**NIST Special Publication 1500-2**

# DRAFT NIST Big Data Interoperability Framework:
# Volume 2, Big Data Taxonomies

NIST Big Data Public Working Group
Definitions and Taxonomies Subgroup

Draft Version 1
April 6, 2015

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

NIST Special Publication 1500-2
Information Technology Laboratory

# DRAFT NIST Big Data Interoperability Framework:
# Volume 2, Big Data Taxonomies

## 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

**National Institute of Standards and Technology Special Publication 1500-2**
32 pages (April 6, 2015)

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

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

# 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 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

# Acknowledgements

# Notice to Readers

NIST is seeking feedback on the proposed working draft of the *NIST Big Data Interoperability Framework: Volume 2, Big Data Taxonomies*. 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), 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 with the volume number in the subject line (e.g., Edits for Volume 2.)

Please contact Wo Chang (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

## FIGURES

# Executive Summary

This *NIST Big Data Interoperability Framework: Volume 2, 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* consists of seven volumes, each of which addresses a specific key topic, resulting from the work of the NBD-PWG. The seven volumes are as follows:

- Volume 1, Definitions
- Volume 2, Taxonomies
- Volume 3, Use Cases and General Requirements
- Volume 4, Security and Privacy
- Volume 5, Architectures White Paper Survey
- Volume 6, Reference Architecture
- Volume 7, Standards Roadmap

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

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

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

# 1 INTRODUCTION

## 1.1 BACKGROUND

There is broad agreement among commercial, academic, and government leaders about the remarkable 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:

- What attributes define Big Data solutions?
- 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 Big Data solutions?

Within this context, on March 29, 2012, the White House announced the Big Data Research and Development Initiative.[1] 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 the 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 Interoperability Framework. Forum participants noted that this framework should define and prioritize Big Data requirements, including interoperability, portability, reusability, extensibility, data usage, analytics, and technology infrastructure. In doing so, the framework would accelerate the adoption of the most secure and effective Big Data techniques and technology.

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

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

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

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

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

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

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

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

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

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

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

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

## 1.3  REPORT PRODUCTION

113  This document derives from discussions in the NBD-PWG Definitions and Taxonomy Subgroup and with
114  interested parties. This volume provides the taxonomy of the components of the NBDRA. This taxonomy

115    was developed using a mind map representation, which provided a mechanism for multiple inputs and
116    easy editing.

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

123    In addition to the reference architecture taxonomy, the Subgroup began the development of a data
124    hierarchy

## 1.4  REPORT STRUCTURE

126    This document provides multiple hierarchical presentations related to Big Data.

127    The first presentation is the taxonomy for the NBDRA. This taxonomy provides the terminology and
128    definitions for the components of technical systems that implement technologies for Big Data. Section 2
129    introduces the NBDRA using concepts of actors and roles and the activities each performs. In the
130    NBDRA presented in *NIST Big Data Interoperability Framework Volume 6: Reference Architecture*,
131    there are two roles that span the activities within the other roles: Management, and Security and Privacy.
132    These two topic areas will be addressed further in future versions of this document. The NBDRA
133    components are more fully described in the *NIST Big Data Interoperability Framework: Volume 6,
134    Reference Architecture* and the *NIST Big Data Interoperability Framework: Volume 4, Security and
135    Privacy* documents. Comparing the related sections in these two documents will give the reader a more
136    complete picture of the consensus of the working groups.

137    The second presentation is a hierarchical description about the data itself. For clarity, a strict taxonomy is
138    not followed; rather, data is examined at different groupings to better describe what is new with Big Data.
139    The grouping-based description presents data elements, data records, datasets, and multiple datasets.  This
140    examination at different groupings provides a way to easily identify the data characteristics that have
141    driven the development of Big Data engineering technologies, as described in the *NIST Big Data
142    Interoperability Framework: Volume 1, Definitions.*

143    Within the following sections, illustrative examples are given to facilitate understanding of the role/actor
144    and activity of the NBDRA. There is no expectation of completeness in the components; the intent is to
145    provide enough context to understand the specific areas that have changed because of the new Big Data
146    paradigm. Likewise, the data hierarchy only expresses the broad overview of data at different levels of
147    granularity to highlight the properties that drive the need for Big Data architectures.

148    For descriptions of the future of Big Data and opportunities to use Big Data technologies, the reader is
149    referred to the *NIST Big Data Interoperability Framework: Volume 7, Standards Roadmap*. Finally, to
150    understand how these systems are architected to meet users' needs, the reader is referred to *NIST Big
151    Data Interoperability Framework: Volume 3, Use Cases and General Requirements.*

## 1.5  FUTURE WORK ON THIS VOLUME

153    As mentioned in the previous section, the Subgroup is continuing to explore the changes in both
154    Management and in Security and Privacy. As changes in the activities within these roles are clarified, the
155    taxonomy will be developed further. In addition, a fuller understanding of Big Data and its technologies
156    should consider the interactions between the characteristics of the data and the desired methods in both
157    technique and time window for performance. These characteristics drive the application and the choice of
158    tools to meet system requirements. Investigation of the interfaces between data characteristics and

159 technologies is a continuing task for the NBD-PWG Definitions and Taxonomy Subgroup and the NBD-
160 PWG Reference Architecture Subgroup. Finally, societal impact issues have not yet been fully explored.
161 There are a number of overarching issues in the implications of Big Data, such as data ownership and data
162 governance, which need more examination. Big Data is a rapidly evolving field, and the initial discussion
163 presented in this volume must be considered a work in progress.

164

## 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, application. For elucidating the taxonomy, a hierarchy has been chosen to allow placing the new technologies within the context 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 independents 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 *NIST Big Data Interoperability Framework: Volume 3, Use Cases and General Requirements* and *NIST Big Data Interoperability Framework: 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 1.

*Figure 1: 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 lifecycle 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
- **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 as follows:

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

Figure 2 outlines potential actors for the seven items listed above. The five central roles are explained in greater detail in the following subsections.

215                                    *Figure 2: Roles and a Sampling of Actors in the NBDRA Taxonomy*

## 2.2 SYSTEM ORCHESTRATOR

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

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

224 Figure 3 lists the actors and activities associated with the System Orchestrator, which are further
225 described below.

226       *Figure 3: System Orchestrator Actors and Activities*

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

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

230         •   Business goals
231         •   Targeted business action
232         •   Data Provider contracts and service-level agreements (SLAs)
233         •   Data Consumer contracts and SLAs
234         •   Negotiation with capabilities provider
235         •   Make/buy cost analysis

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

240    **B. Governance Requirements and Monitoring**

241    The System Orchestrator establishes all policies and regulations to be followed throughout the data
242    lifecycle, including the following:

243         •   Policy compliance requirements and monitoring
244         •   Change management process definition and requirements
245         •   Data stewardship and ownership

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

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

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

251         •   Data source selection (e.g., identifying descriptions, location, file types, and provenance)
252         •   Data collection and storage requirements and monitoring
253         •   Data preparation requirements and monitoring
254         •   Data analysis requirements and monitoring

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

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

### D. System Architecture Requirements and Monitoring

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

267    • Data process requirements
268    • Software requirements
269    • Hardware requirements
270    • Logical data modeling and partitioning
271    • Data export requirements
272    • Scaling requirements

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

## 2.3  DATA PROVIDER

279    A Data Provider makes data available to itself or to others. The actor fulfilling this role can be part of the
280    Big Data system, from another system, or internal or external to the organization orchestrating the system.
281    Once the data is within the local system, requests to retrieve the needed data will be made by the Big Data
282    Application Provider and routed to the Big Data Framework Provider. Data Provider actors include those
283    shown in Figure 4.

284 *Figure 4: Data Provider Actors and Activities*

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

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

289 **A. Data Capture from Sources**

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

293 • Web browsers
294 • Sensors
295 • Deep packet inspection devices (e.g., bridge, router, border controller)
296 • Mobile devices

297 Offline sources can include the following:

298 • Public records
299 • Internal records

300 While perhaps not theoretically different from what has been in use before, data capture from sources is
301 an area that is exploding in the new Big Data paradigm. New forms of sensors are now providing not only
302 a number of sources of data, but also data in large quantities. Smartphones and personal wearable devices
303 (e.g., exercise monitors, household electric meters) can all be used as sensors. In addition, technologies
304 such as radio frequency identification (RFID) chips are sources of data for the location of shipped items.
305 Collectively, all the data-producing sensors are known as the Internet of Things (IoT). The subset of
306 personal information devices are often referred to as "wearable tech", with the resulting data sometimes
307 referred to as "digital exhaust".

308 **B. Data Persistence**

309 The Data Provider stores the data in a repository from which the data can be extracted and made available
310 to others. The stored data is subject to a data retention policy. The data can be stored (i.e., persisted) in the
311 following ways:

312       • Internal hosting
313       • External hosting
314       • Cloud hosting (a different hosting model whether internal or external)

315   Hosting models have expanded through the use of cloud computing. In addition, the data persistence is
316   often accessed through mechanisms such as web services that hide the specifics of the underlying storage.
317   DaaS is a term used for this kind of data persistence that is accessed through specific interfaces.

318   **C. Data Scrubbing**

319   Some datasets contain sensitive data elements that are naturally collected as part of the data production
320   process. Whether for regulatory compliance or sensitivity, such data elements may be altered or removed.
321   As one example of data scrubbing for Personally Identifiable Information (PII), the Data Provider can:

322       • Remove PII
323       • Perform data randomization

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

328   **D. Data Annotation and Metadata Creation**

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

333       • In an ontology: a semantic description of the elements of the data
334       • Within a data file: in any number of formats

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

340   **E. Access Rights Management**

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

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

346   **F. Access Policy Contracts**

347   The Data Provider defines policy for others' use of the accessed data, as well as what data will be made
348   available. These contracts specify:

349       • Policies for primary and secondary rights
350       • Agreements

351   To expand this description, the contracts specify acceptable use policies and any specific restrictions on
352   the use of the data, as well as ownership of the original data and any derivative works from the data.

353 **G. Data Distribution Application Programming Interfaces**
354 Technical protocols are defined for different types of data access from data distribution application
355 programming interfaces (APIs), which can include:

356 • File Transfer Protocol (FTP) or streaming
357 • Compression techniques (e.g., single compressed file, split compressed file, )
358 • Authentication methods
359 • Authorization

360 **H. Capabilities Hosting**
361 In addition to offering data downloads, the Data Provider offers several capabilities to access the data,
362 including the following:

363 • Providing query access without transferring the data
364 • Allowing analytic tools to be sent to operate on the data sets

365 For large volumes of data, it may become impractical to move the data to another location for processing.
366 This is often described as moving the processing to the data, rather than the data to the processing.

367 **I. Data Availability Publication**
368 The Data Provider makes available the information needed to know what data or data services they offer.
369 Such publication can consist of the following:

370 • Web description
371 • Services catalog
372 • Data dictionaries
373 • Advertising

374 A number of third-party locations also currently publish a list of links to available datasets (e.g., U.S.
375 Government's Open Data Initiative[2].)

## 2.4 BIG DATA APPLICATION PROVIDER

377 The Big Data Application Provider executes the manipulations of the data lifecycle to meet requirements
378 established by the System Orchestrator, as well as meeting security and privacy requirements. This is
379 where the general capabilities within the Big Data framework are combined to produce the specific data
380 system. Figure 5 lists the actors and activities associated with the Big Data Application Provider.



381 *Figure 5: Big Data Application Provider Actors and Activities*

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

385 **A. Collection**

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

388 • Transport protocol and security
389 • Data format
390 • Metadata

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

393 **B. Preparation**

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

397 • Data validation (e.g., checksums/hashes, format checks)
398 • Data cleansing (e.g., eliminating bad records/fields, deduplication)
399 • Outlier removal
400 • Data conversion (e.g., standardization, reformatting, and encapsulating)
401 • Calculated field creation and indexing
402 • Data aggregation and summarization
403 • Data partition implementation
404 • Data storage preparation
405 • Data virtualization layer

406 Just as data collection may require a number of resources to handle the load, data preparation may also
407 require new resources or new techniques. For large data volumes, data collection is often followed by
408 storage of the data in its raw form. Data preparation processes then occur after the storage and are handled
409 by the application code. This technique of storing raw data first and applying a schema upon interaction
410 with the data is commonly called "schema on read", and is a new area of emphasis in Big Data due to the
411 size of the datasets. When storing a new cleansed copy of the data is prohibitive, the data is stored in its
412 raw form and only prepared for a specific purpose when requested.

413 Data summarization is a second area of added emphasis due to Big Data. With very large datasets, it is
414 difficult to render all the data for visualization. Proper sampling would need some *a priori* understanding
415 of the distribution of the entire dataset. Summarization techniques can characterize local subsets of the
416 data, and then provide these characterizations for visualization as the data is browsed.

417 **C. Analytics**

418 The term data science is used in many ways. While it can refer to the end-to-end data lifecycle, the most
419 common usage focuses on the steps of discovery (i.e., rapid hypothesis-test cycle) for finding value in big
420 volume datasets. This rapid analytics cycle (also described as agile analytics) starts with quick correlation
421 or trending analysis, with greater effort spent on hypotheses that appear most promising.

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

426 While analytic methods have not changed with Big Data, their implementation has changed to
427 accommodate parallel data distribution across a cluster of independent nodes and data access methods.
428 For example, the overall data analytic task may be broken into subtasks that are assigned to the
429 independent data nodes. The results from each subtask are collected and compiled to achieve the final full

430  dataset analysis. Furthermore, data often resided in simple tables or relational databases. With the
431  introduction of new storage paradigms, analytics techniques should be modified for different types of data
432  access.

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

434  • Metadata matching processes
435  • Analysis complexity considerations (e.g., computational, machine learning, data extent, data
436    location)
437  • Analytics latency considerations (e.g., real-time or streaming, near real-time or interactive, batch
438    or offline)
439  • Human-in-the-loop analytics lifecycle (e.g., discovery, hypothesis, hypothesis testing)

440  While these considerations are not new to Big Data, implementing them can be tightly coupled with the
441  specifics of the data storage and the preparation step. In addition, some of the preparation tasks are done
442  during the analytics phase (the schema-on-read discussed above).

## D. Visualization

444  While visualization (or the human in the loop) is often placed under analytics, the added emphasis due to
445  Big Data warrants separate consideration of visualization. The following are three general categories of
446  data visualization:

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

452  Data science relies on the full dataset type of discovery or exploration visualization from which the data
453  scientist would form a hypothesis. While clearly predating Big Data, a greater emphasis now exists on
454  exploratory visualization, as it is immensely helpful in understanding large volumes of repurposed data
455  because the size of the datasets requires new techniques.

456  Explanatory visualization is the creation of a simplified, digestible visual representation of the results,
457  suitable for assisting a decision or communicating the knowledge gained. Again, while this technique has
458  long been in use, there is now greater emphasis to "tell the story". Often this is done through simple
459  visuals or "infographics". Given the large volumes and varieties of data, and the data's potentially
460  complex relationships, the communication of the analytics to a non-analyst audience requires careful
461  visual representation to communicate the results in a way that can be easily consumed.

## E. Access

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

465  • Data export API processes (e.g., protocol, query language)
466  • Data charging mechanisms
467  • Consumer analytics hosting
468  • Analytics as a service hosting

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

473 **2.5 BIG DATA FRAMEWORK PROVIDER**

474 The Big Data Framework Provider has general resources or services to be used by the Big Data
475 Application Provider in the creation of the specific application. There are many new technologies from
476 which the Big Data Application Provider can choose in using these resources and the network to build the
477 specific system. Figure 6 lists the actors and activities associated with the Big Data Framework Provider.



478 *Figure 6: Big Data Framework Provider Actors and Activities*

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

488 **A. Infrastructure Frameworks**

489 Infrastructure frameworks can be grouped as follows:

490 • Networking: These are the resources that transfer data from one resource to another (e.g.,
491 physical, virtual, software defined)
492 • Computing: These are the physical processors and memory that execute and hold the software of
493 the other Big Data system components (e.g., physical resources, operating system, virtual
494 implementation, logical distribution)
495 • Storage: These are resources which provide persistence of the data in a Big Data system (e.g., in-
496 memory, local disk, hardware/software [HW/SW] redundant array of independent disks [RAID],
497 Storage Area Networks [SAN], network-attached storage [NAS])
498 • Environmental: These are the physical plant resources (e.g., power, cooling) that must be
499 accounted for when establishing an instance of a Big Data system

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

502 **B. Data Platform Frameworks**

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

507 • Physical storage (e.g., distributed and non-distributed file systems and object stores)
508 • File systems (e.g., centralized, distributed)

509    • Logical storage

510         o  Simple tuple (e.g., relational, non-relational or not only SQL [NoSQL] tables both row and
511            column)
512         o  Complex tuple (e.g., indexed document store, non-indexed key-value or queues)
513         o  Graph (e.g., property, hyper-graph, triple stores)

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

### C.  Processing Frameworks

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

523    • Data type processing services (e.g., numeric, textual, spatial, images, video)
524    • Schema information or metadata (e.g., on demand, pre-knowledge)
525    • Query frameworks (e.g., relational, arrays)
526    • Temporal frameworks

527         o  Batch (e.g., dense linear algebra, sparse linear algebra, spectral, N-body, structured grids,
528            unstructured grids, Map/Reduce, Bulk Synchronous Parallel [BSP])
529         o  Interactive
530         o  Real-time/streaming (e.g., event ordering, state management, partitioning)

531    • Application frameworks (e.g., automation, test, hosting, workflow)
532    • Messaging/communications frameworks
533    • Resource management frameworks (e.g., cloud/virtualization, intra-framework, inter-framework)

534    Both the Big Data Application Provider activities and the Big Data Framework Provider activities have
535    changed significantly due to Big Data engineering. Currently, the interchange between these two roles
536    operates over a set of independent, yet coupled, resources. It is in this interchange that the new methods
537    for data distribution over a cluster have developed. Just as simulations went through a process of
538    parallelization (or horizontal scaling) to harness massive numbers of independent process to coordinate
539    them to a single analysis, Big Data services now perform the orchestration of data processes over parallel
540    resources.

## 2.6  DATA CONSUMER

542    The Data Consumer receives the value output of the Big Data system. In many respects, the Data
543    Consumer receives the same functionality that the Data Provider brings to the Big Data Application
544    Provider. After the system adds value to the original data sources, the Big Data Application Provider then
545    offers that same functionality to the Data Consumer. There is less change in this role due to Big Data,
546    except, of course, in the desire for Consumers to extract extensive datasets from the Big Data Application
547    Provider. Figure 7 lists the actors and activities associated with the Data Consumer.

548     *Figure 7: Data Consumer Actors and Activities*

549     The activities listed in Figure 7 are explicit to the Data Consumer role within a data system. If the Data
550     Consumer is in fact a follow-on application, then the Data Consumer would look to the Big Data
551     Application Provider for the activities of any other Data Provider. The follow-on application's System
552     Orchestrator would negotiate with this application's System Orchestrator for the types of data wanted,
553     access rights, and other requirements. The Big Data Application Provider would thus serve as the Data
554     Provider, from the perspective of the follow-on application.

555     **A.  Search and Retrieve**

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

558     **B.  Download**

559     All the data from the Data Provider could be exported to the Data Consumer for download.

560     **C.  Analyze Locally**

561     The Data Provider could allow the Data Consumer to run their own application on the data.

562     **D.  Reporting**

563     The data can be presented according to the chosen filters, values, and formatting.

564     **E.  Visualization**

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

566     ## 2.7  MANAGEMENT FABRIC

567     The Big Data characteristics of volume, velocity, variety, and variability demand a versatile management
568     platform for storing, processing, and managing complex data. Management of Big Data systems should
569     handle both system and data related aspects of the Big Data environment. The Management Fabric of the
570     NBDRA encompasses two general groups of activities: system management and Big Data lifecycle
571     management. System management includes activities such as provisioning, configuration, package
572     management, software management, backup management, capability management, resources
573     management, and performance management. Big Data lifecycle management involves activities
574     surrounding the data lifecycle of collection, preparation/curation, analytics, visualization, and access.
575     More discussion about the Management Fabric is needed, particularly with respect to new issues in the
576     management of Big Data and Big Data engineering. This section will be developed in Version 2of this
577     document.

578     Figure 8 lists an initial set of activities associated with the Management role of the NBDRA.
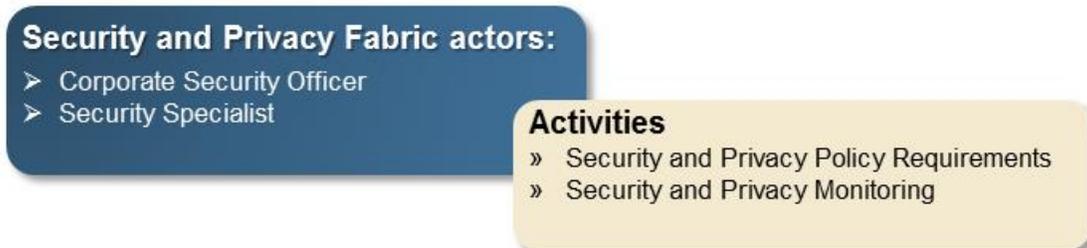
579    *Figure 8: Big Data Management Actors and Activities*

## 2.8  SECURITY AND PRIVACY FABRIC

581    Security and privacy issues affect all other components of the NBDRA, as depicted by the encompassing
582    Security and Privacy box in Figure 1. A Security and Privacy Fabric could interact with the System
583    Orchestrator for policy, requirements, and auditing and also with both the Big Data Application Provider
584    and the Big Data Framework Provider for development, deployment, and operation. These ubiquitous
585    security and privacy activities are described in the *NIST Big Data Interoperability Framework: Volume 4,*
586    *Security and Privacy* document. Figure 9 lists representative actors and activities associated with the
587    Security and Privacy Fabric of the NBDRA. Security and privacy actors and activities will be further
588    developed in Version 2 of *NIST Big Data Interoperability Framework: Volume 4, Security and Privacy*
589    document and summarized in this volume.

590



591    *Figure 9: Big Data Security and Privacy Actors and Activities*
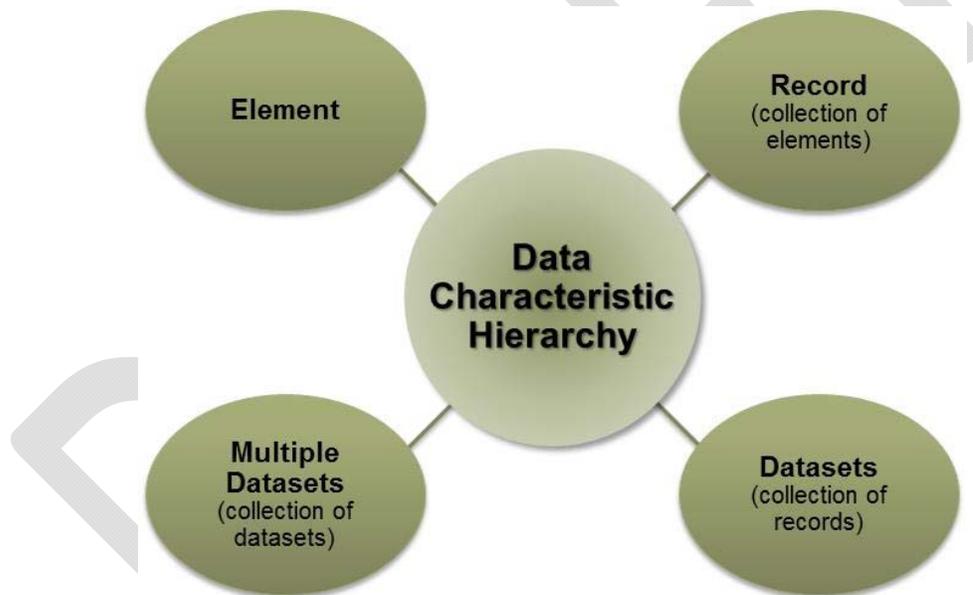
592

# 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. In Section 2 of this document, a taxonomy was presented for the NBDRA, which is described in *NIST Big Data Interoperability Framework: 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 have changed with Big Data can best be done by examining the data scales 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 10. 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, attributes for those objects, and characteristics and subcharacteristics of the attributes. The framework described will help illuminate areas where the driving characteristics for Big Data can be understood in the context of the characteristics of all data.

*Figure 10: 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.

### A. Data Format

Data formats are well characterized through International Organization for Standardization (ISO) standards such as ISO 8601: 2004 Data elements and interchange formats—Information interchange— Representation of dates and times.[3] The data formats have not changed for Big Data.

## B. Data Values and Vocabulary

The data element is characterized 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 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. [a]

## D. Quality And Veracity

Veracity is one of the Big Data 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.[b]

## 3.2 RECORDS

Data elements are grouped into records that describe a specific entity or event or transaction. At the level of records, new emphasis for Big Data begins to be seen.

## 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 XML or HTML, where the contents within the markup can be free text.

These categories again predate Big Data, but two notable changes have occurred with Big Data. First, structured and unstructured data can be stored in one of the new non-relational formats, such as a key-value record structure, a key-document record, or a graph. Second, a greater emphasis is placed on unstructured data due to increasing amounts on the Web (e.g., online images and video.)

## 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 genes and proteins.) Complexity is not new to Big Data.

## C. Volume

Records themselves 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.

## 3.3 DATASETS

Records can be grouped to form datasets. This grouping of records can reveal changes due to Big Data.

[a] Further information about metadata and semantics can be found in: ISO/IEC 11179 Information Technology–Metadata registries; W3C's work on the Semantic Web.

[b] Further information about data quality can be found in ISO 8000 Data Quality.

**Quality and Consistency**

655
656 A new aspect of data quality for records focuses on the characteristic of consistency. As records are
657 distributed horizontally across a collection of data nodes, consistency becomes an issue. In relational
658 databases, consistency was maintained by assuring that all operations in a transaction were completed
659 successfully, otherwise the operations were rolled back. This assured that the database maintained its
660 internal consistency.[c] For Big Data, with multiple nodes and backup nodes, new data is sent in turn to the
661 appropriate nodes. However, constraints may or may not exist to confirm that all nodes have been updated
662 when the query is sent. The time delay in replicating data across nodes can cause an inconsistency. The
663 methods used to update nodes are one of the main areas in which specific implementations of non-
664 relational data storage methods differ. A description of these patterns is a future focus area for this NBD-
665 PWG.

## 3.4 MULTIPLE DATASETS

667 The primary focus on multiple datasets concerns the need to integrate or fuse multiple datasets. The focus
668 here is on the variety characteristic of Big Data. Extensive datasets cannot always be converted into one
669 structure (e.g., all weather data being reported on the same spatio-temporal grid). Since large volume
670 datasets cannot be easily copied into a normalized structure, new techniques are being developed to
671 integrate data as needed.

**Personally Identifiable Information**

673 An area of increasing concern with Big Data is the identification of individuals from the integration of
674 multiple datasets, even when the individual datasets would not allow the identification. For additional
675 discussion, the reader is referred to *NIST Big Data Interoperability Framework: Volume 4, Security and*
676 *Privacy.*

677

678

---

[c] For additional information on this concept the reader is referred to the literature on ACID properties of databases.

# 4 SUMMARY

680 Big Data and data science represent a rapidly changing field due to the recent emergence of new
681 technologies and rapid advancements in methods and perspectives. This document presents a taxonomy
682 for the NBDRA, which is presented in *NIST Big Data Interoperability Framework: Volume 6, Reference*
683 *Architecture*. This taxonomy is a first attempt at providing a hierarchy for categorizing the new
684 components and activities of Big Data systems. This initial version does not incorporate a breakdown of
685 either the Management or the Security and Privacy roles within the NBDRA as those areas need further
686 discussion within the NBD-PWG. In addition, a description of data at different scales was provided to
687 place concepts being ascribed to Big Data into their context. The NBD-PWG will further develop the data
688 characteristics and attributes in the future, in particular determining whether additional characteristics
689 related to data at rest or in-motion should be described. The Big Data patterns related to transactional
690 constraints such as ACID (Atomicity, Consistency, Isolation, Durability—a set of properties guaranteeing
691 reliable processing of database transactions) have not been described here, and are left to future work as
692 the interfaces between resources is an important area for discussion. This document constitutes a first
693 presentation of these descriptions, and future enhancements should provide additional understanding of
694 what is new in Big Data and in specific technology implementations.

695

696 # Appendix A: Acronyms

697 ACID        Atomicity, Consistency, Isolation, Durability

698 APIs        application programming interfaces

699 BSP         Bulk Synchronous Parallel

700 DaaS        Data as a Service

701 FTP         File Transfer Protocol

702 HW/SW RAID  hardware/software redundant array of independent disks

703 IoT         Internet of Things

704 ISO         International Organization for Standardization

705 ITL         Information Technology Laboratory

706 NARA        National Archives and Records Administration

707 NAS         network-attached storage

708 NASA        National Aeronautics and Space Administration

709 NBD-PWG     NIST Big Data Public Working Group

710 NBDRA       NIST Big Data Reference Architecture

711 NIST        National Institute of Standards and Technology

712 NoSQL       not only SQL

713 NSF         National Science Foundation

714 PII         Personally Identifiable Information

715 RFID        radio frequency identification

716 SAN         Storage Area Networks

717 SLAs        service-level agreements

718

# 719    **Appendix B: References**

## 720    DOCUMENT REFERENCES

[1] Tom Kalil, "Big Data is a Big Deal", *The White House, Office of Science and Technology Policy*, accessed February 21, 2014, http://www.whitehouse.gov/blog/2012/03/29/big-data-big-deal.

[2] General Services Administration, "The home of the U.S. Government's open data," *Data.gov,* http://www.data.gov/

[3] ISO 8601: 2004, "Data elements and interchange formats -- Information interchange -- Representation of dates and times," *International Organization of Standardization*, http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=40874