).
- <sup>18</sup> Phil Archer, "W3C Data Activity Building the Web of Data," *W3C*, <http://www.w3.org/2013/data/>.
- <sup>19</sup> Dan Brickley and Ivan Herman, "Semantic Web Interest Group," *W3C*, June 16, 2012, <http://www.w3.org/2001/sw/interest/>.



<sup>20</sup> EMC2, “Digital Universe,” *EMC*, accessed February 21, 2014, <http://www.emc.com/leadership/programs/digital-universe.htm>.

<sup>21</sup> Lee Badger, David Bernstein, Robert Bohn, Frederic de Vault, Mike Hogan, Michaela Iorga, Jian Mao, John Messina, Kevin Mills, Eric Simmon, Annie Sokol, Jin Tong, Fred Whiteside, and Dawn Leaf, “US Government Cloud Computing Technology Roadmap Volume I: High-Priority Requirements to Further USG Agency Cloud Computing Adoption; and Volume II: Useful Information for Cloud Adopters,” *National Institute of Standards and Technology*, October 21, 2014, <http://dx.doi.org/10.6028/NIST.SP.500-293>.

<sup>22</sup> Lee Badger, Tim Grance, Robert Patt-Corner, and Jeff Voas, “Cloud Computing Synopsis and Recommendations,” *National Institute of Standards and Technology*, May 2012, <http://csrc.nist.gov/publications/nistpubs/800-146/sp800-146.pdf>.

DRAFT