# DRAFT NIST Big Data Interoperability Framework:

# Volume 4, Security and Privacy

NIST Big Data Public Working Group
Security and Privacy Subgroup

NIST

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

# DRAFT NIST Big Data Interoperability Framework:
# Volume 4, Security and Privacy

## Draft Version 1

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

April 2015

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

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

**National Institute of Standards and Technology Special Publication 1500-4**
71 pages (April 6, 2015)

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

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

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

## 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 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 4, contains an exploration of security and privacy topics with respect to Big Data. This volume considers new aspects of security and privacy with respect to Big Data, reviews security and privacy use cases, proposes security and privacy taxonomies, presents details of the Security and Privacy Fabric of the NIST Big Data Reference Architecture (NBDRA), and begins mapping the security and privacy use cases to the NBDRA.

## Keywords

Big Data security, Big Data privacy, Big Data taxonomy, use cases, Big Data characteristics, security and privacy fabric, Big Data risk management, cybersecurity, computer security, information assurance, information security frameworks, encryption standards, role-based access controls, Big Data forensics, Big Data audit

# Acknowledgements

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

The document contains input from members of the NBD-PWG Security and Privacy Subgroup, led by Arnab Roy (Fujitsu), Mark Underwood (Krypton Brothers), and Akhil Manchanda (GE); and the Reference Architecture Subgroup, led by Orit Levin (Microsoft), Don Krapohl (Augmented Intelligence), and James Ketner (AT&T).

NIST SP1500-4, Version 1 has been collaboratively authored by the NBD-PWG. As of the date of this publication, there are over six hundred NBD-PWG participants from industry, academia, and government. Federal agency participants include the National Archives and Records Administration (NARA), National Aeronautics and Space Administration (NASA), National Science Foundation (NSF), and the U.S. Departments of Agriculture, Commerce, Defense, Energy, Health and Human Services, Homeland Security, Transportation, Treasury, and Veterans Affairs.

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

Pw Carey
*Compliance Partners, LLC*

Wo Chang
*NIST*

Brent Comstock
*Cox Communications*

Michele Drgon
*Data Probity*

Roy D'Souza
*AlephCloud Systems, Inc.*

Eddie Garcia
*Gazzang, Inc.*

David Harper
*Johns Hopkins University/ Applied Physics Laboratory*

Pavithra Kenjige
*PK Technologies*

Orit Levin
*Microsoft*

Yale Li
*Microsoft*

Akhil Manchanda
*General Electric*

Marcia Mangold
*General Electric*

Serge Mankovski
*CA Technologies*

Robert Marcus
*ET-Strategies*

Lisa Martinez
*Northbound Transportation and Infrastructure, US*

William Miller
*MaCT USA*

Sanjay Mishra
*Verizon*

Ann Racuya-Robbins
*World Knowledge Bank*

Arnab Roy
*Fujitsu*

Anh-Hong Rucker
*Jet Propulsion Laboratory*

Paul Savitz
*ATIS*

John Schiel
*CenturyLink, Inc.*

Mark Underwood
*Krypton Brothers LLC*

Alicia Zuniga-Alvarado
*Consultant*

The editors for this document were Arnab Roy, Mark Underwood, and Wo Chang.

---

[a] "Contributors" are members of the NIST Big Data Public Working Group who dedicated great effort to prepare and substantial time on a regular basis to research and development in support of this document.

# Notice to Readers

NIST is seeking feedback on the proposed working draft of the *NIST Big Data Interoperability Framework: Volume 4, Security and Privacy*. 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 4.)

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

# Tables

# Executive Summary

This *NIST Big Data Interoperability Framework: Volume 4, Security and Privacy* document was prepared by the NIST Big Data Public Working Group (NBD-PWG) Security and Privacy Subgroup to identify security and privacy issues that are specific to Big Data.

Big Data application domains include health care, drug discovery, insurance, finance, retail and many others from both the private and public sectors. Among the scenarios within these application domains are health exchanges, clinical trials, mergers and acquisitions, device telemetry, targeted marketing and international anti-piracy. Security technology domains include identity, authorization, audit, network and device security, and federation across trust boundaries.

Clearly, the advent of Big Data has necessitated paradigm shifts in the understanding and enforcement of security and privacy requirements. Significant changes are evolving, notably in scaling existing solutions to meet the volume, variety, velocity, and variability of Big Data and retargeting security solutions amid shifts in technology infrastructure, e.g., distributed computing systems and non-relational data storage. In addition, diverse datasets are becoming easier to access and increasingly contain personal content. A new set of emerging issues must be addressed, including balancing privacy and utility, enabling analytics and governance on encrypted data, and reconciling authentication and anonymity.

With the key Big Data characteristics of variety, volume, velocity, and variability in mind, the Subgroup gathered use cases from volunteers, developed a consensus-based security and privacy taxonomy, related the taxonomy to the NIST Big Data Reference Architecture (NBDRA), and validated the NBDRA by mapping the use cases to the NBDRA.

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 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.

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

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

93    • Volume 1, Definitions
94    • Volume 2, Taxonomies
95    • Volume 3, Use Cases and General Requirements
96    • Volume 4, Security and Privacy
97    • Volume 5, Architectures White Paper Survey
98    • Volume 6, Reference Architecture
99    • Volume 7, Standards Roadmap

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

102   Stage 1: Identify the high-level Big Data reference architecture key components, which are
103         technology, infrastructure, and vendor agnostic
104   Stage 2: Define general interfaces between the NIST Big Data Reference Architecture (NBDRA)
105         components
106   Stage 3: Validate the NBDRA by building Big Data general applications through the general interfaces

107   The NBDRA, created in Stage 1 and further developed in Stages 2 and 3, is a high-level conceptual model
108   designed to serve as a tool to facilitate open discussion of the requirements, structures, and operations
109   inherent in Big Data. It is discussed in detail in *NIST Big Data Interoperability Framework: Volume 6,*
110   *Reference Architecture*. Potential areas of future work for the Subgroup during stage 2 are highlighted in
111   Section 1.5 of this volume. The current effort documented in this volume reflects concepts developed
112   within the rapidly evolving field of Big Data.

## 113  1.2  SCOPE AND OBJECTIVES OF THE SECURITY AND PRIVACY SUBGROUP

114   The focus of the NBD-PWG Security and Privacy Subgroup is to form a community of interest from
115   industry, academia, and government with the goal of developing consensus on a reference architecture to
116   handle security and privacy issues across all stakeholders. This includes understanding what standards are
117   available or under development, as well as identifying which key organizations are working on these
118   standards.

119   The scope of the Subgroup's work includes the following topics, some of which will be addressed in
120   future versions of this Volume:

121   • Provide a context from which to begin Big Data-specific security and privacy discussions
122   • Gather input from all stakeholders regarding security and privacy concerns in Big Data
123         processing, storage, and services
124   • Analyze/prioritize a list of challenging security and privacy requirements that may delay or
125         prevent adoption of Big Data deployment
126   • Develop a Security and Privacy Reference Architecture that supplements the NBDRA

127 • Produce a working draft of this Big Data Security and Privacy document
128 • Develop Big Data security and privacy taxonomies
129 • Explore mapping between the Big Data security and privacy taxonomies and the NBDRA
130 • Explore mapping between the use cases and the NBDRA

131 While there are many issues surrounding Big Data security and privacy, the focus of this Subgroup is on
132 the technology aspects of security and privacy with respect to Big Data.

## 1.3 REPORT PRODUCTION

134 The NBD-PWG Security and Privacy Subgroup explored various facets of Big Data security and privacy
135 to develop this document. The major steps involved in this effort included:

136 • Announce that the NBD-PWG Security and Privacy Subgroup is open to the public in order to
137   attract and solicit a wide array of subject matter experts and stakeholders in government, industry,
138   and academia
139 • Identify use cases specific to Big Data security and privacy
140 • Develop a detailed security and privacy taxonomy
141 • Expand the security and privacy fabric of the NBDRA and identify specific topics related to
142   NBDRA components
143 • Begin mapping of identified security and privacy use cases to the NBDRA

144 This report is a compilation of contributions from the PWG. Since this is a community effort, there are
145 several topics covered that are related to security and privacy. While an effort has been made to connect
146 the topics, gaps may come to light that could be addressed in Version 2 of this document.

## 1.4 REPORT STRUCTURE

148 Following this introductory section, the remainder of this document is organized as follows:

149 • Section 2 discusses security and privacy issues particular to Big Data
150 • Section 3 presents examples of security and privacy related use cases
151 • Section 4 offers a preliminary taxonomy for security and privacy
152 • Section 5 introduces the details of a draft NIST Big Data security and privacy reference
153   architecture in relation to the overall NBDRA
154 • Section 6 maps the use cases presented in Section 3 to the NBDRA
155 • Appendix A discusses special security and privacy topics
156 • Appendix B contains information about cloud technology
157 • Appendix C lists the terms and definitions appearing in the taxonomy
158 • Appendix D contains the acronyms used in this document
159 • Appendix E lists the references used in the document

## 1.5 FUTURE WORK ON THIS VOLUME

161 The NBD-PWG Security and Privacy Subgroup plans to further develop several topics for the subsequent
162 version (i.e., Version 2) of this document. These topics include the following:

163    •     Examining closely other existing templates[b] in literature: The templates may be adapted to the
164           Big Data security and privacy fabric to address gaps and to bridge the efforts of this Subgroup
165           with the work of others.
166    •     Further developing the security and privacy taxonomy
167    •     Enhancing the connection between the security and privacy taxonomy and the NBDRA
168           components
169    •     Developing the connection between the security and privacy fabric and the NBDRA
170    •     Expanding the privacy discussion within the scope of this volume
171    •     Exploring governance, risk management, data ownership, and valuation with respect to Big Data
172           ecosystem, with a focus on security and privacy
173    •     Mapping the identified security and privacy use cases to the NBDRA
174    •     Contextualizing the content of Appendix B in the NBDRA
175    •     Exploring privacy in actionable terms with respect to the NBDRA

176 Further topics and direction may be added, as warranted, based on future input and contributions to the
177 Subgroup, including those received during the public comments period.

178

---

[b] There are multiple templates developed by others to adapt as part of a Big Data security metadata model. For instance, the subgroup has considered schemes offered in the NIST Preliminary Critical Infrastructure Cybersecurity Framework (CIICF) of October 2013, http://1.usa.gov/1wQuti1 (accessed January 9, 2015.)

## 2  BIG DATA SECURITY AND PRIVACY

The NBD-PWG Security and Privacy Subgroup began this effort by identifying a number of ways that security and Privacy in Big Data projects can be different from traditional implementations. While not all concepts apply all of the time, the following seven principles were considered representative of a larger set of differences:

1. Big Data projects often encompass heterogeneous components in which a single security scheme has not been designed from the outset.
2. Most security and privacy methods have been designed for batch or online transaction processing systems. Big Data projects increasingly involve one or more streamed data sources that are used in conjunction with data at rest, creating unique security and privacy scenarios.
3. The use of multiple Big Data sources not originally intended to be used together can compromise privacy, security, or both. Approaches to de-identify personally identifiable information (PII) that were satisfactory prior to Big Data may no longer be adequate.
4. An increased reliance on sensor streams, such as those anticipated with the Internet of Things (IoT; e.g., smart medical devices, smart cities, smart homes) can create vulnerabilities that were more easily managed before amassed to Big Data scale.
5. Certain types of data thought to be too big for analysis, such as geospatial and video imaging, will become commodity Big Data sources. These uses were not anticipated and/or may not have implemented security and privacy measures.
6. Issues of veracity, provenance, and jurisdiction are greatly magnified in Big Data. Multiple organizations, stakeholders, legal entities, governments, and an increasing amount of citizens will find data about themselves included in Big Data analytics.
7. Volatility is significant because Big Data scenarios envision that data is permanent by default. Security is a fast-moving field with multiple attack vectors and countermeasures. Data may be preserved beyond the lifetime of the security measures designed to protect it.

### 2.1  OVERVIEW

Security and privacy measures are becoming ever more important with the increase of Big Data generation and utilization and increasingly public nature of data storage and availability.

The importance of security and privacy measures is increasing along with the growth in the generation, access, and utilization of Big Data. Data generation is expected to double every two years to about 40,000 exabytes in 2020. It is estimated that over one third of the data in 2020 could be valuable if analyzed.[2] Less than a third of data needed protection in 2010, but more than 40% of data will need protection in 2020.[3]

Security and privacy measures for Big Data involve a different approach than traditional systems. Big Data is increasingly stored on public cloud infrastructure built by employing various hardware, operating systems, and analytical software. Traditional security approaches usually addressed small-scale systems holding static data on firewalled and semi-isolated networks. The surge in streaming cloud technology necessitates extremely rapid responses to security issues and threats.[4]

Big Data system representations that rely on concepts of actors and roles present a different facet to security and privacy. The Big Data systems should be adapted to the emerging Big Data landscape, which is embodied in many commercial and open source access control frameworks. These security approaches will likely persist for some time and may evolve with the emerging Big Data landscape. Appendix C considers actors and roles with respect to Big Data security and privacy.

222 Big Data is increasingly generated and used across diverse industries such as health care, drug discovery,
223 finance, insurance, and marketing of consumer-packaged goods. Effective communication across these
224 diverse industries will require standardization of the terms related to security and privacy. The NBD-
225 PWG Security and Privacy Subgroup aims to encourage participation in the global Big Data discussion
226 with due recognition to the complex and difficult security and privacy requirements particular to Big
227 Data.

228 There is a large body of work in security and privacy spanning decades of academic study and
229 commercial solutions. While much of that work is not conceptually distinct from Big Data, it may have
230 been produced using different assumptions. One of the primary objectives of this document is to
231 understand how Big Data security and privacy requirements arise out of the defining characteristics of
232 Big Data, and how these requirements are differentiated from traditional security and privacy
233 requirements.

234 The following list is a representative—though not exhaustive—list of differences between what is new for
235 Big Data and the requirements that informed previous big system security and privacy.

- **Big Data may be gathered from diverse end points.** Actors include more types than just
  traditional providers and consumers—data owners, such as mobile users and social network users,
  are primary actors in Big Data. Devices that ingest data streams for physically distinct data
  consumers may also be actors. This alone is not new, but the mix of human and device types is on
  a scale that is unprecedented. The resulting combination of threat vectors and potential protection
  mechanisms to mitigate them is new.
- **Data aggregation and dissemination must be secured inside the context of a formal,
  understandable framework.** The availability of data and transparency of its current and past use
  by data consumers is an important aspect of Big Data. However, Big Data systems may be
  operational outside formal, readily understood frameworks, such as those designed by a single
  team of architects with a clearly defined set of objectives. In some settings, where such
  frameworks are absent or have been unsystematically composed, there may be a need for public
  or walled garden portals and ombudsman-like roles for data at rest. These system combinations
  and unforeseen combinations call for a renewed Big Data framework.
- **Data search and selection can lead to privacy or security policy concerns.** There is a lack of
  systematic understanding of the capabilities that should be provided by a data provider in this
  respect.[c] A combination of well-educated users, well-educated architects, and system protections
  may be needed, as well as excluding databases or limiting queries that may be foreseen as
  enabling re-identification. If a key feature of Big Data is, as one analyst called it, "the ability to
  derive differentiated insights from advanced analytics on data at any scale," the search and
  selection aspects of analytics will accentuate security and privacy concerns.[5]
- **Privacy-preserving mechanisms are needed for Big Data, such as for Personally Identifiable
  Information (PII).** Because there may be disparate, potentially unanticipated processing steps
  between the data owner, provider, and data consumer, the privacy and integrity of data coming
  from end points should be protected at every stage. End-to-end information assurance practices
  for Big Data are not dissimilar from other systems but must be designed on a larger scale.
- **Big Data is pushing beyond traditional definitions for information trust, openness, and
  responsibility.** Governance, previously consigned to static roles and typically employed in larger
  organizations, is becoming an increasingly important intrinsic design consideration for Big Data
  systems.
- **Information assurance and disaster recovery for Big Data Systems may require unique and
  emergent practices.** Because of its extreme scalability, Big Data presents challenges for

---

[c] Reference to NBDRA Data Provider.

268         information assurance (IA) and disaster recovery (DR) practices that were not previously
269         addressed in a systematic way. Traditional backup methods may be impractical for Big Data
270         systems. In addition, test, verification, and provenance assurance for Big Data replicas may not
271         complete in time to meet temporal requirements that were readily accommodated in smaller
272         systems.

273     •   **Big Data creates potential targets of increased value.** The effort required to consummate
274         system attacks will be scaled to meet the opportunity value. Big Data systems will present
275         concentrated, high value targets to adversaries. As Big Data becomes ubiquitous, such targets are
276         becoming more numerous—a new information technology scenario in itself.

277     •   **Risks have increased for de-anonymization and transfer of PII without consent traceability**.
278         Security and privacy can be compromised through unintentional lapses or malicious attacks on
279         data integrity. Managing data integrity for Big Data presents additional challenges related to all
280         the Big Data characteristics, but especially for PII. While there are technologies available to
281         develop methods for de-identification, some experts caution that equally powerful methods can
282         leverage Big Data to re-identify personal information. For example, the availability of
283         unanticipated data sets could make re-identification possible. Even when technology is able to
284         preserve privacy, proper consent and use may not follow the path of the data through various
285         custodians.

286     •   **Emerging Risks in Open Data and Big Science.** Data identification, metadata tagging,
287         aggregation, and segmentation—widely anticipated for data science and open datasets—if not
288         properly managed, may have degraded veracity because they are derived and not primary
289         information sources. Retractions of peer-reviewed research due to inappropriate data
290         interpretations may become more commonplace as researchers leverage third party Big Data.
291

## 2.2 EFFECTS OF BIG DATA CHARACTERISTICS ON SECURITY AND PRIVACY

293 Variety, volume, velocity, and variability are key characteristics of Big Data and commonly referred to as
294 the Vs of Big Data. Where appropriate, these characteristics shaped discussions within the NBD-PWG
295 Security and Privacy Subgroup. While the Vs provide a useful shorthand description, used in the public
296 discourse about Big Data, there are other important characteristics of Big Data that affect security and
297 privacy, such as veracity, validity, and volatility. These elements are discussed below with respect to their
298 impact on Big Data security and privacy.

### 2.2.1 VARIETY

300 Variety describes the organization of the data—whether the data is structured, semi-structured, or
301 unstructured. Retargeting traditional relational database security to non-relational databases has been a
302 challenge[6]. These systems were not designed with security and privacy in mind, and these functions are
303 usually relegated to middleware. Traditional encryption technology also hinders organization of data
304 based on semantics. The aim of standard encryption is to provide semantic security, which means that the
305 encryption of any value is indistinguishable from the encryption of any other value. Therefore, once
306 encryption is applied, any organization of the data that depends on any property of the data values
307 themselves are rendered ineffective, whereas organization of the metadata, which may be unencrypted,
308 may still be effective.

309 An emergent phenomenon introduced by Big Data variety that has gained considerable importance is the
310 ability to infer identity from anonymized datasets by correlating with apparently innocuous public
311 databases. While several formal models to address privacy preserving data disclosure have been
312 proposed,[7] [8] in practice, sensitive data is shared after sufficient removal of apparently unique identifiers
313 by the processes of anonymization and aggregation. This is an ad hoc process that is often based on

314    empirical evidence[9] and has led to many instances of de-anonymization in conjunction with publicly
315    available data.[10]

## 2.2.2 VOLUME

317    The volume of Big Data describes how much data is coming in. In Big Data parlance, this typically
318    ranges from gigabytes to exabytes. As a result, the volume of Big Data has necessitated storage in multi-
319    tiered storage media. The movement of data between tiers has led to a requirement of cataloging threat
320    models and a surveying of novel techniques. The threat model for network-based, distributed, auto-tier
321    systems includes the following major scenarios: confidentiality and integrity, provenance, availability,
322    consistency, collusion attacks, roll-back attacks and recordkeeping disputes. [11]

323    A flip side of having volumes of data is that analytics can be performed to help detect security breach
324    events. This is an instance where Big Data technologies can fortify security. This document addresses
325    both facets of Big Data security.

## 2.2.3 VELOCITY

327    Velocity describes the speed at which data is processed. The data usually arrives in batches or is streamed
328    continuously. As with certain other non-relational databases, distributed programming frameworks were
329    not developed with security and privacy in mind.[12] Malfunctioning computing nodes might leak
330    confidential data. Partial infrastructure attacks could compromise a significantly large fraction of the
331    system due to high levels of connectivity and dependency. If the system does not enforce strong
332    authentication among geographically distributed nodes, rogue nodes can be added that can eavesdrop on
333    confidential data.

## 2.2.4 VERACITY

335    Big Data veracity and validity encompass several subcharacteristics:

336    **Provenance**—or what some have called veracity in keeping with the V theme—is important for both data
337    quality and for protecting security and maintaining privacy policies. Big Data frequently moves across
338    individual boundaries to groups and communities of interest, and across state, national, and international
339    boundaries. Provenance addresses the problem of understanding the data's original source, such as
340    through metadata, though the problem extends beyond metadata maintenance. Various approaches have
341    been tried, such as for glycoproteomics,[13] but no clear guidelines yet exist.

342    A common understanding holds that provenance data is metadata establishing pedigree and chain of
343    custody, including calibration, errors, missing data (e.g., time stamp, location, equipment serial number,
344    transaction number, and authority.)

345    Some experts consider the challenge of defining and maintaining metadata to be the overarching
346    principle, rather than provenance. The two concepts, though, are clearly interrelated.

347    **Veracity** (in some circles also called Provenance, though the two terms are not identical) also
348    encompasses information assurance for the methods through which information was collected. For
349    example, when sensors are used, traceability, calibration, version, sampling, and device configuration is
350    needed.

351    Curation is an integral concept which binds veracity and provenance to principles of governance as well
352    as to data quality assurance, Curation, for example, may improve raw data by fixing errors, filling in gaps,
353    modeling, calibrating values, ordering data collection.

354    **Validity** refers to the accuracy and correctness of data. Traditionally this is referred to data quality. In the
355    Big Data security scenario, validity refers to a host of assumptions about data from which analytics are
356    being applied. For example, continuous and discrete measurements have different properties. The field
357    "gender" can be coded as 1=Male, 2=Female, but 1.5 does not mean halfway between male and female.

358 In the absence of such constraints, an analytical tool can make inappropriate conclusions. There are many
359 types of validity whose constraints are far more complex. By definition, Big Data allows for aggregation
360 and collection across disparate data sets in ways not envisioned by system designers.

361 Several examples of "invalid" uses for Big Data have been cited. Click fraud, conducted on a Big Data
362 scale, but which can be detected using Big Data techniques, has been cited as the cause of perhaps $11.6
363 billion in wasted advertisement spending. A software executive listed seven different types of online ad
364 fraud, including non-human generated impressions, non-human generated clicks, hidden ads,
365 misrepresented sources, all-advertising sites, malicious ad injections, and policy-violating content such as
366 pornography or privacy violations.[14] Each of these can be conducted at Big Data scale and may require
367 Big Data solutions to detect and combat.

368 Despite initial enthusiasm, some trend producing applications that use social media to predict the
369 incidence of flu have been called into question. A study by Lazer et al.[15] suggested that one application
370 overestimated the prevalence of flu for 100 of 108 weeks studied. Careless interpretation of social media
371 is possible when attempts are made to characterize or even predict consumer behavior using imprecise
372 meanings and intentions for "like" and "follow."

373 These examples show that what passes for "valid" Big Data can be innocuously lost in translation,
374 interpretation or intentionally corrupted to malicious intent.

### 2.2.5 VOLATILITY

376 Volatility of data—how data management changes over time—directly affects provenance. Big Data is
377 transformational in part because systems may produce indefinitely persisting data—data that outlives the
378 instruments on which it was collected; the architects who designed the software that acquired, processed,
379 aggregated, and stored it; and the sponsors who originally identified the project's data consumers.

380 Roles are time-dependent in nature. Security and privacy requirements can shift accordingly. Governance
381 can shift as responsible organizations merge or even disappear.

382 While research has been conducted into how to manage temporal data (e.g., in e-science for satellite
383 instrument data),[16] there are few standards beyond simplistic timestamps and even fewer common
384 practices available as guidance. To manage security and privacy for long-lived Big Data, data temporality
385 should be taken into consideration.

## 2.3 RELATION TO CLOUD

387 Many Big Data systems will be designed using cloud architectures. Any strategy to achieve proper access
388 control and security risk management within a Big Data cloud ecosystem enterprise architecture for
389 industry must address the complexities associated with cloud-specific security requirements triggered by
390 cloud characteristics, including, but not limited to, the following:

391 - Broad network access
392 - Decreased visibility and control by consumer
393 - Dynamic system boundaries and commingled roles and responsibilities between consumers and
394   providers
395 - Multi-tenancy
396 - Data residency
397 - Measured service
398 - Order-of-magnitude increases in scale (on demand), dynamics (elasticity and cost optimization),
399   and complexity (automation and virtualization)

400 These cloud computing characteristics often present different security risks to an organization than the
401 traditional information technology solutions, altering the organization's security posture.

402    To preserve security when migrating data to the cloud, organizations need to identify all cloud-specific,
403    risk-adjusted security controls or components in advance. It may be necessary in some situations to
404    requests from the cloud service providers through contractual means and service-level agreements that all
405    require security components and controls to be fully and accurately implemented.

406    A further discussion of internal security considerations within cloud ecosystems can be found in
407    Appendix B. Future versions of this document will contextualize the content of Appendix B in the
408    NBDRA.

409

# 3 EXAMPLE USE CASES FOR SECURITY AND PRIVACY

There are significant Big Data challenges in science and engineering. Many of these are described in the use cases in *NIST Big Data Interoperability Framework: Volume 3, Use Cases and General Requirements*. However, these use cases focused primarily on science and engineering applications for which security and privacy were secondary concerns—if the latter had any impact on system architecture at all. Consequently, a different set of use cases was developed in the preparation of this document specifically to discover security and privacy issues. Some of these use cases represent inactive or legacy applications, but were selected because they demonstrate characteristic security / privacy design patterns.

The use cases selected for security and privacy are presented in the following subsections. The use cases included are grouped to organize this presentation, as follows: retail/marketing, healthcare, cybersecurity, government, industrial, aviation, and transportation. However, these groups do not represent the entire spectrum of industries affected by Big Data security and privacy.

The use cases were collected when the reference architecture was not mature. The use cases were collected to identify representative security and privacy scenarios thought to be suitably classified as particular to Big Data. An effort was made to map the use cases to the NBDRA. In Version 2, additional mapping of the use cases to the NBDRA and taxonomy will be developed. Parts of this document were developed in parallel and the connections will be strengthened in Version 2.

## 3.1 RETAIL/MARKETING

### 3.1.1 CONSUMER DIGITAL MEDIA USAGE

**Scenario Description:** Consumers, with the help of smart devices, have become very conscious of price, convenience, and access before they decide on a purchase. Content owners license data for use by consumers through presentation portals, such as Netflix, iTunes, and others.

Comparative pricing from different retailers, store location and/or delivery options, and crowd-sourced rating have become common factors for selection. To compete, retailers are keeping a close watch on consumer locations, interests, and spending patterns to dynamically create marketing strategies and sell products that consumers do not yet know they want.

**Current Security and Privacy:** Individual data is collected by several means, including smartphone GPS (global positioning system) or location, browser use, social media, and applications (apps) on smart devices.

- Privacy:
  - o Most data collection means described above offer weak privacy controls. In addition, consumer unawareness and oversight allow third parties to legitimately capture information. Consumers can have limited to no expectation of privacy in this scenario.
- Security:
  - o Controls are inconsistent and/or not established appropriately to achieve the following:
    - Isolation, containerization, and encryption of data
    - Monitoring and detection of threats
    - Identification of users and devices for data feed
    - Interfacing with other data sources
    - Anonymization of users: while some data collection and aggregation uses anonymization techniques, individual users can be re-identified by leveraging other public Big Data pools
    - Original digital rights management (DRM) techniques were not built to scale to meet demand for the forecasted use for the data. "DRM refers to a broad category of access

453          control technologies aimed at restricting the use and copy of digital content on a wide
454          range of devices."[17] DRM can be compromised, diverted to unanticipated purposes,
455          defeated, or fail to operate in environments with Big Data characteristics—especially
456          velocity and aggregated volume
457 **Current Research:** There is limited research on enabling privacy and security controls that protect
458 individual data (whether anonymized or non-anonymized).

### 3.1.2 NIELSEN HOMESCAN: PROJECT APOLLO

460 **Scenario Description:** Nielsen Homescan is a subsidiary of Nielsen that collects family-level retail
461 transactions. Project Apollo was a project designed to better unite advertising content exposure to
462 purchase behavior among Nielsen panelists. Project Apollo did not proceed beyond a limited trial, but
463 reflects a Big Data intent. The description is a best-effort general description and is not an official
464 perspective from Nielsen, Arbitron or the various contractors involved in the project. The information
465 provided here should be taken as illustrative rather than as a historical record.

466 A general retail transaction has a checkout receipt that contains all SKUs (stock keeping units) purchased,
467 time, date, store location, etc. Nielsen Homescan collected purchase transaction data using a statistically
468 randomized national sample. As of 2005, this data warehouse was already a multi-terabyte data set. The
469 warehouse was built using structured technologies but was built to scale many terabytes. Data was
470 maintained in house by Homescan but shared with customers who were given partial access through a
471 private web portal using a columnar database. Additional analytics were possible using third party
472 software. Other customers would only receive reports that include aggregated data, but greater granularity
473 could be purchased for a fee.

474 **Then Current (2005-2006) Security and Privacy:**

475 - Privacy: There was a considerable amount of PII data. Survey participants are compensated in
476      exchange for giving up segmentation data, demographics, and other information.
477 - Security: There was traditional access security with group policy, implemented at the field level
478      using the database engine, component-level application security and physical access controls.
479 - There were audit methods in place, but were only available to in-house staff. Opt-out data
480      scrubbing was minimal.

### 3.1.3 WEB TRAFFIC ANALYTICS

482 **Scenario Description:** Visit-level webserver logs are high-granularity and voluminous. To be useful, log
483 data must be correlated with other (potentially Big Data) data sources, including page content (buttons,
484 text, navigation events), and marketing-level events such as campaigns, media classification, etc. There
485 are discussions—if not deployment—of plans for traffic analytics using complex event processing (CEP)
486 in real time. One nontrivial problem is segregating traffic types, including internal user communities, for
487 which collection policies and security are different.

488 **Current Security and Privacy:**

489 - Non-European Union (EU): Opt-in defaults are relied upon to gain visitor consent for tracking.
490      Internet Protocol (IP) address logging enables some analysts to identify visitors down to the level
491      of a city block
492 - Media access control (MAC) address tracking enables analysts to identify IP devices, which is a
493      form of PII
494 - Some companies allow for purging of data on demand, but most are unlikely to expunge
495      previously collected web server traffic
496 - The EU has stricter regulations regarding collection of such data, which is treated as PII. Such
497      web traffic is to be scrubbed (anonymized) or reported only in aggregate, even for multinationals
498      operating in the EU but based in the United States

## 499  3.2 HEALTHCARE

### 500  3.2.1 HEALTH INFORMATION EXCHANGE

501  **Scenario Description:** Health Information Exchanges (HIEs) facilitate sharing of healthcare information
502  that might include electronic health records (EHRs) so that the information is accessible to relevant
503  covered entities, but in a manner that enables patient consent.

504  HIEs tend to be federated, where the respective covered entity retains custodianship of its data. This poses
505  problems for many scenarios, such as emergencies, for a variety of reasons that include technical (such as
506  interoperability), business, and security concerns.

507  Cloud enablement of HIEs, through strong cryptography and key management, that meets the Health
508  Insurance Portability and Accountability Act (HIPAA) requirements for protected health information
509  (PHI)—ideally without requiring the cloud service operator to sign a business associate agreement
510  (BAA)—would provide several benefits, including patient safety, lowered healthcare costs, and regulated
511  accesses during emergencies that might include break-the-glass and Centers for Disease Control and
512  Prevention (CDC) scenarios.

513  The following are some preliminary scenarios that have been proposed by the NBD PWG:

514  - **Break-the-Glass:** There could be situations where the patient is not able to provide consent due
515    to a medical situation, or a guardian is not accessible, but an authorized party needs immediate
516    access to relevant patient records. Cryptographically enhanced key life cycle management can
517    provide a sufficient level of visibility and nonrepudiation that would enable tracking violations
518    after the fact
519  - **Informed Consent:** When there is a transfer of EHRs between covered entities and business
520    associates, it would be desirable and necessary for patients to be able to convey their approval, as
521    well as to specify what components of their EHR can be transferred (e.g., their dentist would not
522    need to see their psychiatric records.) Through cryptographic techniques, one could leverage the
523    ability to specify the fine-grain cipher text policy that would be conveyed. (For related standards
524    efforts regarding consent, see NIST 800-53, Appendix J, Section IP-1), US DHS Health IT Policy
525    Committee, Privacy and Security Workgroup) and Health Level Seven (HL7) International
526    Version 3 standards for Data Access Consent, Consent Directives)
527  - **Pandemic Assistance:** There will be situations when public health entities, such as the CDC and
528    perhaps other nongovernmental organizations that require this information to facilitate public
529    safety, will require controlled access to this information, perhaps in situations where services and
530    infrastructures are inaccessible. A cloud HIE with the right cryptographic controls could release
531    essential information to authorized entities through authorization and audits in a manner that
532    facilitates the scenario requirement

533  **Project Current and/or Proposed Security and Privacy:**

534  - Security:
535    - Lightweight but secure off-cloud encryption: There is a need for the ability to perform
536      lightweight but secure off-cloud encryption of an EHR that can reside in any container that
537      ranges from a browser to an enterprise server, and that leverages strong symmetric
538      cryptography
539    - Homomorphic encryption
540    - Applied cryptography: Tight reductions, realistic threat models, and efficient techniques
541  - Privacy:
542    - Differential privacy: Techniques for guaranteeing against inappropriate leakage of PII
543    - HIPAA

### 3.2.2 GENETIC PRIVACY

**Scenario Description:** A consortium of policy makers, advocacy organizations, individuals, academic centers, and industry has formed an initiative, **Free the Data!**, to fill the public information gap caused by the lack of available genetic information for the BRCA1 and BRCA2 genes. The consortium also plans to expand to provide other types of genetic information in open, searchable databases, including the National Center for Biotechnology Information's database, ClinVar. The primary founders of this project include Genetic Alliance, the University of California San Francisco, InVitae Corporation, and patient advocates.

This initiative invites individuals to share their genetic variation on their own terms and with appropriate privacy settings in a public database so that their family, friends, and clinicians can better understand what the mutation means. Working together to build this resource means working toward a better understanding of disease, higher-quality patient care, and improved human health.

**Current Security and Privacy:**

- Security:
    - Secure Sockets Layer (SSL)-based authentication and access control. Basic user registration with low attestation level
    - Concerns over data ownership and custody upon user death
    - Site administrators may have access to data—strong encryption and key escrow are recommended
- Privacy:
    - Transparent, logged, policy-governed controls over access to genetic information
    - Full lifecycle data ownership and custody controls

### 3.2.3 PHARMA CLINICAL TRIAL DATA SHARING[18]

**Scenario Description:** Companies routinely publish their clinical research, collaborate with academic researchers, and share clinical trial information on public websites, atypically at three different stages: the time of patient recruitment, after new drug approval, and when investigational research programs have been discontinued. Access to clinical trial data is limited, even to researchers and governments, and no uniform standards exist.

The Pharmaceutical Research and Manufacturers of America (PhRMA) represents the country's leading biopharmaceutical researchers and biotechnology companies. In July 2013, PhRMA joined with the European Federation of Pharmaceutical Industries and Associations (EFPIA) in adopting joint Principles for Responsible Clinical Trial Data Sharing. According to the agreement, companies will apply these Principles as a common baseline on a voluntary basis, and PhRMA encouraged all medical researchers, including those in academia and government, to promote medical and scientific advancement by adopting and implementing the following commitments:

- Enhancing data sharing with researchers
- Enhancing public access to clinical study information
- Sharing results with patients who participate in clinical trials
- Certifying procedures for sharing trial information
- Reaffirming commitments to publish clinical trial results

**Current and Proposed Security and Privacy:**

PhRMA does not directly address security and privacy, but these issues were identified either by PhRMA or by reviewers of the proposal.

- Security:
    - Longitudinal custody beyond trial disposition is unclear, especially after firms merge or dissolve
    - Standards for data sharing are unclear

590  o   There is a need for usage audit and security
591  o   Publication restrictions: Additional security will be required to protect the rights of
592      publishers; for example, Elsevier or Wiley
593  •   Privacy:
594  o   Patient-level data disclosure—elective, per company
595  o   The PhRMA mentions anonymization (re-identification), but mentions issues with small
596      sample sizes
597  o   Study-level data disclosure—elective, per company

## 3.3  CYBERSECURITY

### 3.3.1  NETWORK PROTECTION

600  **Scenario Description:** Network protection includes a variety of data collection and monitoring. Existing
601  network security packages monitor high-volume data sets, such as event logs, across thousands of
602  workstations and servers, but they are not yet able to scale to Big Data. Improved security software will
603  include physical data correlates (e.g., access card usage for devices as well as building entrance/exit) and
604  likely be more tightly integrated with applications, which will generate logs and audit records of
605  previously undetermined types or sizes. Big Data analytics systems will be required to process and
606  analyze this data to deliver meaningful results. These systems could also be multi-tenant, catering to more
607  than one distinct company.

608  This scenario highlights two subscenarios:

609  •   Security for Big Data
610  •   Big Data for security
611  **Current Security and Privacy:**

612  •   Security in this area is mature; privacy concepts less so.
613  o   Traditional policy-type security prevails, though temporal dimension and monitoring of
614      policy modification events tends to be nonstandard or unaudited
615  o   Cybersecurity apps run at high levels of security and thus require separate audit and security
616      measures
617  o   No cross-industry standards exist for aggregating data beyond operating system collection
618      methods
619  o   Implementing Big Data cybersecurity should include data governance, encryption/key
620      management, and tenant data isolation/containerization
621  o   Volatility should be considered in the design of backup and disaster recovery for Big Data
622      cybersecurity. The useful life of logs may extend beyond the lifetime of the devices which
623      created them
624  •   Privacy:
625  o   Enterprise authorization for data release to state/national organizations
626  o   Protection of PII data
627  Currently vendors are adopting Big Data analytics for mass-scale log correlation and incident response,
628  such as for security information and event management (SIEM).

## 3.4  GOVERNMENT

### 3.4.1  MILITARY: UNMANNED VEHICLE SENSOR DATA

631  **Scenario Description**: Unmanned vehicles (or drones) and their onboard sensors (e.g., streamed video)
632  can produce petabytes of data that should be stored in nonstandard formats. These streams are often not
633  processed in real time, but the U.S. Department of Defense (DOD) is buying technology to make this

634  possible. Because correlation is key, GPS, time, and other data streams must be co-collected. The Bradley
635  Manning leak situation is one security breach use case.

636  **Current Security and Privacy:**

637  - Separate regulations for agency responsibility apply.
638    - For domestic surveillance: The U.S. Federal Bureau of Investigation (FBI)
639    - For overseas surveillance: Multiple agencies, including the U.S. Central Intelligence Agency
640      (CIA) and various DOD agencies
641  - Not all uses will be military; for example, the National Oceanic and Atmospheric Administration
642  - Military security classifications are moderately complex and determined on need to know basis
643  - Information assurance practices are rigorously followed, unlike in some commercial settings

644  **Current Research:**

645  - Usage is audited where audit means are provided, software is not installed/deployed until
646    'certified,' and development cycles have considerable oversight; for example, the U.S. Army's
647    Army Regulation 25-2[19]
648  - Insider threats (e.g., Edward Snowden, Bradley Manning, and spies) are being addressed in
649    programs such as the Defense Advanced Research Projects Agency's (DARPA) Cyber-Insider
650    Threat (CINDER) program. This research and some of the unfunded proposals made by industry
651    may be of interest

## 3.4.2 EDUCATION: COMMON CORE STUDENT PERFORMANCE REPORTING

653  **Scenario Description:** Forty-five states have decided to unify standards for K–12 student performance
654  measurement. Outcomes are used for many purposes, and the program is incipient, but it will obtain
655  longitudinal Big Data status. The data sets envisioned include student-level performance across students'
656  entire school history and across schools and states, as well as taking into account variations in test stimuli.

657  **Current Security and Privacy:**

658  - Data is scored by private firms and forwarded to state agencies for aggregation. Classroom,
659    school, and district identifiers remain with the scored results. The status of student PII is
660    unknown; however, it is known that teachers receive classroom-level performance feedback. The
661    extent of student/parent access to test results is unclear
662  - Privacy-related disputes surrounding education Big Data are illustrated by the reluctance of states
663    to participate in the InBloom initiative[20]
664  - According to some reports, parents can opt students out of state tests, so opt-out records must also
665    be collected and used to purge ineligible student records. [21]

666  **Current Research:**

667  - Longitudinal performance data would have value for program evaluators if data scales up
668  - Data-driven learning[22] will involve access to students' performance data, probably more often
669    than at test time, and at higher granularity, thus requiring more data. One example enterprise is
670    Civitas Learning's[23] predictive analytics for student decision making

## 3.5 INDUSTRIAL: AVIATION

## 3.5.1 SENSOR DATA STORAGE AND ANALYTICS

673  **Scenario Description**: Most commercial airlines are equipped with hundreds of sensors to constantly
674  capture engine and/or aircraft health information during a flight. For a single flight, the sensors may
675  collect multiple gigabytes of data and transfer this data stream to Big Data analytics systems. Several
676  companies manage these Big Data analytics systems, such as parts/engine manufacturers, airlines, and
677  plane manufacturers, and data may be shared across these companies. The aggregated data is analyzed for

678  maintenance scheduling, flight routines, etc. One common request from airline companies is to secure and
679  isolate their data from competitors, even when data is being streamed to the same analytics system.
680  Airline companies also prefer to control how, when, and with whom the data is shared, even for analytics
681  purposes. Most of these analytics systems are now being moved to infrastructure cloud providers.

682  **Current and Proposed Security and Privacy:**

683  • Encryption at rest: Big Data systems should encrypt data stored at the infrastructure layer so that
684     cloud storage administrators cannot access the data
685  • Key management: The encryption key management should be architected so that end customers
686     (e.g., airliners) have sole/shared control on the release of keys for data decryption
687  • Encryption in motion: Big Data systems should verify that data in transit at the cloud provider is
688     also encrypted
689  • Encryption in use: Big Data systems will desire complete obfuscation/encryption when
690     processing data in memory (especially at a cloud provider)
691  • Sensor validation and unique identification (e.g., device identity management)
692  Researchers are currently investigating the following security enhancements:

693  • Virtualized infrastructure layer mapping on a cloud provider
694  • Homomorphic encryption
695  • Quorum-based encryption
696  • Multi-party computational capability
697  • Device public key infrastructure (PKI)

698  ## 3.6 TRANSPORTATION

699  ### 3.6.1 CARGO SHIPPING
700  The following use case outlines how the shipping industry (e.g., FedEx, UPS, DHL) regularly uses Big
701  Data. Big Data is used in the identification, transport, and handling of items in the supply chain. The
702  identification of an item is important to the sender, the recipient, and all those in between with a need to
703  know the location of the item while in transport and the time of arrival. Currently, the status of shipped
704  items is not relayed through the entire information chain. This will be provided by sensor information,
705  GPS coordinates, and a unique identification schema based on the new International Organization for
706  Standardization (ISO) 29161 standards under development within the ISO technical committee ISO JTC1
707  SC31 WG2. The data is updated in near real time when a truck arrives at a depot or when an item is
708  delivered to a recipient. Intermediate conditions are not currently known, the location is not updated in
709  real-time, and items lost in a warehouse or while in shipment represent a potential problem for homeland
710  security. The records are retained in an archive and can be accessed for system-determined number of
711  days.

712                                   *Figure 1: Cargo Shipping Scenario*

713

714

## 715 4 TAXONOMY OF SECURITY AND PRIVACY TOPICS

716 A candidate set of topics from the Cloud Security Alliance Big Data Working Group (CSA BDWG)
717 article, *Top Ten Challenges in Big Data Security and Privacy Challenges*, was used in developing these
718 security and privacy taxonomies.[24] Candidate topics and related material used in preparing this section are
719 provided for reference in Appendix A.

720 A taxonomy for Big Data security and privacy should encompass the aims of existing, useful taxonomies.
721 While many concepts surrounding security and privacy exist, the objective in the taxonomies contained
722 herein is to highlight and refine new or emerging principles specific to Big Data.

723 The following subsections present an overview of each security and privacy taxonomy, along with lists of
724 topics encompassed by the taxonomy elements. These lists are the results of preliminary discussions of
725 the Subgroup and may be developed further in Version 2.

### 726 4.1 CONCEPTUAL TAXONOMY OF SECURITY AND PRIVACY TOPICS

727 The conceptual security and privacy taxonomy, presented in Figure 2, contains four main groups: data
728 confidentiality; data provenance; system health; and public policy, social, and cross-organizational topics.
729 The first three topics broadly correspond with the traditional classification of confidentiality, integrity,
730 and availability (CIA), reoriented to parallel Big Data considerations.



731 *Figure 2: Security and Privacy Conceptual Taxonomy*

### 732 4.1.1 DATA CONFIDENTIALITY
733    • Confidentiality of data in transit: For example, enforced by using Transport Layer Security (TLS)
734    • Confidentiality of data at rest
735       o Policies to access data based on credentials
736          ▪ Systems: Policy enforcement by using systems constructs such as Access Control Lists
737            (ACLs) and Virtual Machine (VM) boundaries
738          ▪ Crypto-enforced: Policy enforcement by using cryptographic mechanisms, such as PKI
739            and identity/attribute-based encryption
740    • Computing on encrypted data

741    o Searching and reporting: Cryptographic protocols that support searching and reporting on
742     encrypted data—any information about the plain text not deducible from the search criteria is
743     guaranteed to be hidden
744    o Homomorphic encryption: Cryptographic protocols that support operations on the underlying
745     plain text of an encryption—any information about the plain text is guaranteed to be hidden
746   • Secure data aggregation: Aggregating data without compromising privacy
747   • Data anonymization
748    o De-identification of records to protect privacy
749   • Key management
750    o As noted by Chandramouli and Iorga, cloud security for cryptographic keys, an essential
751     building block for security and privacy, takes on "additional complexity," which can be
752     rephrased for Big Data settings: (1) greater variety due to more cloud consumer-provider
753     relationships, and (2) greater demands and variety of infrastructures "on which both the Key
754     Management System and protected resources are located." [25]
755    o Big Data systems are not purely cloud systems, but as is noted elsewhere in this document,
756     the two are closely related. One possibility is to retarget the key management framework that
757     Chandramouli and Iorga developed for cloud service models to the NBDRA security and
758     privacy fabric. Cloud models would correspond to the NBDRA and cloud security concepts
759     to the proposed fabric. NIST 800-145 provides definitions for cloud computing concepts,
760     including infrastructure as a service (IaaS), platform as a service (PaaS), and software as a
761     service (SaaS) cloud service models [26]
762    o Challenges for Big Data key management systems (KMS) reflect demands imposed by Big
763     Data characteristics (i.e., volume, velocity, variety, and variability). For example, leisurely
764     key creation and workflow associated with legacy—and often fastidious—data warehouse
765     key creation is insufficient for Big Data systems deployed quickly and scaled up using
766     massive resources. The lifetime for a Big Data KMS will likely outlive the period of
767     employment of the Big Data system architects who designed it. Designs for location, scale,
768     ownership, custody, provenance, and audit for Big Data key management is an aspect of a
769     security and privacy fabric

770 **4.1.2 PROVENANCE**
771   • End-point input validation: A mechanism to validate whether input data is coming from an
772    authenticated source, such as digital signatures
773    o Syntactic: Validation at a syntactic level
774    o Semantic: Semantic validation is an important concern. Generally, semantic validation would
775     validate typical business rules such as a due date. Intentional or unintentional violation of
776     semantic rules can lock up an application. This could also happen when using data translators
777     that do not recognize the particular variant. Protocols and data formats may be altered by a
778     vendor using, for example, a reserved data field that will allow their products to have
779     capabilities that differentiate them from other products. This problem can also arise in
780     differences in versions of systems for consumer devices, including mobile devices. The
781     semantics of a message and the data to be transported should be validated to verify, at a
782     minimum, conformity with any applicable standards. The use of digital signatures will be
783     important to provide assurance that the data from a sensor or data provider has been verified
784     using a validator or data checker and is, therefore, valid. This capability is important,
785     particularly if the data is to be transformed or involved in the curation of the data. If the data
786     fails to meet the requirements, it may be discarded, and if the data continues to present a
787     problem, the source may be restricted in its ability to submit the data. These types of errors
788     would be logged and prevented from being disseminated to consumers
789    o Digital signatures will be very important in the Big Data system

| | | |
|---|---|---|
| 790 | • | Communication integrity: Integrity of data in transit, enforced, for example, by using TLS |
| 791 | • | Authenticated computations on data: Ensuring that computations taking place on critical |
| 792 | | fragments of data are indeed the expected computations |
| 793 | | o Trusted platforms: Enforcement through the use of trusted platforms, such as Trusted |
| 794 | | Platform Modules (TPMs) |
| 795 | | o Crypto-enforced: Enforcement through the use of cryptographic mechanisms |
| 796 | • | Granular audits: Enabling audit at high granularity |
| 797 | • | Control of valuable assets |
| 798 | | o Life cycle management |
| 799 | | o Retention and disposition |
| 800 | | o DRM |

### 801 4.1.3 SYSTEM HEALTH

| | | |
|---|---|---|
| 802 | • | Security against denial-of-service (DoS) |
| 803 | | o Construction of cryptographic protocols proactively resistant to DoS |
| 804 | • | Big Data for Security |
| 805 | | o Analytics for security intelligence |
| 806 | | o Data-driven abuse detection |
| 807 | | o Big Data analytics on logs, cyberphysical events, intelligent agents |
| 808 | | o Security breach event detection |
| 809 | | o Forensics |
| 810 | | o Big Data in support of resilience |

### 811 4.1.4 PUBLIC POLICY, SOCIAL AND CROSS-ORGANIZATIONAL TOPICS

812 The following set of topics is drawn from an Association for Computing Machinery (ACM) grouping.[27]
813 Each of these topics has Big Data security and privacy dimensions that could affect how a fabric overlay
814 is implemented for a specific Big Data project. For instance, a medical devices project might need to
815 address human safety risks, whereas a banking project would be concerned with different regulations
816 applying to Big Data crossing borders. Further work to develop these concepts for Big Data is anticipated
817 by the Subgroup.

| | | |
|---|---|---|
| 818 | • | Abuse and crime involving computers |
| 819 | • | Computer-related public / private health systems |
| 820 | • | Ethics (within data science, but also across professions) |
| 821 | • | Human safety |
| 822 | • | Intellectual property rights and associated information management[d] |
| 823 | • | Regulation |
| 824 | • | Transborder data flows |
| 825 | • | Use/abuse of power |
| 826 | • | Assistive technologies for persons with disabilities (e.g., added or different security / privacy |
| 827 | | measures may be needed for subgroups within the population) |
| 828 | • | Employment (e.g., regulations applicable to workplace law may govern proper use of Big Data |
| 829 | | produced or managed by employees) |
| 830 | • | Social aspects of ecommerce |
| 831 | • | Legal: Censorship, taxation, contract enforcement, forensics for law enforcement |

---

[d] For further information, see the frameworks suggested by the Association for Information and Image Management (AIIM; http://www.aiim.org/) and the MIKE 2.0 Information Governance Association (http://mike2.openmethodology.org/wiki/MIKE2.0_Governance_Association)

## 4.2  OPERATIONAL TAXONOMY OF SECURITY AND PRIVACY TOPICS

Current practice for securing Big Data systems is diverse, employing widely disparate approaches that often are not part of a unified conceptual framework. The elements of the operational taxonomy, shown in Figure 3, represent groupings of practical methodologies. These elements are classified as "operational" because they address specific vulnerabilities or risk management challenges to the operation of Big Data systems. At this point in the standards development process, these methodologies have not been incorporated as part of a cohesive security fabric. They are potentially valuable checklist-style elements that can solve specific security or privacy needs. Future work must better integrate these methodologies with risk management guidelines developed by others (e.g., NIST Special Publication 800-37 Guide for Applying the Risk Management Framework to Federal Information Systems[28] and COBIT Risk IT Framework[29].)

In the proposed operational taxonomy, broad considerations of the conceptual taxonomy appear as recurring features. For example, confidentiality of communications can apply to governance of data at rest and access management, but it is also part of a security metadata model.[30]

The operational taxonomy will overlap with small data taxonomies while drawing attention to specific issues with Big Data.[31] [32]



*Figure 3: Security and Privacy Operational Taxonomy*

### 4.2.1  DEVICE AND APPLICATION REGISTRATION

- Device, User, Asset, Services, and Applications Registration: Includes registration of devices in machine to machine (M2M) and IoT networks, DRM-managed assets, services, applications, and user roles
- Security Metadata Model

854         o   The metadata model maintains relationships across all elements of a secured system. It
855              maintains linkages across all underlying repositories. Big Data often needs this added
856              complexity due to its longer life cycle, broader user community, or other aspects
857         o   A Big Data model must address aspects such as data velocity, as well as temporal aspects of
858              both data and the life cycle of components in the security model

859     •  Policy Enforcement
860         o   Environment build
861         o   Deployment policy enforcement
862         o   Governance model
863         o   Granular policy audit
864         o   Role-specific behavioral profiling

## 4.2.2  IDENTITY AND ACCESS MANAGEMENT

866     •  Virtualization layer identity (e.g., cloud console, platform as a service [PaaS])
867         o   Trusted platforms
868     •  Application layer Identity
869     •  End-user layer identity management
870         o   Roles
871     •  Identity provider (IdP)
872         o   An IdP is defined in the Security Assertion Markup Language (SAML). [33] In a Big Data
873              ecosystem of data providers, orchestrators, resource providers, framework providers, and data
874              consumers, a scheme such as the SAML/Security Token Service (STS) or eXtensible Access
875              Control Markup Language (XACML) is seen as a helpful—but not proscriptive—way to
876              decompose the elements in the security taxonomy
877         o   Big Data may have multiple IdPs. An IdP may issue identities (and roles) to access data from
878              a resource provider. In the SAML framework, trust is shared via SAML/web services
879              mechanisms at the registration phase
880         o   In Big Data, due to the density of the data, the user "roams" to data (whereas in conventional
881              virtual private network [VPN]-style scenarios, users roam across trust boundaries). Therefore,
882              the conventional authentication/authorization (authn/authz) model needs to be extended
883              because the relying party is no longer fully trusted—they are custodians of somebody else's
884              data. Data is potentially aggregated from multiple resource providers
885         o   One approach is to extend the claims-based methods of SAML to add security and privacy
886              guarantees
887     •  Additional XACML Concepts
888         o   XACML introduces additional concepts that may be useful for Big Data security. In Big
889              Data, parties are not just sharing claims, but also sharing policies about what is authorized.
890              There is a policy access point at every data ownership and authoring location, and a policy
891              enforcement point at the data access. A policy enforcement point calls a designated policy
892              decision point for an auditable decision. In this way, the usual meaning of non-repudiation
893              and trusted third parties is extended in XACML. Big Data presumes an abundance of
894              policies, "points," and identity issuers, as well as data
895              ▪   Policy authoring points
896              ▪   Policy decision points
897              ▪   Policy enforcement point
898              ▪   Policy access points

## 4.2.3  DATA GOVERNANCE

900  However large and complex Big Data becomes in terms of data volume, velocity, variety, and variability,
901  Big Data governance will, in some important conceptual and actual dimensions, be much larger. Big Data

902 without Big Data governance may become less useful to its stakeholders. To stimulate positive change,
903 data governance will need to persist across the data lifecycle—at rest, in motion, in incomplete stages,
904 and transactions—while serving the security and privacy of the young, the old, individuals as
905 organizations, and organizations as organizations. It will need to cultivate economic benefits and
906 innovation but also enable freedom of action and foster individual and public welfare. It will need to rely
907 on standards governing technologies and practices not fully understood while integrating the human
908 element. Big Data governance will require new perspectives yet accept the slowness or inefficacy of some
909 current techniques. Some data governance considerations are listed below.

910 **Big Data Apps to Support Governance:** The development of new applications employing Big Data
911 principles and designed to enhance governance may be among the most useful Big Data applications on
912 the horizon.

913 - Encryption and key management
914   - At rest
915   - In memory
916   - In transit
917 - Isolation/containerization
918 - Storage security
919 - Data loss prevention and detection
920 - Web services gateway
921 - Data transformation
922   - Aggregated data management
923   - Authenticated computations
924   - Computations on encrypted data
925 - Data life cycle management
926   - Disposition, migration, and retention policies
927   - PII microdata as "hazardous" [34]
928   - De-identification and anonymization
929   - Re-identification risk management
930 - End-point validation
931 - DRM
932 - Trust
933 - Openness
934 - Fairness and information ethics [35]

935 ### 4.2.4 INFRASTRUCTURE MANAGEMENT
936 Infrastructure management involves security and privacy considerations related to hardware operation and
937 maintenance. Some topics related to infrastructure management are listed below.

938 - Threat and vulnerability management
939   - DoS-resistant cryptographic protocols
940 - Monitoring and alerting
941   - As noted in the Critical Infrastructure Cybersecurity Framework (CIICF), Big Data affords
942     new opportunities for large-scale security intelligence, complex event fusion, analytics, and
943     monitoring
944 - Mitigation
945   - Breach mitigation planning for Big Data may be qualitatively or quantitatively different
946 - Configuration Management
947   - Configuration management is one aspect of preserving system and data integrity. It can
948     include the following:

949        o   Patch management
950        o   Upgrades
951    •   Logging
952        o   Big Data must produce and manage more logs of greater diversity and velocity. For example,
953            profiling and statistical sampling may be required on an ongoing basis
954    •   Malware surveillance and remediation
955        o   This is a well-understood domain, but Big Data can cross traditional system ownership
956            boundaries. Review of NIST's "Identify, Protect, Detect, Respond, and Recover" framework
957            may uncover planning unique to Big Data
958    •   Network boundary control
959        o   Establishes a data-agnostic connection for a secure channel
960            ▪   Shared services network architecture, such as those specified as "secure channel use cases
961                and requirements" in the European Telecommunications Standards Institute (ETSI) TS
962                102 484 Smart Card specifications [36]
963            ▪   Zones/cloud network design (including connectivity)
964    •   Resilience, Redundancy, and Recovery
965        o   Resilience
966            ▪   The security apparatus for a Big Data system may be comparatively fragile in comparison
967                to other systems. A given security and privacy fabric may be required to consider this.
968                Resilience demands are domain-specific, but could entail geometric increases in Big Data
969                system scale
970        o   Redundancy
971            ▪   Redundancy within Big Data systems presents challenges at different levels. Replication
972                to maintain intentional redundancy within a Big Data system takes place at one software
973                level. At another level, entirely redundant systems designed to support failover, resilience
974                or reduced data center latency may be more difficult due to velocity, volume or other
975                aspects of Big Data
976        o   Recovery
977            ▪   Recovery for Big Data security failures may require considerable advance provisioning
978                beyond that required for small data. Response planning and communications with users
979                may be on a similarly large scale

### 980   4.2.5  RISK AND ACCOUNTABILITY

981 Risk and accountability encompass the following topics:

982    •   Accountability
983        o   Information, process, and role behavior accountability can be achieved through various
984            means, including:
985            ▪   Transparency portals and inspection points
986            ▪   Forward- and reverse-provenance inspection
987    •   Compliance
988        o   Big Data compliance spans multiple aspects of the security and privacy taxonomy, including
989            privacy, reporting, and nation-specific law
990    •   Forensics
991        o   Forensics techniques enabled by Big Data
992        o   Forensics used in Big Data security failure scenarios
993    •   Business risk level
994        o   Big Data risk assessments should be mapped to each element of the taxonomy.[37] Business
995            risk models can incorporate privacy considerations

## 4.3 ROLES RELATED TO SECURITY AND PRIVACY TOPICS

Discussions of Big Data security and privacy should be accessible to a diverse audience, including individuals who specialize in cryptography, security, compliance, or information technology. In addition, there are domain experts and corporate decision makers who should understand the costs and impact of these controls. Ideally, these documents would be prefaced by information that would help specialists find the content relevant to them. The specialists could then provide feedback on those sections.

Organizations typically contain diverse roles and workflows for participating in a Big Data ecosystem. Therefore, this document proposes a pattern to help identify the "axis" of an individual's roles and responsibilities, as well as classify the security controls in a similar manner to make these more accessible to each class.

### 4.3.1 INFRASTRUCTURE MANAGEMENT

Typically, the individual role axis contains individuals and groups who are responsible for technical reviews before their organization is on-boarded in a data ecosystem. After the on-boarding, they are usually responsible for addressing defects and security issues.

When infrastructure technology personnel work across organizational boundaries, they accommodate diverse technologies, infrastructures, and workflows and the integration of these three elements. For Big Data security, these include identity, authorization, access control, and log aggregation.

Their backgrounds and practices, as well as the terminologies they use, tend to be uniform, and they face similar pressures within their organizations to constantly do more with less. "Save money" is the underlying theme, and infrastructure technology usually faces pressure when problems arise.

### 4.3.2 GOVERNANCE, RISK MANAGEMENT, AND COMPLIANCE

Data governance is a fundamental element in the management of data and data systems. Data governance refers to administering, or formalizing, discipline (e.g., behavior patterns) around the management of data. Risk management involves the evaluation of positive and negative risks resulting from the handling of Big Data. Compliance encompasses adherence to laws, regulations, protocols, and other guiding rules for operations related to Big Data. Typically, governance, risk management, and compliance (GRC) is a function that draws participation from multiple areas of the organization, such as legal, human resources (HR), information technology (IT), and compliance. In some industries and agencies, there may be a strong focus on compliance, often in isolation from disciplines.

Professionals working in GRC tend to have similar backgrounds, share a common terminology, and employ similar processes and workflows, which typically influence other organizations within the corresponding vertical market or sector.

Within an organization, GRC professionals aim to protect the organization from negative outcomes that might arise from loss of intellectual property, liability due to actions by individuals within the organization, and compliance risks specific to its vertical market.

In larger enterprises and government agencies, GRC professionals are usually assigned to legal, marketing, or accounting departments or staff positions connected to the CIO. Internal and external auditors are often involved.

Smaller organizations may create, own, or process Big Data, yet may not have GRC systems and practices in place, due to the newness of the Big Data scenario to the organization, a lack of resources, or other factors specific to small organizations. Prior to Big Data, GRC roles in smaller organizations received little attention.

A one-person company can easily construct a Big Data application and inherit numerous unanticipated related GRC responsibilities. This is a new GRC scenario.

1040     A security and privacy fabric entails additional data and process workflow in support of GRC, which is
1041     most likely under the control of the System Orchestrator component of the NBDRA, as explained in
1042     Section 5.

1043     ### 4.3.3 INFORMATION WORKER
1044     Information workers are individuals and groups who work on the generation, transformation, and
1045     consumption of content. Due to the nascent nature of the technologies and related businesses in which
1046     they work, they tend to use common terms at a technical level within a specialty. However, their roles and
1047     responsibilities and the related workflows do not always align across organizational boundaries. For
1048     example, a data scientist has deep specialization in the content and its transformation, but may not focus
1049     on security or privacy until it adds effort, cost, risk, or compliance responsibilities to the process of
1050     accessing domain-specific data or analytical tools.

1051     Information workers may serve as data curators. Some may be research librarians, operate in quality
1052     management roles, or be involved in information management roles such as content editing, search
1053     indexing, or performing forensic duties as part of legal proceedings.

1054     Information workers are exposed to a great number of products and services. They are under pressure
1055     from their organizations to deliver concrete business value from these new Big Data analytics capabilities
1056     by monetizing available data, monetizing the capability to transform data by becoming a service provider,
1057     or optimizing and enhancing business by consuming third-party data.

1058     ## 4.4 RELATION OF ROLES TO THE SECURITY AND PRIVACY CONCEPTUAL
1059     TAXONOMY

1060     The next sections cover the four components of the conceptual taxonomy: data confidentiality, data
1061     provenance, system health, and public policy, social and cross-organizational topics. To leverage these
1062     three axes and to facilitate collaboration and education, a stakeholder can be defined as an individual or
1063     group within an organization who is directly affected by the selection and deployment of a Big Data
1064     solution. A ratifier is defined as an individual or group within an organization who is tasked with
1065     assessing the candidate solution before it is selected and deployed. For example, a third-party security
1066     consultant may be deployed by an organization as a ratifier, and an internal security specialist with an
1067     organization's IT department might serve as both a ratifier and a stakeholder if tasked with ongoing
1068     monitoring, maintenance, and audits of the security.

1069     The upcoming sections also explore potential gaps that would be of interest to the anticipated
1070     stakeholders and ratifiers who reside on these three new conceptual axes.

1071     ### 4.4.1 DATA CONFIDENTIALITY
1072     IT specialists who address cryptography should understand the relevant definitions, threat models,
1073     assumptions, security guarantees, and core algorithms and protocols. These individuals will likely be
1074     ratifiers, rather than stakeholders. IT specialists who address end-to-end security should have an
1075     abbreviated view of the cryptography, as well as a deep understanding of how the cryptography would be
1076     integrated into their existing security infrastructures and controls.

1077     GRC should reconcile the vertical requirements (e.g., HIPAA requirements related to EHRs) and the
1078     assessments by the ratifiers that address cryptography and security. GRC managers would in turn be
1079     ratifiers to communicate their interpretation of the needs of their vertical. Persons in these roles also serve
1080     as stakeholders due to their participation in internal and external audits and other workflows.

1081     ### 4.4.2 PROVENANCE
1082     Provenance (or veracity) is related in some ways to data privacy, but it might introduce information
1083     workers as ratifiers because businesses may need to protect their intellectual property from direct leakage

1084 or from indirect exposure during subsequent Big Data analytics. IWs would need to work with the
1085 ratifiers from cryptography and security to convey the business need, as well as understand how the
1086 available controls may apply.

1087 Similarly, when an organization is obtaining and consuming data, information workers may need to
1088 confirm that the data provenance guarantees some degree of information integrity and address incorrect,
1089 fabricated, or cloned data before it is presented to an organization.

1090 Additional risks to an organization could arise if one of its data suppliers does not demonstrate the
1091 appropriate degree of care in filtering or labeling its data. As noted in the U.S. Department of Health and
1092 Human Services (HHS) press release announcing the HIPAA final omnibus rule:

1093 *"The changes announced today expand many of the requirements to business associates*
1094 *of these entities that receive protected health information, such as contractors and*
1095 *subcontractors. Some of the largest breaches reported to HHS have involved business*
1096 *associates. Penalties are increased for noncompliance based on the level of negligence*
1097 *with a maximum penalty of $1.5 million per violation."[38]*

1098 Organizations using or sharing health data among ecosystem partners, including mobile apps and SaaS
1099 providers, will need to verify that the proper legal agreements are in place to require data veracity and
1100 provenance.

### 4.4.3 SYSTEM HEALTH MANAGEMENT
1101
1102 System health is typically the domain of IT, and IT managers will be ratifiers and stakeholders of
1103 technologies, protocols, and products that are used for system health. IT managers will also design how
1104 the responsibilities to maintain system health would be shared across the organizations that provide data,
1105 analytics, or services—an area commonly known as operations support systems (OSS) in the telecom
1106 industry, which has significant experience in syndication of services.

1107 Security and cryptography specialists should scrutinize the system health to spot potential gaps in the
1108 operational architectures. The likelihood of gaps increases when a system infrastructure includes diverse
1109 technologies and products.

1110 System health is an umbrella concept that emerges at the intersection of information worker and
1111 infrastructure management. As with human health, monitoring nominal conditions for Big Data systems
1112 may produce Big Data volume and velocity—two of the Big Data characteristics. Following the human
1113 health analogy, some of those potential signals reflect defensive measures such as white cell count. Others
1114 could reflect compromised health, such as high blood pressure. Similarly, Big Data systems may employ
1115 applications like Security Information and Event Management (SIEM) or Big Data analytics more
1116 generally to monitor system health.

1117 Volume, velocity, variety, and variability of Big Data systems health make it different from small data
1118 system health. Health tools and design patterns for existing systems are likely insufficient to handle Big
1119 Data—including Big Data security and privacy. At least one commercial web services provider has
1120 reported that its internal accounting and systems management tool uses more resources than any other
1121 single application. The volume of system events and the complexity of event interactions is a challenge
1122 that demands Big Data solutions to defend Big Data systems. Managing systems health—including
1123 security—will require roles defined as much by the tools needed to manage as by the organizational
1124 context. Stated differently, Big Data is transforming the role of the Computer Security Officer.

1125 For example, one aspect motivated by the DevOps movement (i.e., move toward blending tasks
1126 performed by applications development and systems operations teams) is the rapid launch,
1127 reconfiguration, redeployment and distribution of Big Data systems. Tracking intended vs. accidental or
1128 malicious configuration changes is increasingly a Big Data challenge.

### 4.4.4 PUBLIC POLICY, SOCIAL, AND CROSS-ORGANIZATIONAL TOPICS

Roles in setting public policy related to security and privacy are established in the U.S. by federal agencies such as the Federal Trade Commission, the Food and Drug Administration or the DHHS Office of National Coordinator. DHS is responsible for aspects of domestic U.S. computer security through the activities of US-CERT. Social roles include the influence of NGO's, interest groups, professional organizations and standards development organizations. Cross-organizational roles include design patterns employed across or within certain industries such as pharmaceuticals, logistics, manufacturing, distribution to facilitate data sharing, curation, and even orchestration. Big Data frameworks will impact, and are impacted by cross-organizational considerations, possibly industry-by-industry. Further work to develop these concepts for Big Data is anticipated by the Subgroup.

## 4.5 ADDITIONAL TAXONOMY TOPICS

Additional areas have been identified but not carefully scrutinized, and it is not yet clear whether these would fold into existing categories or if new categories for security and privacy concerns would need to be identified and developed. Some candidate topics are briefly described below.

### 4.5.1 PROVISIONING, METERING, AND BILLING

Provisioning, metering and billing are elements in typically commercial systems used to manage assets, meter their use and invoice clients for that usage. Commercial pipelines for Big Data can be constructed and monetized more readily if these systems are agile in offering services, metering access suitably, and integrating with billing systems. While this process can be manual for a small number of participants, it can become complex very quickly when there are many suppliers, consumers, and service providers. Information workers and IT professionals who are involved with existing business processes would be candidate ratifiers and stakeholders. Assuring privacy and security of provisioning and metering data may or may not have already been designed into these systems. The scope of metering and billing data will explode, so potential uses and risks have likely not been fully explored.

There are both veracity and validity concerns with these systems. GRC considerations, such as audit and recovery, may overlap with provisioning and metering.

### 4.5.2 DATA SYNDICATION

A feature of Big Data systems is that data is bought and sold as a valuable asset. That Google Search is free relies on users giving up information about their search terms on a Big Data scale. Google and Facebook can choose to repackage and syndicate that information for use by others for a fee.

Similar to service syndication, a data ecosystem is most valuable if any participant can have multiple roles, which could include supplying, transforming, or consuming Big Data. Therefore, a need exists to consider what types of data syndication models should be enabled; again, information workers and IT professionals are candidate ratifiers and stakeholders, For some domains, more complex models may be required to accommodate PII, provenance and governance. Syndication involves transfer of risk and responsibility for security and privacy.

# 1166   5  SECURITY AND PRIVACY FABRIC

1167   Security and privacy considerations are a fundamental aspect of the NBDRA. Using the material gathered
1168   for this volume and extensive brainstorming among the NBD-PWG Security and Privacy Subgroup
1169   members and others, the following proposal for a security and privacy fabric was developed. [e]

1170   **Security and Privacy Fabric:** Security and privacy considerations form a fundamental aspect of the
1171   NBDRA. This is geometrically depicted in Figure 4 by the Security and Privacy Fabric surrounding the
1172   five main components, since all components are affected by security and privacy considerations. Thus,
1173   the role of security and privacy is correctly depicted in relation to the components but does not expand
1174   into finer details, which may be more accurate but are best relegated to a more detailed security and
1175   privacy reference architecture. The Data Provider and Data Consumer are included in the Security and
1176   Privacy Fabric since, at the least, they should agree on the security protocols and mechanisms in place.
1177   The Security and Privacy Fabric is an approximate representation that alludes to the intricate
1178   interconnected nature and ubiquity of security and privacy throughout the NBDRA.

1179   This pervasive dimension is depicted in Figure 4 by the presence of the security and privacy fabric
1180   surrounding all of the functional components., NBD-PWG decided to include the Data Provider and Data
1181   Consumer as well as the Big Data Application and Framework Providers in the Security and Privacy
1182   Fabric because these entities should agree on the security protocols and mechanisms in place. The *NIST*
1183   *Big Data Interoperability Framework: Volume 6, Reference Architecture* document discusses in detail the
1184   other components of the NBDRA.

1185   At this time, explanations as to how the proposed fabric concept is implemented across each NBDRA
1186   component are cursory—more suggestive than prescriptive. However, it is believed that, in time, a
1187   template will evolve and form a sound basis for more detailed iterations.

1188

---

[e] The concept of a "fabric" for security and privacy has precedent in the hardware world, where the notion of a
fabric of interconnected nodes in a distributed computing environment was introduced. Computing fabrics were
invoked as part of cloud and grid computing, as well as for commercial offerings from both hardware and software
manufacturers.

*Figure 4: NIST Big Data Reference Architecture*

Figure 4 introduces two new concepts that are particularly important to security and privacy considerations: information value chain and IT value chain.

**Information value chain:** While it does not apply to all domains, there may be an implied processing progression through which information value is increased, decreased, refined, defined, or otherwise transformed. Application of provenance-preservation and other security mechanisms at each stage may be conditioned by the state-specific contributions to information value.

**IT value chain** Platform-specific considerations apply to Big Data systems when scaled-up or -out. In the process of scaling, specific security, privacy, or GRC mechanism or practices may need to be invoked.

## 5.1 SECURITY AND PRIVACY FABRIC IN THE NBDRA

Figure 5 provides an overview of several security and privacy topics with respect to some key NBDRA components and interfaces. The figure represents a beginning characterization of the interwoven nature of the Security and Privacy Fabric with the NBDRA components.

It is not anticipated that Figure 5 will be further developed for Version 2 of this document. However, the relationships between the Security and Privacy Fabric and the NBDRA and the Security and Privacy Taxonomy and the NBDRA will be investigated for Version 2 of this document.
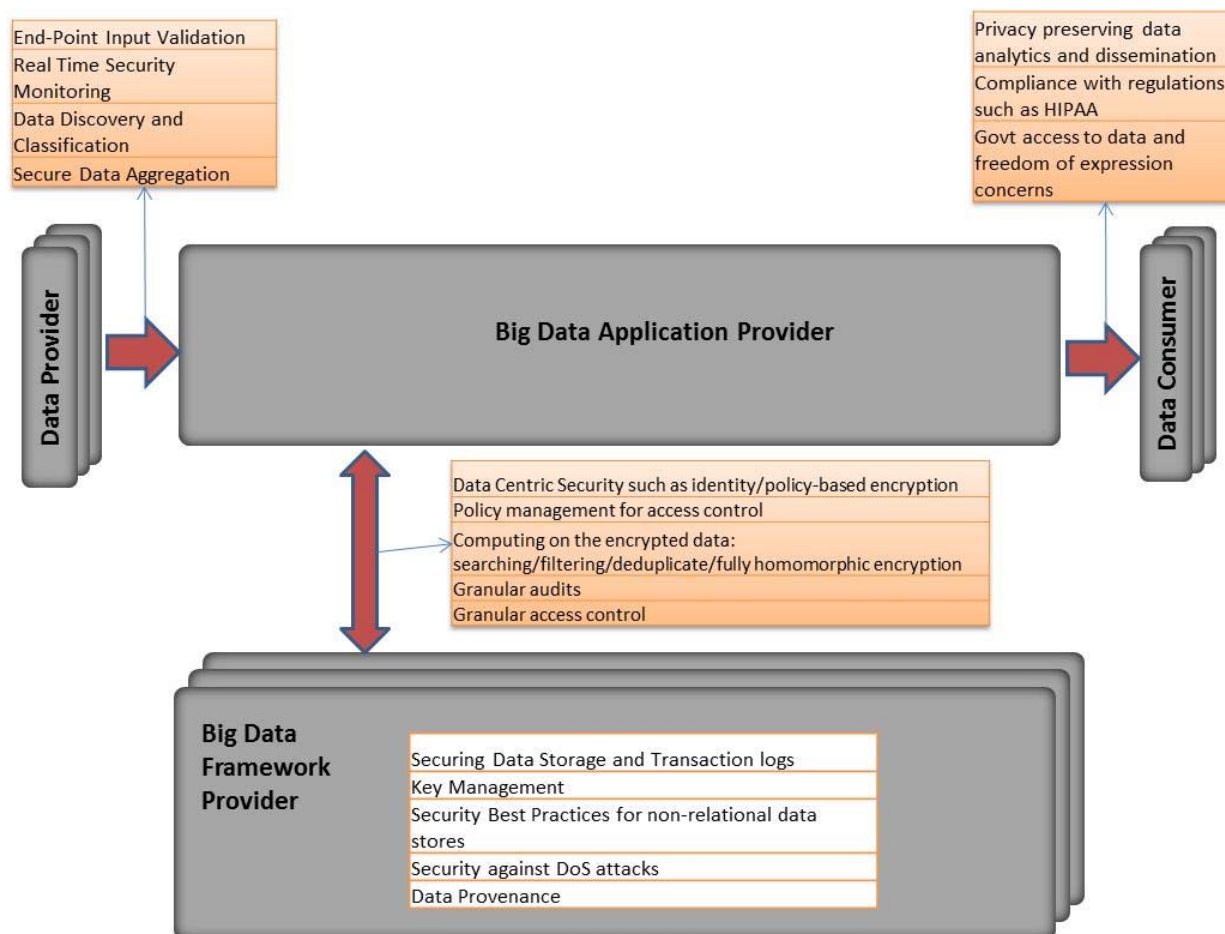
**Figure 5: Notional Security and Privacy Fabric Overlay to the NBDRA**

1207    The groups and interfaces depicted in Figure 5 are described below.

1208    **A. INTERFACE BETWEEN DATA PROVIDERS → BIG DATA APPLICATION PROVIDER**

1209    Data coming in from data providers may have to be validated for integrity and authenticity. Incoming
1210    traffic may be maliciously used for launching DoS attacks or for exploiting software vulnerabilities on
1211    premise. Therefore, real-time security monitoring is useful. Data discovery and classification should be
1212    performed in a manner that respects privacy.

1213    **B. INTERFACE BETWEEN BIG DATA APPLICATION PROVIDER → DATA CONSUMER**

1214    Data, including aggregate results delivered to data consumers, must preserve privacy. Data accessed by
1215    third parties or other entities should follow legal regulations such as HIPAA. Concerns include access to
1216    sensitive data by the government.

1217    **C. INTERFACE BETWEEN APPLICATION PROVIDER ↔ BIG DATA FRAMEWORK PROVIDER**

1218    Data can be stored and retrieved under encryption. Access control policies should be in place to assure
1219    that data is only accessed at the required granularity with proper credentials. Sophisticated encryption
1220    techniques can allow applications to have rich policy-based access to the data as well as enable searching,
1221    filtering on the encrypted data, and computations on the underlying plaintext.

1222    **D. INTERNAL TO BIG DATA FRAMEWORK PROVIDER**

1223    Data at rest and transaction logs should be kept secured. Key management is essential to control access
1224    and keep track of keys. Non-relational databases should have a layer of security measures. Data

1225 provenance is essential to having proper context for security and function of the data at every stage. DoS
1226 attacks should be mitigated to assure availability of the data.

1227 **E. SYSTEM ORCHESTRATOR**

1228 A System Orchestrator may play a critical role in identifying, managing, auditing, and sequencing Big
1229 Data processes across the components. For example, a workflow that moves data from a collection stage
1230 to further preparation may implement aspects of security or privacy.

1231 System Orchestrators present an additional, attractive attack surface for adversaries. System Orchestrators
1232 often require permanent or transitory elevated permissions. System Orchestrators present opportunities to
1233 implement security mechanisms, monitor provenance, access systems management tools, provide audit
1234 points, and inadvertently subjugate privacy or other information assurance measures.

## 5.2 PRIVACY ENGINEERING PRINCIPLES

1236 Big Data security and privacy should leverage existing standards and practices. In the privacy arena, a
1237 systems approach that considers privacy throughout the process is a useful guideline to consider when
1238 adapting security and privacy practices to Big Data scenarios. The Organization for the Advancement of
1239 Structured Information Standards (OASIS) Privacy Management Reference Model (PMRM), consisting
1240 of seven foundational principles, provides appropriate basic guidance for Big System architects. [39,40]
1241 When working with any personal data, privacy should be an integral element in the design of a Big Data
1242 system.

1243 Other privacy engineering frameworks are also under consideration.[41 42 43 44 45 46]

1244 Related principles include identity management frameworks such as proposed in the National Strategy for
1245 Trusted Identities in Cyberspace (NSTIC)[47] and considered in the NIST Cloud Computing Security
1246 Reference Architecture.[48] Aspects of identity management that contribute to a security and privacy fabric
1247 will be addressed in future versions of this document.

1248 Big Data frameworks can also be used for strengthening security. Big Data analytics can be used for
1249 detecting privacy breaches through security intelligence, event detection, and forensics.

## 5.3 RELATION OF THE BIG DATA SECURITY OPERATIONAL TAXONOMY TO THE NBDRA

1252 Table 1 represents a preliminary mapping of the operational taxonomy to the NBDRA components. The
1253 topics and activities listed for each operational taxonomy element (Section 4.2) have been allocated to a
1254 NBDRA component under the Activities column in Table 1. The description column provides additional
1255 information about the security and privacy aspects of each NBDRA component.

1256 *Table 1: Draft Security Operational Taxonomy Mapping to the NBDRA Components*

| Activities | Description |
|---|---|
| **System Orchestrator** | |
| <ul><li>Policy Enforcement</li><li>Security Metadata Model</li><li>Data Loss Prevention, Detection</li><li>Data Lifecycle Management</li><li>Threat and Vulnerability Management</li><li>Mitigation</li><li>Configuration Management</li><li>Monitoring, Alerting</li><li>Malware Surveillance and Remediation</li></ul> | Several security functions have been mapped to the System Orchestrator block, as they require architectural level decisions and awareness. Aspects of these functionalities are strongly related to the Security Fabric and thus touch the entire architecture at various points in different forms of operational details. Such security functions include nation-specific compliance requirements, vastly expanded demand for forensics, and domain-specific, privacy-aware business |

| Activities | Description |
|---|---|
| • Resiliency, Redundancy and Recovery<br>• Accountability<br>• Compliance<br>• Forensics<br>• Business Risk Model | risk models. |
| **Data Provider** | |
| • Device, User, Asset, Services, Applications Registration<br>• Application Layer Identity<br>• End User Layer Identity Management<br>• End Point Input Validation<br>• Digital Rights Management<br>• Monitoring, Alerting | Data Providers are subject to guaranteeing authenticity of data and in turn require that sensitive, copyrighted, or valuable data be adequately protected. This leads to operational aspects of entity registration and identity ecosystems. |
| **Data Consumer** | |
| • Application Layer Identity<br>• End User Layer Identity Management<br>• Web Services Gateway<br>• Digital Rights Management<br>• Monitoring, Alerting | Data Consumers exhibit a duality with Data Providers in terms of obligations and requirements – only they face the access/visualization aspects of the Application Provider. |
| **Application Provider** | |
| • Application Layer Identity<br>• Web Services Gateway<br>• Data Transformation<br>• Digital Rights Management<br>• Monitoring, Alerting | Application Provider interfaces between the Data Provider and Data Consumer. It takes part in all the secure interface protocols with these blocks as well as maintains secure interaction with the Framework Provider. |
| **Framework Provider** | |
| • Virtualization Layer Identity<br>• Identity Provider<br>• Encryption and Key Management<br>• Isolation/Containerization<br>• Storage Security<br>• Network Boundary Control<br>• Monitoring, Alerting | Framework Provider is responsible for the security of data/computations for a significant portion of the lifecycle of the data. This includes security of data at rest through encryption and access control; security of computations via isolation/virtualization; and security of communication with the Application Provider. |

1257

1258

# 6 MAPPING USE CASES TO NBDRA

In this section, the security and privacy related use cases presented in Section 3 are mapped to the NBDRA components and interfaces explored in Figure 5, Notional Security and Privacy Fabric Overlay to the NBDRA.

## 6.1 CONSUMER DIGITAL MEDIA USE

Content owners license data for use by consumers through presentation portals. The use of consumer digital media generates Big Data, including both demographics at the user level and patterns of use such as play sequence, recommendations, and content navigation.

*Table 2: Mapping Consumer Digital Media Usage to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Varies and is vendor dependent. Spoofing is possible. For example, protections afforded by securing Microsoft Rights Management Services. [49] Secure/Multipurpose Internet Mail Extensions (S/MIME) |
| | Real-time security monitoring | Content creation security |
| | Data discovery and classification | Discovery/classification is possible across media, populations, and channels |
| | Secure data aggregation | Vendor-supplied aggregation services—security practices are opaque |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Aggregate reporting to content owners |
| | Compliance with regulations | PII disclosure issues abound |
| | Government access to data and freedom of expression concerns | Various issues; for example, playing terrorist podcast and illegal playback |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Unknown |
| | Policy management for access control | User, playback administrator, library maintenance, and auditor |
| | Computing on the encrypted data: searching/ filtering/ deduplicate/ fully homomorphic encryption | Unknown |
| | Audits | Audit DRM usage for royalties |
| Framework Provider | Securing data storage and transaction logs | Unknown |
| | Key management | Unknown |
| | Security best practices for non-relational data stores | Unknown |
| | Security against DoS attacks | N/A |
| | Data provenance | Traceability to data owners, producers, consumers is preserved |
| Fabric | Analytics for security intelligence | Machine intelligence for unsanctioned use/access |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | Event detection | "Playback" granularity defined |
| | Forensics | Subpoena of playback records in legal disputes |

## 6.2 NIELSEN HOMESCAN: PROJECT APOLLO

1268

1269 Nielsen Homescan involves family-level retail transactions and associated media exposure using a
1270 statistically valid national sample. A general description[50] is provided by the vendor. This project
1271 description is based on a 2006 Project Apollo architecture. (Project Apollo did not emerge from its
1272 prototype status.)

1273
*Table 3: Mapping Nielsen Homescan to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Device-specific keys from digital sources; receipt sources scanned internally and reconciled to family ID (Role issues) |
| | Real-time security monitoring | None |
| | Data discovery and classification | Classifications based on data sources (e.g., retail outlets, devices, and paper sources) |
| | Secure data aggregation | Aggregated into demographic crosstabs. Internal analysts had access to PII |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Aggregated to (sometimes) product-specific, statistically valid independent variables |
| | Compliance with regulations | Panel data rights secured in advance and enforced through organizational controls |
| | Government access to data and freedom of expression concerns | N/A |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Encryption not employed in place; only for data-center-to-data-center transfers. XML (Extensible Markup Language) cube security mapped to Sybase IQ and reporting tools |
| | Policy management for access control | Extensive role-based controls |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption | N/A |
| | Audits | Schematron and process step audits |
| Framework Provider | Securing data storage and transaction logs | Project-specific audits secured by infrastructure team |
| | Key management | Managed by project chief security officer (CSO). Separate key pairs issued for customers and internal users |
| | Security best practices for non-relational data stores | Regular data integrity checks via XML schema validation |
| | Security against DoS attacks | Industry-standard webhost protection provided for query subsystem |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | Data provenance | Unique |
| Fabric | Analytics for security intelligence | No project-specific initiatives |
| | Event detection | N/A |
| | Forensics | Usage, cube-creation, and device merge audit records were retained for forensics and billing |

## 6.3 WEB TRAFFIC ANALYTICS

Visit-level webserver logs are of high-granularity and voluminous. Web logs are correlated with other sources, including page content (buttons, text, and navigation events) and marketing events such as campaigns and media classification.

*Table 4: Mapping Web Traffic Analytics to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Device-dependent. Spoofing is often easy |
| | Real-time security monitoring | Web server monitoring |
| | Data discovery and classification | Some geospatial attribution |
| | Secure data aggregation | Aggregation to device, visitor, button, web event, and others |
| Application Provider → Data Consumer | Privacy-preserving data analytics | IP anonymizing and timestamp degrading. Content-specific opt-out |
| | Compliance with regulations | Anonymization may be required for EU compliance. Opt-out honoring |
| | Government access to data and freedom of expression concerns | Yes |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Varies depending on archivist |
| | Policy management for access control | System- and application-level access controls |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption | Unknown |
| | Audits | Customer audits for accuracy and integrity are supported |
| Framework Provider | Securing data storage and transaction logs | Storage archiving—this is a big issue |
| | Key management | CSO and applications |
| | Security best practices for non-relational data stores | Unknown |
| | Security against DoS attacks | Standard |
| | Data provenance | Server, application, IP-like identity, page point-in-time Document Object Model (DOM), and point-in-time marketing events |
| Fabric | Analytics for security intelligence | Access to web logs often requires privilege elevation |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | Event detection | Can infer; for example, numerous sales, marketing, and overall web health events |
| | Forensics | See the SIEM use case |

## 6.4 HEALTH INFORMATION EXCHANGE

Health information exchange (HIE) data is aggregated from various data providers, which might include covered entities such as hospitals and contract research organizations (CROs) identifying participation in clinical trials. The data consumers would include emergency room personnel, the CDC, and other authorized health (or other) organizations. Because any city or region might implement its own HIE, these exchanges might also serve as data consumers and data providers for each other.

*Table 5: Mapping HIE to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Strong authentication, perhaps through X.509v3 certificates, potential leverage of SAFE (Signatures & Authentication for Everything[51]) bridge in lieu of general PKI |
| | Real-time security monitoring | Validation of incoming records to assure integrity through signature validation and to assure HIPAA privacy through ensuring PHI is encrypted. May need to check for evidence of informed consent |
| | Data discovery and classification | Leverage Health Level Seven (HL7) and other standard formats opportunistically, but avoid attempts at schema normalization. Some columns will be strongly encrypted while others will be specially encrypted (or associated with cryptographic metadata) for enabling discovery and classification. May need to perform column filtering based on the policies of the data source or the HIE service provider |
| | Secure data aggregation | Clear text columns can be deduplicated, perhaps columns with deterministic encryption. Other columns may have cryptographic metadata for facilitating aggregation and deduplication. Retention rules are assumed, but disposition rules are not assumed in the related areas of compliance |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Searching on encrypted data and proofs of data possession. Identification of potential adverse experience due to clinical trial participation. Identification of potential professional patients. Trends and epidemics, and co-relations of these to environmental and other effects. Determination of whether the drug to be |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | Compliance with regulations | administered will generate an adverse reaction, without breaking the double blind. Patients will need to be provided with detailed accounting of accesses to, and uses of, their EHR data<br>HIPAA security and privacy will require detailed accounting of access to EHR data. Facilitating this, and the logging and alerts, will require federated identity integration with data consumers |
| | Government access to data and freedom of expression concerns | CDC, law enforcement, subpoenas and warrants. Access may be toggled based on occurrence of a pandemic (e.g., CDC) or receipt of a warrant (e.g., law enforcement) |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Row-level and column-level access control |
| | Policy management for access control | Role-based and claim-based. Defined for PHI cells |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption | Privacy-preserving access to relevant events, anomalies, and trends for CDC and other relevant health organizations |
| | Audits | Facilitate HIPAA readiness and HHS audits |
| Framework Provider | Securing data storage and transaction logs | Need to be protected for integrity and privacy, but also for establishing completeness, with an emphasis on availability |
| | Key management | Federated across covered entities, with the need to manage key life cycles across multiple covered entities that are data sources |
| | Security best practices for non-relational data stores | End-to-end encryption, with scenario-specific schemes that respect min-entropy to provide richer query operations without compromising patient privacy |
| | Security against distributed denial of Service (DDoS) attacks | A mandatory requirement: systems must survive DDoS attacks |
| | Data provenance | Completeness and integrity of data with records of all accesses and modifications. This information could be as sensitive as the data and is subject to commensurate access policies |
| Fabric | Analytics for security intelligence | Monitoring of informed patient consent, authorized and unauthorized transfers, and accesses and modifications |
| | Event detection | Transfer of record custody, addition/modification of record (or cell), authorized queries, unauthorized queries, and modification attempts |
| | Forensics | Tamper-resistant logs, with evidence of tampering events. Ability to identify record- |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | | level transfers of custody and cell-level access or modification |

## 6.5 GENETIC PRIVACY

1286

1287 Mapping of genetic privacy is under development and will be included in future versions of this
1288 document.

## 6.6 PHARMACEUTICAL CLINICAL TRIAL DATA SHARING

1289

1290 Under an industry trade group proposal, clinical trial data for new drugs will be shared outside intra-
1291 enterprise warehouses.

1292 *Table 6: Mapping Pharmaceutical Clinical Trial Data Sharing to the Reference Architecture*

| NBDRA Component and Interfaces | Security & Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation<br>Real-time security monitoring<br>Data discovery and classification<br>Secure data aggregation | Opaque—company-specific<br>None<br>Opaque—company-specific<br>Third-party aggregator |
| Application Provider → Data Consumer | Privacy-preserving data analytics<br><br>Compliance with regulations<br>Government access to data and freedom of expression concerns | Data to be reported in aggregate but preserving potentially small-cell demographics<br>Responsible developer and third-party custodian<br>Limited use in research community, but there are possible future public health data concerns. Clinical study reports only, but possibly selectively at the study- and patient-levels |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption<br>Policy management for access control<br><br>Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption<br>Audits | TBD<br><br>Internal roles; third-party custodian roles; researcher roles; participating patients' physicians<br>TBD<br><br><br>Release audit by a third party |
| Framework Provider | Securing data storage and transaction logs<br>Key management<br>Security best practices for non-relational data stores<br>Security against DoS attacks<br>Data provenance | TBD<br><br>Internal varies by firm; external TBD<br>TBD<br><br>Unlikely to become public<br>TBD—critical issue |
| Fabric | Analytics for security intelligence<br>Event detection<br>Forensics | TBD<br>TBD |

## 6.7 NETWORK PROTECTION

1294    SIEM is a family of tools used to defend and maintain networks.

1295    *Table 7: Mapping Network Protection to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Software-supplier specific; refer to commercially available end point validation[52] |
| | Real-time security monitoring | --- |
| | Data discovery and classification | Varies by tool, but classified based on security semantics and sources |
| | Secure data aggregation | Aggregates by subnet, workstation, and server |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Platform-specific |
| | Compliance with regulations | Applicable, but regulated events are not readily visible to analysts |
| | Government access to data and freedom of expression concerns | NSA and FBI have access on demand |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Usually a feature of the operating system |
| | Policy management for access control | For example, a group policy for an event log |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption | Vendor and platform-specific |
| | Audits | Complex—audits are possible throughout |
| Framework Provider | Securing data storage and transaction logs | Vendor and platform-specific |
| | Key management | Chief Security Officer and SIEM product keys TBD |
| | Security best practices for non-relational data stores | |
| | Security against DDoS attacks | Big Data application layer DDoS attacks can be mitigated using combinations of traffic analytics, correlation analysis |
| | Data provenance | For example, how to know an intrusion record was actually associated with a specific workstation |
| Fabric | Analytics for security intelligence | Feature of current SIEMs |
| | Event detection | Feature of current SIEMs |
| | Forensics | Feature of current SIEMs |

## 6.8 MILITARY: UNMANNED VEHICLE SENSOR DATA

1297    Unmanned vehicles (drones) and their onboard sensors (e.g., streamed video) can produce petabytes of
1298    data that should be stored in nonstandard formats. The U.S. government is pursuing capabilities to expand
1299    storage capabilities for Big Data such as streamed video. For more information, refer to the Defense
1300    Information Systems Agency (DISA) large data object contract for exabytes in the DOD private cloud.[53]

1301    *Table 8: Mapping Military Unmanned Vehicle Sensor Data to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation | Need to secure the sensor (e.g., camera) to prevent spoofing/stolen sensor streams. There are new transceivers and protocols in the DOD pipeline. Sensor streams will include smartphone and tablet sources |
| | Real-time security monitoring | Onboard and control station secondary sensor security monitoring |
| | Data discovery and classification | Varies from media-specific encoding to sophisticated situation-awareness enhancing fusion schemes |
| | Secure data aggregation | Fusion challenges range from simple to complex. Video streams may be used[54] unsecured or unaggregated |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Geospatial constraints: cannot surveil beyond Universal Transverse Mercator (UTM). Military secrecy: target and point of origin privacy |
| | Compliance with regulations | Numerous. There are also standards issues |
| | Government access to data and freedom of expression concerns | For example, the Google lawsuit over Street View |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption | Policy-based encryption, often dictated by legacy channel capacity/type |
| | Policy management for access control | Transformations tend to be made within DOD/contractor-devised system schemes |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption | Sometimes performed within vendor-supplied architectures, or by image-processing parallel architectures |
| | Audits | CSO and Inspector General (IG) audits |
| Framework Provider | Securing data storage and transaction logs | The usual, plus data center security levels are tightly managed (e.g., field vs. battalion vs. headquarters) |
| | Key management | CSO—chain of command |
| | Security best practices for non-relational data stores | Not handled differently at present; this is changing in DOD |
| | Security against DoS attacks | DOD anti-jamming e-measures |
| | Data provenance | Must track to sensor point in time configuration and metadata |
| Fabric | Analytics for security intelligence | DOD develops specific field of battle security software intelligence—event driven and monitoring—that is often remote |
| | Event detection | For example, target identification in a video stream infers height of target from shadow. Fuse data from satellite infrared with separate sensor stream |
| | Forensics | Used for after action review (AAR)—desirable to have full playback of sensor streams |

## 6.9  EDUCATION: COMMON CORE STUDENT PERFORMANCE REPORTING

Cradle-to-grave student performance metrics for every student are now possible—at least within the K-12 community, and probably beyond. This could include every test result ever administered.

*Table 9: Mapping Common Core K–12 Student Reporting to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation<br>Real-time security monitoring<br><br>Data discovery and classification<br>Secure data aggregation | Application-dependent. Spoofing is possible<br>Vendor-specific monitoring of tests, test-takers, administrators, and data<br>Unknown<br>Typical: Classroom-level |
| Application Provider → Data Consumer | Privacy-preserving data analytics<br><br>Compliance with regulations<br><br>Government access to data and freedom of expression concerns | Various: For example, teacher-level analytics across all same-grade classrooms<br>Parent, student, and taxpayer disclosure and privacy rules apply<br>Yes. May be required for grants, funding, performance metrics for teachers, administrators, and districts |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption<br>Policy management for access control<br><br>Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption<br>Audits | Support both individual access (student) and partitioned aggregate<br>Vendor (e.g., Pearson) controls, state-level policies, federal-level policies; probably 20-50 different roles are spelled out at present<br>Proposed [55]<br><br>Support both internal and third-party audits by unions, state agencies, responses to subpoenas |
| Framework Provider | Securing data storage and transaction logs<br>Key management<br><br>Security best practices for non-relational data stores<br>Security against DDoS attacks<br>Data provenance | Large enterprise security, transaction level controls—classroom to the federal government<br>CSOs from the classroom level to the national level<br>---<br><br>Standard<br>Traceability to measurement event requires capturing tests at a point in time, which may itself require a Big Data platform |
| Fabric | Analytics for security intelligence<br>Event detection<br>Forensics | Various commercial security applications<br>Various commercial security applications<br>Various commercial security applications |

## 6.10 SENSOR DATA STORAGE AND ANALYTICS

Mapping of sensor data storage and analytics is under development and will be included in future versions of this document.

1309 **6.11 CARGO SHIPPING**

1310 This use case provides an overview of a Big Data application related to the shipping industry for which
1311 standards may emerge in the near future.

1312 *Table 10: Mapping Cargo Shipping to the Reference Architecture*

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| Data Provider → Application Provider | End-point input validation Real-time security monitoring | Ensuring integrity of data collected from sensors Sensors can detect abnormal temperature/environmental conditions for packages with special requirements. They can also detect leaks/radiation |
| | Data discovery and classification Secure data aggregation | --- Securely aggregating data from sensors |
| Application Provider → Data Consumer | Privacy-preserving data analytics | Sensor-collected data can be private and can reveal information about the package and geo-information. The revealing of such information needs to preserve privacy |
| | Compliance with regulations Government access to data and freedom of expression concerns | --- The U.S. Department of Homeland Security may monitor suspicious packages moving into/out of the country |
| Data Provider ↔ Framework Provider | Data-centric security such as identity/policy-based encryption Policy management for access control | --- Private, sensitive sensor data and package data should only be available to authorized individuals. Third-party commercial offerings may implement low-level access to the data |
| | Computing on the encrypted data: searching/filtering/deduplicate/fully homomorphic encryption Audits | See above section on "Transformation" --- |
| Framework Provider | Securing data storage and transaction logs | Logging sensor data is essential for tracking packages. Sensor data at rest should be kept in secure data stores |
| | Key management Security best practices for non-relational data stores Security against DoS attacks Data provenance | For encrypted data The diversity of sensor types and data types may necessitate the use of non-relational data stores --- Metadata should be cryptographically attached to the collected data so that the integrity of origin and progress can be assured. Complete preservation of provenance will sometimes mandate a separate Big Data application |
| Fabric | Analytics for security intelligence | Anomalies in sensor data can indicate tampering/fraudulent insertion of data traffic |

| NBDRA Component and Interfaces | Security and Privacy Topic | Use Case Mapping |
|---|---|---|
| | Event detection | Abnormal events such as cargo moving out of the way or being stationary for unwarranted periods can be detected |
| | Forensics | Analysis of logged data can reveal details of incidents after they occur |

1313

1314

1315

# Appendix A: Candidate Security and Privacy Topics for Big Data Adaptation

The following set of topics was initially adapted from the scope of the CSA BDWG charter and organized according to the classification in CSA BDWG's *Top 10 Challenges in Big Data Security and Privacy*.[56] Security and privacy concerns are classified in four categories:

- Infrastructure Security
- Data Privacy
- Data Management
- Integrity and Reactive Security

NBD-PWG Security and Privacy Subgroup identified the Big Data topics below for possible inspection during the preparation of Version 2 of this document. A complete rework of these topics is beyond the scope of this document. This material may be refined and organized if needed in future versions of this document.

### Infrastructure Security

- Review of technologies and frameworks that have been primarily developed for performance, scalability, and availability, massively parallel processing (MPP) databases, and others.
- High-availability
  - o Use of Big Data to enhance defenses against DDoS attacks.
- DevOps Security

### Data Privacy

- System architects should consider the impact of the social data revolution on the security and privacy of Big Data implementations. Some systems not designed to include social data could be connected to social data systems by third parties, or by other project sponsors within an organization.
  - o Unknowns of innovation: When a perpetrator, abuser, or stalker misuses technology to target and harm a victim, there are various criminal and civil charges that might be applied to ensure accountability and promote victim safety. A number of U.S. federal and state, territory, or tribal laws might apply. To support the safety and privacy of victims, it is important to take technology-facilitated abuse and stalking seriously. This includes assessing all ways that technology is being misused to perpetrate harm, and considering all charges that could or should be applied.
  - o Identify laws that address violence and abuse
    - Stalking and cyberstalking (e.g., felony menacing by, via electronic surveillance)
    - Harassment, threats, and assault
    - Domestic violence, dating violence, sexual violence, and sexual exploitation
    - Sexting and child pornography: electronic transmission of harmful information to minors, providing obscene material to a minor, inappropriate images of minors, and lascivious intent
    - Bullying and cyberbullying
    - Child abuse
  - o Identify possible criminal or civil laws applicable related to Big Data technology, communications, privacy, and confidentiality

- Unauthorized access, unauthorized recording/taping, illegal interception of electronic communications, illegal monitoring of communications, surveillance, eavesdropping, wiretapping, and unlawful party to call
- Computer and internet crimes: fraud and network intrusion
- Identity theft, impersonation, and pretexting
- Financial fraud and telecommunications fraud
- Privacy violations
- Consumer protection laws
- Violation of no contact, protection, and restraining orders
- Technology misuse: Defamatory libel, slander, economic or reputational harms, and privacy torts
- Burglary, criminal trespass, reckless endangerment, disorderly conduct, mischief, and obstruction of justice
- Data-centric security may be needed to protect certain types of data no matter where it is stored or accessed (e.g., attribute-based encryption and format-preserving encryption). There are domain-specific particulars that should be considered when addressing encryption tools available to system users.
- Big data privacy and governance
  - Data discovery and classification
  - Policy management for accessing and controlling Big Data
    - Are new policy language frameworks specific to Big Data architectures needed?
  - Data masking technologies: Anonymization, rounding, truncation, hashing, and differential privacy
    - It is important to consider how these approaches degrade performance or hinder delivery all together—*for Big Data systems in particular*. Often these solutions are proposed and then cause an outage at the time of the release, forcing the removal of the option.
  - Compliance with regulations such as the Health Insurance Portability and Accountability Act (HIPAA), European Union (EU) data protection regulations, Asia-Pacific Economic Cooperation (APEC) Cross-Border Privacy Rules (CBPR) requirements, and country-specific regulations
    - Regional data stores enable regional laws to be enforced
      - Cybersecurity Executive Order 1998—assumed data and information would remain within the region
    - People-centered design makes the assumption that private-sector stakeholders are operating ethically and respecting the freedoms and liberties of all Americans.
      - Litigation, including class action suits, could follow increased threats to Big Data security, when compared to other systems
        - People before profit must be revisited to understand the large number of Executive Orders overlooked
        - People before profit must be revisited to understand the large number of domestic laws overlooked
      - Indigenous and aboriginal people and the privacy of all associated vectors and variables must be excluded from any Big Data store in any case in which a person must opt in
        - All tribal land is an exclusion from any image capture and video streaming or capture
        - Human rights
  - Government access to data and freedom of expression concerns

1406            ▪   Polls show that U.S. citizens are less concerned about the loss of privacy than
1407            Europeans are, but both are concerned about data misuse and their inability to
1408            govern private- and public-sector use

1409     o   Potentially unintended/unwanted consequences or uses
1410       ▪   Appropriate uses of data collected or data aggregation and problem management
1411       capabilities must be enabled
1412       ▪   Mechanisms for the appropriate secondary or subsequent data uses, such as filtered upon
1413       entry processed and presented in the inbound framework
1414     o   Issues surrounding permission to collect data, consent, and privacy
1415       ▪   Differences between where the privacy settings are applied in web services and the user's
1416       perception of the privacy setting application
1417       ▪   Permission based on clear language and not forced by preventing users to access their
1418       online services
1419       ▪   People do not believe the government would allow businesses to take advantage of their
1420       rights
1421     o   Data deletion: Responsibility to purge data based on certain criteria and/or events
1422       ▪   Examples include legal rulings that affect an external data source. For example, if
1423       Facebook were to lose a legal challenge and required to purge its databases of certain
1424       private information. Is there then a responsibility for downstream data stores to follow suit
1425       and purge their copies of the same data? The provider, producer, collector or social media
1426       supplier, or host absolutely must inform and remove all versions. Enforcement?
1427       Verification?
1428     o   Computing on encrypted data
1429       ▪   Deduplication of encrypted data
1430       ▪   Searching and reporting on the encrypted data
1431       ▪   Fully homomorphic encryption
1432       ▪   Anonymization of data (no linking fields to reverse identify)
1433       ▪   De-identification of data (individual centric)
1434       ▪   Non-identifying data (individual and context centric)
1435     o   Secure data aggregation
1436     o   Data loss prevention
1437     o   Fault tolerance—recovery for zero data loss
1438       ▪   Aggregation in end-to-end scale of resilience, record, and operational scope for integrity
1439       and privacy in a secure or better risk management strategy
1440       ▪   Fewer applications will require fault tolerance with clear distinction around risk and scope
1441       of the risk

## Data Management

1443   •   Securing data stores
1444     o   Communication protocols
1445       ▪   Database links
1446       ▪   Access control list (ACL)
1447       ▪   Application programming interface (API)
1448       ▪   Channel segmentation
1449     o   Attack surface reduction
1450   •   Key management and ownership of data
1451     o   Providing full control of the keys to the data owner
1452     o   Transparency of data life cycle process: Acquisition, uses, transfers, dissemination, and
1453       destruction

1454          o   Maps to aid non-technical people determine who is using their data and how their data is
1455              being used, including custody over time

### *Integrity and Reactive Security*

1456

1457    •   Big Data analytics for security intelligence (identifying malicious activity) and situational
1458        awareness (understanding the health of the system)
1459        o   Large-scale analytics
1460            ▪   Need assessment of the public sector
1461        o   Streaming data analytics
1462            ▪   This could require, for example, segregated virtual machines and secure channels
1463            ▪   This is a low-level requirement
1464            ▪   Roadmap
1465            ▪   Priority of security and return on investment must be done to move to this degree of
1466                maturity
1467    •   Event detection
1468        o   Respond to data risk events trigger by application-specific analysis of user and system
1469            behavior patterns
1470        o   Data-driven abuse detection
1471    •   Forensics
1472    •   Security of analytics results
1473

1474

# 1475 Appendix B: Internal Security Considerations within
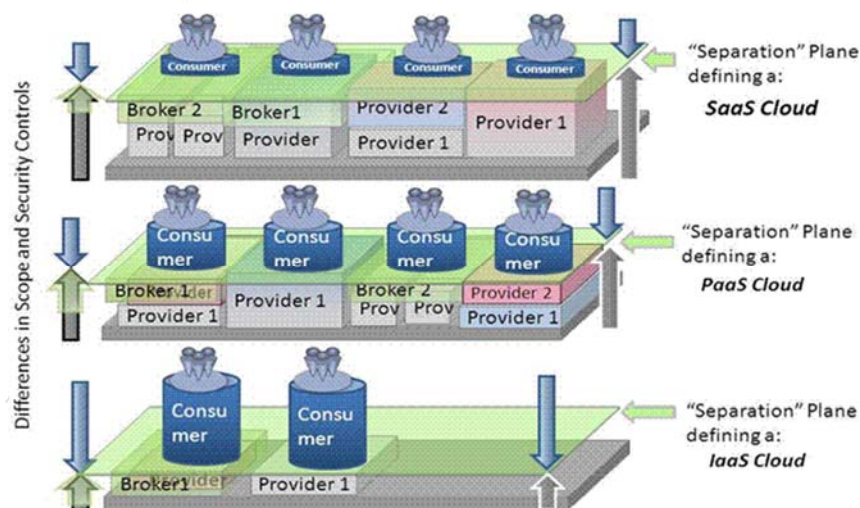# 1476 Cloud Ecosystems

1477 Many Big Data systems will be designed using cloud architectures. Any strategy to implement a mature
1478 security and privacy framework within a Big Data cloud ecosystem enterprise architecture must address
1479 the complexities associated with cloud-specific security requirements triggered by the cloud
1480 characteristics. These requirements could include the following:

1481 • Broad network access
1482 • Decreased visibility and control by consumer
1483 • Dynamic system boundaries and comingled roles/responsibilities between consumers and
1484 providers
1485 • Multi-tenancy
1486 • Data residency
1487 • Measured service
1488 • Order-of-magnitude increases in scale (on demand), dynamics (elasticity and cost optimization),
1489 and complexity (automation and virtualization)

1490 These cloud computing characteristics often present different security risks to an agency than the
1491 traditional information technology solutions, thereby altering the agency's security posture.

1492 To preserve the security-level after the migration of their data to the cloud, organizations need to identify
1493 all cloud-specific, risk-adjusted security controls or components in advance. The organizations must also
1494 request from the cloud service providers, through contractual means and service-level agreements, to have
1495 all identified security components and controls fully and accurately implemented.

1496 The complexity of multiple interdependencies is best illustrated by Figure B-1.



1497
1498 *Figure B-1: Composite Cloud Ecosystem Security Architecture[57]*

1499 When unraveling the complexity of multiple interdependencies, it is important to note that enterprise-
1500 wide access controls fall within the purview of a well thought out Big Data and cloud ecosystem risk
1501 management strategy for end-to-end enterprise access control and security (AC&S), via the following five
1502 constructs:

1503 1. Categorize the data value and criticality of information systems and the data custodian's duties and
1504     responsibilities to the organization, demonstrated by the data custodian's choice of either a
1505     discretionary access control policy or a mandatory access control policy that is more restrictive. The
1506     choice is determined by addressing the specific organizational requirements, such as, but not limited
1507     to the following:
1508     a. GRC
1509     b. Directives, policy guidelines, strategic goals and objectives, information security requirements,
1510        priorities, and resources available (filling in any gaps)
1511 2. Select the appropriate level of security controls required to protect data and to defend information
1512     systems
1513 3. Implement access security controls and modify them upon analysis assessments
1514 4. Authorize appropriate information systems
1515 5. Monitor access security controls at a minimum of once a year

1516 To meet GRC and confidentiality, integrity, and availability regulatory obligations required from the
1517 responsible data custodians—which are directly tied to demonstrating a valid, current, and up-to-date
1518 AC&S policy—one of the better strategies is to implement a layered approach to AC&S, comprised of
1519 multiple access control gates, including, but not limited to, the following infrastructure AC&S via:

1520     • Physical security/facility security, equipment location, power redundancy, barriers, security
1521       patrols, electronic surveillance, and physical authentication
1522     • Information Security and residual risk management
1523     • Human resources (HR security, including, but not limited to, employee codes of conduct, roles
1524       and responsibilities, job descriptions, and employee terminations
1525     • Database, end point, and cloud monitoring
1526     • Authentication services management/monitoring
1527     • Privilege usage management/monitoring
1528     • Identify management/monitoring
1529     • Security management/monitoring
1530     • Asset management/monitoring

1531 The following section revisits the traditional access control framework. The traditional framework
1532 identifies a standard set of attack surfaces, roles, and tradeoffs. These principles appear in some existing
1533 best practices guidelines. For instance, they are an important part of the Certified Information Systems
1534 Security Professional (CISSP) body of knowledge. [f] This framework for Big Data may be adopted during
1535 the future work of the NBD-PWG.

1536 ### *Access Control*

1537 Access control is one of the most important areas of Big Data. There are multiple factors, such as
1538 mandates, policies, and laws that govern the access of data. One overarching rule is that the highest
1539 classification of any data element or string governs the protection of the data. In addition, access should
1540 only be granted on a need-to-know/-use basis that is reviewed periodically in order to control the access.

1541 Access control for Big Data covers more than accessing data. Data can be accessed via multiple channels,
1542 networks, and platforms—including laptops, cell phones, smart phones, tablets, and even fax machines—
1543 that are connected to internal networks, mobile devices, the internet, or all of the above. With this reality
1544 in mind, the same data may be accessed by a user, administrator, another system, etc., and it may be
1545 accessed via a remote connection/access point as well as internally. Therefore, visibility as to who is

---

[f] CISSP is a professional computer security certification administered by (ISC)[2].
(https://www.isc2.org/cissp/default.aspx)

1546 accessing the data is critical in protecting the data. The trade-offs between strict data access control versus
1547 conducting business requires answers to questions such as the following.

1548 • How important/critical is the data to the lifeblood and sustainability of the organization?
1549 • What is the organization responsible for (e.g., all nodes, components, boxes, and machines within
1550 the Big Data/cloud ecosystem)?
1551 • Where are the resources and data located?
1552 • Who should have access to the resources and data?
1553 • Have GRC considerations been given due attention?

1554 Very restrictive measures to control accounts are difficult to implement, so this strategy can be considered
1555 impractical in most cases. However, there are best practices, such as protection based on classification of
1556 the data, least privilege[58], and separation of duties that can help reduce the risks.

1557 The following measures are often included in Best Practices lists for security and privacy. Some, and
1558 perhaps all, of the measure require adaptation or expansion for Big Data systems.

1559 • Least privilege—access to data within a Big Data/cloud ecosystem environment should be based
1560 on providing an individual with the minimum access rights and privileges to perform his/her job
1561 • If one of the data elements is protected because of its classification (e.g., PII, HIPAA, payment
1562 card industry [PCI]), then all of the data that it is sent with it inherits that classification, retaining
1563 the original data's security classification. If the data is joined to and/or associated with other data
1564 that may cause a privacy issue, then all data should be protected. This requires due diligence on
1565 the part of the data custodian(s) to ensure that this secure and protected state remains throughout
1566 the entire end-to-end data flow. Variations on this theme may be required for domain-specific
1567 combinations of public and private data hosted by Big Data applications.
1568 • If data is accessed from, transferred to, or transmitted to the cloud, internet, or another external
1569 entity, then the data should be protected based on its classification.
1570 • There should be an indicator/disclaimer on the display of the user if private or sensitive data is
1571 being accessed or viewed. Openness, trust, and transparency considerations may require more
1572 specific actions, depending on GRC or other broad considerations of how the Big Data system is
1573 being used
1574 • All system roles ("accounts") should be subjected to periodic meaningful audits to check that they
1575 are still required
1576 • All accounts (except for system-related accounts) that have not been used within 180 days should
1577 be deactivated
1578 • Access to PII data should be logged. Role-based access to Big Data should be enforced. Each role
1579 should be assigned the fewest privileges needed to perform the functions of that role
1580 • Roles should be reviewed periodically to check that they are still valid and that the accounts
1581 assigned to them are still appropriate

## User Access Controls

1583 • Each user should have his or her personal account. Shared accounts should not be the default
1584 practice in most settings
1585 • A user role should match the system capabilities for which it was intended. For example, a user
1586 account intended only for information access or to manage an Orchestrator should not be used as
1587 an administrative account or to run unrelated production jobs

## System Access Controls

1589 • There should not be shared accounts in cases of system-to-system access. "Meta-accounts" that
1590 operate across systems may be an emerging Big Data concern

1591
1592
1593

- Access for a system that contains Big Data needs to be approved by the data owner or his/her representative. The representative should not be infrastructure support personnel (e.g., a system administrator), because that may cause a separation of duties issue.

1594
1595
1596
1597

- Ideally, the same type of data stored on different systems should use the same classifications and rules for access controls to provide the same level of protection. In practice, Big Data systems may not follow this practice, and different techniques may be needed to map roles across related but dissimilar components or even across Big Data systems

1598

## *Administrative Account Controls*

1599
1600

- System administrators should maintain a separate user account that is not used for administrative purposes. In addition, an administrative account should not be used as a user account

1601
1602

- The same administrative account should not be used for access to the production and non-production (e.g., test, development, and quality assurance) systems

1603

# Appendix C: Big Data Actors and Roles: Adaptation to Big Data Scenarios

Service-oriented architectures (SOA) were a widely discussed paradigm through the early 2000's. While the concept is employed less often, SOA has influenced systems analysis processes, and perhaps to a lesser extent, systems design. As noted by Patig and Lopez-Sanz et al., actors and roles were incorporated into Unified Modeling Language so that these concepts could be represented within and well as across services. [59] [60] Big Data calls for further adaptation of these concepts. While actor/role concepts have not been fully integrated into the proposed security fabric, the Subgroup felt it important to emphasize to Big Data system designers how these concepts may need to be adapted from legacy and SOA usage.

Similar adaptations from Business Process Execution Language, Business Process Model and Notation frameworks offer additional patterns for Big Data security and privacy fabric standards. Ardagna et al. [61] suggest how adaptations might proceed from SOA, but Big Data systems offer somewhat different challenges.

Big Data systems can comprise simple machine-to-machine actors, or complex combinations of persons and machines that are systems of systems.

A common meaning of actor assigns roles to a person in a system. From a citizen's perspective, a person can have relationships with many applications and sources of information in a Big Data system.

The following list describes a number of roles as well as how roles can shift over time. For some systems, roles are only valid for a specified point in time. Reconsidering temporal aspects of actor security is salient for Big Data systems, as some will be architected without explicit archive or deletion policies.

- A retail organization refers to a person as a consumer or prospect before a purchase; afterwards, the consumer becomes a customer
- A person has a customer relationship with a financial organization for banking services
- A person may have a car loan with a different organization or the same financial institution
- A person may have a home loan with a different bank or the same bank
- A person may be "the insured" on health, life, auto, homeowners, or renters insurance
- A person may be the beneficiary or future insured person by a payroll deduction in the private sector, or via the employment development department in the public sector
- A person may have attended one or more public or private schools
- A person may be an employee, temporary worker, contractor, or third-party employee for one or more private or public enterprises
- A person may be underage and have special legal or other protections
- One or more of these roles may apply concurrently

For each of these roles, system owners should ask themselves whether users could achieve the following:

- Identify which systems their PII has entered
- Identify how, when, and what type of de-identification process was applied
- Verify integrity of their own data and correct errors, omissions, and inaccuracies
- Request to have information purged and have an automated mechanism to report and verify removal
- Participate in multilevel opt-out systems, such as will occur when Big Data systems are federated
- Verify that data has not crossed regulatory (e.g., age-related), governmental (e.g., a state or nation), or expired ("I am no longer a customer") boundaries

## OPT-IN REVISITED

1647   While standards organizations grapple with frameworks such as the one developed here, and until an
1648   individual's privacy and security can be fully protected using such a framework, some observers believe
1649   that the following two simple "protocols" ought to govern PII Big Data collection in the meantime.

1650   **Suggested Protocol one**: An individual can only decide to opt-in for inclusion of their personal data
1651   manually, and it is a decision that they can revoke at any time.

1652   **Suggested Protocol number two:** The individual's privacy and security opt-in process should enable
1653   each individual to modify their choice at any time, to access and review log files and reports and establish
1654   a self-destruct timeline (similar to the EU's "right to be forgotten".)

1655

# 1656 **Appendix D: Acronyms**

| 1657 | AC&S | access control and security |
|------|------|------|
| 1658 | ACLs | Access Control Lists |
| 1659 | AuthN/AuthZ | Authentication/Authorization |
| 1660 | BAA | business associate agreement |
| 1661 | CDC | U.S. Centers for Disease Control and Prevention |
| 1662 | CEP | complex event processing |
| 1663 | CIA | U.S. Central Intelligence Agency |
| 1664 | CIICF | Critical Infrastructure Cybersecurity Framework |
| 1665 | CINDER | DARPA Cyber-Insider Threat |
| 1666 | CMS | U.S. Centers for Medicare & Medicaid Services |
| 1667 | CoP | communities of practice |
| 1668 | CSA | Cloud Security Alliance |
| 1669 | CSA BDWG | Cloud Security Alliance Big Data Working Group |
| 1670 | CSP | Cloud Service Provider |
| 1671 | DARPA | Defense Advanced Research Projects Agency's |
| 1672 | DDoS | distributed denial of Service |
| 1673 | DOD | U.S. Department of Defense |
| 1674 | DoS | denial of service |
| 1675 | DRM | digital rights management |
| 1676 | EFPIA | European Federation of Pharmaceutical Industries and Associations |
| 1677 | EHRs | electronic health records |
| 1678 | EU | European Union |
| 1679 | FBI | U.S. Federal Bureau of Investigation |
| 1680 | FTC | Federal Trade Commission |
| 1681 | GPS | global positioning system |
| 1682 | GRC | governance, risk management, and compliance |
| 1683 | HIEs | Health Information Exchanges |
| 1684 | HIPAA | Health Insurance Portability and Accountability Act |
| 1685 | HITECH Act | Health Information Technology for Economic and Clinical Health Act |
| 1686 | HR | human resources |
| 1687 | IdP | Identity Provider |
| 1688 | IoT | internet of things |
| 1689 | IP | Internet Protocol |

| 1690 | IT | information technology |
|------|------|------|
| 1691 | LHNCBC | Lister Hill National Center for Biomedical Communications |
| 1692 | M2M | machine to machine |
| 1693 | MAC | media access control |
| 1694 | NBD-PWG | NIST Big Data Public Working Group |
| 1695 | NBDRA | NIST Big Data Reference Architecture |
| 1696 | NBDRA-SP | NIST Big Data Security and Privacy Reference Architecture |
| 1697 | NIEM | National Information Exchange Model |
| 1698 | NIST | National Institute of Standards and Technology |
| 1699 | NSA | U.S. National Security Agency |
| 1700 | OSS | operations systems support |
| 1701 | PaaS | platform as a service |
| 1702 | PHI | protected health information |
| 1703 | PII | personally identifiable information |
| 1704 | PKI | public key infrastructure |
| 1705 | SAML | Security Assertion Markup Language |
| 1706 | SIEM | Security Information and Event Management |
| 1707 | SKUs | stock keeping units |
| 1708 | SLAs | Service Level Agreements |
| 1709 | STS | Security Token Service |
| 1710 | TLS | Transport Layer Security |
| 1711 | VM | virtual machine |
| 1712 | VPN | virtual private network |
| 1713 | WS | web services |
| 1714 | XACML | eXtensible Access Control Markup Language |
| 1715 | | |

1716 # Appendix E: References

1717 ## GENERAL RESOURCES

1718 Luciano, Floridi (ed.), *The Cambridge Handbook of Information and Computer Ethics* (New York, NY:
1719 Cambridge University Press, 2010).

1720 Julie Lane, Victoria Stodden, Stefen Bender, and Helen Nissenbaum (eds.), *Privacy, Big Data and the*
1721 *Public Good: Frameworks for Engagement* (New York, NY: Cambridge University Press, 2014).

1722 Martha Nussbaum, *Creating Capabilities: The Human Development Approach* (Cambridge, MA:
1723 Belknap Press, 2011).

1724 John Rawls, *A Theory of Justice* (Cambridge, MA: Belknap Press, 1971).

1725 Martin Rost and Kirsten Bock, "Privacy by Design and the New Protection Goals," English translation of
1726 Privacy By Design und die Neuen Schutzziele, *Datenschutz und Datensicherheit*, Volume 35, Issue 1
1727 (2011), pages 30-35.

1728

1729 ## DOCUMENT REFERENCES

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

[2] EMC², "Digital Universe," *EMC,* accessed February 21, 2014, http://www.emc.com/leadership/programs/digital-universe.htm.

[3] EMC², "Digital Universe," *EMC,* accessed February 21, 2014, http://www.emc.com/leadership/programs/digital-universe.htm.

[4] Big Data Working Group, "Expanded Top Ten Big Data Security and Privacy Challenges*," Cloud Security Alliance,* April 2013, https://downloads.cloudsecurityalliance.org/initiatives/bdwg/Expanded_Top_Ten_Big_Data_Security_and_Privacy_Challenges.pdf.

[5] Subgroup correspondence with James G Kobielus (IBM), August 28, 2014.

[6] Big Data Working Group, "Top 10 Challenges in Big Data Security and Privacy," *Cloud Security Alliance*, November 2012, http://www.isaca.org/Groups/Professional-English/big-data/GroupDocuments/Big_Data_Top_Ten_v1.pdf.

[7] Benjamin Fung, Ke Wang, Rui Chen, and Philip S. Yu. "Privacy-preserving data publishing: A survey of recent developments", ACM Computing Surveys (CSUR), 42(4):14, 2010.

[8] Cynthia Dwork. "Differential privacy", In Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener, editors, ICALP 2006: 33rd International Colloquium on Automata, Languages and Programming, Part II, volume 4052 of Lecture Notes in Computer Science, pages 1-12, Venice, Italy, July 10-14, 2006. Springer, Berlin, Germany.

[9] Latanya Sweeney. "k-anonymity: A model for protecting privacy", International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, 10(05):557-570, 2002.

[10] Arvind Narayanan and Vitaly Shmatikov. "Robust de-anonymization of large sparse datasets", In 2008 IEEE Symposium on Security and Privacy, pages 111-125, Oakland, California, USA, May 18-21, 2008. IEEE Computer Society Press.

[11] Big Data Working Group, "Top 10 Challenges in Big Data Security and Privacy," *Cloud Security Alliance*, November 2012, http://www.isaca.org/Groups/Professional-English/big-data/GroupDocuments/Big_Data_Top_Ten_v1.pdf.

[12] Big Data Working Group, "Top 10 Challenges in Big Data Security and Privacy," *Cloud Security Alliance*, November 2012, http://www.isaca.org/Groups/Professional-English/big-data/GroupDocuments/Big_Data_Top_Ten_v1.pdf.

[13] S. S. Sahoo, A. Sheth, and C. Henson, "Semantic provenance for eScience: Managing the deluge of scientific data," *Internet Computing, IEEE,* Volume 12, Issue 4 (2008), pages 46–54, http://dx.doi.org/10.1109/MIC.2008.86.

[14] Ronan Shields, "AppNexus CTO on the fight against ad fraud," *Exchange Wire*, October 29, 2014, https://www.exchangewire.com/blog/2014/10/29/appnexus-cto-on-the-fight-against-ad-fraud/.

[15] David Lazer, Ryan Kennedy, Gary King, Alessandro Vespignani, "The parable of google flu: Traps in big data analysis" *Science* Volume 343, Issue 6176 (2014), pages 1203-1205, http://dx.doi.org/10.1126/science.1248506.

[16] Peng Chen, Beth Plale, and Mehmet Aktas, "Temporal representation for mining scientific data provenance," *Future Generation Computer Systems*, Volume 36, Special Issue (2014), pages 363-378, http://dx.doi.org/10.1016/j.future.2013.09.032.

[17] Xiao Zhang, edited by Raj Jain, "A survey of digital rights management technologies," *Washington University in Saint Louis,* accessed January 9, 2015, http://bit.ly/1y3Y1P1.

[18] PhRMA, "Principles for Responsible Clinical Trial Data Sharing," *European Federation of Pharmaceutical Industries and Associations,* July 18, 2013, http://phrma.org/sites/default/files/pdf/PhRMAPrinciplesForResponsibleClinicalTrialDataSharing.pdf.

[19] U.S. Army, "Army Regulation 25-2," *U.S. Army Publishing Directorate*, October 27, 2007, www.apd.army.mil/jw2/xmldemo/r25_2/main.asp.

[20] Jon Campbell, "Cuomo panel: State should cut ties with inBloom," *Albany Bureau,* March 11, 2014, http://lohud.us/1mV9U2U.

[21] Lisa Fleisher, "Before Tougher State Tests, Officials Prepare Parents," *Wall Street Journal*, April 15, 2013, http://blogs.wsj.com/metropolis/2013/04/15/before-tougher-state-tests-officials-prepare-parents/.

[22] Debra Donston-Miller, "Common Core Meets Aging Education Technology," *InformationWeek*, July 22, 2013, www.informationweek.com/big-data/news/common-core-meets-aging-education-techno/240158684.

[23] Civitas Learning, "About," *Civitas Learning*, www.civitaslearning.com/about/.

[24] Big Data Working Group, "Top 10 Challenges in Big Data Security and Privacy," *Cloud Security Alliance*, November 2012, http://www.isaca.org/Groups/Professional-English/big-data/GroupDocuments/Big_Data_Top_Ten_v1.pdf.

[25] R. Chandramouli, M. Iorga, and S. Chokhani, "Cryptographic key management issues & challenges in cloud services," *National Institute of Standards and Technology*, September 2013, http://dx.doi.org/10.6028/NIST.IR.7956.

[26] Peter Mell and Timothy Grance, "The NIST Definition of Cloud Computing: Recommendations of the National Institute of Standards and Technology," *National Institute of Standards and Technology,* September 2011, http://csrc.nist.gov/publications/nistpubs/800-145/SP800-145.pdf.

[27] ACM, Inc., "The ACM Computing Classification System" *Association for Computing Machinery, Inc.,* 1998, http://www.acm.org/about/class/ccs98-html#K.4.

[28] Computer Security Division, Information Technology Laboratory, "Guide for Applying the Risk Management Framework to Federal Information Systems: A Security Life Cycle Approach," *National Institute for Standards and Technology,* February 2010, http://csrc.nist.gov/publications/nistpubs/800-37-rev1/sp800-37-rev1-final.pdf.

[29] ISACA, "The Risk IT Framework," *www.isaca.org,* 2009, http://www.isaca.org/Knowledge-Center/Research/ResearchDeliverables/Pages/The-Risk-IT-Framework.aspx.

[30] Cybersecurity Framework, "Framework for Improving Critical Infrastructure Cybersecurity" *National Institute for Standards and Technology,* accessed January 9, 2015, http://1.usa.gov/1wQuti1.

[31] OASIS "SAML V2.0 Standard," *SAML Wiki*, accessed January 9, 2015, http://bit.ly/1wQByit.

[32] James Cebula and Lisa Young, "A taxonomy of operational cyber security risks," *Carnegie Melon University*, December 2010, http://resources.sei.cmu.edu/asset_files/TechnicalNote/2010_004_001_15200.pdf.

[33] OASIS "SAML V2.0 Standard," *SAML Wiki*, accessed January 9, 2015, http://bit.ly/1wQByit.

[34] H. C. Kum and S. Ahalt, "Privacy-by-Design: Understanding Data Access Models for Secondary Data," *AMIA Summits on Translational Science Proceedings*, *2013*, pages 126–130.

[35] John Rawls, "Justice as Fairness," *A Theory of Justice*, 1985.

[36] ETSI, "Smart Cards; Secure channel between a UICC and an end-point terminal," *etsy.org,* December 2007, http://bit.ly/1x2HSUe.

[37] James Cebula and Lisa Young, *"*Taxonomy of Operational Cyber Security Risks," (Pittsburgh, PA: Carnegie Mellon University, Software Engineering Institute, 2010).

[38] HHS Press Office, "New rule protects patient privacy, secures health information," *U.S. Department of Health and Human Services,* January 17, 2013, http://www.hhs.gov/news/press/2013pres/01/20130117b.html.

[39] John Sabo, Michael Willet, Peter Brown, and Dawn Jutla, "Privacy Management Reference Model and Methodology (PMRM) Version 1.0," *OASIS,* March 26, 2012, http://docs.oasis-open.org/pmrm/PMRM/v1.0/csd01/PMRM-v1.0-csd01.pdf.

[40] NIST, "National Strategy for Trusted Identities in Cyberspace (NSTIC)," *National Institute for Standards and Technology,* 2015, http://www.nist.gov/nstic/.

[41] Wayne Jansen and Timothy Grance, SP800-144, "Guidelines on Security and Privacy in Public Cloud Computing," *National Institute for Standards and Technology,* December 2011, http://csrc.nist.gov/publications/nistpubs/800-144/SP800-144.pdf.

[42] Wayne Jansen and Timothy Grance, SP 800-144, "Guidelines on Security and Privacy in Public Cloud Computing," *National Institute for Standards and Technology,* December 2011, http://csrc.nist.gov/publications/nistpubs/800-144/SP800-144.pdf.

[43] Carolyn Brodie, Clare-Marie Karat, John Karat, and Jinjuan Feng, "Usable security and privacy: A case study of developing privacy management tools," *Proceedings of the 2005 Symposium on Usable Privacy and Security,* 2005, http://doi.acm.org/10.1145/1073001.1073005.

[44] W. Knox Carey, Jarl Nilsson, and Steve Mitchell, "Persistent security, privacy, and governance for healthcare information," *Proceedings of the 2nd USENIX Conference on Health Security and Privacy,* 2011, http://dl.acm.org/citation.cfm?id=2028026.2028029.

[45] Paul Dunphy, John Vines, Lizzie Coles-Kemp, Rachel Clarke, Vasilis Vlachokyriakos, Peter Wright, John McCarthy, and Patrick Olivier, "Understanding the Experience-Centeredness of privacy and security technologies," *Proceedings of the 2014 Workshop on New Security Paradigms Workshop,* 2014, http://doi.acm.org/10.1145/2683467.2683475.

[46] Ebenezer Oladimeji, Lawrence Chung, Hyo Taeg Jung, and Jaehyoun Kim, "Managing security and privacy in ubiquitous eHealth information interchange," *Proceedings of the 5th International Conference on Ubiquitous Information Management and Communication,"* 2011, http://doi.acm.org/10.1145/1968613.1968645.

[47] NIST, "National Strategy for Trusted Identities in Cyberspace (NSTIC)," *National Institute for Standards and Technology,* 2015, http://www.nist.gov/nstic/.

[48] NIST Cloud Computing Security Working Group, "NIST Cloud Computing Security Reference Architecture," *National Institute for Standards and Technology,* May 15, 2013, http://collaborate.nist.gov/twiki-cloud-computing/pub/CloudComputing/CloudSecurity/NIST_Security_Reference_Architecture_2013.05.15_v1.0.pdf.

[49] Microsoft, "Deploying Windows Rights Management Services at Microsoft," *Microsoft*, 2015, http://technet.microsoft.com/en-us/library/dd277323.aspx.

[50] The Nielsen Company, "Consumer Panel and Retail Measurement," *Nielsen*, 2015, www.nielsen.com/us/en/nielsen-solutions/nielsen-measurement/nielsen-retail-measurement.html.

[51] SAFE-BioPharma, "Welcome to SAFE-BioPharma," *SAFE-BioPharma Association,* accessed March 3, 2015, http://www.safe-biopharma.org/.

[52] Microsoft, "How to set event log security locally or by using Group Policy in Windows Server 2003," *Microsoft*, http://support.microsoft.com/kb/323076.

[53] Kathleen Hickey, "DISA plans for exabytes of drone, satellite data," *GCN*, April 12, 2013, http://gcn.com/articles/2013/04/12/disa-plans-exabytes-large-data-objects.aspx.

[54] DefenseSystems, "UAV video encryption remains unfinished job," *DefenseSystems*, October 31, 2012, http://defensesystems.com/articles/2012/10/31/agg-drone-video-encryption-lags.aspx.

[55] K. A. G. Fisher, A. Broadbent, L. K. Shalm, Z. Yan, J. Lavoie, R. Prevedel, T. Jennewein, and K. J. Resch, "Quantum computing on encrypted data 5," *Nature Communications,* January 2015, http://www.nature.com/ncomms/2014/140121/ncomms4074/full/ncomms4074.html.

[56] Big Data Working Group, "Top 10 Challenges in Big Data Security and Privacy," Cloud Security Alliance, November 2012, http://www.isaca.org/Groups/Professional-English/big-data/GroupDocuments/Big_Data_Top_Ten_v1.pdf.

[57] Fang Liu, Jin Tong, Jian Mao, Robert Bohn, John Messina, Lee Badger, and Dawn Leaf, SP500-292, "NIST Cloud Computing Reference Architecture," *National Institute of Standards and Technology,* September 2011, http://www.nist.gov/customcf/get_pdf.cfm?pub_id=909505.

[58] John Mutch and Brian Anderson, "Preventing Good People From Doing Bad Things: Implementing Least Privilege," (Berkeley, CA: Apress, 2011).

[59] S. Patig, "Model-Driven development of composite applications," *Communications in Computer and Information Science*, 2008, http://dx.doi.org/10.1007/978-3-540-78999-4_8.

[60] M. López-Sanz, C. J. Acuña, C. E. Cuesta, and E. Marcos, "Modelling of Service-Oriented Architectures with UML," *Theoretical Computer Science*, Volume 194, Issue 4 (2008), pages 23–37.

[61] D. Ardagna, L. Baresi, S. Comai, M. Comuzzi, and B. Pernici, "A Service-Based framework for flexible business processes," *IEEE*, March 2011, http://dx.doi.org/10.1109/ms.2011.28.