

**NIST Special Publication 1500-2r2**

---

**NIST Big Data Interoperability  
Framework:  
Volume 2, Big Data Taxonomies**

---

**Version 3**

NIST Big Data Public Working Group  
Definitions and Taxonomies Subgroup

This publication is available free of charge from:  
<https://doi.org/10.6028/NIST.SP.1500-2r2>

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

**NIST Special Publication 1500-2r2**

**NIST Big Data Interoperability  
Framework:  
Volume 2, Big Data Taxonomies**

**Version 3**

NIST Big Data Public Working Group  
Definitions and Taxonomies Subgroup  
*Information Technology Laboratory*  
*National Institute of Standards and Technology*  
*Gaithersburg, MD 20899*

This publication is available free of charge from:  
<https://doi.org/10.6028/NIST.SP.1500-2r2>

November 2019



U.S. Department of Commerce  
*Wilbur L. Ross, Jr., Secretary*

National Institute of Standards and Technology  
*Walter Copan, NIST Director and Undersecretary of Commerce for Standards and Technology*

**National Institute of Standards and Technology (NIST) Special Publication 1500-2r2**  
36 pages (November 2019)

NIST Special Publication series 1500 is intended to capture external perspectives related to NIST standards, measurement, and testing-related efforts. These external perspectives can come from industry, academia, government, and others. These reports are intended to document external perspectives and do not represent official NIST positions.

Certain commercial entities, equipment, or materials may be identified in this document 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 publications during public comment periods and provide feedback to NIST. All NIST publications are available at <http://www.nist.gov/publication-portal.cfm>.

### **Copyrights and Permissions**

Official publications of the National Institute of Standards and Technology are not subject to copyright in the United States. Foreign rights are reserved. Questions concerning the possibility of copyrights in foreign countries should be referred to the Office of Chief Counsel at NIST via email to [nistcounsel@nist.gov](mailto:nistcounsel@nist.gov).

### **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 large amount of data in the networked, digitized, sensor-laden, information-driven world. While opportunities exist with Big Data, the data can overwhelm traditional technical approaches and the growth of data 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 concepts related to Big Data. The results are reported in the *NIST Big Data Interoperability Framework* (NBDIF) series of volumes. This volume, Volume 2, contains the Big Data taxonomies developed by the NBD-PWG. These taxonomies organize the reference architecture components, fabrics, and other topics to lay the groundwork for discussions surrounding Big Data.

### Keywords

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

## Acknowledgements

This document reflects the contributions and discussions by the membership of the NBD-PWG, cochaired by Wo Chang (NIST ITL), Bob Marcus (ET-Strategies), and Chaitan Baru (San Diego Supercomputer Center; National Science Foundation). For all versions, the Subgroups were led by the following people: Nancy Grady (SAIC), Natasha Balac (SDSC), and Eugene Luster (R2AD) for the Definitions and Taxonomies Subgroup; Geoffrey Fox (Indiana University) and Tsegereda Beyene (Cisco Systems) for the Use Cases and Requirements Subgroup; Arnab Roy (Fujitsu), Mark Underwood (Krypton Brothers; Synchrony Financial), and Akhil Manchanda (GE) for the Security and Privacy Subgroup; David Boyd (InCadence Strategic Solutions), Orit Levin (Microsoft), Don Krapohl (Augmented Intelligence), and James Ketner (AT&T) for the Reference Architecture Subgroup; and Russell Reinsch (Center for Government Interoperability), David Boyd (InCadence Strategic Solutions), Carl Buffington (Vistrionix), and Dan McClary (Oracle), for the Standards Roadmap Subgroup.

The editors for this document were the following:

- **Version 1:** Nancy Grady (SAIC) and Wo Chang (NIST)
- **Version 2:** Nancy Grady (SAIC) and Wo Chang (NIST)
- **Version 3:** Nancy Grady (SAIC) and Wo Chang (NIST)

Laurie Aldape (Energetics Incorporated) and Elizabeth Lennon (NIST) provided editorial assistance across all NBDIF volumes.

NIST SP1500-2, Version 3 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<sup>a</sup> to this volume, during Version 1, Version 2, and/or Version 3 activities, by the following NBD-PWG members:

**Natasha Balac**  
*University of California, San Diego,  
Supercomputer Center*

**Chaitan Baru**  
*University of California, San Diego,  
Supercomputer Center*

**Peter Baumann**  
*rasdaman GmbH; Jacobs University  
Bremen*

**Deborah Blackstock**  
*MITRE Corporation*

**Pw McKenna Carey, III**  
*Compliance Partners,  
LLC/RALAND*

**Karen Guertler**  
*Consultant*

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

**Pavithra Kenjige**  
*PK Technologies*

**Orit Levin**  
*Microsoft*

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

**Bill Mandrick**  
*Data Tactics*

**Robert Marcus**  
*ET-Strategies*

**William Miller**  
*MaCT USA*

**Sanjay Mishra**  
*Verizon*

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

**Russell Reinsch**  
*Center for Government  
Interoperability*

**John Rogers**  
*HP*

**William Vorhies**  
*Predictive Modeling LLC*

**Mark Underwood**  
*Krypton Brothers; Synchrony  
Financial*

<sup>a</sup> “Contributors” are members of the NIST Big Data Public Working Group who dedicated great effort to prepare and gave substantial time on a regular basis to research and development in support of this document.

**Wo Chang**

*National Institute of Standards and  
Technology*

**Gary Mazzaferro**

*AlloyCloud, Inc.*

**Alicia Zuniga-Alvarado**

*Consultant*

**Yuri Demchenko**

*University of Amsterdam*

**Nancy Grady**

*SAIC*

# 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.....	3
1.3 REPORT PRODUCTION .....	3
1.4 REPORT STRUCTURE .....	3
<b>2 REFERENCE ARCHITECTURE TAXONOMY .....</b>	<b>6</b>
2.1 ACTORS AND ROLES.....	6
2.2 SYSTEM ORCHESTRATOR.....	8
2.3 DATA PROVIDER .....	10
2.4 BIG DATA APPLICATION PROVIDER .....	13
2.5 BIG DATA FRAMEWORK PROVIDER .....	16
2.6 DATA CONSUMER .....	18
2.7 MANAGEMENT FABRIC.....	19
2.8 SECURITY AND PRIVACY FABRIC.....	19
<b>3 DATA CHARACTERISTIC HIERARCHY.....</b>	<b>20</b>
3.1 DATA ELEMENTS .....	20
3.2 RECORDS .....	21
3.2.1 <i>Record Characteristics</i> .....	21
3.2.2 <i>Inter-Record Structure</i> .....	22
3.3 DATASETS.....	23
3.4 MULTIPLE DATASETS .....	23
<b>4 SUMMARY .....</b>	<b>24</b>
<b>APPENDIX A: ACRONYMS .....</b>	<b>25</b>
<b>APPENDIX B: BIBLIOGRAPHY .....</b>	<b>26</b>

## FIGURES

FIGURE 1: NBDIF DOCUMENTS NAVIGATION DIAGRAM PROVIDES CONTENT FLOW BETWEEN VOLUMES .....	5
FIGURE 2: NIST BIG DATA REFERENCE ARCHITECTURE .....	7
FIGURE 3: ROLES AND A SAMPLING OF ACTORS IN THE NBDRA TAXONOMY.....	8
FIGURE 4: SYSTEM ORCHESTRATOR ACTORS AND ACTIVITIES.....	9
FIGURE 5: DATA PROVIDER ACTORS AND ACTIVITIES .....	11
FIGURE 6: BIG DATA APPLICATION PROVIDER ACTORS AND ACTIVITIES .....	14
FIGURE 7: BIG DATA FRAMEWORK PROVIDER ACTORS AND ACTIVITIES .....	16
FIGURE 8: DATA CONSUMER ACTORS AND ACTIVITIES .....	18
FIGURE 9: BIG DATA MANAGEMENT ACTORS AND ACTIVITIES.....	19
FIGURE 10: BIG DATA SECURITY AND PRIVACY ACTORS AND ACTIVITIES .....	19
FIGURE 11: DATA CHARACTERISTIC HIERARCHY.....	20
FIGURE 12: CLASSIFICATION OF COMMON DATA STRUCTURES (RECORDS REPRESENTED ABSTRACTLY AS YELLOW ELLIPSES) .....	22

# Executive Summary

---

This *NIST Big Data Interoperability Framework (NBDIF): Volume 2, Big Data Taxonomies* was prepared by the NIST Big Data Public Working Group (NBD-PWG) Definitions and Taxonomy Subgroup to facilitate communication and improve understanding across Big Data stakeholders by describing the functional components of the NIST Big Data Reference Architecture (NBDRA). The top-level roles of the taxonomy are System Orchestrator, Data Provider, Big Data Application Provider, Big Data Framework Provider, Data Consumer, Security and Privacy, and Management. The actors and activities for each of the top-level roles are outlined in this document as well. The NBDRA taxonomy aims to describe new issues in Big Data systems but is not an exhaustive list. In some cases, exploration of new Big Data topics includes current practices and technologies to provide needed context.

The *NIST Big Data Interoperability Framework (NBDIF)* was released in three versions, which correspond to the three stages of the NBD-PWG work. Version 3 (current version) of the NBDIF volumes resulted from Stage 3 work with major emphasis on the validation of the NBDRA Interfaces and content enhancement. Stage 3 work built upon the foundation created during Stage 2 and Stage 1. The current effort documented in this volume reflects concepts developed within the rapidly evolving field of Big Data. The three stages (in reverse order) aim to achieve the following with respect to the NIST Big Data Reference Architecture (NBDRA).

- Stage 3: Validate the NBDRA by building Big Data general applications through the general interfaces;
- Stage 2: Define general interfaces between the NBDRA components; and
- Stage 1: Identify the high-level Big Data reference architecture key components, which are technology-, infrastructure-, and vendor-agnostic.

The *NBDIF* consists of nine volumes, each of which addresses a specific key topic, resulting from the work of the NBD-PWG. The nine volumes are as follows:

- Volume 1, Definitions [1]
- Volume 2, Taxonomies (this volume)
- Volume 3, Use Cases and General Requirements [2]
- Volume 4, Security and Privacy [3]
- Volume 5, Architectures White Paper Survey [4]
- Volume 6, Reference Architecture [5]
- Volume 7, Standards Roadmap [6]
- Volume 8, Reference Architecture Interfaces [7]
- Volume 9, Adoption and Modernization [8]

During Stage 1, Volumes 1 through 7 were conceptualized, organized, and written. The finalized Version 1 documents can be downloaded from the V1.0 Final Version page of the NBD-PWG website ([https://bigdatawg.nist.gov/V1\\_output\\_docs.php](https://bigdatawg.nist.gov/V1_output_docs.php)).

During Stage 2, the NBD-PWG developed Version 2 of the NBDIF Version 1 volumes, with the exception of Volume 5, which contained the completed architecture survey work that was used to inform Stage 1 work of the NBD-PWG. The goals of Stage 2 were to enhance the Version 1 content, define general interfaces between the NBDRA components by aggregating low-level interactions into high-level general interfaces, and demonstrate how the NBDRA can be used. As a result of the Stage 2 work, the need for NBDIF Volume 8 and NBDIF Volume 9 was identified and the two new volumes were created.



Version 2 of the NBDIF volumes, resulting from Stage 2 work, can be downloaded from the V2.0 Final Version page of the NBD-PWG website ([https://bigdatawg.nist.gov/V2\\_output\\_docs.php](https://bigdatawg.nist.gov/V2_output_docs.php)).

# 1 INTRODUCTION

---

## 1.1 BACKGROUND

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

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

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

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

- How is Big Data defined?
- What attributes define Big Data solutions?
- What is new in Big Data?
- What is the difference between Big Data and *bigger data* that has been collected for years?
- How is Big Data different from traditional data environments and related applications?
- What are the essential characteristics of Big Data environments?
- How do these environments integrate with currently deployed architectures?
- What are the central scientific, technological, and standardization challenges that need to be addressed to accelerate the deployment of robust, secure Big Data solutions?

Within this context, on March 29, 2012, the White House announced the Big Data Research and Development Initiative [9]. The initiative's goals include helping to accelerate the pace of discovery in science and engineering, strengthening national security, and transforming teaching and learning by improving analysts' ability to extract knowledge and insights from large and complex collections of digital data.

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

Motivated by the White House initiative and public suggestions, the National Institute of Standards and Technology (NIST) has accepted the challenge to stimulate collaboration among industry professionals to further the secure and effective adoption of Big Data. As one result of NIST's Cloud and Big Data Forum held on January 15–17, 2013, there was strong encouragement for NIST to create a public working group for the development of a Big Data Standards Roadmap. Forum participants noted that this roadmap

should define and prioritize Big Data requirements, including interoperability, portability, reusability, extensibility, data usage, analytics, and technology infrastructure. In doing so, the roadmap would accelerate the adoption of the most secure and effective Big Data techniques and technology.

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

The *NIST Big Data Interoperability Framework* (NBDIF) was released in three versions, which correspond to the three stages of the NBD-PWG work. Version 3 (current version) of the NBDIF volumes resulted from Stage 3 work with major emphasis on the validation of the NBDRA Interfaces and content enhancement. Stage 3 work built upon the foundation created during Stage 2 and Stage 1. The current effort documented in this volume reflects concepts developed within the rapidly evolving field of Big Data. The three stages (in reverse order) aim to achieve the following with respect to the NIST Big Data Reference Architecture (NBDRA).

- Stage 3: Validate the NBDRA by building Big Data general applications through the general interfaces;
- Stage 2: Define general interfaces between the NBDRA components; and
- Stage 1: Identify the high-level Big Data reference architecture key components, which are technology-, infrastructure-, and vendor-agnostic.

The *NBDIF* consists of nine volumes, each of which addresses a specific key topic, resulting from the work of the NBD-PWG. The nine volumes are as follows:

- Volume 1, Definitions [1]
- Volume 2, Taxonomies (this volume)
- Volume 3, Use Cases and General Requirements [2]
- Volume 4, Security and Privacy [3]
- Volume 5, Architectures White Paper Survey [4]
- Volume 6, Reference Architecture [10]
- Volume 7, Standards Roadmap [6]
- Volume 8, Reference Architecture Interfaces [7]
- Volume 9, Adoption and Modernization [8]

During Stage 1, Volumes 1 through 7 were conceptualized, organized, and written. The finalized Version 1 documents can be downloaded from the V1.0 Final Version page of the NBD-PWG website ([https://bigdatawg.nist.gov/V1\\_output\\_docs.php](https://bigdatawg.nist.gov/V1_output_docs.php)).

During Stage 2, the NBD-PWG developed Version 2 of the NBDIF Version 1 volumes, with the exception of Volume 5, which contained the completed architecture survey work that was used to inform Stage 1 work of the NBD-PWG. The goals of Stage 2 were to enhance the Version 1 content, define general interfaces between the NBDRA components by aggregating low-level interactions into high-level general interfaces, and demonstrate how the NBDRA can be used. As a result of the Stage 2 work, the need for NBDIF Volume 8 and NBDIF Volume 9 was identified and the two new volumes were created. Version 2 of the NBDIF volumes, resulting from Stage 2 work, can be downloaded from the V2.0 Final Version page of the NBD-PWG website ([https://bigdatawg.nist.gov/V2\\_output\\_docs.php](https://bigdatawg.nist.gov/V2_output_docs.php)).

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

The NBD-PWG Definitions and Taxonomy Subgroup focused on identifying Big Data concepts, defining terms needed to describe this new paradigm, and defining reference architecture terms. This taxonomy provides a hierarchy of the components of the reference architecture. It is designed to meet the needs of specific user groups, as follows:

For **managers**, the terms will distinguish the categorization of techniques needed to understand this changing field.

For **procurement officers**, it will provide the framework for discussing organizational needs and distinguishing among offered approaches.

For **marketers**, it will provide the means to promote Big Data solutions and innovations.

For the **technical community**, it will provide a common language to better differentiate Big Data's specific offerings.

## 1.3 REPORT PRODUCTION

This document derives from discussions in the NBD-PWG Definitions and Taxonomy Subgroup and with interested parties. This volume provides the taxonomy of the components of the NBDRA. This taxonomy was developed using a mind map representation, which provided a mechanism for multiple inputs and easy editing.

It is difficult to describe the new components of Big Data systems without fully describing the context in which they reside. The Subgroup attempted to describe only what has changed in the shift to the new Big Data paradigm, and only the components needed to clarify this shift. For example, there is no attempt to create a taxonomy of analytics techniques as these predate Big Data. This taxonomy will be a work in progress to mature as new technologies are developed and the patterns within data and system architectures are better understood.

To achieve technical and high-quality document content, this document will go through a public comments period along with NIST internal review.

## 1.4 REPORT STRUCTURE

This document provides multiple hierarchical presentations related to Big Data.

The first presentation is the taxonomy for the NBDRA, providing the terminology and definitions for the components of technical systems that implement technologies for Big Data. Section 2 introduces the NBDRA using concepts of actors and roles and the activities each performs. The NBDRA components are more fully described in the *NBDIF: Volume 6, Reference Architecture* and the *NBDIF: Volume 4, Security and Privacy* documents. Comparing the related sections in these two documents will give the reader a more complete picture of the consensus of the working groups. Illustrative examples are given to facilitate understanding of the role/actor and activity of the NBDRA. There is no expectation of completeness in the components; the intent is to provide enough context to understand the specific areas that have changed because of the new Big Data paradigm.

The second presentation in Section 3 is a hierarchical description about the data itself to understand the relationships of the new concepts related to Big Data. For clarity, a strict taxonomy is not followed; rather, data is examined at different groupings to better describe what is new with Big Data. The grouping-based description presents data elements, data records, datasets, and multiple datasets. This

examination at different groupings provides a way to easily identify the data characteristics that have driven the development of Big Data engineering technologies, as described in the *NBDIF: Volume 1, Definitions*. Note that the data hierarchy only expresses the broad overview of data at different levels of granularity to highlight the properties that drive the need for Big Data architectures.

Section 4 provides a brief summary of the document. For descriptions of the future of Big Data and opportunities to use Big Data technologies, the reader is referred to the *NBDIF: Volume 7, Standards Roadmap*. Finally, to understand how these systems are organized and integrated to meet users' needs, the reader is referred to *NBDIF: Volume 3, Use Cases and General Requirements*.

While each NBDIF volume was created with a specific focus within Big Data, all volumes are interconnected. During the creation of the volumes, information from some volumes was used as input for other volumes. Broad topics (e.g., definition, architecture) may be discussed in several volumes with each discussion circumscribed by the volume's particular focus. Arrows shown in Figure 1 indicate the main flow of information input and/or output from the volumes. Volumes 2, 3, and 5 (blue circles) are essentially standalone documents that provide output to other volumes (e.g., to Volume 6). These volumes contain the initial situational awareness research. During the creation of Volumes 4, 7, 8, and 9 (green circles), input from other volumes was used. The development of these volumes took into account work on the other volumes. Volumes 1 and 6 (red circles) were developed using the initial situational awareness research and continued to be modified based on work in other volumes. The information from these volumes was also used as input to the volumes in the green circles.



Figure 1: NBDIF Documents Navigation Diagram Provides Content Flow Between Volumes

## 2 REFERENCE ARCHITECTURE TAXONOMY

---

This section focuses on a taxonomy for the NBDRA, and is intended to describe the hierarchy of actors and roles and the activities the actors perform in those roles. There are a number of models for describing the technologies needed for an application, such as a layer model of network, hardware, operating system, and application. For elucidating the taxonomy, a hierarchy has been chosen to allow placing the new technologies within the context of previous technologies. As this taxonomy is not definitive, it is expected that the taxonomy will mature as new technologies emerge and increase understanding of how to best categorize the different methods for building data systems.

### 2.1 ACTORS AND ROLES

In system development, actors and roles have the same relationship as in the movies. The roles are the parts the actors play in the overall system. One actor can perform multiple roles. Likewise, a role can be played by multiple actors, in the sense that a team of independent entities—perhaps from independent organizations—may be used to satisfy end-to-end system requirements. System development actors can represent individuals, organizations, software, or hardware. Each activity in the taxonomy can be executed by a different actor. Examples of actors include the following:

- Sensors
- Applications
- Software agents
- Individuals
- Organizations
- Hardware resources
- Service abstractions

In the past, data systems tended to be hosted, developed, and deployed with the resources of only one organization. Currently, roles may be distributed, analogous to the diversity of actors within a given cloud-based application. Actors in Big Data systems can likewise come from multiple organizations.

Developing the reference architecture taxonomy began with a review of the NBD-PWG analyses of the use cases and reference architecture survey provided in *NBDIF: Volume 3, Use Cases and General Requirements* and *NBDIF: Volume 5, Reference Architecture Survey*, respectively. From these analyses, several commonalities between Big Data architectures were identified and formulated into five general architecture components, and two fabrics interwoven in the five components, as shown in Figure 2.



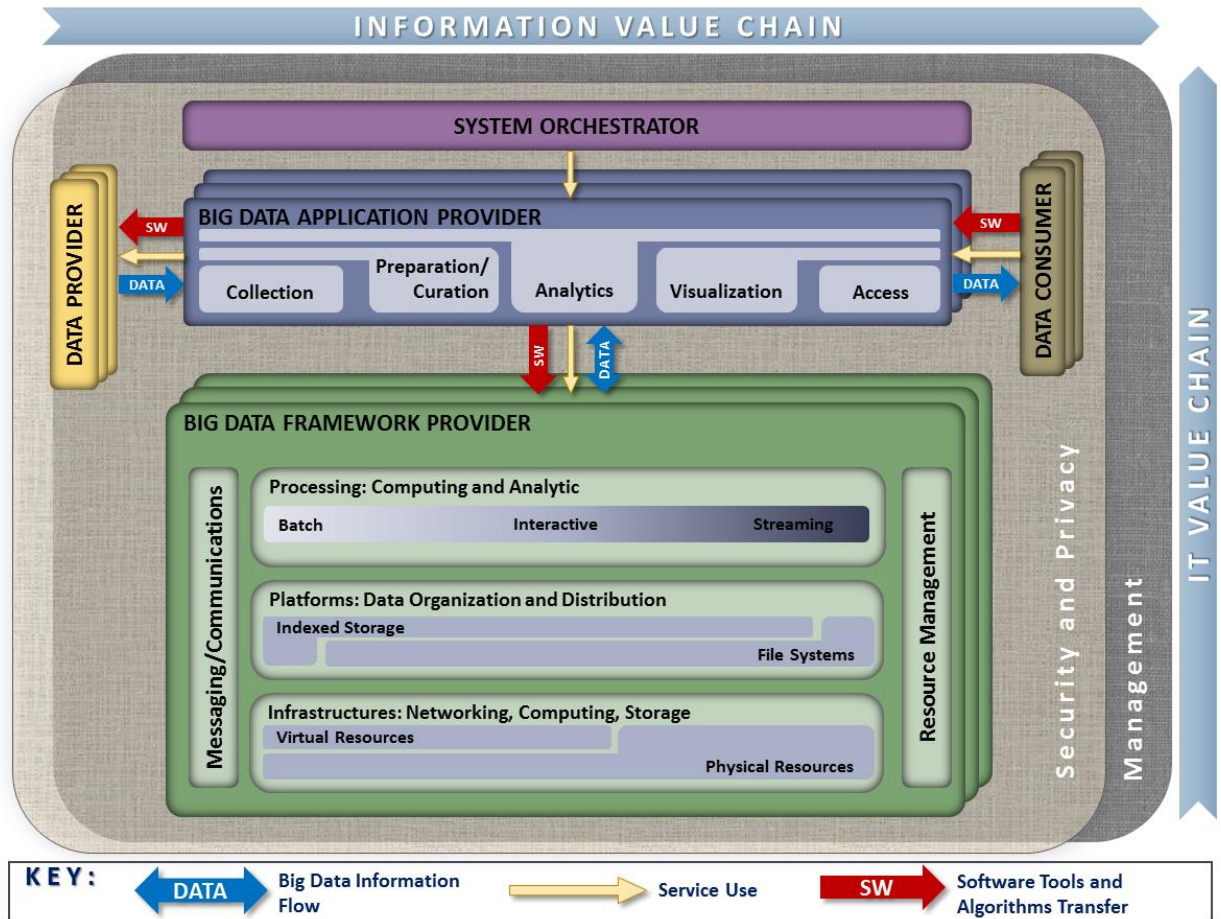


Figure 2: NIST Big Data Reference Architecture

These seven items—five main architecture components and two fabrics interwoven in them—form the foundation of the reference architecture taxonomy.

The five main components, which represent the central roles, are summarized below and discussed in this section (Section 2).

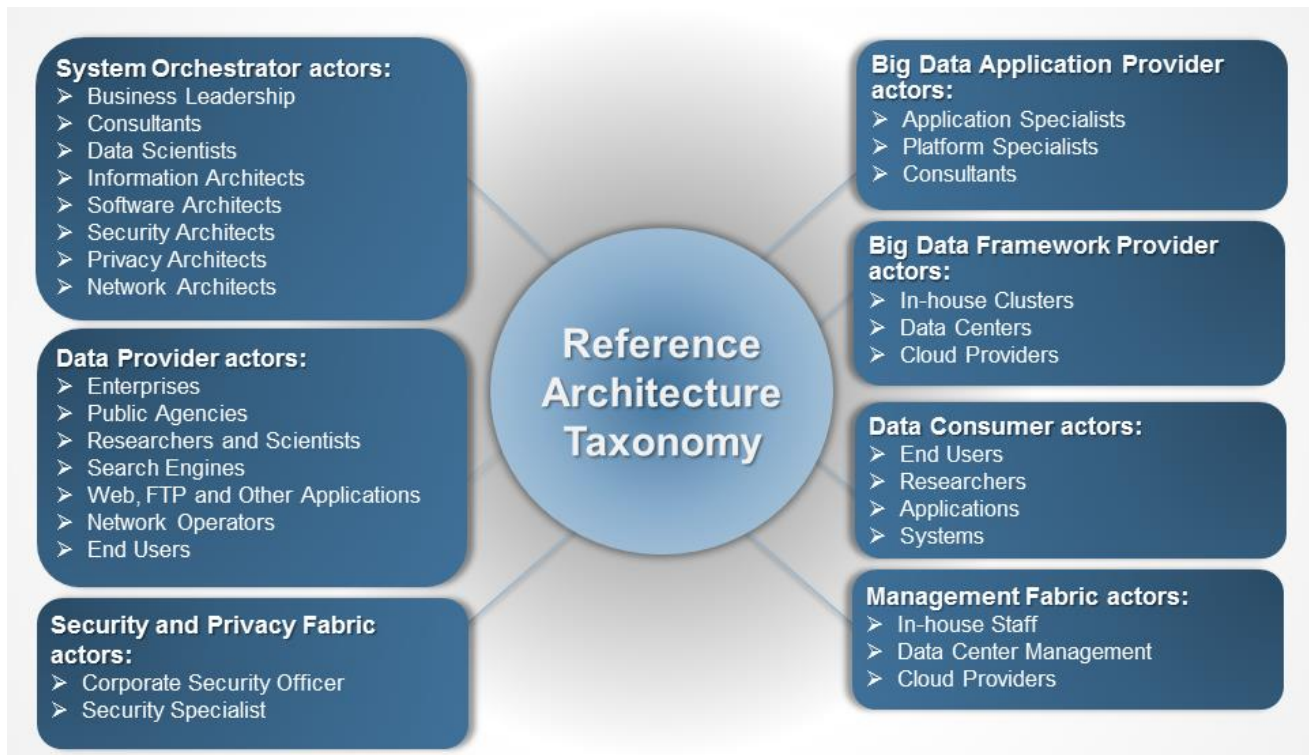
- **System Orchestrator:** Defines and integrates the required data application activities into an operational vertical system;
- **Data Provider:** Introduces new data or information feeds into the Big Data system;
- **Big Data Application Provider:** Executes a life cycle to meet security and privacy requirements as well as System Orchestrator-defined requirements;
- **Big Data Framework Provider:** Establishes a computing framework in which to execute certain transformation applications while protecting the privacy and integrity of data; and
- **Data Consumer:** Includes end users or other systems who use the results of the Big Data Application Provider.

The two fabrics, which are discussed separately in Sections 3 and 4, are:

- **Security and Privacy Fabric**
- **Management Fabric**

Figure 3 outlines potential actors for the seven items listed above. The five central roles are explained in greater detail in the following subsections.





*Figure 3: Roles and a Sampling of Actors in the NBDRA Taxonomy*

## 2.2 SYSTEM ORCHESTRATOR

The System Orchestrator provides the overarching requirements that the system must fulfill, including policy, governance, architecture, resources, and business requirements, as well as monitoring or auditing activities to ensure that the system complies with those requirements.

The System Orchestrator role includes defining and integrating the required data application activities into an operational vertical system. The System Orchestrator role provides system requirements, high-level design, and monitoring for the data system. While the role predates Big Data systems, some related design activities have changed within the Big Data paradigm.

Figure 4 lists the actors and activities associated with the System Orchestrator, which are further described below.



*Figure 4: System Orchestrator Actors and Activities*

### **A. Business Ownership Requirements and Monitoring**

As the business owner of the system, the System Orchestrator oversees the business context within which the system operates, including specifying the following:

- Business goals
- Targeted business action
- Data Provider contracts and service-level agreements (SLAs)
- Data Consumer contracts and SLAs
- Negotiation with capabilities provider
- Make/buy cost analysis

A number of new business models have been created for Big Data systems, including Data as a Service (DaaS), where a business provides the Big Data Application Provider role as a service to other actors. In this case, the business model is to process data received from a Data Provider and provide the transformed data to the contracted Data Consumer.

### **B. Governance Requirements and Monitoring**

The System Orchestrator establishes all policies and regulations to be followed throughout the data life cycle, including the following:

- Policy compliance requirements and monitoring
- Change management process definition and requirements
- Data stewardship and ownership

Big Data systems potentially interact with processes and data being provided by other organizations, requiring more detailed governance and monitoring between the components of the overall system.

### **C. Data Science Requirements and Monitoring**

The System Orchestrator establishes detailed requirements for functional performance of the analytics for the end-to-end system, translating the business goal into data and analytics design, including:

- Data source selection (e.g., identifying descriptions, location, file types, and provenance)
- Data collection and storage requirements and monitoring
- Data preparation requirements and monitoring
- Data analysis requirements and monitoring

- Analytical model choice (e.g., search, aggregation, correlation and statistics, and causal modeling)
- Data visualization requirements and monitoring
- Application type specification (e.g., streaming, real-time, and batch)

A number of the design activities have changed in the new paradigm. In particular, a greater choice of data models now exists beyond the relational model. Choosing a non-relational model will depend on the data type. Choosing the data fields that are used to decide how to distribute the data across multiple nodes will depend on the organization's data analysis needs, and on the ability to use those fields to distribute the data evenly across resources.

#### **D. System Architecture Requirements and Monitoring**

The System Orchestrator establishes detailed architectural requirements for the data system, including the following:

- Data process requirements
- Software requirements
- Hardware requirements
- Logical data modeling and partitioning
- Data import and export requirements
- Scaling requirements

The system architecture has changed in the Big Data paradigm due to the potential interplay between the different actors. The coordination between the five functional NBDRA components is more complex, with additional communications and interconnectivity requirements among the independently operated component activities. Maintaining the needed performance can lead to a very different architecture from that used prior to the new distribution of data across system data nodes.

## **2.3 DATA PROVIDER**

A Data Provider makes data available to itself or to others. The actor fulfilling this role can be part of the Big Data system, from another system, or internal or external to the organization orchestrating the system. Once the data is within the local system, requests to retrieve the needed data will be made by the Big Data Application Provider and routed to the Big Data Framework Provider. Data Provider actors include those shown in Figure 5.

**Data Provider actors:**

- Enterprises
- Public Agencies
- Researchers and Scientists
- Search Engines
- Web, FTP and Other Applications
- Network Operators
- End Users

**Activities**

- » Data Collection from Sources
- » Data Persistence
- » Data Scrubbing
- » Data Annotation/ Metadata Creation
- » Access Rights Management
- » Access Policy Contracts
- » Data Distribution APIs
- » Capabilities Hosting
- » Data Availability Publication

*Figure 5: Data Provider Actors and Activities*

While the concept of a Data Provider is not new, the greater data collection and analytics capabilities have opened new possibilities for providing valuable data. The U.S. government's Open Data Initiative advocates that federal agencies which are stewards of public data also serve the role of Data Provider.

The nine possible Data Provider activities outlined in Figure 4 are discussed further below.

**A. Data Collection from Sources**

The Data Provider captures data from its own sources or others. This activity could be described as the capture from a data producer, whether it is a sensor or an organizational process. Aspects of the data sources activity include both online (available through the internet) and offline sources. Among possible online sources are the following:

- Web servers
- Sensors
- Deep packet inspection devices (e.g., bridge, router, border controller)
- Mobile devices

Offline sources can include the following:

- Public records
- Internal records

While perhaps not theoretically different from what has been in use before, data capture from sources is an area that is exploding in the new Big Data paradigm. New forms of sensors are now providing not only a number of sources of data, but also data in large quantities. Smartphones and personal wearable devices (e.g., exercise monitors, household electric meters) can all be used as sensors. In addition, technologies such as radio frequency identification (RFID) chips can be tracked for the location of shipped items. Collectively, the interaction with data-producing physical devices are known as the internet of things (IoT), or are described as cyber-physical systems (CPS). An example is the set of personal information sensors (watches, health monitors, trackers) which are often referred to as "wearable tech". The recorded digital data about daily activities is sometimes referred to as "digital exhaust."

## B. Data Persistence

The Data Provider stores the data in a repository from which the data can be extracted and made available to others. The stored data is subject to a data retention policy. The data may be stored (i.e., persisted) on an internally hosted system or an externally hosted system. The storage system, may be on a single computer, spread across multiple computers, or on a cloud computing server (dynamically allocated storage spread across multiple computers).

- Internal hosting
- External hosting
- Cloud hosting (a different hosting model whether internal or external)

Hosting models have expanded through the use of cloud computing which are a hybrid of physically being stored externally while under control of the organization. In addition, the data persistence is often accessed through mechanisms such as web services that hide the specifics of the underlying storage. DaaS is a term used for this kind of data persistence that is accessed through specific interfaces.

## C. Data Scrubbing

Some datasets contain sensitive data elements that are naturally collected as part of the data production process. Whether for regulatory compliance or sensitivity, such data elements may need to be altered or removed. As one example of data scrubbing of personally identifiable information (PII), the Data Provider might:

- Remove PII
- Perform data randomization

The latter obscures the PII to remove the possibility of directly tracing the data back to an individual, while maintaining the value distributions within the data. In the era of Big Data, data scrubbing requires greater diligence. While individual sources may not contain PII, when combined with other data sources, the risk arises that individuals may be identified from the integrated datasets.

## D. Data Annotation and Metadata Creation

The Data Provider maintains information about the data and its processing, called metadata, in their repository, and maintains the data itself. The metadata, or data about data, would provide information about the origins and history of the data, in sufficient detail to enable proper use and interpretation of the data. The following approaches can be used to encode the metadata:

- Using an ontology (a semantic description of the elements of the data) pointed to from the data.
- Within a data file: in any number of formats

With the push for open data, where data is repurposed to draw out additional value beyond the initial reason for which it was generated, it has become even more critical that information about the data be encoded to clarify the data's origins and processing. While the actors that collected the data will have a clear understanding of the data history, repurposing data for other uses is open to misinterpretations when other actors use the data at a later time.

In addition, the ability to merge disparate datasets can lead to the identification of individuals even when the original datasets were anonymous. This issue is discussed in Volume 4, *Security and Privacy*.

## E. Access Rights Management

The Data Provider determines the different mechanisms that will be used to define the rights of access, which may be specified individually or by groupings such as the following:

- Data sources: the collection of datasets from a specific source
- Data producer: the collection of datasets from a given producer
- PII access rights: as an example of restrictions on data elements

**F. Access Policy Contracts**

The Data Provider defines policy for others' use of the accessed data, as well as what data will be made available. To expand this description, the contracts specify acceptable use policies and any specific restrictions on the use of the data, as well as ownership of the original data and any derivative works from the data. These contracts specify:

- Policies for primary and secondary rights
- Agreements

**G. Data Distribution Application Programming Interfaces**

Technical protocols are defined for different types of data access from data distribution application programming interfaces (API), which may include:

- File Transfer Protocol (FTP) or streaming
- Compression techniques (e.g., single compressed file, split compressed file)
- Authentication methods
- Authorization

**H. Capabilities Hosting**

For large volumes of data, it may become impractical to move the data to another location for processing. So, in addition to offering data downloads, the Data Provider may offer capabilities to access and manipulate the data while it's on the server, including the following:

- Providing query access without transferring the data
- Allowing analytic tools to be uploaded to operate on the datasets

This is often described as moving the processing to the data, rather than the data to the processing.

**I. Data Availability Publication**

The Data Provider should make available the information that is needed to understand what data or data services they offer. Such publication may consist of the following:

- Web description
- Services and API catalog
- Data dictionaries
- Advertising

A number of third-party locations currently publish a list of links to available datasets (e.g., U.S. government's Open Data Initiative [11]).

**2.4 BIG DATA APPLICATION PROVIDER**

The Big Data Application Provider executes the processes and transformations of the data life cycle to meet requirements established by the System Orchestrator—including meeting the security and privacy requirements. This is where the general capabilities within the Big Data framework are combined to produce the specific data system. Figure 6 lists the actors and activities associated with the Big Data Application Provider.

**Big Data Application Provider actors:**

- Application Specialists
- Platform Specialists
- Consultants

**Activities**

- » Collection
- » Preparation
- » Analytics
- » Visualization
- » Access

*Figure 6: Big Data Application Provider Actors and Activities*

While the activities of an application provider are the same whether the solution being built concerns Big Data or not, the methods and techniques have changed for Big Data because the data and data processing is distributed across resources.

### A. Collection

The Big Data Application Provider must establish the mechanisms to capture data from the Data Provider. These mechanisms include the following:

- Transport protocol and security
- Data format
- Metadata

While the foregoing transport mechanisms predate Big Data, the resources to handle the large volumes or velocities do result in changes in the way the processes are resourced.

### B. Preparation

Whether processes are involved before or after the storage of raw data, a number of them are used in the data preparation activity, analogous to current data system activities. Preparation processes include the following:

- Data validation (e.g., checksums/hashes, format checks)
- Data cleaning (e.g., eliminating bad records/fields, deduplication)
- Outlier removal
- Data conversion (e.g., standardization, reformatting, and encapsulating)
- Calculated field creation and indexing
- Data aggregation and summarization
- Data partition implementation
- Data storage preparation
- Data virtualization layer

Just as data collection may require a number of resources to handle the load, data preparation may also require new resources and/or new techniques. For large data volumes, data collection is often followed by storage of the data in its raw form. Data preparation processes then occur after the storage and are handled by the application code. This technique of storing raw data first and applying a schema upon interaction with the data is commonly called “schema on read.” This dynamic preparation is a new area of interest for Big Data. When storing a new cleaned copy of the data is prohibitive, the data is stored in its raw form and only prepared dynamically for a specific purpose when requested.

Data summarization is a second area of expanded interest due to Big Data. With very large datasets, it is difficult to render all the data for visualization. Proper sampling would need some *a priori* understanding



of the distribution of the entire dataset. Summarization techniques can characterize local subsets of the data, and then provide these characterizations for visualization as the data is browsed.

### C. Analytics

Data science is about the collective set of activities to make sense of large amounts of data. While it can refer to the end-to-end data analytics life cycle, the most common usage focuses on the steps of discovery (i.e., rapid hypothesis-test cycle) for finding value in big volume datasets. This rapid hypothesis-testing analytics cycle (also described as agile analytics) starts with quick correlation or trending analysis, with greater effort spent on hypotheses that appear most promising—along with all activities needed in the handling of the big volumes of data.

Analytics processes for structured and unstructured data have been maturing for many years. There is now more emphasis on the analytics of unstructured data because of the greater quantities now available. The knowledge that valuable information resides in unstructured data promotes a greater attention to the analysis of this type of data.

While analytic methods have not changed with Big Data, their implementation has changed to accommodate parallel data distribution across a cluster of independent nodes and data access methods. For example, the overall data analytic task may be broken into subtasks that are assigned to the independent data nodes. The results from each subtask are collected and compiled to achieve the final full dataset analysis. With the introduction of new storage paradigms, analytics techniques have to be modified for different types of data access.

Some considerations for analytical processes used for Big Data or small data are the following:

- Metadata matching processes
- Analysis complexity considerations (e.g., computational, machine learning, data extent, data location)
- Analytics latency considerations (e.g., real-time or streaming, near real-time or interactive, batch or offline)
- Human-in-the-loop analytics life cycle (e.g., discovery, hypothesis, hypothesis testing)

While these considerations are not new to Big Data, implementing them can be tightly coupled with the specifics of the data storage and the preparation step.

### D. Visualization

While visualization (for the human in the loop) can be considered a type of analytics, the importance warrants special consideration. The following are three general categories of data visualization:

- Exploratory data visualization for data understanding (e.g., browsing, outlier detection, boundary conditions)
- Explicatory visualization for analytical results (e.g., confirmation, near real-time presentation of analytics, interpreting analytic results)
- Explanatory visualization to “tell the story” (e.g., reports, business intelligence, summarization)

Data science relies on the full dataset type of discovery or exploration visualization from which the data scientist would form a hypothesis. While clearly predating Big Data, a greater emphasis now exists on exploratory visualization. It is more critical in understanding large volumes of repurposed data because the size of the datasets can make statistical profiles of the entire datasets quite difficult.

Explanatory visualization is the creation of a simplified, digestible visual representation of the results, suitable for assisting a decision or communicating the knowledge gained. Again, while this technique has long been in use, there is now greater emphasis to “tell the story.” Often this is done through simple visuals or “infographics.” Given the large volumes and varieties of data, and the data’s potentially



complex analysis, the communication of the analytics to a non-analyst audience requires careful visual representation to communicate the results in a way that can be easily consumed.

### E. Access

The Big Data Application Provider gives the Data Consumer access to the results of the data system, including the following:

- Data export API processes (e.g., protocol, query language)
- Data charging mechanisms
- Consumer analytics hosting where the application hosts the Consumer's code
- Analytics as a service hosting where the Consumer accesses the analytics application (such as a web-based business intelligence application)

The access activity of the Big Data Application Provider should mirror all actions of the Data Provider, since the Data Consumer may view this system as the Data Provider for their follow-on tasks. Many of the access-related tasks have changed with Big Data, as algorithms have been rewritten to accommodate for and optimize the distributed resources.

## 2.5 BIG DATA FRAMEWORK PROVIDER

The Big Data Framework Provider has general resources or services to be used by the Big Data Application Provider in the creation of the specific application. There are many new technologies from which the Big Data Application Provider can choose in using these resources and the network to build the specific system. Figure 7 lists the actors and activities associated with the Big Data Framework Provider.



*Figure 7: Big Data Framework Provider Actors and Activities*

The Big Data Framework Provider role has seen the most significant changes with the introduction of Big Data. The Big Data Framework Provider makes available tools and techniques within the three activities: infrastructure frameworks, data platform frameworks, and processing frameworks. There is no requirement that all activities at a given level in the hierarchy leverage the same technology and, in fact, most Big Data implementations are hybrids combining multiple technology approaches. These provide flexibility and can meet the complete range of requirements that are driven from the Big Data Application Provider. Due to the rapid emergence of new techniques, this is an area that will continue to need discussion. As the Subgroup continues its discussion into patterns within these techniques, different orderings will no doubt be more representative and understandable.

### A. Infrastructure Frameworks

Infrastructure frameworks can be grouped as follows:

- Networking: These are the components that transfer data from one resource to another (e.g., physical, virtual, software-defined).

- **Computing:** These are the physical processors and memory that execute and hold the software of the other Big Data system components (e.g., physical resources, operating system, virtual implementation, logical distribution).
- **Storage:** These are components which provide persistence of the data in a Big Data system (e.g., in-memory, local disk, hardware/software redundant array of independent disks [RAID], storage area networks [SANs], network-attached storage [NAS]).
- **Environmental:** These are the facility components (e.g., power, cooling) that must be accounted for when establishing an instance of a Big Data system.

The biggest change under the Big Data paradigm is the cooperation of horizontally scaled, independent resources to implement the desired component performance.

## **B. Data Platform Frameworks**

This is the most recognized area for changes in Big Data engineering, and given the rapid changes, the hierarchy in this area will likely change in the future to better represent the patterns within the techniques. The data platform frameworks activity was expanded into the following logical data organization and distribution approaches to provide additional clarity needed for the new approaches of Big Data.

- Logical File systems (e.g., centralized, distributed)
- Logical data repositories
  - Simple tuple (e.g., relational, non-relational or not only or no Structured Query Language [NoSQL] tables both row and column)
  - Complex tuple (e.g., indexed document store, non-indexed key-value or queues)
  - Graph (e.g., property, hyper-graph, triple stores)

The logical storage paradigm has expanded beyond the “flat file” and relational model paradigms to develop new non-relational models. This has implications for the concurrency of the data across nodes within the non-relational model. Transaction support in this context refers to the completion of an entire data update sequence and the maintenance of eventual consistency across data nodes. This is an area that needs more exploration and categorization.

## **C. Processing Frameworks**

Processing frameworks provide the software support for applications which can deal with the volume, velocity, variety, and variability of data. Some aspects related to processing frameworks are the following:

- Data type processing services (e.g., numeric, textual, spatial, images, video)
- Schema information or metadata management (e.g., on demand, pre-knowledge)
- Query frameworks (e.g., relational, arrays)
- Temporal frameworks
  - Batch (e.g., dense linear algebra, sparse linear algebra, spectral, N-body, structured grids, unstructured grids, MapReduce, bulk synchronous parallel [BSP])
  - Interactive
  - Real-time/streaming (e.g., event ordering, state management, partitioning)
- Application frameworks (e.g., automation, test, hosting, workflow)
- Messaging/communications frameworks
- Resource management frameworks (e.g., cloud/virtualization, intra-framework, inter-framework)

Both the Big Data Application Provider activities and the Big Data Framework Provider activities have changed significantly due to Big Data engineering. Currently, the interchange between these two roles operates over a set of independent, yet coupled, resources. It is in this interchange that the new methods for data distribution over a cluster have developed. Just as simulations went through a process of parallelization (or horizontal scaling) to harness massive numbers of independent process to coordinate

them to a single analysis, Big Data services now perform the orchestration of data processes over parallel or distributed resources.

## 2.6 DATA CONSUMER

The Data Consumer receives the value output of the Big Data system. In many respects, the Data Consumer receives the same functionality that the Data Provider brings to the Big Data Application Provider. After the system adds value to the original data sources, the Big Data Application Provider then offers the additional value to the Data Consumer. There is less change in this role due to Big Data, except, of course, in the desire for Consumers to extract extensive datasets from the Big Data Application Provider. Figure 8 lists the actors and activities associated with the Data Consumer.



*Figure 8: Data Consumer Actors and Activities*

The activities listed in Figure 8 are explicit to the Data Consumer role within a data system. If the Data Consumer is in fact a follow-on application, then the Data Consumer would look to the Big Data Application Provider for the activities of any other Data Provider. The follow-on application's System Orchestrator would negotiate with this application's System Orchestrator for the types of data wanted, access rights, and other requirements. The Big Data Application Provider would thus serve as the Consumer's Data Provider, from the perspective of the follow-on application.

### A. Search and Retrieve

The Big Data Application Provider could allow the Data Consumer to search across the data, and query and retrieve data for its own usage.

### B. Download

Data from the Data Provider could be exported to the Data Consumer for download to the Consumer's environment. This is the same process the Application Provider follows to download data from the Data Provider.

### C. Analyze Locally

The Application Provider could allow the Data Consumer to run their own application on the data. This would imply that the application provided hosting capability to allow the consumer's code to run directly in the application environment.

### D. Reporting

The data could be presented according to the chosen filters, values, and formatting as a reporting-as-a-service application.

### E. Visualization

The Data Consumer could be allowed to browse the raw data, or the data output from the analytics.

## 2.7 MANAGEMENT FABRIC

The Big Data characteristics of volume, velocity, variety, and variability demand a versatile management platform for storing, processing, and managing complex data. Management of Big Data systems should handle both system- and data-related aspects of the Big Data environment. The Management Fabric of the NBDRA encompasses two general groups of activities: system management and Big Data life cycle management. System management includes activities such as provisioning, configuration, package management, software management, backup management, capability management, resources management, and performance management. Big Data life cycle management involves activities surrounding the data life cycle of collection, preparation/curation, analytics, visualization, and access. More discussion about the Management Fabric is needed, particularly with respect to new issues in the management of Big Data and Big Data engineering.

Figure 9 lists an initial set of activities associated with the Management role of the NBDRA.



Figure 9: Big Data Management Actors and Activities

## 2.8 SECURITY AND PRIVACY FABRIC

Security and privacy issues affect all other components of the NBDRA, as depicted by the encompassing Security and Privacy box in Figure 2. A Security and Privacy Fabric could interact with the System Orchestrator for policy, requirements, and auditing, and also with both the Big Data Application Provider and the Big Data Framework Provider for development, deployment, and operation. These ubiquitous security and privacy activities are described in the *NBDIF: Volume 4, Security and Privacy* document. Figure 10 lists representative actors and activities associated with the Security and Privacy Fabric of the NBDRA.



Figure 10: Big Data Security and Privacy Actors and Activities

### 3 DATA CHARACTERISTIC HIERARCHY

Equally important to understanding the new Big Data engineering that has emerged in the last ten years, is the need to understand what data characteristics have driven the need for the new technologies—and what data characteristics are affected. In Section 2 of this document, a taxonomy was presented for the NBDRA, which is described in *NBDIF: Volume 6, Reference Architecture*. The NBDRA taxonomy has a hierarchy of roles/actors, and activities. To understand the characteristics of data and how they have changed with the new Big Data Paradigm, it is illustrative to look at the data characteristics at different levels of granularity. Understanding what characteristics are affected by Big Data can best be done by examining the granularity of data elements, of related data elements grouped into a record that represents a specific entity or event, of records collected into a dataset, and of multiple datasets—all in turn, as shown in Figure 11. Therefore, this section does not present a strict taxonomy, breaking down each element into parts, but provides a description of data objects at a specific granularity along with attributes for those objects, and characteristics and sub-characteristics of the attributes. The framework described will help illuminate areas where the effects of Big Data can be understood in the context of the characteristics of all data. For easier presentation, the hierarchy will be traversed from the smallest level of data elements, through elements grouped into records, records into datasets, and then finally the consideration of multiple datasets.

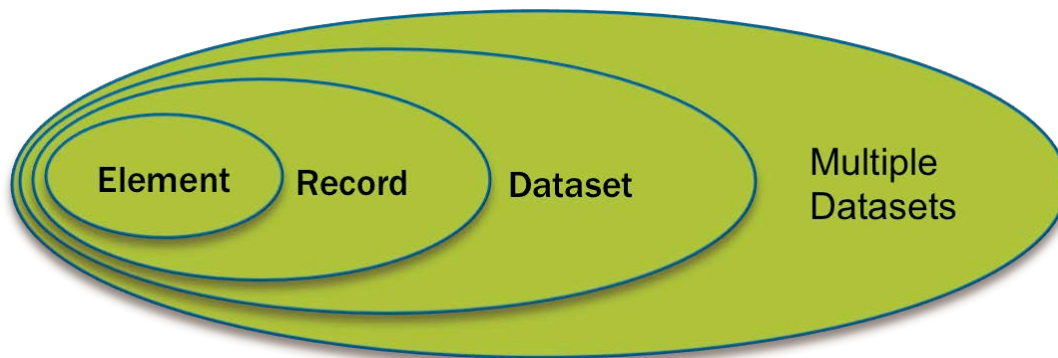


Figure 11: Data Characteristic Hierarchy

#### 3.1 DATA ELEMENTS

Individual data elements have naturally not changed in the new Big Data paradigm. Data elements are understood by their data type and additional contextual data, or metadata, which provides history or additional understanding about the data. For example, in the context of unstructured text, a data element would refer to a single token such as a word.

##### A. Data Format

Data formats are well characterized through Standards Development Organizations (SDOs) including International Organization for Standardization (ISO) standards such as ISO 8601: 2004 Data elements and interchange formats—Information interchange—Representation of dates and times [12]. Data element formats have not changed for Big Data.

**B. Data Values and Vocabulary**

A data element is populated by its actual value. This value is restricted to its defined data type (e.g., numeric, string, date) and chosen data format. Sometimes the value is restricted to a specific standard vocabulary for interoperability with others in the field, or to a set of allowed values.

**C. Metadata and Semantics**

Metadata is data about objects, sometimes simplistically described as “data about data.” Metadata can refer to a number of categories of contextual information, including the origins and history of the data, the processing times, the software versions, and other information. In addition, data can be described semantically to better understand what the value represents, and to make the data machine-operable. Both metadata and semantic data are not specific to Big Data. Further information about metadata and semantics can be found in: ISO/IEC 11179 Information Technology–Metadata registries [13]; and W3C’s work on the Semantic Web [14].

**D. Quality and Veracity**

Data Quality and Veracity are characteristics used in describing Big Data, but the accuracy of the data is not a new concern. Data quality is another name for the consideration of the reliability of the data. Again, this topic predated Big Data and is beyond the scope of this volume. Further information about data quality can be found in ISO/TS 8000 Data Quality [15].

**3.2 RECORDS**

Data elements are grouped into records that describe a specific entity or event or transaction. At the level of records, new techniques for Big Data have been developed. For example, in the context of unstructured text, a data record could refer to a phrase or sentence or entire document. Note that a greater emphasis is placed now on unstructured data due to increasing amounts of web and mobile data (e.g., online text, images and video).

**3.2.1 RECORD CHARACTERISTICS****A. Record Format**

Records have structure and formats. Record structures are commonly grouped as structured, semi-structured, and unstructured. Structured data was traditionally described through formats such as comma-separated values, or as a row in a relational database. Unstructured refers to free text, such as in a document or a video stream. An example of semi-structured is a record wrapped with a markup language such as eXtensible Markup Language (XML) or HyperText Markup Language (HTML), where the contents within the markup can be free text.

These categories again predate Big Data, but a notable change has occurred with Big Data in terms of storage in new non-relational formats.

**B. Complexity**

Complexity refers to the interrelationship between data elements in a record, or between records (e.g., in the interrelationships in genomic data between sequences, genes and proteins). Complexity is not new to Big Data.

**C. Volume**

Records themselves can have an aspect of volume in the emerging data sources, such as considering an entire DNA on an organism as a record.

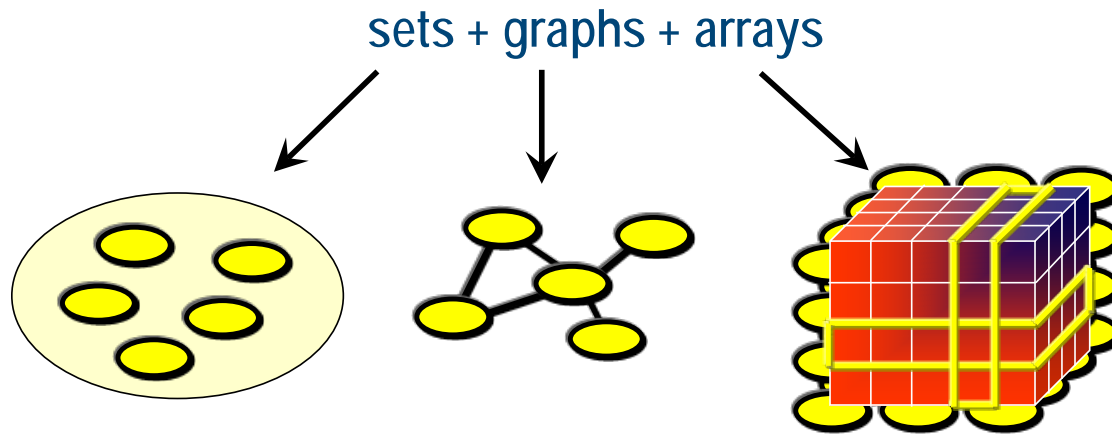
**D. Metadata and Semantics**

The same metadata categories described for data elements can be applied to records. In addition, relationships between data elements can be described semantically in terms of an ontology.



### 3.2.2 INTER-RECORD STRUCTURE

Data structures appearing in science and engineering, business, finance, and statistics are all based on a remarkably small slate of core categories. Among the most frequent and important basic data structuring categories are sets, hierarchies, graphs, and arrays.



Adapted from original figure by Peter Baumann

**Figure 12: Classification of common data structures (records represented abstractly as yellow ellipses)**

Sets are traditionally managed in relational databases where the structuring centers around tables which can be seen as sets of records, typically called tuples. One example of such record sets is employee data in a human resources database. Records in a table are uncorrelated in the first place and, hence, form unordered sets, which is why the relational model is called set-oriented. MapReduce-based techniques [16] on data sets where the individual records / tuples are loosely connected, or not at all. Structured Query Language (SQL) [17] databases—along with some NoSQL and NewSQL databases—operate on this concept of a set of independent records, where each record describes a specific event or entity. In SQL databases each record is typically a set of values for specific fields, essentially the structure of a table. There is potentially significant structure within such records, such as the nested structures often seen in XML [18] or JSON [19] databases, or possibly the unstructured text in a record representing a document.

There are two types of data structures where there is an explicit relationship between the records. Inter-record relationships cannot be represented in traditional set-oriented database systems. The first such structure is given by general graphs, where each record (node) can have a relationship with other records (links). This type of data is represented in the NoSQL graph databases. Examples of graphical data are social networks where users are connected on the basis based on “friend” relationships; the representation of genes and proteins in the human genome; and the concept graphs in Semantic Web ontologies [20]. Query languages and analytics techniques differ in these graphical models to take advantage of this inherent structure.

For relationships such as taxonomies, where there are only parent-child relationships between records, a hierarchical graph representation is used. Here the relationships between nodes are limited to cases of going from the general to the specific as the specific being a “part-of” the more general parent (for example in the taxonomic characterization of species). As a graph with restricted relationships, hierarchical graphs can be seen as a special case of the more general concept of graphs.

Arrays comprise another fundamental data structure being a 3D (or higher) table or matrix, instead of the 2D row and column structure of tables. Objects of the same structure (such as <red, green, blue> pixel records) are aligned along some multi-dimensional grid. The grid structure leads to a well-defined

neighborhood among array records (often referred to as cells): any given cell has exactly two neighbors along each dimension, one in positive and one in negative direction – except for the border cells which have only one such neighbor each. This highly regular structure makes arrays very amenable to bulk processing. Examples for arrays include 2-D imagery, 3-D x/y/t image timeseries, 4-D x/y/z/t climate and weather data, or multiple dimension “data cubes” (analogous to the construction of pivot tables in a spreadsheet).

Each of these fundamental data categories comes with typical operations: sets allow filtering, intersection, etc. while hierarchies suggest collecting all direct and indirect parts of some given object. Interesting operations on graphs include finding unconnected graph parts, determining shortest paths between two objects, and many more. Arrays allow extraction of sub-arrays, including slices of lower dimension, but also aggregation, in general: Linear Algebra (in particular on 2D arrays, i.e.: matrices).

### 3.3 DATASETS

Records can be grouped to form datasets. This grouping of records can reveal changes due to Big Data. For example, in the context of unstructured text, a record could be a sentence, paragraph, or section, and the dataset could refer to the complete document.

#### A. Quality and Consistency

A new aspect of data quality for records due to the distribution of data across multiple nodes focuses on the characteristic of consistency. As records are distributed horizontally across a collection of data nodes, consistency becomes an issue. In relational databases, consistency was maintained by assuring that all operations in a transaction were completed successfully; otherwise the operations were rolled back. This assured that the database maintained its internal consistency. For additional information on this concept, the reader is referred to the literature on atomicity, consistency, isolation, durability (ACID) properties of databases. For Big Data, with multiple nodes and backup nodes, new data is sent in turn to the appropriate nodes. However, constraints may or may not exist to confirm that all nodes have been updated when a query is sent. The time delay in replicating data across nodes can cause an inconsistency. The methods used to update nodes are one of the main areas in which specific implementations of non-relational data storage methods differ.

### 3.4 MULTIPLE DATASETS

The primary focus on multiple datasets concerns the need to integrate or fuse the multiple datasets. The focus is on the variety characteristic of Big Data. Extensive datasets cannot always be converted into one structure (e.g., all weather data being reported on the same spatio-temporal grid). Since large volume datasets cannot be easily copied into a normalized structure, new techniques are being developed to integrate data as needed. For example, in the context of unstructured text, multiple datasets could refer to a document collection.

#### A. Personally Identifiable Information (PII)

An area of increasing concern with Big Data is the identification of individuals from the integration of multiple datasets, even when the individual datasets would not allow the identification. For additional discussion, the reader is referred to *NBDIF: Volume 4, Security and Privacy* [3].

#### B. Data Virtualization

Data virtualization, with respect to multiple datasets, are where the disparate datasets continue to reside in their repositories, but are accessed through a single logically fused structure.



## 4 SUMMARY

---

Big Data and data science represent a rapidly changing field due to the recent emergence of new technologies and rapid advancements in methods and perspectives. This document presents a taxonomy for the NBDRA, which is presented in *NBDIF: Volume 6, Reference Architecture* [5]. This taxonomy provides a base hierarchy for categorizing the new components and activities of Big Data systems. In addition, a description of data at different scales was provided to place concepts being ascribed to Big Data into their context.

# Appendix A: Acronyms

---

ACID	atomicity, consistency, isolation, durability
API	application programming interface
BSP	bulk synchronous parallel
CPS	cyber-physical systems
DaaS	Data as a Service
FTP	File Transfer Protocol
HTML	HyperText Markup Language
IoT	internet of things
ISO	International Organization for Standardization
ITL	Information Technology Laboratory (NIST)
JSON	JavaScript Object Notation
NARA	National Archives and Records Administration
NAS	network-attached storage
NASA	National Aeronautics and Space Administration
NBDIF	NIST Big Data Interoperability Framework
NBD-PWG	NIST Big Data Public Working Group
NBDRA	NIST Big Data Reference Architecture
NIST	National Institute of Standards and Technology
NoSQL	not only (or no) Structured Query Language
NSF	National Science Foundation
PII	personally identifiable information
RAID	redundant array of independent disks
RFID	radio frequency identification
SAN	storage area network
SLA	Service-level Agreement
SQL	Structured Query Language
XML	eXtensible Markup Language

## Appendix B: Bibliography

---

- [1] W. L. Chang (Co-Chair), N. Grady (Subgroup Co-chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 1, Big Data Definitions (NIST SP 1500-1 VERSION 3),” Gaithersburg MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-1r2>
- [2] W. L. Chang (Co-Chair), G. Fox (Subgroup Co-chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 3, Big Data Use Cases and General Requirements (NIST SP 1500-3 VERSION 3),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-3r2>
- [3] W. L. Chang (Co-Chair), A. Roy (Subgroup Co-chair), M. Underwood (Subgroup Co-chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 4, Big Data Security and Privacy (NIST SP 1500-4 VERSION 3),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-4r2>
- [4] W. Chang and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 5, Architectures White Paper Survey (SP1500-5),” 2015 [Online]. Available: <https://www.nist.gov/publications/nist-big-data-interoperability-framework-volume-5-architectures-white-paper-survey>
- [5] W. L. Chang (Co-Chair), D. Boyd (Subgroup Co-chair), O. Levin (Version 1 Subgroup Co-Chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 6, Big Data Reference Architecture (NIST SP 1500-6 VERSION 3),” Gaithersburg MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-6r2>
- [6] W. L. Chang (Co-Chair), R. Reinsch (Subgroup Co-chair), D. Boyd (Version 1 Subgroup Co-chair), C. Buffington (Version 1 Subgroup Co-chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 7, Big Data Standards Roadmap (NIST SP 1500-7 VERSION 3),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-7r2>
- [7] W. L. Chang (Co-Chair), G. von Laszewski (Editor), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 8, Big Data Reference Architecture Interfaces (NIST SP 1500-9 VERSION 2),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-9r1>
- [8] W. L. Chang (Co-Chair), R. Reinsch (Subgroup Co-chair), C. Austin (Editor), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 9, Adoption and Modernization (NIST SP 1500-10 VERSION 2),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-10r1>
- [9] T. White House Office of Science and Technology Policy, “Big Data is a Big Deal,” *OSTP Blog*,

2012. [Online]. Available: <http://www.whitehouse.gov/blog/2012/03/29/big-data-big-deal>. [Accessed: 21-Feb-2014]
- [10] W. L. Chang (Co-Chair), D. Boyd (Subgroup Co-chair), O. Levin (Version 1 Subgroup Co-Chair), and NIST Big Data Public Working Group, “NIST Big Data Interoperability Framework: Volume 6, Big Data Reference Architecture (NIST SP 1500-6 VERSION 3),” Gaithersburg, MD, Sep. 2019 [Online]. Available: <https://doi.org/10.6028/NIST.SP.1500-6r2>
- [11] U.S. Federal Government, “The home of the U.S. Government’s open data,” *Online*, 2014. [Online]. Available: <https://www.data.gov/>
- [12] International Standards Organization, “ISO 8601:2004(E) Data elements and interchange formats - Information interchange - Representation of dates and times,” *Reference number ISO*, vol. 2004. p. 40, 2004.
- [13] *ISO/IEC 11179-1:2015, Information technology - Metadata registries (MDR) – Part 1: Framework*. International Organization for Standardization / International Electrotechnical Commission, 2015 [Online]. Available: <https://www.iso.org/standard/61932.html>
- [14] Dan Brickley and I. Herman, “Semantic Web Interest Group,” W3C. [Online]. Available: <https://www.w3.org/2001/sw/interest/>
- [15] *ISO/TS 8000-1:2011 Data Quality*. International Organization for Standardization, 2011 [Online]. Available: <https://www.iso.org/standard/50798.html>
- [16] The Apache Software Foundation, “MapReduce Tutorial,” 22-Aug-2019. [Online]. Available: [https://hadoop.apache.org/docs/r1.2.1/mapred\\_tutorial.html](https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html)
- [17] *ISO/IEC 9075-2:2016 Information technology -- Database languages -- SQL -- Part 2: Foundation (SQL/Foundation)*. International Organization for Standardization / International Electrotechnical Commission, 2016 [Online]. Available: <https://www.iso.org/standard/63556.html>
- [18] W3C, “W3C: Extensible Markup Language (XML),” 2019. [Online]. Available: <https://www.w3.org/XML/>
- [19] “Introducing JSON.” [Online]. Available: <https://www.json.org/>
- [20] W3C, “Semantic Web,” W3C, 2015. [Online]. Available: <https://www.w3.org/standards/semanticweb/>