

## Table of Contents

### 1. Story

1. [Data Science for Cyber Physical Systems-Internet of Things](#)
2. [Training Students to Extract Value from Big Data](#)
  1. [Slide 1 An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library](#)
  2. [Slide 2 Outline](#)
  3. [Slide 3 Semantic Community: Spotfire Cloud Library MindTouch](#)
  4. [Slide 4 Semantic Community: Spotfire Cloud Library Recent Analyses](#)
  5. [Slide 5 Semantic Community: Spotfire Cloud Library Browse Library](#)
  6. [Slide 6 Semantic Community: Spotfire Cloud Library Shared Folders](#)
  7. [Slide 7 Semantic Community: Spotfire Web Player](#)
  8. [Slide 8 Semantic Community: Spotfire Web Player Edit](#)
  9. [Slide 9 Semantic Community: Spotfire Web Player Menu](#)
  10. [Slide 10 Semantic Community: Spotfire Web Player Menu and New Visualizations](#)
  11. [Slide 11 TIBCO Spotfire: Cloud User's Guide](#)
  12. [Slide 12 TIBCO: Spotfire Cloud User's Guide Getting Started](#)
  13. [Slide 13 TIBCO: Spotfire Cloud User's Guide Setting Up Analyses](#)
  14. [Slide 14 TIBCO: Spotfire Cloud User's Guide Data Preparation in Microsoft Excel](#)
  15. [Slide 15 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 1](#)
  16. [Slide 16 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 2](#)
  17. [Slide 17 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 3](#)
  18. [Slide 18 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Markers](#)
  19. [Slide 19 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Colored Regions](#)
  20. [Slide 20 TIBCO: Spotfire Cloud User's Guide Exporting Your Analysis](#)
  21. [Slide 21 TIBCO: Spotfire Cloud User's Guide Sharing Your Analysis](#)
  22. [Slide 22 TIBCO Spotfire: Education](#)
3. [Summit Highlights](#)
  1. [Slide 1 Opening Remarks 1](#)
  2. [Slide 2 Opening Remarks 2](#)
  3. [Slide 3 Dr. Shoumen Data 1](#)
  4. [Slide 4 Dr. Shoumen Data 2](#)
  5. [Slide 5 Dr. Shoumen Data 3](#)
  6. [Slide 6 Dr. Shoumen Data 4](#)
  7. [Slide 7 Track A 1](#)
  8. [Slide 8 Track A 2](#)
  9. [Slide 9 Track B 1](#)
  10. [Slide 10 Track B 2](#)
  11. [Slide 11 Keith Mazullo 1](#)
  12. [Slide 12 Keith Mazullo 2](#)
  13. [Slide 13 Keith Mazullo 3](#)
  14. [Slide 14 Keith Mazullo 4](#)
  15. [Slide 15 Track C 1](#)
  16. [Slide 16 Track C 2](#)
  17. [Slide 17 Track D 1](#)
  18. [Slide 18 Track D 2](#)
  19. [Slide 19 Dr. Bradford Hess 1](#)
  20. [Slide 20 Dr. Bradford Hess 2](#)
  21. [Slide 21 Dr. Harry Foxwell 1](#)
  22. [Slide 22 Dr. Harry Foxwell 2](#)
  23. [Slide 23 Eric Simmon 1](#)
  24. [Slide 24 Eric Simmon 2](#)
  25. [Slide 25 William Miller 1](#)
  26. [Slide 26 William Miller 2](#)

## 2. [Slides](#)

1. [Slide 1 Data Science for Big Data](#)
2. [Slide 2 Overview](#)
3. [Slide 3 The Profit and Data Enterprises](#)
4. [Slide 4 Federal Big Data Working Group Meetup](#)
5. [Slide 5 Silicon Valley to Washington](#)
6. [Slide 6 First White House Data Chief Discusses His Top Priorities](#)
7. [Slide 7 Precision Medicine and Natural Medicine](#)
8. [Slide 8 Tech Meetup at White House](#)
9. [Slide 9 USDA Data Science MOOC](#)
10. [Slide 10 Upcoming Meetups](#)
11. [Slide 11 Summary 1](#)
12. [Slide 12 Summary 2](#)

## 3. [Slides](#)

1. [Slide 1 Semantic Data Discovery: Proof of Concept for DHS](#)
  2. [Slide 2 Information Sharing at DHS](#)
  3. [Slide 3 NIEM as Big Data in a Network with Data Science](#)
  4. [Slide 4 NIEM 3.0 Alpha 2 Release and Thetus Savanna Review](#)
  5. [Slide 5 NIEM and UCore 2.0 Semantic Layer for Information Sharing](#)
  6. [Slide 6 Global Terrorism Database Experience](#)
  7. [Slide 7 A Quint for Cross Information Sharing and Integration in the Intelligence Community](#)
  8. [Slide 8 Dynamic Case Management Pilot for Healthcare.gov](#)
  9. [Slide 9 Proof of Concept Steps](#)
  10. [Slide 10 Semantic Community](#)
4. [Spotfire Dashboard](#)
  5. [Research Notes](#)
  6. [Ontology Summit 2015 Agenda](#)
    1. [Monday, April 13](#)
    2. [Tuesday, April 14](#)
  7. [Ontology Summit 2015 Background](#)
    1. [Prepared Presentation Materials](#)
    2. [Audio Recordings](#)
    3. [Additional Resources](#)
    4. [Abstract](#)
  8. [Ontology Summit 2015 Communique](#)
    1. [Introduction](#)
    2. [The Case for IoT Ontologies](#)
    3. [How Ontologies are Used in IoT](#)
      1. [Ontology Mapping](#)
      2. [Standards Integration](#)
      3. [Decision Support for IoT](#)
    4. [Beyond Semantic Sensor Network Ontologies](#)
    5. [Ontological Issues](#)
      1. [Scalability](#)
      2. [Standards Integration](#)
    6. [Challenges](#)
    7. [Forecasts](#)
    8. [Recommendations](#)
    9. [Terminology](#)
  9. [Training Students to Extract Value from Big Data](#)
    1. [The National Academies Press](#)
      1. [Authors](#)
      2. [Description](#)
      3. [Topics](#)
      4. [Publication Info](#)

5. [Copyright Information](#)
2. [Cover Page](#)
3. [The National Academies Press](#)
4. [Planning Committee on Training Students](#)
5. [Committee on Applied and Theoretical Statistics](#)
6. [Board of Mathematical Sciences and Their Applications](#)
7. [Acknowledgment of Reviewers](#)
8. **1 INTRODUCTION**
  1. [Workshop Overview](#)
    1. [BOX 1.1 Statement of Task](#)
  2. [National Efforts in Big Data](#)
    1. [Suzanne Iacono, National Science Foundation](#)
  3. [Organization of This Report](#)
9. **2 THE NEED FOR TRAINING: EXPERIENCES AND CASE STUDIES**
  1. [Training Students to Do Good with Big Data](#)
    1. [Rayid Ghani, University of Chicago](#)
  2. [The Need for Training in Big Data: Experiences and Case Studies](#)
    1. [Guy Lebanon, Amazon Corporation](#)
10. **3 PRINCIPLES FOR WORKING WITH BIG DATA 13**
  1. [Teaching about MapReduce](#)
    1. [Jeffrey Ullman, Stanford University](#)
  2. [Big Data Machine Learning—Principles for Industry](#)
    1. [Alexander Gray, Skytree Corporation](#)
  3. [Principles for the Data Science Process](#)
    1. [Duncan Temple Lang, University of California, Davis](#)
  4. [Principles for Working with Big Data](#)
    1. [Juliana Freire, New York University](#)
    2. [FIGURE 3.1 Simplified schematic of the big data analysis pipeline](#)
11. **4 COURSES, CURRICULA, AND INTERDISCIPLINARY PROGRAMS**
  1. [Computational Training and Data Literacy for Domain Scientists](#)
    1. [Joshua Bloom, University of California, Berkeley](#)
  2. [Data Science and Analytics Curriculum Development at Rensselaer \(and the Tetherless World Constellation\)](#)
    1. [Peter Fox, Rensselaer Polytechnic Institute](#)
    2. [FIGURE 4.1 Framework for modern informatics](#)
    3. [FIGURE 4.2 Generations of mediation](#)
  3. [Experience with a First Massive Online Open Course on Data Science](#)
    1. [William Howe, University of Washington](#)
12. **5 SHARED RESOURCES**
  1. [Can Knowledge Bases Help Accelerate Science?](#)
    1. [Christopher Ré, Stanford University](#)
  2. [Divide and Recombine for Large, Complex Data](#)
    1. [Bill Cleveland, Purdue University](#)
  3. [Yahoo's Webscope Data Sharing Program](#)
    1. [Ron Brachman, Yahoo Labs](#)
  4. [Resource Sharing](#)
    1. [Mark Ryland, Amazon Corporation](#)
13. **6 WORKSHOP LESSONS**
  1. [Whom to Teach: Types of Students to Target in Teaching Big Data](#)
  2. [How to Teach: The Structure of Teaching Big Data](#)
  3. [What to Teach: Content in Teaching Big Data](#)
  4. [Parallels in Other Disciplines](#)
14. **Footnotes**
  1. [1](#)
  2. [2](#)
  3. [3](#)

4. [4](#)
5. [1](#)
6. [1](#)
7. [2](#)
8. [3](#)
9. [4](#)
10. [5](#)
11. [6](#)
12. [7](#)
13. [1](#)
14. [2](#)
15. [3](#)
16. [4](#)
17. [5](#)
18. [6](#)
19. [7](#)
20. [8](#)
21. [9](#)
22. [1](#)
23. [2](#)
24. [3](#)
25. [4](#)
26. [5](#)
27. [6](#)
28. [7](#)
29. [8](#)
30. [9](#)
31. [10](#)
32. [11](#)

## 15. REFERENCES

## 16. APPENDIXES

1. [A Registered Workshop Participants](#)
2. [B Workshop Agenda](#)
  1. [APRIL 11, 2014](#)
  2. [APRIL 12, 2014](#)
3. [C Acronyms](#)

## 10. Getting Data Right: Tackling The Challenges of Big Data Volume and Variety

1. [CHAPTER 1 The Solution: Data Curation at Scale](#)
  1. [Three Generations of Data Integration Systems](#)
    1. [Table 1-1. Evolution of Three Generations of Data Integration Systems](#)
  2. [Five Tenets for Success](#)
    1. [Tenet 1: Data curation is never done](#)
    2. [Tenet 2: A PhD in AI can't be a requirement for success](#)
    3. [Tenet 3: Fully automatic data curation is not likely to be successful](#)
    4. [Tenet 4: Data curation must fit into the enterprise ecosystem](#)
    5. [Tenet 5: A scheme for "finding" data sources must be present](#)

## 11. NEXT

1. [Story](#)
  1. [Data Science for Cyber Physical Systems-Internet of Things](#)
  2. [Training Students to Extract Value from Big Data](#)
    1. [Slide 1 An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library](#)
    2. [Slide 2 Outline](#)
    3. [Slide 3 Semantic Community: Spotfire Cloud Library MindTouch](#)
    4. [Slide 4 Semantic Community: Spotfire Cloud Library Recent Analyses](#)

5. [Slide 5 Semantic Community: Spotfire Cloud Library Browse Library](#)
6. [Slide 6 Semantic Community: Spotfire Cloud Library Shared Folders](#)
7. [Slide 7 Semantic Community: Spotfire Web Player](#)
8. [Slide 8 Semantic Community: Spotfire Web Player Edit](#)
9. [Slide 9 Semantic Community: Spotfire Web Player Menu](#)
10. [Slide 10 Semantic Community: Spotfire Web Player Menu and New Visualizations](#)
11. [Slide 11 TIBCO Spotfire: Cloud User's Guide](#)
12. [Slide 12 TIBCO: Spotfire Cloud User's Guide Getting Started](#)
13. [Slide 13 TIBCO: Spotfire Cloud User's Guide Setting Up Analyses](#)
14. [Slide 14 TIBCO: Spotfire Cloud User's Guide Data Preparation in Microsoft Excel](#)
15. [Slide 15 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 1](#)
16. [Slide 16 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 2](#)
17. [Slide 17 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 3](#)
18. [Slide 18 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Markers](#)
19. [Slide 19 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Colored Regions](#)
20. [Slide 20 TIBCO: Spotfire Cloud User's Guide Exporting Your Analysis](#)
21. [Slide 21 TIBCO: Spotfire Cloud User's Guide Sharing Your Analysis](#)
22. [Slide 22 TIBCO Spotfire: Education](#)
3. [Summit Highlights](#)
  1. [Slide 1 Opening Remarks 1](#)
  2. [Slide 2 Opening Remarks 2](#)
  3. [Slide 3 Dr. Shoumen Data 1](#)
  4. [Slide 4 Dr. Shoumen Data 2](#)
  5. [Slide 5 Dr. Shoumen Data 3](#)
  6. [Slide 6 Dr. Shoumen Data 4](#)
  7. [Slide 7 Track A 1](#)
  8. [Slide 8 Track A 2](#)
  9. [Slide 9 Track B 1](#)
  10. [Slide 10 Track B 2](#)
  11. [Slide 11 Keith Mazullo 1](#)
  12. [Slide 12 Keith Mazullo 2](#)
  13. [Slide 13 Keith Mazullo 3](#)
  14. [Slide 14 Keith Mazullo 4](#)
  15. [Slide 15 Track C 1](#)
  16. [Slide 16 Track C 2](#)
  17. [Slide 17 Track D 1](#)
  18. [Slide 18 Track D 2](#)
  19. [Slide 19 Dr. Bradford Hess 1](#)
  20. [Slide 20 Dr. Bradford Hess 2](#)
  21. [Slide 21 Dr. Harry Foxwell 1](#)
  22. [Slide 22 Dr. Harry Foxwell 2](#)
  23. [Slide 23 Eric Simmon 1](#)
  24. [Slide 24 Eric Simmon 2](#)
  25. [Slide 25 William Miller 1](#)
  26. [Slide 26 William Miller 2](#)
2. [Slides](#)
  1. [Slide 1 Data Science for Big Data](#)
  2. [Slide 2 Overview](#)
  3. [Slide 3 The Profit and Data Enterprises](#)
  4. [Slide 4 Federal Big Data Working Group Meetup](#)
  5. [Slide 5 Silicon Valley to Washington](#)
  6. [Slide 6 First White House Data Chief Discusses His Top Priorities](#)
  7. [Slide 7 Precision Medicine and Natural Medicine](#)
  8. [Slide 8 Tech Meetup at White House](#)
  9. [Slide 9 USDA Data Science MOOC](#)

10. [Slide 10 Upcoming Meetups](#)
  11. [Slide 11 Summary 1](#)
  12. [Slide 12 Summary 2](#)
3. [Slides](#)
1. [Slide 1 Semantic Data Discovery: Proof of Concept for DHS](#)
  2. [Slide 2 Information Sharing at DHS](#)
  3. [Slide 3 NIEM as Big Data in a Network with Data Science](#)
  4. [Slide 4 NIEM 3.0 Alpha 2 Release and Thetus Savanna Review](#)
  5. [Slide 5 NIEM and UCore 2.0 Semantic Layer for Information Sharing](#)
  6. [Slide 6 Global Terrorism Database Experience](#)
  7. [Slide 7 A Quint for Cross Information Sharing and Integration in the Intelligence Community](#)
  8. [Slide 8 Dynamic Case Management Pilot for Healthcare.gov](#)
  9. [Slide 9 Proof of Concept Steps](#)
  10. [Slide 10 Semantic Community](#)
4. [Spotfire Dashboard](#)
5. [Research Notes](#)
6. [Ontology Summit 2015 Agenda](#)
  1. [Monday, April 13](#)
  2. [Tuesday, April 14](#)
7. [Ontology Summit 2015 Background](#)
  1. [Prepared Presentation Materials](#)
  2. [Audio Recordings](#)
  3. [Additional Resources](#)
  4. [Abstract](#)
8. [Ontology Summit 2015 Communique](#)
  1. [Introduction](#)
  2. [The Case for IoT Ontologies](#)
  3. [How Ontologies are Used in IoT](#)
    1. [Ontology Mapping](#)
    2. [Standards Integration](#)
    3. [Decision Support for IoT](#)
  4. [Beyond Semantic Sensor Network Ontologies](#)
  5. [Ontological Issues](#)
    1. [Scalability](#)
    2. [Standards Integration](#)
  6. [Challenges](#)
  7. [Forecasts](#)
  8. [Recommendations](#)
  9. [Terminology](#)
9. [Training Students to Extract Value from Big Data](#)
  1. [The National Academies Press](#)
    1. [Authors](#)
    2. [Description](#)
    3. [Topics](#)
    4. [Publication Info](#)
    5. [Copyright Information](#)
  2. [Cover Page](#)
  3. [The National Academies Press](#)
  4. [Planning Committee on Training Students](#)
  5. [Committee on Applied and Theoretical Statistics](#)
  6. [Board of Mathematical Sciences and Their Applications](#)
  7. [Acknowledgment of Reviewers](#)
  8. [1 INTRODUCTION](#)
    1. [Workshop Overview](#)
      1. [BOX 1.1 Statement of Task](#)

2. [National Efforts in Big Data](#)
  1. [Suzanne Iacono, National Science Foundation](#)
  3. [Organization of This Report](#)
9. [2 THE NEED FOR TRAINING: EXPERIENCES AND CASE STUDIES](#)
  1. [Training Students to Do Good with Big Data](#)
    1. [Rayid Ghani, University of Chicago](#)
  2. [The Need for Training in Big Data: Experiences and Case Studies](#)
    1. [Guy Lebanon, Amazon Corporation](#)
10. [3 PRINCIPLES FOR WORKING WITH BIG DATA](#)
  13. [Teaching about MapReduce](#)
    1. [Jeffrey Ullman, Stanford University](#)
  2. [Big Data Machine Learning—Principles for Industry](#)
    1. [Alexander Gray, Skytree Corporation](#)
  3. [Principles for the Data Science Process](#)
    1. [Duncan Temple Lang, University of California, Davis](#)
  4. [Principles for Working with Big Data](#)
    1. [Juliana Freire, New York University](#)
    2. [FIGURE 3.1 Simplified schematic of the big data analysis pipeline](#)
11. [4 COURSES, CURRICULA, AND INTERDISCIPLINARY PROGRAMS](#)
  1. [Computational Training and Data Literacy for Domain Scientists](#)
    1. [Joshua Bloom, University of California, Berkeley](#)
  2. [Data Science and Analytics Curriculum Development at Rensselaer \(and the Tetherless World Constellation\)](#)
    1. [Peter Fox, Rensselaer Polytechnic Institute](#)
    2. [FIGURE 4.1 Framework for modern informatics](#)
    3. [FIGURE 4.2 Generations of mediation](#)
  3. [Experience with a First Massive Online Open Course on Data Science](#)
    1. [William Howe, University of Washington](#)
12. [5 SHARED RESOURCES](#)
  1. [Can Knowledge Bases Help Accelerate Science?](#)
    1. [Christopher Ré, Stanford University](#)
  2. [Divide and Recombine for Large, Complex Data](#)
    1. [Bill Cleveland, Purdue University](#)
  3. [Yahoo's WebScope Data Sharing Program](#)
    1. [Ron Brachman, Yahoo Labs](#)
  4. [Resource Sharing](#)
    1. [Mark Ryland, Amazon Corporation](#)
13. [6 WORKSHOP LESSONS](#)
  1. [Whom to Teach: Types of Students to Target in Teaching Big Data](#)
  2. [How to Teach: The Structure of Teaching Big Data](#)
  3. [What to Teach: Content in Teaching Big Data](#)
  4. [Parallels in Other Disciplines](#)
14. [Footnotes](#)
  1. [1](#)
  2. [2](#)
  3. [3](#)
  4. [4](#)
  5. [1](#)
  6. [1](#)
  7. [2](#)
  8. [3](#)
  9. [4](#)
  10. [5](#)
  11. [6](#)
  12. [7](#)
  13. [1](#)

- 14. [2](#)
- 15. [3](#)
- 16. [4](#)
- 17. [5](#)
- 18. [6](#)
- 19. [7](#)
- 20. [8](#)
- 21. [9](#)
- 22. [1](#)
- 23. [2](#)
- 24. [3](#)
- 25. [4](#)
- 26. [5](#)
- 27. [6](#)
- 28. [7](#)
- 29. [8](#)
- 30. [9](#)
- 31. [10](#)
- 32. [11](#)

15. [REFERENCES](#)

16. [APPENDIXES](#)

- 1. [A Registered Workshop Participants](#)
- 2. [B Workshop Agenda](#)
  - 1. [APRIL 11, 2014](#)
  - 2. [APRIL 12, 2014](#)
- 3. [C Acronyms](#)

10. [Getting Data Right: Tackling The Challenges of Big Data Volume and Variety](#)

- 1. [CHAPTER 1 The Solution: Data Curation at Scale](#)
  - 1. [Three Generations of Data Integration Systems](#)
    - 1. [Table 1-1. Evolution of Three Generations of Data Integration Systems](#)
  - 2. [Five Tenets for Success](#)
    - 1. [Tenet 1: Data curation is never done](#)
    - 2. [Tenet 2: A PhD in AI can't be a requirement for success](#)
    - 3. [Tenet 3: Fully automatic data curation is not likely to be successful](#)
    - 4. [Tenet 4: Data curation must fit into the enterprise ecosystem](#)
    - 5. [Tenet 5: A scheme for "finding" data sources must be present](#)

11. [NEXT](#)

---

## Story

---

### Data Science for Cyber Physical Systems-Internet of Things

Dr. Leo Obrst, The MITRE Corporation, Information Semantics Cognitive Science & Artificial Intelligence, CCG, said recently:

Thanks to you all for your participation and Michael Gruninger and Mark Underwood for organizing the Ontology Summit 2015 as co-chairs. Also thanks to our Organizing Committee, Co-Champions, Ram Sriram and Fauod Ramia of NITRD for organizing the Symposium, and all the invited speakers during the many virtual sessions. In addition, thanks to Ken Baclawski and Peter Yim for background infrastructure support, and to Joe Kopena for web page creativity and support. It's been a very interesting Ontology Summit 2015!

The [Ontology Summit Wiki](#) says:

The Ontology Summit is an annual series of events (first started by Ontolog and NIST in 2006) that involves the ontology community and communities related to each year's theme chosen for the summit. The Ontology Summit program is now co-organized by Ontolog, NIST, NCOR, NCBO, IAOA, NCO\_NITRD along with the co-sponsorship of other organizations that are supportive of the Summit goals and objectives.

we are witnessing the emergence of “intelligent devices,” such as smart meters, smart cars, etc., with considerable sensing and networking capabilities. Hence, these devices – and the network -- will be constantly sensing, monitoring, and interpreting the environment – this is sometimes referred to as the Internet of Things. And as local and wide area networks became almost secondary to the WWW (World-Wide Web), users and their usage patterns will become increasingly visible. This will have significant implications for both the market for advanced computing and communication infrastructure and the future markets – for nearly 4.5 billion people -- that net-centric societies will create. Ontologies will play a significant role in the realization of Smart networked systems and societies (SNSS).

In computer science and information science, [an ontology](#) is a formal naming and definition of the types, properties, and interrelationships of the entities that really or fundamentally exist for a particular domain of discourse. Many fields all create ontologies to limit complexity and to organize information. The ontology can then be applied to problem solving.

So the Internet has evolved from the Web of documents, to the Web of data, now big data, and to the Web of things like [cyber physical systems](#), a system of collaborating computational elements controlling physical entities. An ontology is part of [Data Science](#), the extraction of knowledge from data.

A familiar example is MIT's ongoing CarTel project where a fleet of taxis collecting real-time traffic information in the Boston area. Together with historical data, this information is then used for calculating fastest routes for a given time of the day. The [Federal Big Data Working Group Meetup](#) is familiar with this from its participation in the [MIT Big Data Initiative](#) and [Tackling the Challenges of Big Data](#).

This is all about Turning Data Into Value: Organizations do not fear a shortage of data. Systems, applications, and devices all produce and consume exponentially increasing amounts of it. The Federal Big Data Working Group Meetup Data Scientists help organizations use their data to best advantage by integrating, analyzing, and acting on it via event processing as follows:

- Integration provides the right data to the right system or person in real time.
- Analytics lets users develop insights using vast amounts of data to understand the past and anticipate the future.
- Event processing combines the knowledge gained from analytics with real-time information to identify patterns of events and act to bring about the best outcomes.

My comments during the Ontology Summit were:

- <https://www.linkedin.com/in/shoumendatta>
- Proposal to digitize the semantic web using a numerical nomenclature as a sub-layer to URI (utilizing a IPv6 type format).
- Unified Theory of Relativistic Identification of Information in a Systems Age: Proposed Convergence of Unique Identification with Syntax and Semantics through Internet Protocol version 6:  
<http://dspace.mit.edu/handle/1721.1/41902>
- Tim Berners-Lee on the Future of His Invention: <http://www.emc.com/leadership/featur...erners-lee.htm>
- Sounds like the MITRE-VA work might benefit from Michael Stonebraker's Tamr that we have had presented at our Federal Big Data Working Group Meetup on March 16th: <http://www.meetup.com/Federal-Big-Da...nts/220121871/>
- To turn data (IoT) into value you need: Integration, Analytics, and Event Processing
- Semantic Big Data Science is more about putting semantics in the data (especially big data, IoT, etc.) than putting semantics in the technology (like the past emphasis)!
- Is the Data Science for USGS Minerals Big Data an IoT? [http://semanticcommunity.info/Data\\_Sc...erals\\_Big\\_Data](http://semanticcommunity.info/Data_Sc...erals_Big_Data) I have been working on this will listening all day so I would have something to show as an example. Thank you for an excellent day of listening!
- Glad to hear good medical data science from Kaiser and other leading providers. My Kaiser medical records are available anywhere I can get on the Web and even downloadable to a memory stick.
- See Data Science for Natural Medicines and Epigenetics Meetup on May 4th: <http://www.meetup.com/Federal-Big-Da...nts/221457330/> There is a history of failed medical theories like all disease was thought to be genetic. There is a new book, Epigenetics, The Death of The Genetic Theory of Disease Transmission, which says all of the diseases that you have been told are genetic, are indeed not. "Epigenetics" explains various diseases and cites the nutrients that are missing. Genes need "trace" nutrients such as minerals in order to make them function optimally. So we are in the Age of Biochemistry, Nutrition, and Epigenetics
- DavidBlevins: @Brand: sure you have plenty to say about interoperability, given your history and participation Thank you, I do, namely data science simplifies the interoperability problem and one should start with the science and not try to boil the ocean of the general interoperability problem - e.g. Alan Turing's code breaking work

## Training Students to Extract Value from Big Data

A recent NAS Workshop report on [Training Students to Extract Value from Big Data](#) said:

Data sets—whether in science and engineering, economics, health care, public policy, or business—have been growing rapidly; the recent National Research Council (NRC) report [Frontiers in Massive Data Analysis](#) documented the rise of “big data,” as systems are routinely returning terabytes, petabytes, or more of information (National Research Council, 2013). Big data has become pervasive because of the availability of high-throughput data collection technologies, such as information-sensing mobile devices, remote sensing, radiofrequency identification readers, Internet log records, and wireless sensor networks. Science, engineering, and business have rapidly transitioned from the longstanding state of striving to develop information from scant data to a situation in which the challenge is now that the amount of information exceeds a human’s ability to examine, let alone absorb, it. Web companies—such as Yahoo, Google, and Amazon—commonly work with data sets that consist of billions of items, and they are likely to increase by an order of magnitude or more as the Internet of Things [1](#) matures. In other words, the size and scale of data, which can be overwhelming today, are only increasing. In addition, data sets are increasingly complex, and this potentially increases the problems associated with such concerns as missing information and other quality concerns, data heterogeneity, and differing data formats.

Advances in technology have made it easier to assemble and access large amounts of data. Now, a key challenge is to develop the experts needed to draw reliable inferences from all that information. The nation’s ability to make use of the data depends heavily on the availability of a workforce that is properly trained and ready to tackle these high-need areas. A report from McKinsey & Company (Manyika et al., 2011) has predicted shortfalls of 150,000 data analysts and 1.5 million managers who are knowledgeable about data and their relevance. It is becoming increasingly important to increase the pool of qualified scientists and engineers who can extract value from big data. Training students to be capable in exploiting big data requires experience with statistical analysis, machine learning, and computational infrastructure that permits the real problems associated with massive data to be revealed and, ultimately, addressed. The availability of repositories (of both data and software) and computational infrastructure will be necessary to train the next generation of data scientists. Analysis of big data requires cross-disciplinary skills, including the ability to make modeling decisions while balancing trade-offs between optimization and approximation, all while being attentive to useful metrics and system robustness. To develop those skills in students, it is important to identify whom to teach, that is, the [EDUCATIONAL](#) background, experience, and characteristics of a prospective data science student; what to teach, that is, the technical and practical content that should be taught to the student; and how to teach, that is, the structure and organization of a data science [program](#).

The topic of training students in big data is timely, as universities are already experimenting with courses and [PROGRAMS](#) tailored to the needs of students who will work with big data. Eight university [programs](#) have been or will be launched in 2014 alone. [2](#) The workshop that is the subject of this report was designed to enable participants to learn and benefit from emerging insights while innovation in [EDUCATION](#) is still ongoing.

Suzanne Iacono, of NSF, set the stage for the workshop by speaking about [national efforts in big data](#), current challenges, and NSF’s motivations for sponsoring the workshop. She explained that the workshop was an outgrowth of the national big data research and development (R&D) initiative. The federal government is interested in big data for three reasons:

- To stimulate commerce and the economy.
- To accelerate the pace of discovery and enable new activities.
- To address pressing national challenges in [EDUCATION](#), health care, and public safety.

The [Federal Big Data Working Group Meetup](#) is training its members, including students, to extract value from big data in support of the national efforts in big data.

The simple tutorial below for An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library, by Dr. Brand L. Nleemann, Founder and Co-Organizer of the Federal Big Data Working Group Meetup, is meant to supplement the recent NAS Workshop and support the NSF National efforts in big data.

[Slides](#)

Slide 1 An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library  
<http://semanticcommunity.info/>

# An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library

Dr. Brand Niemann

Director and Senior Data Scientist/Data Journalist

Semantic Community

<http://semanticcommunity.info/>

June 16, 2015

1

Slide 2 Outline

# Outline

- Semantic Community Spotfire Cloud Library:
  - MindTouch
  - Recent Analyses
  - Browse Library
  - Shared Folders
- Spotfire Web Player
  - Normal Mode
  - Edit Mode
  - Menu (details) and New Visualizations (details)
- TIBCO: Spotfire Cloud User's Guide:
  - Cloud User's Guide
  - Getting Started
  - Setting Up Analyses
  - Data Preparation in Microsoft Excel
  - Creating a Visualization
  - Creating a Map Chart with Markers
  - Creating a Map Chart with Colored Regions
  - Exporting Your Analysis
  - Sharing Your Analysis
- TIBCO Spotfire:
  - Education

2

Slide 3 Semantic Community: Spotfire Cloud Library MindTouch

[http://semanticcommunity.info/#Spotfire\\_Cloud\\_Library](http://semanticcommunity.info/#Spotfire_Cloud_Library)

## Semantic Community: Spotfire Cloud Library MindTouch

Semantic Community - S... x HDP2015DataLab-Spotfire x Learn Spotfire x SpotfireCloudInventory-S... x TIBCO Spotfire® Cloud U... x

semanticcommunity.info/#Spotfire\_Cloud\_Library

Apps ssid\_ss\_x\_cmmc vc... Yahoo! Twitter / Home Chili's To Go Menu LEGO.com Star Wars... Cox Communicatio... Incredible Ball Boy C...

Other bookmarks

## Most Recent Migrated Spotfire Files

Title	Link	Story	Comments
Global Database of Events, Language, and Tone (GDELT)	<a href="#">Web Player</a> <a href="#">Web Player</a> & <a href="#">Web Player</a>	<a href="#">Story</a>	<a href="#">Hackathon</a>
AmericasDataFest	<a href="#">Web Player</a>	<a href="#">Story</a>	<a href="#">CONTEST</a>
Hubway Data Visualization Challenge	<a href="#">Web Player</a>	<a href="#">Slides</a>	<a href="#">CONTEST</a>
Global Terrorism Database	<a href="#">Web Player</a>	<a href="#">Story</a>	Meetup
Data Science Central Meteors	<a href="#">Web Player</a>	<a href="#">Story</a>	Training

[Spreadsheet](#)

For [INTERNET EXPLORER](#) Users and Those Wanting Full Screen Display Use: [Web Player](#) Get [Spotfire](#) [iPad App](#)

Cover Page | Spotfire Silver to Spotfire Cloud Migration | Spotfire Learning | [Edit](#) | [Twitter](#) [Facebook](#) [Filter](#) [Print](#) [Close](#)

Navigation and Metadata

**Data Science for Big Data: Internet of Things by Brand Niemann, April 29, 2015.**

[Spreadsheet](#)

Instructions: Use Filter to Right to Select Category, and/or Comment. Then [Reset All Filters](#) and [WS](#)

Distribution – Organization  
Data table:  
  
(RowCount)

Distribution – Topic  
Data table:  
  
(RowCount)

Filters   
Type to search filters

Spotfire Cloud Library   
Name   
Type to search in list   
(All) 343 values  
2008SApartEventsPUF...  
2008SAPartIDEEventsPUF...  
2010Census-Spotfire  
2010EIISymposium-Spotf...  
2010ReporttoCongress...  
2010StatAbstract-Spotfire

[http://semanticommunity.info/#Spotfire\\_Cloud\\_Library](http://semanticommunity.info/#Spotfire_Cloud_Library)

3

Slide 4 Semantic Community: Spotfire Cloud Library Recent Analyses

# Semantic Community:

## Spotfire Cloud Library Recent Analyses

The screenshot shows the TIBCO Spotfire Cloud Library interface. At the top, there's a navigation bar with tabs for 'New analysis', 'Recent analyses' (which is selected), 'Browse library', and 'Shared folders'. Below the navigation bar, there's a 'Live Chat' button on the left. The main area displays a grid of eight analysis cards, each with a thumbnail, title, version, and a brief description. The cards are arranged in two rows of four. The titles of the analyses include 'SpotfireCloudInventory-Spotfire', 'HDP2015DataLab-Spotfire', 'ESRIGISHealth-Spotfire', 'HubwayDataChapter9-Spotfire', 'NetworkAnalytics-Spotfire', 'FederalIDashboard07151010...', 'EPAGreenVehiclesRatings20...', and 'EPASubstanceRegistrySyste...'. The descriptions provide details about the data science applications and versions.

4

### Slide 5 Semantic Community: Spotfire Cloud Library Browse Library

# Semantic Community:

## Spotfire Cloud Library Browse Library

The screenshot shows a web browser window with multiple tabs open. The active tab is 'Library - TIBCO Spotfire' at <https://spotfire.cloud.tibco.com/spotfire/wp/startPage#/libraryBrowser?id=01c5fbef-7c04-4ea7-945c-97e374512467>. The page displays a list of analyses in the 'users/bniemann/Public' folder. The analyses listed are:

Name	Modified
2008BSAPartDEventsPUF1-Spotfire Version 2.0 Test of Web Player	12/12/13 1:33:12 AM
2008BSAPartDEventsPUFData-Spotfire Version 2.0 Test of Web Player	8/3/12 1:46:32 PM
2010Census-Spotfire Version 2.0 Test of Web Player	7/25/12 7:19:19 PM
2010OEISymposium-Spotfire Version 2.0 Test of Web Player	7/29/12 1:02:32 PM
2010ReporttoCongressonWhiteHouseStaffSalaries-Spotfire Version 2.0 Test of Web Player	7/29/12 1:10:43 AM
2010StatAbstract-Spotfire Version 2.0 Test of Web Player	7/25/12 7:30:35 PM
2011DataCenterConsolidations-Spotfire Version 2.0 Test of Web Player	7/26/12 4:00:25 AM

A green vertical bar on the left is labeled 'Live Chat'. The bottom status bar shows the Windows taskbar with icons for Google, Chrome, and the Spotfire application, along with system status like battery level and time (2:47 PM).

5

## Slide 6 Semantic Community: Spotfire Cloud Library Shared Folders

# Semantic Community: Spotfire Cloud Library Shared Folders

The screenshot shows a web browser window with multiple tabs open. The active tab is 'Library - TIBCO Spotfire' at <https://spotfire.cloud.tibco.com/spotfire/wp/startPage#/share>. The page displays a table of shared folders:

Name	Type	Modified	Owner
bniemann Workgroup	Work Group	9/8/2014	bniemann workgroup
Public	Public	6/13/2015	bniemann

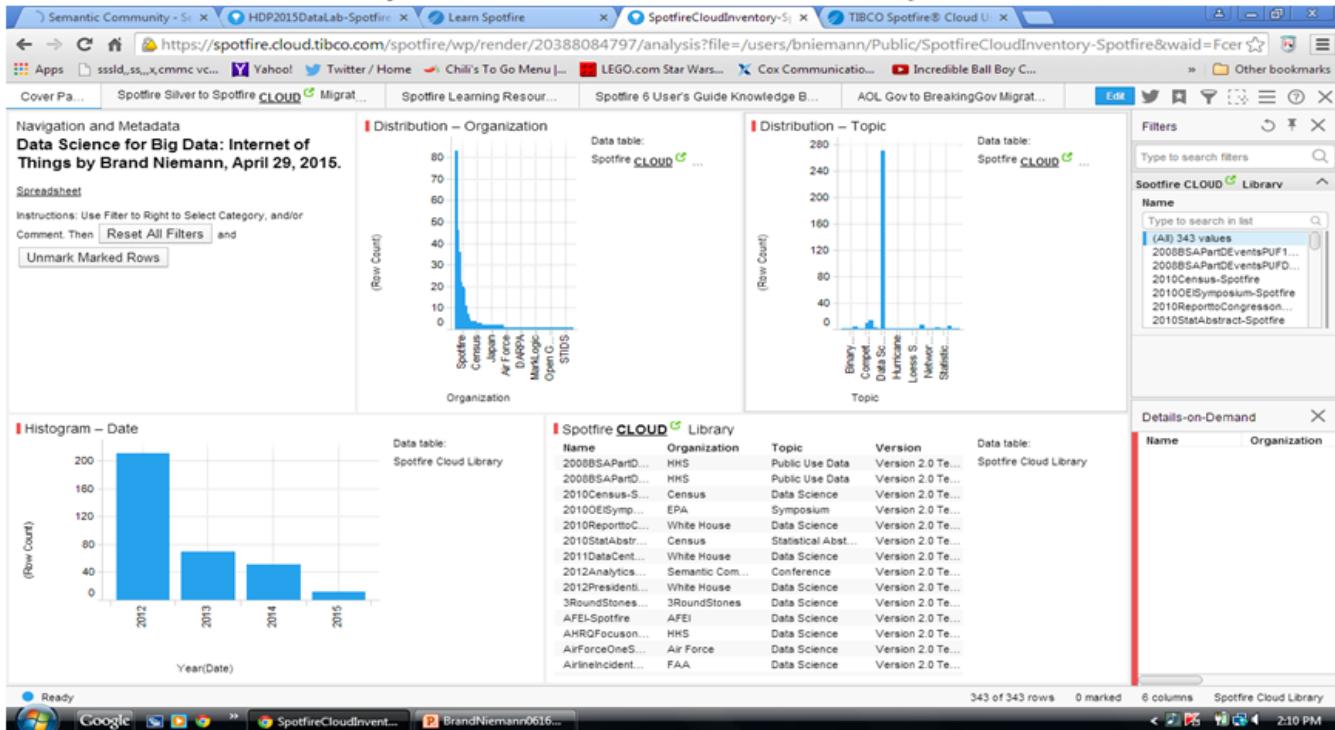
A green vertical bar on the left contains the text 'Live Chat'. The bottom of the screen shows the Windows taskbar with icons for Google, Chrome, and the Spotfire application.

6

Slide 7 Semantic Community: Spotfire Web Player

# Semantic Community:

## Spotfire Web Player



### Web Player

7

### Slide 8 Semantic Community: Spotfire Web Player Edit

# Semantic Community:

## Spotfire Web Player Edit

The screenshot shows the Spotfire Web Player Edit interface with the following components:

- Top Bar:** Shows multiple tabs including "Data Science for Cyber", "HDP2015DataLab-Spotfi", "Learn Spotfire", "SpotfireCloudInventory", "TIBCO Spotfire® Cloud", and "TIBCO Spotfire® Cloud".
- Left Panel:** A "Cover Page" section titled "Data Science for Big Things by Brand Niemann, April 20, 2010." It includes a "Spreadsheet" section with instructions and a "Unmark Marked Rows" button.
- Middle Panel:**
  - A histogram titled "Histogram – Date" showing the distribution of dates from 2012 to 2015.
  - A bar chart titled "Distribution – Organization" showing the count of rows for various organizations.
  - A bar chart titled "Distribution – Topic" showing the count of rows for various topics.
- Right Panel:**
  - A "Filters" panel containing a search bar and a list of filters: "2008BSAPartD...", "2008BSAPartD...", "2010Census-S...", "2010OEISymp...", "2010ReporttoC...", "2010StatAbstr...", "2011DataCent...", "2012Analytics...", "2012Presidenti...", "3RoundStones...", "AFEI-Spotfire", "AHRQFocuson...", "AirForceOneS...".
  - A "Details-on-Demand" panel showing a table with columns "Name" and "Organization".
- Bottom:** A status bar showing "343 of 343 rows", "0 marked", "6 columns", and "Spotfire Cloud Library".

8

### Slide 9 Semantic Community: Spotfire Web Player Menu

# Semantic Community:

## Spotfire Web Player Menu

- Save
- Undo
- Add Data
- Recommend Visualizations
- Add Annotation
- Properties
- Filters
- Data
- Details-on-Demand
- Bookmarks
- Visual Themes
- Arrange Visualizations
- Share to Twitter
- Menu
- Help
- Close Analysis

9

Slide 10 Semantic Community: Spotfire Web Player Menu and New Visualizations

# Semantic Community:

## Spotfire Web Player Menu and New Visualizations

- Menu (details):
  - Add new page
  - Rename page
  - Go to page
  - Page navigation
  - Edit data tables
  - Edit markings
  - Analysis information
  - Export
  - Share
  - Print
  - Other tools
  - Download as DXP file
  - About TIBCO Spotfire Web PLayer
- New Visualizations(details):
  - Table
  - Cross Table
  - Bar Chart
  - Line Chart
  - Combination Chart
  - Pie Chart
  - Scatter Plot
  - Map Chart (Markers)
  - Map Chart (Features)
  - Treemap
  - Parallel Coordinate Plot
  - Recommended Visualizations

10

Slide 11 TIBCO Spotfire: Cloud User's Guide

# TIBCO Spotfire: Cloud User's Guide

The image shows two screenshots side-by-side. The left screenshot is a web browser displaying the 'TIBCO Spotfire® Cloud User's Guide' on the docs.tibco.com website. It features a navigation bar with links like 'Home', 'CONTENTS', and 'Search'. The main content area is titled 'TIBCO Spotfire® Cloud User's Guide' and contains sections such as 'Getting started', 'Using Recommended visualizations', and 'Learning more'. The right screenshot shows the TIBCO Spotfire software interface with multiple tabs open, including 'Intro', 'Sales Performance', 'Territory Analysis', and 'Effect of Promotions'. Below the tabs, there are several visualization windows, such as a 'Performance Matrix - Year 1 vs Year 2' and a 'Brand Share Change per Promotions' chart.

## Help

11

### Slide 12 TIBCO: Spotfire Cloud User's Guide Getting Started

# TIBCO: Spotfire Cloud User's Guide

## Getting Started

- **Getting Started**
  - To get a quick start, see:
    - Creating a new analysis and Opening an existing analysis to learn how to bring data into an analysis
    - Creating a visualization to learn how to set up various types of visualizations
    - Sharing your analysis to learn how make your analyses available to others

12

Slide 13 TIBCO: Spotfire Cloud User's Guide Setting Up Analyses

# TIBCO: Spotfire Cloud User's Guide

## Setting Up Analyses

- **Setting Up Analyses**

- Before loading the data, it is important that its structure is correct to avoid misinterpretations.

- **Data Preparation in Microsoft Excel:**

- Before loading a Microsoft Excel file into an analysis, it is important that the data spreadsheet is free from irrelevant information and has a good structure to prevent misinterpretation. Possible actions that can be done before loading data are removing contextual information and combining columns into one.

13

Slide 14 TIBCO: Spotfire Cloud User's Guide Data Preparation in Microsoft Excel

# TIBCO: Spotfire Cloud User's Guide

## Data Preparation in Microsoft Excel

- **Data Preparation in Microsoft Excel:**

- The tabular format of the data in an Excel spreadsheet will be represented as a data table in your analysis. The first row with data in the spreadsheet will be interpreted as names of the data columns in the table, and the following rows will be interpreted as data rows.

- **Remove contextual information:**

- If your spreadsheet containing some contextual information above the actual data table, it needs to be removed otherwise this will cause misinterpretation of the data.

14

Slide 15 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 1

# TIBCO: Spotfire Cloud User's Guide

## Creating a Visualization 1

- **Creating a Visualization:**

- Data can be of many different kinds. Therefore, several types of visualizations are offered to present it in a relevant way. To create visualizations, either use **Recommended visualizations** or create the visualizations from scratch.
- **Using Recommended Visualizations**
  - When adding data and selecting data columns of interest, you are presented with recommended visualizations. Simply choose visualizations you find suitable, and use them as they are, or adjust them to suit your needs.
- **Creating Visualization from Scratch**
  - Another option is to create visualizations from the very beginning and make your own settings. To learn how to set up the different visualizations from scratch, see the next slides.

15

Slide 16 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 2

# TIBCO: Spotfire Cloud User's Guide

## Creating a Visualization 2

- **Creating a table:**  
The table visualization presents the details of the data. The individual values are arranged in columns and rows.
- **Creating a cross table:**  
Cross tables are used to summarize large amounts of data, and then present the result in a structured table format.
- **Creating a bar chart:**  
In a bar chart, you can compare values for different categories in your data.
- **Creating a line chart:**  
A line chart is used for showing trends, and in most cases trends over time. It can also be used for discerning certain patterns.
- **Creating a combination chart:**  
In a combination chart, you have the option to display both bars and lines in a single visualization. Because of the overlay effect, lines are drawn on top of the bars, it is easy to compare values for different columns or categories in your data. Trends can be identified, and you can spot deviations directly.
- **Creating a pie chart:**  
A pie chart is a circle graph that is divided into sectors. It is used to compare values for different categories in your data on a relative basis. Each pie sector represents a specific category, and its size the category's contribution to the whole value, expressed as a percentage. The values are usually sums.
- **Creating a scatter plot:**  
In a scatter plot, markers are displayed in a two-dimensional coordinate system. It is useful for getting an overview of how your data is distributed across two dimensions.

16

Slide 17 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 3

# TIBCO: Spotfire Cloud User's Guide

## Creating a Visualization 3

- **Creating a map chart with markers:**  
In a map chart with markers, data associated with locations is represented by markers placed in a geographical context.
- **Creating a map chart with colored regions:**  
In a map chart with colored regions, colors represent data associated with geographical regions. The regions, also called *features*, can be for example countries, provinces or postal code areas. The entire shape of a region is colored, and it is automatically placed in its geographical context.
- **Creating a treemap:**  
A treemap is used for displaying huge amounts of data that can be structured hierarchically (tree-structured). It presents the data using differently sized and colored rectangles.
- **Creating a parallel coordinate plot:**  
A parallel coordinate plot is used to compare data values which are of completely different types or magnitudes within a single visualization. The values are normalized and then presented as points on a line with one point per data column. The visualization is useful also for examining patterns.
- **Creating a histogram:**  
In a histogram, you can show the distribution of numerical data. The entire range of the numerical values is divided into equal intervals on the Category axis, and for each interval, it is indicated on the Value axis how many individual data values that fall within it.
- **Creating a trellised visualization:**  
A visualization that is trellised is split into a number of panels, where each panel represents a subset of the data. Using trellised visualizations, you can spot similarities and differences between the subsets of data, or within the subsets.

17

Slide 18 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Markers

# TIBCO: Spotfire Cloud User's Guide

## Creating a Map Chart with Markers

- **Creating a Map Chart with Markers:**
  - In a map chart with markers, data associated with locations is represented by markers placed in a geographical context.
    - **Note:** Geographical names are not unique. Names can be identical but their positions different. You can avoid misinterpretations by specifying a distinct geographical hierarchy, if possible. For example, above you have access to both Country and City columns. To prevent ambiguity, we recommend that you select not only City but both Country and City on the **Marker** axis.

18

Slide 19 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Colored Regions

# TIBCO: Spotfire Cloud User's Guide

## Creating a Map Chart with Colored Regions

- **Creating a Map Chart with Colored Regions:**

- In a map chart with colored regions, colors represent data associated with geographical regions. The regions, also called *features*, can be for example countries, provinces or postal code areas. The entire shape of a region is colored, and it is automatically placed in its geographical context.
  - The map chart consists of layers, where the background map forms one layer, and the colored regions (or features) representing the actual data form another layer. Navigation controls for zooming and panning are located to the right in the visualization. The geographical regions in your data are automatically positioned on top of the background map. Their colors reflect values from a certain data column, usually a numerical column containing, for example, sales figures. The values may represent aggregated data or not aggregated data for the particular region. Examples of aggregated values are sums and averages.

19

Slide 20 TIBCO: Spotfire Cloud User's Guide Exporting Your Analysis

# TIBCO: Spotfire Cloud User's Guide

## Exporting Your Analysis

- **Exporting Your Analysis:**

- You can export an analysis to different file formats to share it with others.
- When you want to share your discoveries in an analysis with others, you can export the entire analysis, or parts of it, as a PDF or Microsoft PowerPoint document.
  - **Exporting to PDF**  
You can export an entire analysis, or parts of it, to a PDF document.
  - **Exporting to Microsoft PowerPoint**  
You can export an entire analysis or parts of it to a Microsoft PowerPoint document.

20

Slide 21 TIBCO: Spotfire Cloud User's Guide Sharing Your Analysis

# TIBCO: Spotfire Cloud User's Guide

## Sharing Your Analysis

- **Sharing Your Analysis:**

- Your analysis can be shared with others. You can share it publicly, or share it with a selected group of people.
  - To share an analysis, you must first save it to a folder in the library. Then you make the folder available to others. The folders can be set as **Private**, **Public**, or **Shared**. A new folder can be created anytime, and initially it is **Private**. Analyses in a **Private** folder are available only to you.
  - Analyses in a **Public** folder are available to anyone on the internet. You can, for example, make them available through social media, or by posting a link on your website.
  - Analyses in a **Shared** folder are available to a selected group of people. The individuals in the group must have consumer licenses to get access to the folder.

21

Slide 22 TIBCO Spotfire: Education

<http://learn.spotfire.tibco.com/>

# TIBCO Spotfire: Education

<http://learn.spotfire.tibco.com/>

22

## Summit Highlights

An IoT from the Ontology Summit 2015

Slide 1 Opening Remarks 1

# Perspectives on Ontology Summit 2015

Mark Underwood and Michael Grüninger

Ontology Summit 2014 Symposium

April 13, 2015

Underwood and Grüninger ()

Ontology Summit Symposium

April 13, 2015

1 / 5

Slide 2 Opening Remarks 2

# What Have We Learned?

- IoT ontologies need to deal with dynamic time varying data vs. the often static Semantic Web. In particular, more work is needed on the development of event ontologies for targeted domains, building from core ontologies.
- We lack tools for a wide range of tasks, including for semantic annotation, ontology validation and integration.
- A more coordinated effort is required to compile IoT case studies which can serve as the basis for ontology reuse and the design of new ontologies. Key areas included Sensor integration, Smart Grid, and Smart Healthcare.

# Ontology Summit • NSF-NITRD • April 13, 2015

IoS

# Internet of Systems

Dr Shoumen Palit Austin Datta  
SVP, IIC • Research Affiliate, MIT

- [shoumen@mit.edu](mailto:shoumen@mit.edu) • [datta@iiconsortium.org](mailto:datta@iiconsortium.org) • [www.iiconsortium.org](http://www.iiconsortium.org) • <http://bit.ly/MIT-IOT> •

Slide 4 Dr. Shoumen Data 2

# MIT News

[Browse](#)

or

[Search](#)

## Michael Stonebraker wins \$1 million Turing Award

CSAIL researcher invented core database concepts, turned many into companies.

Adam Conner-Simons | CSAIL

March 25, 2015

▼ Press Inquiries

PRESS MENTIONS

Michael Stonebraker, a researcher at MIT's Computer Science and Artificial Intelligence Laboratory (CSAIL) who has revolutionized the field of database management systems (DBMSs) and founded multiple successful database companies, has won the Association for Computing Machinery's (ACM) **A.M. Turing Award**, often referred to as "the Nobel Prize of computing." This year marks the first time that the Turing Award comes with a Google-funded \$1 million prize.

In its announcement today, ACM said that Stonebraker "invented many of the concepts that are used in almost all modern database systems ... and founded numerous companies successfully commercializing his pioneering database technology work."

An adjunct professor of computer science and engineering at MIT and a principal investigator at CSAIL, Stonebraker sometimes jokes that he didn't know what he was researching for more than 30 years. "But then, out of nowhere, some marketing guys started talking about 'big data,'" he says. "That's when I realized that I'd been studying this thing for the better part of my academic life."

Stonebraker's work over the past four decades has helped spur a multibillion-dollar "big data" industry that he himself has participated in, creating and leading nine separate companies,

The ACM has awarded the A.M. Turing Award, widely regarded as the "Nobel Prize in Computing," to CSAIL researcher and adjunct professor Michael Stonebraker, reports Barb Darrow for *Fortune*. Stonebraker is "famous for arguing that database is not a one-size-fits-all category."

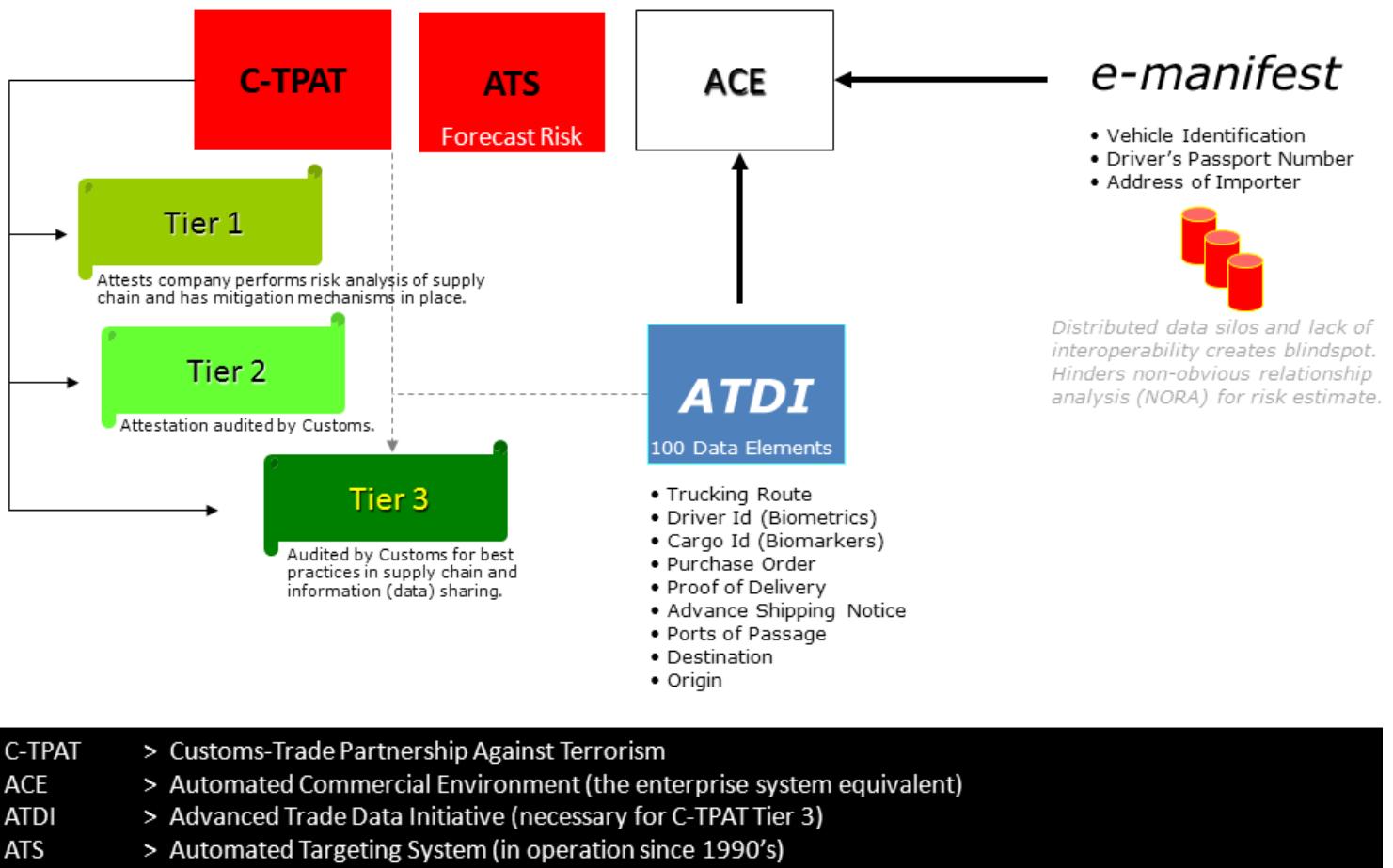
**FORTUNE**

Michael Stonebraker, a principal investigator at the MIT Computer Science and Artificial Intelligence Lab and an adjunct professor at MIT, has won the A.M. Turing Award for his work with database management systems, reports Nidhi Subbaraman for *BetaBoston*. "This is every computer scientist's lifetime dream, and it came true for me," said Stonebraker.

**BetaBoston**

Slide 5 Dr. Shoumen Data 3

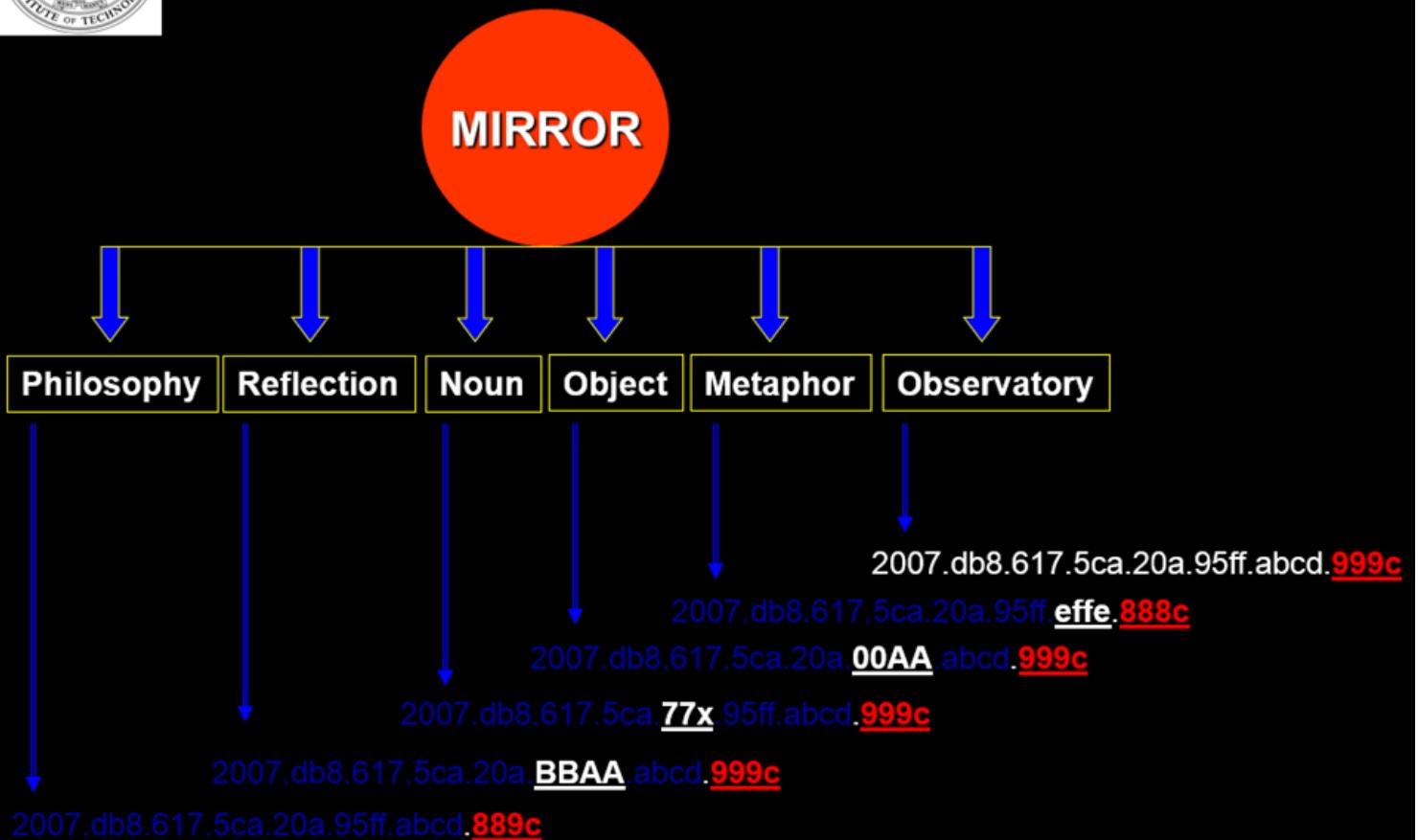
# Department of Homeland Security • Operation Safe Commerce



Slide 6 Dr. Shoumen Data 4



# Digital Semantics – Ontological UID with IPv6 type format



*This is an idea proposed by the author. It is not a fact or form of identification of ontologies, in practice, yet (2006).*

Dr Shoumen Datta <http://dft.ba/-shoumen> (fmr Research Director, Forum for Supply Chain Innovation, School of Engineering, MIT)

Slide 7 Track A 1

Ontology Summit 2015 Symposium:  
Internet of Things: Toward Smart Networked Systems  
and Societies  
National Science Foundation  
April 13-14, 2015

## Track A: Ontology Integration in the Internet of Things

Co-Champions:  
Leo Obrst (MITRE)  
Ram D. Sriram (NIST)

Slide 8 Track A 2

## **Summary: Methods for Integration and Interoperability**

- **Encode dynamic semantics**, which includes incorporating time semantics [Barnaghi]
- **Develop event ontologies** with above dynamic semantics [Barnaghi, Sriram]
- **Explore category theory**, which reveals deep connections between formal logic, computer science and theoretical physics. Approach would involve the following:
  - 1) Begin with two or more models;
  - 2) Identify the overlap between these;
  - 3) Map the overlap into each piece;
  - 4) Define aggregates semantically (using category theory); and
  - 5) Push out, [Breiner and Subrahmanian]
- **Use design patterns** toward ontology virtualization: Given a set of ontology design patterns and their combination into micro-ontologies, one can abstract the underlying axiomatization by dynamically reconfiguring patterns in a plug and play style; bridging between different patterns as micro-theories; providing ontological views and semantic shortcuts that suit particular provide, user, and use case needs by highlighting or hiding certain aspects of the underlying ontological model; and mapping between major modeling styles [Janowicz]
- **Expand on techniques presented by Ray and Hodges.**

13

Slide 9 Track B 1

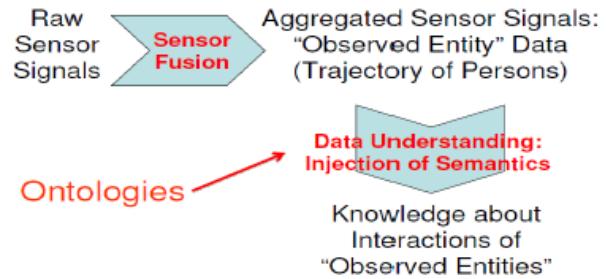
# Synthesis "Beyond Semantic Sensor Network Ontologies"

Ontology Summit 2015 Track B

Internet of Things: Toward Smart Networked Systems and Societies



## Raw Sensor Data to Knowledge



April 13, 2015

- Gary Berg-Cross (SOCoP) & Torsten Hahmann (U of Maine)

2015-04-06

1

Ontology Summit 2015 Track B

Slide 10 Track B 2

# Approach to Synthesis of “Beyond Semantic Sensor Network Ontologies”

The range of work springing from or leveraging SSNO is broad.

- We have leveraged insight from our 10 speakers and the community discussion of approaches, issues and problems.
- We have attempted to distilled the virtual meeting topics to a useful summary for the face-to-face Symposium.
- Our Synthesis is organized into several parts:
  1. Evolving SSN in order to place it into the larger IoT context, including interfacing it with other ontologies etc.
    1. Sensor-network interactions, services etc.
  2. Future visions including two directions for using more semantics in sensor networks
    1. Semantics in the cloud, including sensor registries
    2. More local semantics on the edge, including smart sensors & control entities

5

Slide 11 Keith Mazullo 1

# Cyber Physical Systems and IoT

Keith Marzullo

Division Director, Computer and Networking Systems

Directorate of Computer and Information  
Science and Engineering

National Science Foundation

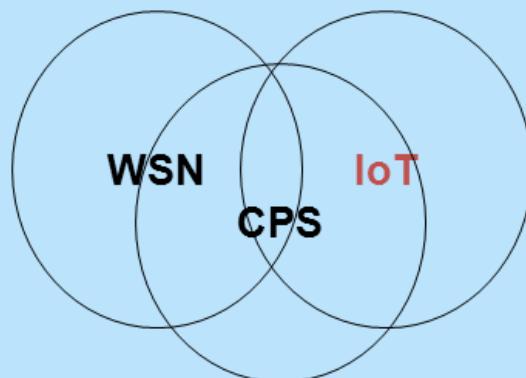
13 April 2015



1

Slide 12 Keith Mazullo 2

# IoT and CPS – Embrace the Synergies



**WSN – Wireless Sensor Networks**

**CPS – Cyber Physical Systems**

**IoT – Internet of Things**



- IoT & CPS share many core technology elements
- Open Programmable Devices and Objects are here
  - By 2020, 26B smart devices: lights, locks, security sensors ...
  - By 2019, 70% homes with IoT devices
- Industry standards promise simple IoT inter-operability
  - Examples: AllJoyn, Thread, SmartThings, Open Interconnect, ...
- Developers and IoT Apps are coming!
- 1,400 global developer survey – 40% working on IoT apps today

[Source: Evans Data Corporation]



Slide 13 Keith Mazullo 3

# Knowledge and Big Data

- IoT sensors will be generating potentially huge streams of data
- How will the data be utilized? How will it be converted into actionable information? How will it be fused with other streams?
- What algorithms will be used for analysis?
- Protocols and formats
- Will it be usable for real-time control?
- How and to where will it be transmitted? How will it be stored – if at all?
- What are the energy / bandwidth trades?
- Privacy / Security / Trust – more detail later



Slide 14 Keith Mazullo 4

# Mission Agency Partners

Core of foundational technologies with direct relevance to application areas of value to partners. Enables accelerated maturation and transition

- DHS S&T Directorate Cyber Security Division – 2014/15
- DOT FHWA & Intelligent Transportation System JPO - 2014/15
- NASA Aeronautics Research Mission Directorate - 2015
- NIH institutes including NIBIB (Imaging and Bioengineering), OBSSR (Behavioral and Social Sciences Research), National Cancer Institute (NCI), and National Center for Advancing Translational Sciences (NCATS) – 2015



53

Slide 15 Track C 1

# Ontology Summit 2015 Internet of Things: Toward Smart Networked Systems and Societies

Symposium – April 13, 2015  
Track C: Decision Making in Different Domains

Co-Champions  
Mike Bennett, Michael Grüninger

13 April 2015

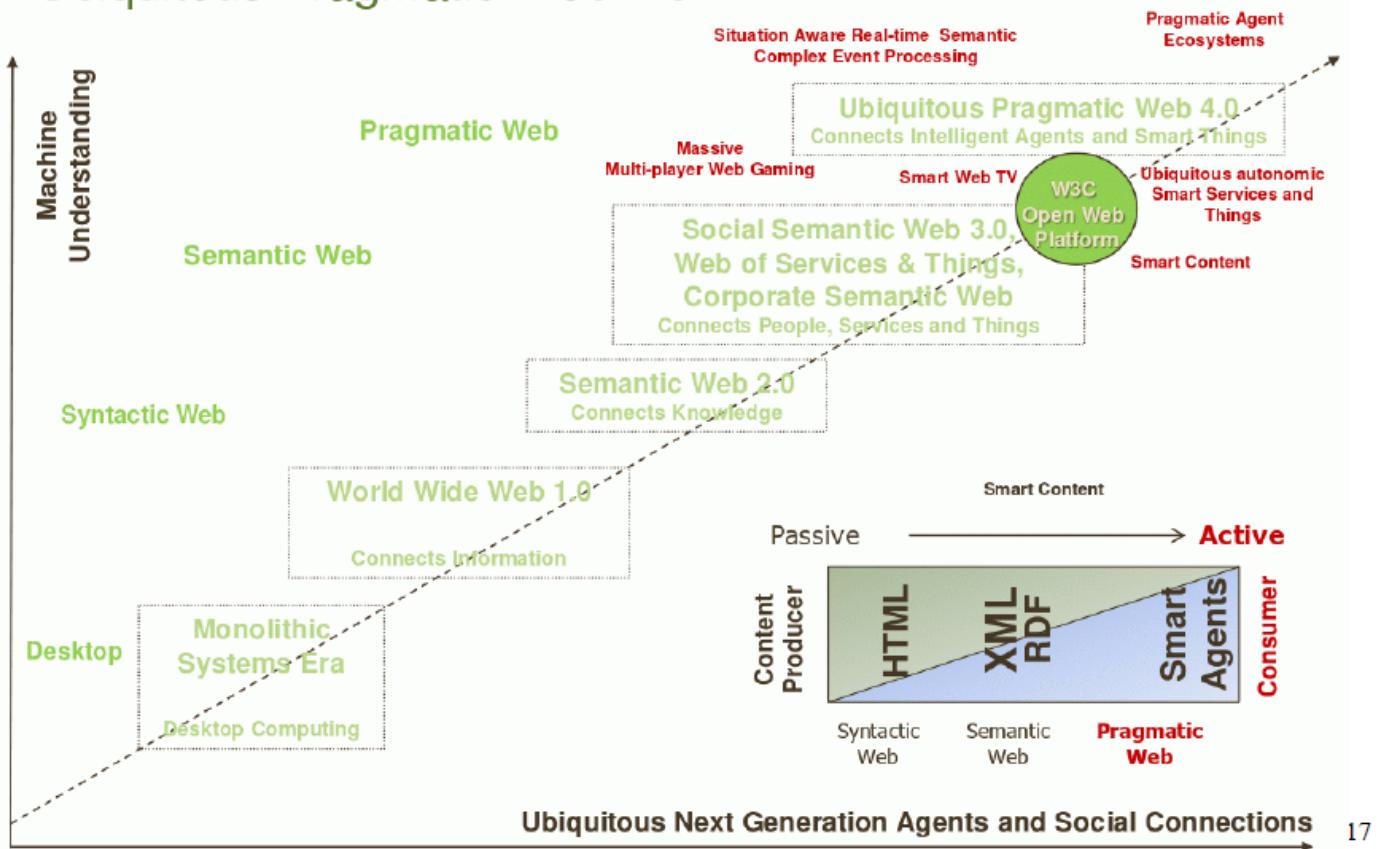
Ontology Summit 2015 Symposium

1

Slide 16 Track C 2

## Vision of the future of the Web (from Paschke)

### Ubiquitous Pragmatic Web 4.0



Slide 17 Track D 1

# Standards: What are Design Patterns for Ontologies in IoT?



MARK UNDERWOOD | KRYPTON BROTHERS

@KNOWLENGR @KRYPTONBROTHERS

CO-CHAIR NIST BIG DATA WORKING GROUP, SECURITY & PRIVACY SUBGROUP

13 APR 2015

TRACK D SYNTHESIS - FINAL - F2F MEETING

1

Slide 18 Track D 2

## Semantic Sensor Network Ontology

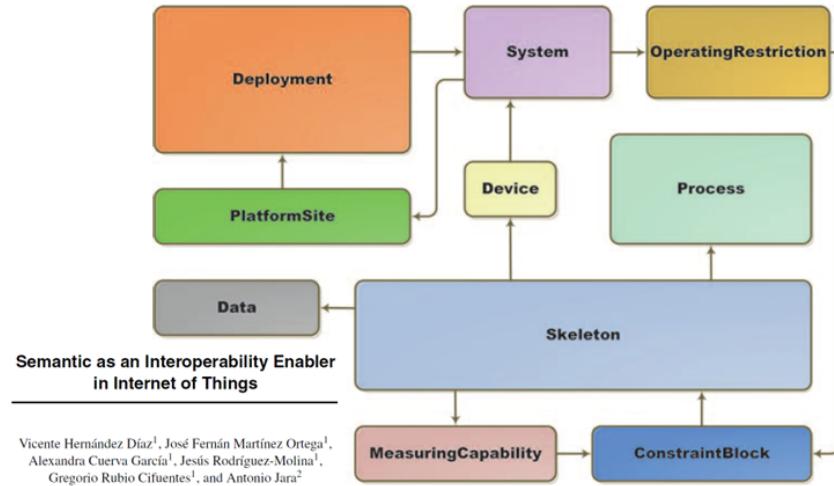


Fig. 9.4 Overview of the Semantic Sensor Network ontology classes and properties.

13 APR 2015

TRACK D SYNTHESIS - FINAL - F2F MEETING

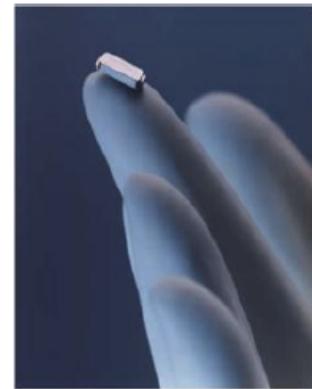
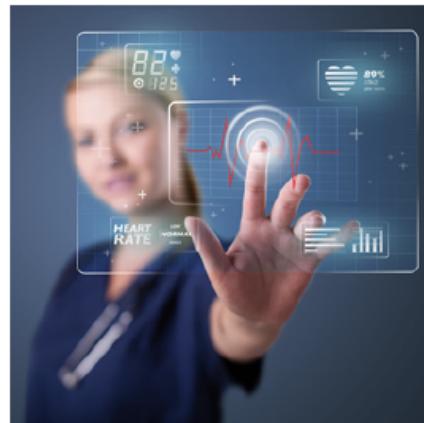
38

Slide 19 Dr. Bradford Hess 1

# Medicine and the “Internet of Things”



*Remote Liquid Biopsy*  
Singapore's A\*STAR Technologies



*Intraaortic Monitoring Device*  
Remon Medical Technologies

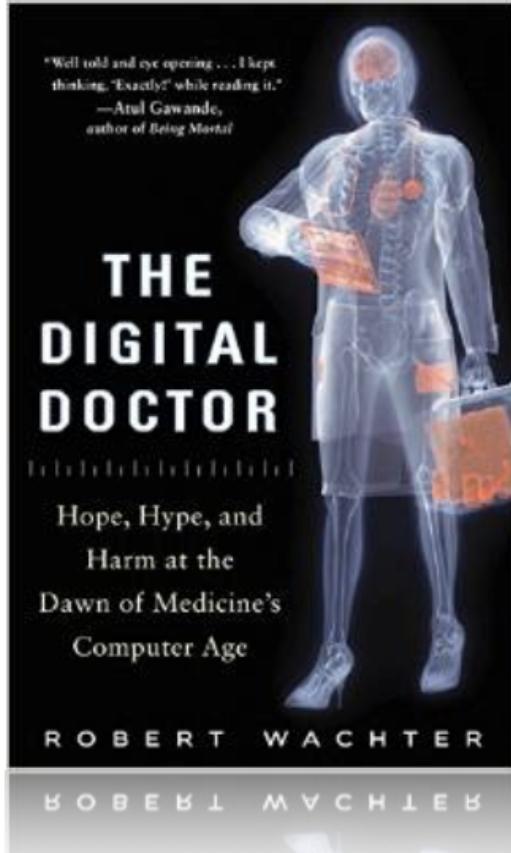
Bradford W. Hesse, PhD  
Chief, Health Communication and  
Informatics Research Branch



Slide 20 Dr. Bradford Hess 2



## Ultimately, the goal is to empower patients



“The real action — and the money — will shift to creating innovative tools to allow patients to stay healthy and manage chronic illness.”

- Robert Wachter

Slide 21 Dr. Harry Foxwell 1



**KEEP  
CALM  
AND  
CODE  
JAVA**

**Java Standards' role  
in the Internet of Things**

**ORACLE®**

## **Standards & Internet of Things: Oracle Perspective**

**Harry J Foxwell, PhD**

**Principal Consultant, Oracle Public Sector**

**harry.foxwell@oracle.com**

Slide 22 Dr. Harry Foxwell 2

BY THE YEAR 2020, THERE WILL BE

**50,000,000,000** connected devices,  
creating and sharing

**40,000,000,000,000 GB**

worth of data across the Internet of Things.

**IoT gains value through data...Big Data,  
which also needs ontology, standards,  
reference architecture(s), use cases, tools, ...**

ORACLE

Slide 23 Eric Simmon 1



# Ontology and Standards in the Internet of Things

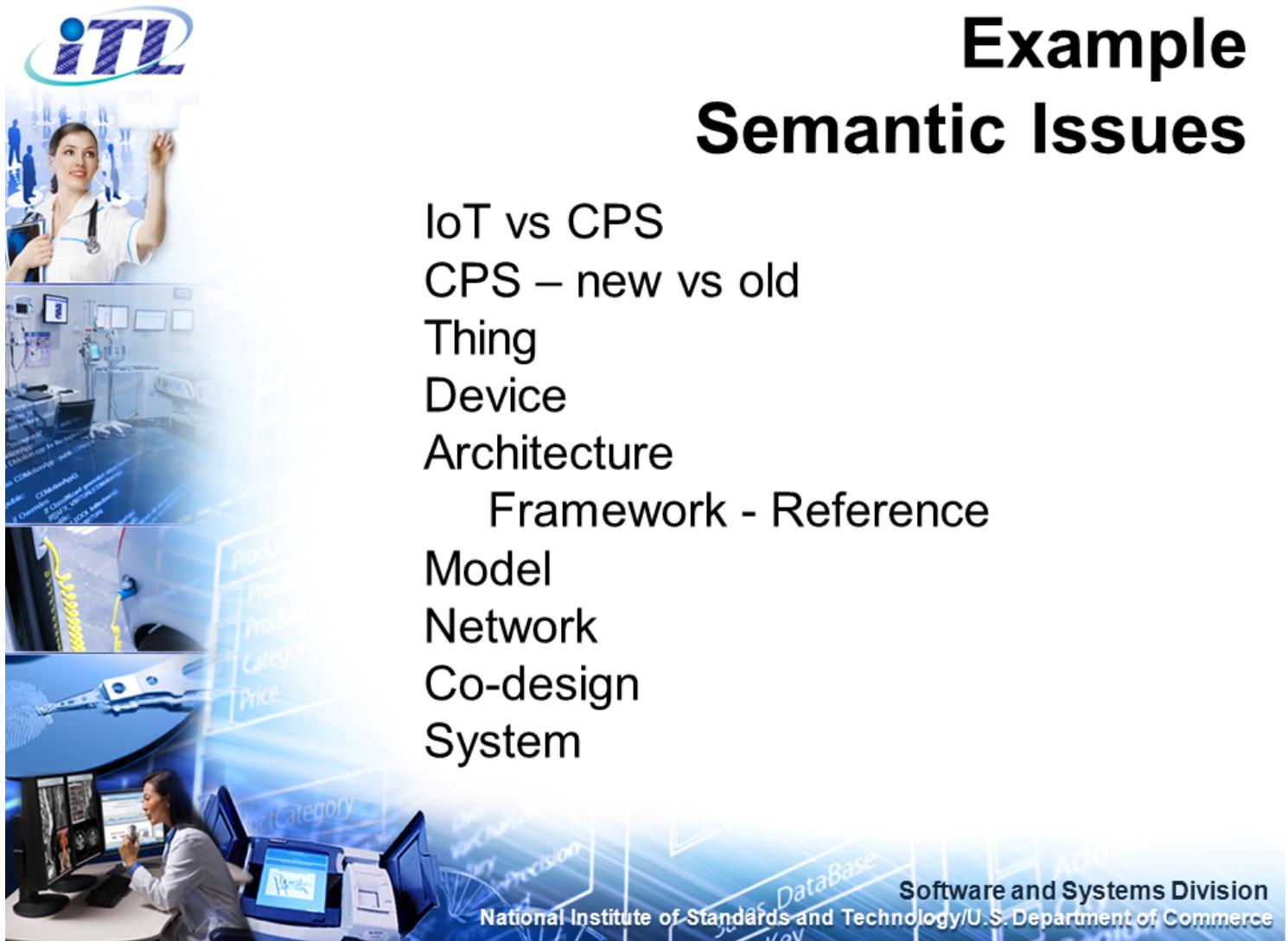
**Eric Simmon**

Systems and Software Division  
Information Technology Laboratory  
[eric.simmon@nist.gov](mailto:eric.simmon@nist.gov)

Software and Systems Division

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

Slide 24 Eric Simmon 2



Slide 25 William Miller 1



# Ontology Summit F2F Meeting

## ISO/IEC/IEEE P21451-1-4

### Sensei/IoT\* XMPP

1<sup>st</sup> International Semantic Web 3.0 Standard  
for the Internet of Things (IoT)  
William J. Miller  
Chairman

1

Slide 26 William Miller 2

## UNIVERAL UNIQUE IDENTIFICATION

- ISO/IEC/IEEE P21451-1-4 will use a JID (EUI-64) which is a Universal Unique IDentifier (UUID), defined in the draft ISO/IEC 29161 Automatic Identification for the Internet of Things developed by ISO/IEC/JTC1/SG31/WG6 Automatic Identification & Data Capture and ISO/IEC/TC122 Packaging and Internet of Things (IoT).
- jid = [ node "@" ] domain [ "/" resource {device} ]
- There are hundreds of ways to identify Things and ISO/IEC 29161 offers a unified approach.
- NOTE - EUI-64 is a IEEE SA 64-bit Global Identification.
- Example: 

12

---

## Slides

These slides were presented at the WhartonDC Innovation Summit 2015 and contain an example of the Internet of Things by Semantic Community for the [Federal Big Data Working Group Meetup](#).

[Slides](#)

---

### Slide 1 Data Science for Big Data

<http://semanticcommunity.info/>  
<http://whartondcinnovation.com/>

# Data Science for Big Data

Dr. Brand Niemann  
Director and Senior Data Scientist/Data Journalist  
Semantic Community  
<http://semanticcommunity.info/>  
<http://whartondcinnovation.com/>  
April 29, 2015

1

---

## Slide 2 Overview

# Overview

- Federal Big Data Working Group Meetup
- Silicon Valley to Washington
- First White House Data Chief Discusses His Top Priorities
- Precision Medicine and Natural Medicine
- Tech Meetup at White House
- USDA Data Science MOOC

2

---

## Slide 3 The Profit and Data Enterprises

[Federal Big Data Working Group Meetup](#) and [Eastern Foundry](#)

# The Profit and Data Enterprises

Marcus Lemonis (born November 16, 1973) is a Lebanese-born American businessman, investor, television personality and philanthropist. He is currently the chairman and CEO of Camping World and Good Sam Enterprises, and the star of The Profit, a CNBC reality show about saving small businesses through People, Process, and Products.

[http://en.wikipedia.org/wiki/Marcus\\_Lemonis](http://en.wikipedia.org/wiki/Marcus_Lemonis)

•The [Federal Big Data Working Group Meetup](#) is also about helping government agencies develop:

- **People – Data Scientists/Chief Data Officers**
- **Process – Data Infrastructure**
- **Products – Data Publications**

•Some examples:

- [EPA](#)
- [FDA](#)
- [NOAA](#)
- [HHS](#)
- [USDA](#)
- [Eastern Foundry](#)

•And provide MOOCs/Meetups for training and networking.

3

---

Slide 4 Federal Big Data Working Group Meetup

# Federal Big Data Working Group Meetup

- **Federal:** Supports the Federal Big Data Initiative, but not endorsed by the Federal Government or its Agencies;
- **Big Data:** Supports the Federal Digital Government Strategy which is "treating all content as data", so big data = all your content;
- **Working Group:** Data Science Teams composed of Federal Government and Non-Federal Government experts producing big data products; and
- **Meetup:** The world's largest network of local groups to revitalize local community and help people around the world self-organize like MOOCs (Massive Open On-line Courses) being considered by the White House.

4

---

Slide 5 Silicon Valley to Washington

# Silicon Valley to Washington

- Crafting Obama Administration Tech Policy:
  - Megan Smith, from Google Inc., to be U.S. Chief Technology Officer (CTO)
  - Alexander Macgillivray, from Twitter Inc., to be Deputy CTO
  - Tony Scott, from VMware, to U.S. Chief Information Officer (CIO)
  - Mikey Dickerson, from QSSI, Google, and Obama for America, to U.S. Digital Service Administrator
  - DJ Patil, from VP of Product at RelateIQ and the Data Scientist in Residence at Greylock Partners, to be Chief Data Scientist
  - David Portnoy, Aginity LLC, to HHS IDEA Lab Fellow, Datalytx, Inc. and Healthbox

5

---

## Slide 6 First White House Data Chief Discusses His Top Priorities

<http://www.scientificamerican.com/article/first-white-house-data-chief-discusses-his-top-priorities/>

# First White House Data Chief Discusses His Top Priorities

- At the top of my list right now is the Precision Medicine Initiative. Science has enabled us to unlock the [human genome](#). Now we want to combine that with the power of data science, which uses new techniques like machine learning as well as the explosion of data now available about individual patients, whether through their phones or other sensors in their environment. The challenge is putting this together to come up with new ways to think about health care and medical treatments.
  - **Semantic Medline and Natural Medicine for Disease and Wellness Meetup**
- My second priority is opening up more data and making it available for people [both the government and general public] to build an ecosystem of research, mobile apps and visualizations on top of that information.
  - **Semantic Community and Federal Big Data Working Group Meetup**
- The third main priority is inserting more data capacity into agencies throughout the government. We're seeing a rise of data scientists and chief data officers at the [National Institutes of Health](#) as well as within [the Department of] [Health and Human Services](#). The Commerce Department [announced its first chief data officer](#) [Ian Kalin] last week. We have to decide how to use the best of what we see in data science and statistics groups throughout the government to develop new services.
  - **Federal Big Data Working Group Meetup and Eastern Foundry**

<http://www.scientificamerican.com/article/first-white-house-data-chief-discusses-his-top-priorities/>

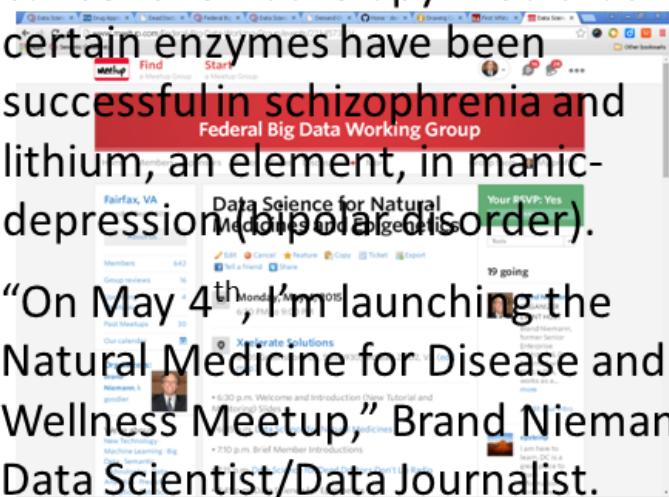
6

---

## Slide 7 Precision Medicine and Natural Medicine

# Precision Medicine and Natural Medicine

- “Tonight, I'm launching a new Precision Medicine Initiative to bring us closer to curing diseases like cancer and diabetes — and to give all of us access to the personalized information we need to keep ourselves and our families healthier.”
- — President Barack Obama, State of the Union Address, January 20, 2015

In early 2013, our Semantic Medline Data Science Team found that cancer immunotherapy was becoming more successful than cancer chemotherapy. Also that certain enzymes have been successful in schizophrenia and lithium, an element, in manic-depression (bipolar disorder).  


“On May 4<sup>th</sup>, I'm launching the Natural Medicine for Disease and Wellness Meetup,” Brand Niemann, Data Scientist/Data Journalist.

7

## Slide 8 Tech Meetup at White House

<https://www.whitehouse.gov/blog/2015/04/14/friday-tech-meetup-white-house-0>

# Tech Meetup at White House

The screenshot shows a Microsoft Internet Explorer browser window displaying the White House website. The main content is a blog post titled "This Friday: Tech Meetup at the White House". The post is dated April 14, 2015, and is written by Megan Smith and Jerry Abramson. It discusses the first-ever White House Tech Meetup, which will bring local leaders from across the country to the White House. The post includes a small image of a laptop screen with "#WHTMeetup" handwritten on it. To the right of the main content, there is a sidebar with links to "YOUR FEDERAL TAXPAYER RECEIPT" and "WHITE HOUSE BLOGS". Below the main content, there is a download bar with several files listed.

<https://www.whitehouse.gov/blog/2015/04/14/friday-tech-meetup-white-house-0>

8

## Slide 9 USDA Data Science MOOC

[http://semanticcommunity.info/Data\\_Sc...a\\_Science\\_MOOC](http://semanticcommunity.info/Data_Sc...a_Science_MOOC)  
<http://www.meetup.com/Federal-Big-Da...nts/221457264/>

# USDA Data Science MOOC

My Note: In Process for May 18<sup>th</sup> Meetup

The screenshot shows a web browser with two main windows. On the left is a page titled "USDA Data Science MOOC" listing various modules and resources. On the right is a "Web Player" interface for a "Spotfire" dashboard titled "Module 3 USDA Open Data Success Stories 1: AMS Spotfire". The dashboard includes a map of the United States with red dots representing farmers' markets, a bar chart titled "Distribution - SNAP" comparing two categories (red and green), and a table of data for "Geo U.S. Farmers Markets".

[http://semanticommunity.info/Data\\_Science/USDA\\_Data\\_Science\\_MOOC](http://semanticommunity.info/Data_Science/USDA_Data_Science_MOOC)  
<http://www.meetup.com/Federal-Big-Data-Working-Group/events/221457264/>

9

## Slide 10 Upcoming Meetups

<http://www.meetup.com/Federal-Big-Data-Working-Group>

# Upcoming Meetups

- President's Chief Data Scientist and EPA Big Data Analytics, April 20, 2015.
- The Wharton DC Alumni Innovation Summit, April 28-29, 2015.
- Natural Medicine for Disease and Wellness Meetup (New Meetup) Data Science for Natural Medicines and Epigenetics, May 4, 2015.
- USDA Data Science MOOC Meetup, May 18, 2015.
- Government Technology & Innovation Incubator for Big Data Analytics III, late May, 2015. DATA Act Challenge Cup.
- Data Science for Homeless Data: QlikView. Tableau, & Spotfire Bakeoff, June 1, 2015.
- Data Science for USGS Minerals Big Data, June 15, 2015.
- Data Science for Cyber Physical Systems-Internet of Things, June 29, 2015.

<http://www.meetup.com/Federal-Big-Data-Working-Group>

10

---

## Slide 11 Summary 1

[YouTube TIBCO How FAST Data Works](#)

1. Web of Documents
2. Web of Data (Semantic Web)
3. Internet of Things (People, Processes, & Products)

1. Data: Data (Statistical, Found, Confidential, & Classified) Information, Knowledge
2. Big Data: Structured & Unstructured (80%)
3. Fast Data: Real-Time Awareness from Big Data

Federal Big Data Working Group

Meetup:

1. Data Mining Data Science Data

Publication:

- How was the Data collected?
- Where is the Data stored?
- What are the Data results?
- Why should we believe the Data results?

2. Massive Open Online Courses (5 + USDA, etc.)

3. Boost employment and entrepreneurship (Like The Profit-Marcus Lemonis with Eastern Foundry)

## Data Mining Process Standard

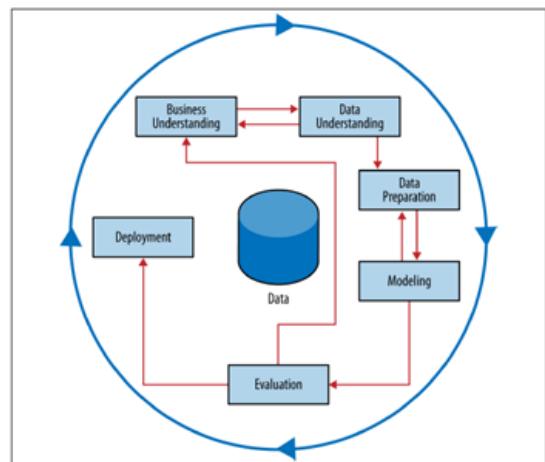
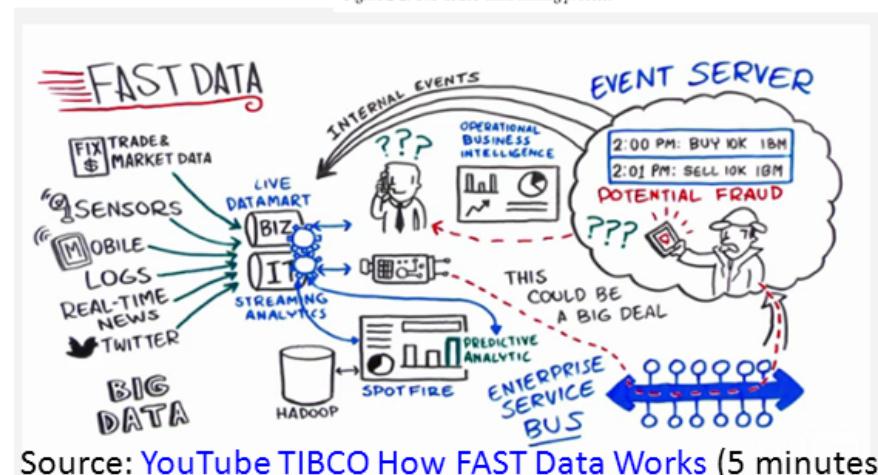


Figure 2-2. The CRISP data mining process.



Source: [YouTube TIBCO How FAST Data Works \(5 minutes\)](#)

11

## Slide 12 Summary 2

[Web Player](#) and [Spotfire Cloud Library](#)

## An Internet of Things: People, Processes, & Products

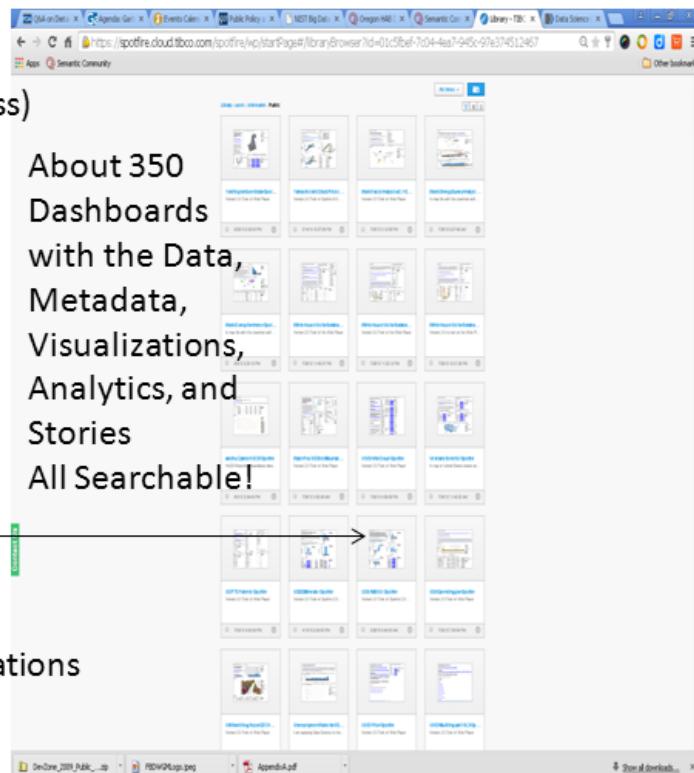
Data Science Data Mining Data Publication  
Part of USDA Data Science MOOC (in process)

### Module 3 USDA Open Data Success Stories 1: AMS Spotfire



Small Adjacent Dynamically Linked Visualizations

Desktop and [Web Player](#)



[Spotfire Cloud Library](#)

12

## Slides

These slides were prepared for a DHS briefing by Semantic Community to show the history of Semantic Communities DHS Information Sharing work and suggested proof of concept steps.

[Slides](#)

## Slide 1 Semantic Data Discovery: Proof of Concept for DHS

<http://semanticcommunity.info/>

<http://www.meetup.com/Virginia-Big-Data-Meetup/>

<http://www.meetup.com/Federal-Big-Data-Working-Group/>

<http://www.meetup.com/Northern-Virg...ic-Web-Meetup/>

[http://semanticcommunity.info/Data\\_Sc...g\\_Group\\_Meetup](http://semanticcommunity.info/Data_Sc...g_Group_Meetup)

# Semantic Data Discovery: Proof of Concept for DHS

Dr. Brand Niemann

Director and Senior Data Scientist/Data Journalist

Semantic Community

<http://semanticcommunity.info/>

<http://www.meetup.com/Virginia-Big-Data-Meetup/>

<http://www.meetup.com/Federal-Big-Data-Working-Group/>

<http://www.meetup.com/Northern-Virginia-Semantic-Web-Meetup/>

[http://semanticcommunity.info/Data\\_Science/Federal\\_Big\\_Data\\_Working\\_Group\\_Meetup](http://semanticcommunity.info/Data_Science/Federal_Big_Data_Working_Group_Meetup)

March 25, 2015

1

---

## Slide 2 Information Sharing at DHS

# Information Sharing at DHS

- NIEM (Michael Daconta):
  - XML Messages
- SOA (Wolf Tombe):
  - XML Messages and XML Data in an ESB
- Semantic Ontology (Barry Smith):
  - RDF/OWL UCore
- Semantic Knowledge Bases (Brand Niemann):
  - NIEM and NIEM and Thetus Savana
- Semantic Search Data Browser (Brand Niemann)
  - Global Terrorism Database Experience
- Semantic Quint Dynamic Case Management (Brand Niemann):
  - XML/RDF/OWL/RML in Be Informed

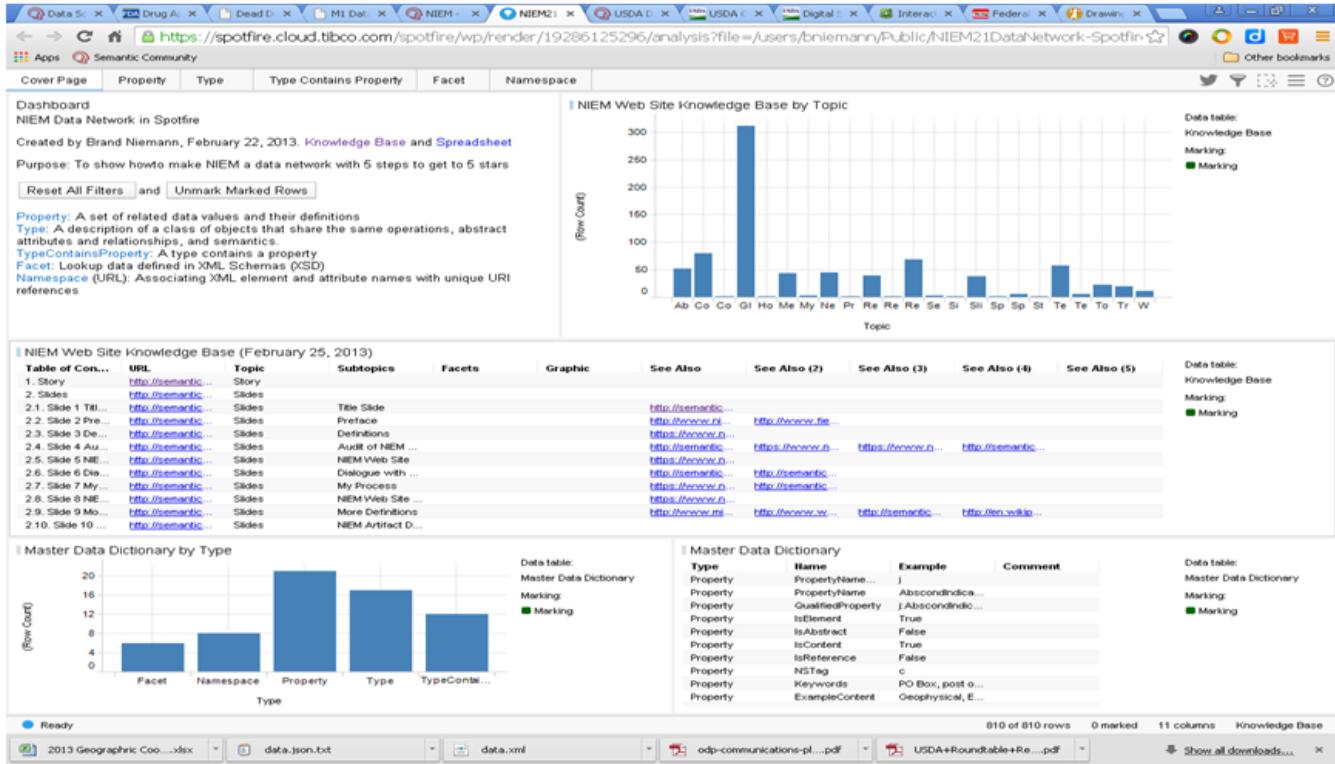
2

---

## Slide 3 NIEM as Big Data in a Network with Data Science

MindTouch Knowledge Base: [NIEM](#) and Spotfire Dashboard: [Web Player](#)

# NIEM as Big Data in a Network with Data Science



MindTouch Knowledge Base: [NIEM](#)

Spotfire Dashboard: [Web Player](#)

3

## Slide 4 NIEM 3.0 Alpha 2 Release and Thetus Savanna Review

MindTouch Knowledge Base: [NIEM and Thetus](#) and Spotfire Dashboard: [Web Player](#)

# NIEM 3.0 Alpha 2 Release and Thetus Savanna Review

**Key Questions!**

How can one review NIEM 3.0 Alpha 2 without some data science analytics?

Does Thetus Savanna do what NIEM is ultimately trying to accomplish without NIEM 3.0?

The dashboard includes the following components:

- Knowledge Base and Spreadsheet:** A table showing the consolidation of NIEM 3.0 across various categories like Core, Activity, and Document.
- Bar Charts:** Two bar charts showing the count of items by type ('Type') and comment ('Comment').
- Table:** A table titled 'Inventory of NIEM' listing various tabs and their corresponding fields.
- Data Table:** A table titled 'Consolidation of NIEM 3.0' showing the relationship between categories, types, and definitions.
- File List:** A list of files including '2013 Geographic Co...xlsx', 'data.json.txt', 'data.xml', 'odp-communications-pl.pdf', and 'USDA+RoundTable+Re....pdf'.

MindTouch Knowledge Base: [NIEM and Thetus](#)

Spotfire Dashboard: [Web Player](#)

4

## Slide 5 NIEM and UCore 2.0 Semantic Layer for Information Sharing

MindTouch Knowledge Base: [Universal Core Semantic Layer](#) and Spotfire Dashboard: [Web Player](#)

# NIEM and UCore 2.0 Semantic Layer for Information Sharing

**NIEM and UCore 2.0 Semantic Layer for Information Sharing**

**Web Sites:** <http://www.niem.gov/> and <https://ucore.gov/>

**Wiki Pages:** [http://semanticcommunity.info/National\\_Information\\_Exchange\\_Model](http://semanticcommunity.info/National_Information_Exchange_Model) and [http://semanticcommunity.info/Universal\\_Core\\_Semantic\\_Layer](http://semanticcommunity.info/Universal_Core_Semantic_Layer)

**Created by Brand Niemann, August 22, 2010 (in process).**

**Contents:**

- UCORE 2.0 Semantic Layer: Axioms
- UCORE Semantic Layer: Properties and Classes
- UCORE 2.0 Semantic Layer: OWL-DL Files
- NIEM 2.1 Spreadsheet (in process)
- NIEM 2.1 Developers Spreadsheet

**Table**

**Column 1**

**Description**

**Ver...**

**Release**

**Format**

**Data table:**

**UCORE2.0: Taxonomy**

**ucability uc:event**

**uc:Cargo uc:AlertEvent**

**uc:CollectionO... uc:Communica...**

**uc:CyberAgent uc:CriminalEvent**

**uc:Document uc:CyberSpat...**

**uc:Environment uc:DisasterEv...**

**uc:Equipment uc:EconomicE...**

**uc:Facility uc:Emergency...**

**uc:Financial... uc:Execution...**

**uc:GroupOrG... uc:ExerciseEv...**

**uc:GroupOrOr... uc:FinancialEv...**

**uc:Information uc:Hazardous...**

**uc:Infrastructure uc:Humanitaria...**

**uc:LivingThing uc:Instructiu...**

**ucore:SL\_Taxonomy-OWL-DL\_1\_0\_06-30-2009.owl**

**UCore-SL\_Relations-OWL-DL\_1\_0\_06-30-2009.owl**

**UCore-SL\_Axioms\_1\_0\_08-26-2009.xlsx**

**UCore-SL\_Class\_Definitions\_1\_0\_07-31-2009.xlsx**

**UCore-SL\_Relations-OWL-1\_0\_06-30-2009.owl**

**UCORE2.0: SL Relations**

**ucore:Re... rdfs:subClassOf**

**ucore:Affiliate... srl:affiliated\_with**

**ucore:CauseOf srl:cause\_of**

**ucore:Controls srl:controls**

**ucore:DistructF... owt:differentFr...**

**ucore:Employee... srl:employed\_by**

**ucore:HasDest... srl:has\_destin...**

**ucore:HasFunct... srl:has\_funct...**

**ucore:HasOrga... srl:has\_orga\_of**

**ucore:InvolvedIn srl:involved\_in**

**ucore:LocatedAt srl:located\_at**

**ucore:OccursAt srl:occurs\_at**

**ucore:SameAs owt:sameAs**

**ucore:Subordi... srl:subordinate...**

**UCORE2.0: SL Relations**

**rdfs:subClassOf**

**srl:affiliated\_with**

**srl:cause\_of**

**srl:controls**

**owt:differentFr...**

**srl:employed\_by**

**srl:has\_destin...**

**srl:has\_funct...**

**srl:has\_orga\_of**

**srl:involved\_in**

**srl:located\_at**

**srl:occurs\_at**

**owt:sameAs**

**srl:subordinate...**

**stability stEvent**

**st:InformationC... st:Act**

**st:Analysis st:ActionComm...**

**st:Objective st:ActionFunda...**

**st:ObjectiveSp... st:ActionObse...**

**st:Opinion st:CriminalAct**

**st:Plan st:Immigration...**

**st:TaskSpecifi... st:LawEnfor...**

**st:TimePointEntity st:TerroristAct**

**st:Agent st:PhysicalSpace...**

**st:Artifact st:Danger**

**st:ArtificialAg... st:Disaster**

**st:Equipment st:Equipment**

**st:Facility st:FinancialEvent**

**st:Sensor st:Environment...**

**Data table:**

**UCore-SL\_File\_Descri...**

**Ver... 1 6/30/2009 OML**

**Format OML**

**Marking:**

**Data table:**

**UCORE2.0: SL\_Taxo...**

**Ver... 1 6/30/2009 OML**

**Format OML**

**Marking:**

**Data table:**

**UCORE2.0: SL Relations**

**Ver... 1 6/30/2009 OML**

**Format OML**

**Marking:**

**Show all downloads...**

MindTouch Knowledge Base: [Universal Core Semantic Layer](#)

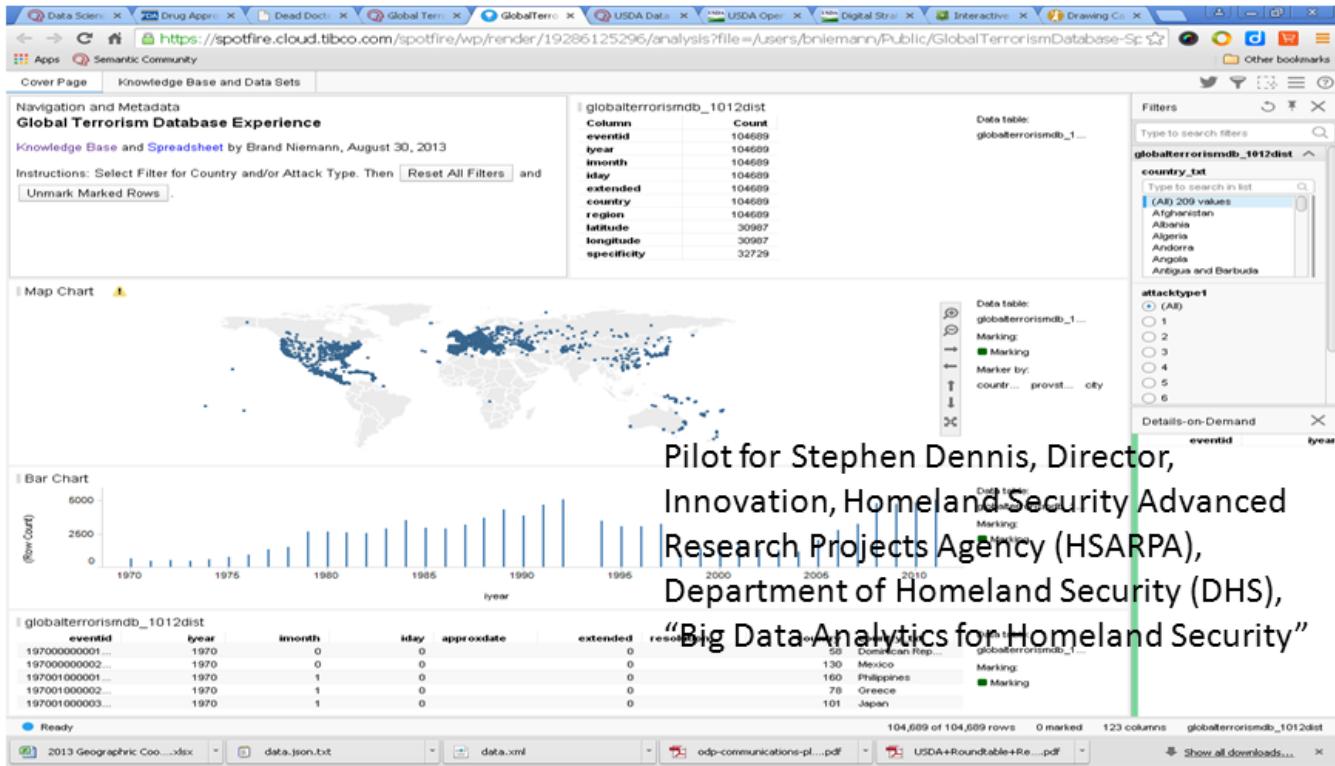
Spotfire Dashboard: [Web Player](#)

5

## Slide 6 Global Terrorism Database Experience

MindTouch Knowledge Base: [Global Terrorism Database](#) and Spotfire Dashboard: [Web Player](#)

# Global Terrorism Database Experience



MindTouch Knowledge Base: [Global Terrorism Database](#)

Spotfire Dashboard: [Web Player](#)

6

Slide 7 A Quint for Cross Information Sharing and Integration in the Intelligence Community

MindTouch Knowledge Base: [A Quint-Cross Information Sharing and Integration](#) and Spotfire Dashboard: [Web Player](#)

# A Quint for Cross Information Sharing and Integration in the Intelligence Community

The screenshot shows a complex data visualization dashboard titled "A Quint - Cross Information Sharing and Integration App". The interface includes:

- Log Table:** Shows a list of entries with columns: Section, URL, Comment, Comment (2), URL (2), Comment (3), and URL (3). The log includes items like "1. Story", "2. Spotfire Das...", "3. Top Secret ...", etc.
- World Map:** A world map where cities are plotted as colored dots, categorized by country. A legend on the right maps colors to countries: Afghanistan (red), Albania (blue), Algeria (green), American Samoa (yellow), Andorra (orange), Angola (purple), and Antigua & Barbuda (pink).
- Filters:** A sidebar with two main sections: "World Cities" and "Country". "World Cities" lists 3325 values including "Adan", "Abdo", "S-Gravenhage", "A Coruna", "Aachen", and "Aba". "Country" lists 200 values including "Afghanistan", "Albania", "Algeria", "American Samoa", "Andorra", and "Angola".
- Navigation:** Includes tabs for "CIA Site Map & World Fact Book Countries in Regions", "CIA Fact Book Appendices A-D", and "NCOIC Geospatial Interoperability Team Working Files".
- Downloads:** A bottom bar showing several PDF files available for download, such as "2013 Geographic Co...xlsx", "data.json.txt", "data.xml", "odp-communications-pl...pdf", and "USDA+RoundTable+Re...pdf".

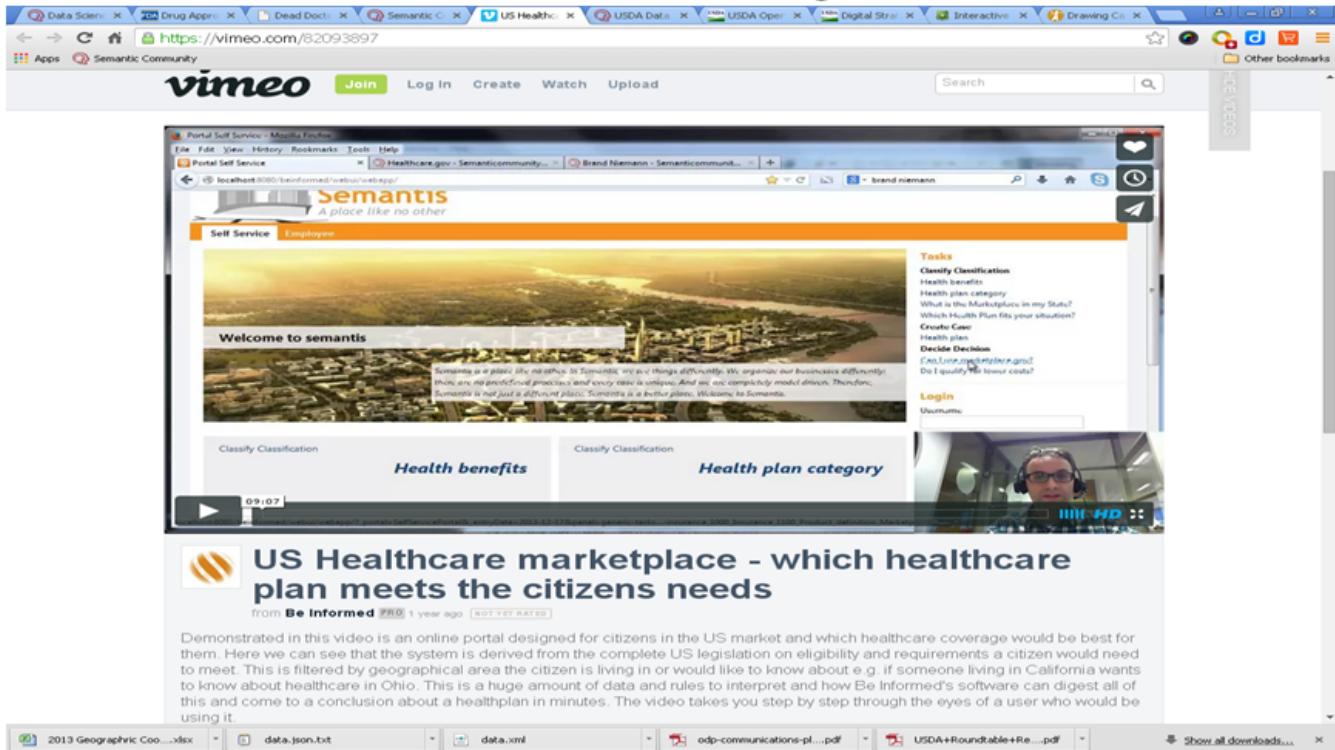
MindTouch Knowledge Base: [A Quint-Cross Information Sharing and Integration](#) Spotfire Dashboard: [Web Player](#)

7

## Slide 8 Dynamic Case Management Pilot for Healthcare.gov

MindTouch Knowledge Base: [Healthcare.gov Data Science](#) and Video Demo: [Vimeo](#)

# Dynamic Case Management Pilot for Healthcare.gov



MindTouch Knowledge Base: [Healthcare.gov Data Science](#)

Video Demo: [Vimeo](#)

8

## Slide 9 Proof of Concept Steps

# Proof of Concept Steps

- Introduction:
  - Best practice system example with data dictionary (DD) and Application Programming Interface (API)
  - Challenge all systems to do that and submit for internal web page
- Phase I:
  - All systems gain experience using internal web page for improved information sharing
  - Initial work on semantic harmonization for multiple DDs
- Phase II:
  - Develop various semantic harmonization methods and tools (e.g. ontology)
  - Pilot those methods and tools (Dynamic Case Management – Be Informed?)
- Phase III:
  - Develop requirements for improved information sharing system based on Phase I and II experience
  - Release RFI for RFQ

9

---

## Slide 10 Semantic Community

<http://semanticcommunity.info/>

<http://www.meetup.com/Federal-Big-Data-Working-Group>

# Semantic Community

- Former Senior Enterprise Architect & Data Scientist with the US EPA.
- Led Federal CIO Council Web Services, SOA (with Mitre), Semantic Interoperability, and Semantic Community Work.
- Founded and Co-organize the Federal Big Data Working Group Meetup to Continue the Above as a Private Citizen.
- Helping Government Agencies (US, Europe, and Japanese) Develop Data Scientists/Chief Data Officers, Data Infrastructure, and Data Publications.
  - <http://semanticcommunity.info/>
- Providing MOOCS/Meetups for Training and Networking.
  - <http://www.meetup.com/Federal-Big-Data-Working-Group>

10

---

## Spotfire Dashboard

For Internet Explorer Users and Those Wanting Full Screen Display Use: [Web Player](#) Get [Spotfire iPad App](#)

[Cover Page](#)[Spotfire Silver to Spotfire Cloud Migration](#)[Spotfire Learning Resources](#)

## Navigation and Metadata

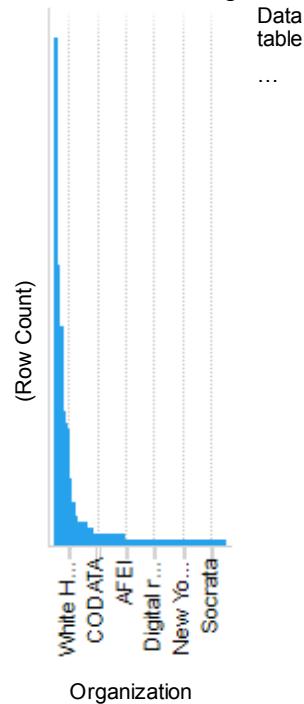
### Data Science for Big Data: Internet of Things by Brand Niemann, April 29, 2015.

[Spreadsheet](#)

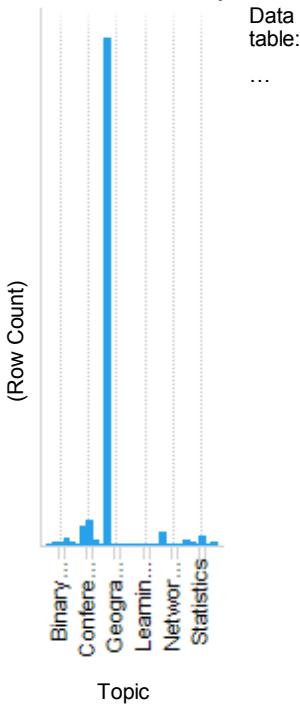
Instructions: Use Filter to Right to Select Category, and/or Comment. Then

[Reset All Filters](#) and[Unmark Marked Rows](#)

#### Distribution – Organization



#### Distribution – Topic



#### Filters

 Type to search filters

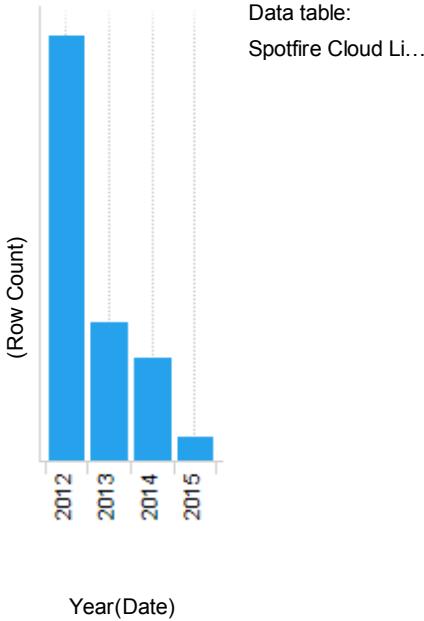
#### Spotfire Cloud Library

##### Name

 Type to search in list

- (All) 343 values
- 2008BSAPartDEventsPUF...
- 2008BSAPartDEventsPUF...
- 2010Census-Spotfire
- 2010OEISymposium-Spotfire
- 2010ReporttoCongression...
- 2010StatAbstract-Spotfire

#### Histogram – Date



#### Spotfire Cloud Library

Name	Data table:
2008BSAPart...	Spotfire Cloud Library
2008BSAPart...	
2010Census-...	
2010OEISymp...	
2010Reportto...	
2010StatAbstr...	
2011DataCent...	
2012Analytics...	
2012Presidenti...	
3RoundStones...	
AFEI-Spotfire	
AHRQFocuso...	
AirForceOneS...	
AirlineIncident...	
AlionScienceW...	
AmazonGover...	
AmericasData...	

#### Details-on-Demand

Name	Organization

● Ready

343 of 343 rows 0 marked 6 columns Spotfire Cloud Library

## Research Notes

### Looking for Big Data

Tim, Excellent email for our work with MOOCs and upcoming Meetup on June 29th. Many thanks, Brand

NAS Workshop Report "Training Students to Extract Value from Big Data: Summary of a Workshop"

Hello, The National Academy of Sciences (NAS) published a report titled "Training Students to Extract Value from Big Data: Summary of a Workshop" in 2014. The workshop was a followup on the NAS report "Frontiers in Massive Data Analysis". It is

available for free download at <http://www.nap.edu/catalog/18981/tra...ata-summary-of>.

The same NAS. One year later. A different group of people. A very different report. This report has actionable information for the NBD-PWG in several subject areas. The panels and participants openly discuss their experiences and plans and aspirations for big data. It is worth reading, except the conclusions.

Traditional statisticians approach big data using their traditional methods based on carefully controlled data sampling, questions known far in advance, assumed statistical model, unique statistically valid answers. Their methods work well for insurance companies, mass production factories, census calculations, epidemiology, classical controlled experiments, social psychology, and much more. Traditional statisticians attempt to generalize from sampled experiences (e.g., insurance claim rates, product failure rates). The Central Limit Theorem applies.

The fourth paradigm emphasizes, relying on a surfeit of data, discovery of patterns within the data w/o precisely framing a priori questions steering the analysis, reliance on multiple data models competing to best match the data. The fourth paradigm attempts to individualize, personalize by connecting the dots into complex models for each person, vehicle, residence, machine. Complex systems and their synthetic models co-evolve.

The workshop turned up many important big data projects and ideas. The academics usually noted the lack of big data sets to work and teach with.

(NIST is very fortunate to have 6 independent data sets covering 6 different types of data.) In their conclusions, it is apparent that the experts and teachers are still learning how best to teach students to work with big data.

Pilots start learning theory of flight, meteorology, radio, navigation, etc.

Pilots progress to light single engine planes under VFR, then IFR. Pilots can progress to multi engine planes, heavy transports, jets, etc. Perhaps, a gradual progression is the best approach to teaching data scientists, programmers, operators, curators, etc. There are many different career paths. Big data expertise has a short half life, about 1 year today. I favor the MOOCs today.

Recommended reading, Tim Zimmerlin

## Ontology Summit 2015 Agenda

Source: <http://ontolog.cim3.net/OntologySumm.../schedule.html>

### Monday, April 13

Time Activity

08:30 Breakfast

09:00 Welcome, Introductions, Opening Remarks

Hosts: [Michael Gruninger](#), [Mark Underwood](#), [Ram Sriram](#), [Leo Obrst](#), [George Strawn](#)

Keynote Speaker: [Dr. Shoumen Palit Austin Datta](#)

09:30 Senior Vice President

[Industrial Internet Consortium](#)

10:30 Break

11:00 Summary Report for Track A---Ontology Integration in IoT

Champions: [Ram Sriram](#), [Leo Obrst](#)

11:45 Summary Report for Track B---Beyond Semantic Sensor Network Ontologies

Champions: [Gary Berg-Cross](#), [Torsten Hahmann](#)

12:30 Lunch

Keynote Speaker: [Dr. Keith Marzullo](#)

14:00 Director, [Computer and Networks Systems Division](#)

[Directorate for Computer & Information Science & Engineering](#)

[National Science Foundation](#)

15:00 Summary Report for Track C---Decision Making in Different Domains

Champions: [Mike Bennett](#), [Michael Gruninger](#), [Ken Baclawski](#)

15:45 Summary Report for Track D---Related Standards And Synergies for Emerging IoT Ontologies

Champion: [Mark Underwood](#)

16:30 Open Discussion

18:30 Group Dinner

---

## Tuesday, April 14

Time Activity

08:30 Breakfast

Keynote Speaker: [Dr. Bradford Hesse](#)

Chief, [Health Communication and Informatics Research Branch](#)

09:00 [Behavioral Research Program](#)

[National Cancer Institute](#)

[National Institutes of Health](#)

10:00 Presentation of the 2015 Ontology Summit Communique

10:45 Affirmation of Communique

11:00 Break

Demo Session

11:15 Champion: [Mark Underwood](#)

12:45 Lunch

Panel Discussion: Ontology and Standardization

[Prof. Eswaran Subrahmanian](#), Carnegie Mellon University (moderator)

Harry Foxwell, Principle Consultant, Oracle

Elaine Newton, Deputy Standards Liaison, ITL at NIST

[Eric Simmons](#), NIST

Coulin Soutar, Technology Executive, Deloitte

[Mark Underwood](#), Krypton Brothers

14:00 The topic of Ontology and Standardization is interesting and timely because the process of interlinking devices, people, and creating an Internet of Things is a challenge of interoperability and compositionality. Information technology standards are all around us but the problem of interoperability persists from hospitals to defense systems. Many have advocated the use of ontologies as a means to mediate between and among devices similar to the efforts in domains such as gene-ontology, parts of medicine, and other domains. However, we need to understand what role ontologies are expected to play, how are standards in information technology created, and how they may shed light on the approaches to addressing the issue of standardization in ontology. Does standardization take place in fragments and is it possible to be put together? If not, what approaches could overcome some of the consequences?

15:30 Conclusion and Next Steps

16:00 Symposium Adjourns

---

## Ontology Summit 2015 Background

Source: [http://ontolog-02.cim3.net/wiki/Onto...2015\\_Symposium](http://ontolog-02.cim3.net/wiki/Onto...2015_Symposium)

- Dates: Monday 2015-04-13 and Tuesday 2015-04-14
  - Summit Theme: Internet of Things: Toward Smart Networked Systems and Societies
  - Session Topic: Ontology Summit 2015 Symposium
  - Ontology Summit Co-chairs: [MarkUnderwood](#) and [MichaelGruninger](#)
  - Symposium Co-chairs: [RamSriram](#) and [LeoObrst](#)
  - Schedule: See above
- 

## Prepared Presentation Materials

Slides can be accessed by clicking on each of the following:

- Chair opening: [Mark Underwood and Michael Grüniger: Introduction](#)

- Keynote: [Shoumen Datta Part 1](#) ; [Shoumen Datta Part 2](#) ;[Shoumen Datta Part 3](#)
- [Leo Obrst and Ram Sriram: Track A Synthesis](#)
- [Gary Berg-Cross and Torsten Hahmann: Track B Synthesis](#)
- [Ken Baclawski, Michael Grüninger and Mike Bennett: Track C Synthesis](#)
- Keynote: [Keith Marzullo: CPS-IoT](#)
- [Mark Underwood: Track D Synthesis](#)
- Keynote: [BrafordHesse: Medicine and the Internet of Things](#)
- Demo: [William Miller: Standard for the Internet of Things](#)
- Panelist: [Eric Simmon: Ontology and Standards in the Internet of Things](#)
- Panelist: [Harry Foxwell: Standards & Internet of Things: Oracle Perspective](#)

## Audio Recordings

- [Morning Session on Monday, 13 April 2015](#)
- [Afternoon Session on Monday, 13 April 2015](#)
- [Morning Session on Tuesday, 14 April 2015](#)
- [Afternoon Session on Tuesday, 14 April 2015](#)

## Additional Resources

- [Commenque](#)

## Abstract

The OntologySummit is an annual series of events (first started by Ontolog and NIST in 2006) that involves the ontology community and communities related to each year's theme chosen for the summit. The Ontology Summit program is now co-organized by Ontolog, NIST, NCOR, NCBO, IAOA, NCO\_NITRD along with the co-sponsorship of other organizations that are supportive of the Summit goals and objectives.

We are witnessing a new revolution in computing and communication. The Internet, which has spanned several networks in a wide variety of domains, is having a significant impact on every aspect of our lives. The next generation of networks will utilize a wide variety of resources with significant sensing capabilities. Such networks will extend beyond physically linked computers to include multimodal information from biological, cognitive, semantic, and social networks. This paradigm shift will involve symbiotic networks of people, intelligent devices, and mobile personal computing and communication devices (mPCDs), which will form net-centric societies or smart networked systems and societies (SNSS). mPCDs are already equipped with a myriad of sensors, with regular updates of additional sensing capabilities. Additionally, we are witnessing the emergence of "intelligent devices," such as smart meters, smart cars, etc., with considerable sensing and networking capabilities. Hence, these devices – and the network -- will be constantly sensing, monitoring, and interpreting the environment – this is sometimes referred to as the Internet of Things. And as local and wide area networks became almost secondary to the WWW (World-Wide Web), users and their usage patterns will become increasingly visible. This will have significant implications for both the market for advanced computing and communication infrastructure and the future markets – for nearly 4.5 billion people -- that net-centric societies will create

Well-designed and constructed net-centric societies will result in better quality of life, reduced threat from external sources, and improved commerce. For example, assume a scenario where people at various locations suffer from flu-like symptoms. In a net-centric society, mPCDs will send vital signs and other associated information to appropriate laboratories and medical centers. These centers will analyze the information, including searching the Internet for potential solutions, and will aid in determining possible causes for this phenomenon. Based on the diagnosis, people will be directed to the nearest clinic for treatment. Here we have several types of information flowing through the net: data from mPCDs; location information; images; video; audio; etc.

Ontologies will play a significant role in the realization of SNSS. For example, a considerable amount of data passes through the network and should be converted into higher abstractions that can be used in appropriate reasoning. This requires the development of standard terminologies which capture objects and events. Creating and testing such terminologies will aid in effective recognition and reaction in a network-centric situation awareness environment. This would involve identifying a methodology for development of terminologies for multimodal data (or ontologies), developing appropriate ontologies,

developing testing methods for these ontologies, demonstrating interoperability for selected domains (e.g., healthcare, situational awareness), and using these ontologies in decision making.

In today's session, we will take inventory of the what has transpired in the [OntologySummit2015](#) proceedings so far, and present the syntheses of the discourse of each of the four content tracks. The co-lead Editors will be presenting a first draft of the Communiqué Outline. An open discussion among the editors, the track co-champions and all the participants will ensue, with an aim towards arriving at a near-final [OntologySummit2015](#) Communiqué Outline, which will frame how this year's Communiqué will get developed by all parties concerned.

## Ontology Summit 2015 Communiqué

Source: <http://ontolog.cim3.net/file/work/On...Communique.pdf> (PDF)

Internet of Things: Toward Smart Networked Systems and Societies

Lead Editors: Mark Underwood and Michael Gruninger

CoEditors: Ken Baclawski, Mike Bennett, Gary BergCross, Torsten Hahmann, Leo Obrst, Ram Sriram

MY NOTE: This is still incomplete

### Introduction

We are witnessing another phase in the evolution in computing and communication. The Internet, which spans networks in a wide variety of domains, is having a significant impact on every aspect of our lives. The next generation of networks will extend beyond physically linked computers to include multimodal information from biological, cognitive, semantic, social, and sensor networks. This paradigm shift will involve symbiotic networks of people, intelligent devices, and mobile personal computing and communication devices (mPCDs), which will form netcentric societies or smart networked systems and societies (SNSS). mPCDs are already equipped with a myriad of sensors, with regular updates of additional sensing capabilities. Additionally, we are witnessing the emergence of "intelligent devices," such as smart meters, smart cars, etc., with considerable sensing and networking capabilities. Hence, these devices – and the network will be constantly sensing, monitoring, and interpreting the environment – this is sometimes referred to as the Internet of Things (IoT). And as local and wide area networks became almost secondary to the WWW (WorldWide Web), users and their usage patterns will become increasingly visible. This will have significant implications for both the market for advanced computing and communication infrastructure and the future markets – for nearly 4.5 billion people that netcentric societies will create.

Smart networked systems and societies will result in better quality of life, reduced threat from external sources, and improved commerce. For example, assume a scenario where people at various locations suffer from flulike symptoms. In a netcentric society, mPCDs will send vital signs and other associated information to appropriate laboratories and medical centers. These centers will analyze the information, including searching the Internet for potential solutions, and will aid in determining possible causes for this phenomenon. Based on the diagnosis, people will be directed to the nearest clinic for treatment. Here we have several types of information flowing through the net: data from mPCDs; location information; images; video; and audio.

The development of a trusted, secure, reliable, and interoperable netcentric computing environment will need technologies that can assure a flexible and scalable system allowing the application of diverse and robust privacy requirements, thus enabling the trusted and meaningful growth of netcentric infrastructures for the benefit of all societies. One such technical challenge is that the network consists of things (both devices and humans) which are heterogeneous, yet need to have seamless interoperability. Devices need to interoperate and data needs to be compatible to be integrated. This requires the development of standard terminologies which capture the meaning and relations of objects and events. Creating and testing such terminologies will aid in effective recognition and reaction in a networkcentric situation awareness environment. The primary goal of this summit to discuss the role of ontologies in the development of smart networked systems and societies.

Several key issues were addressed within the Ontology Summit, especially:

1. Making the case for IoT ontologies
2. How ontologies are used in IoT

- 3. The challenge of scalability
  - 4. Ontologybased standards for IoT
- 

## The Case for IoT Ontologies

Ontologies play a significant role in the realization of SNSS. For example, a considerable amount of data passes through the network and should be converted into higher abstractions that can be used in appropriate reasoning. This requires the development of standard terminologies which capture objects and events. Moreover, such terminologies must align with the intended semantics of generic and domainspecific concepts. Creating and testing such terminologies will aid in effective recognition and reaction in a networkcentric situation awareness environment. This involves identifying a methodology for development of terminologies for multimodal data (or ontologies), developing appropriate ontologies, both foundational (such as time, situation, events) and domain specific, developing testing methods for these ontologies, demonstrating interoperability for selected domains (e.g., healthcare, situational awareness), and using these ontologies in decision making.

Sensors are most closely in touch with the outside world and are thus are a big part of IoT since they provide an observational basis for data about things of interest. Since sensors are a big embedded part of the sensing and processing infrastructure of IoT, this results in many Big Data challenges related to semantic heterogeneity. Data can be hard to use because it is in different formats, uses inconsistent naming conventions, and is often provided at a low level of abstraction that makes it difficult to integrate it with other knowledge bases and software systems. To address these challenges, the Semantic Sensor Network Ontology (SSNO) was developed by W3C SSN-XG (2011) to help process and understand sensor information, and to allow the discovery, understanding, and querying of sensor data. SSNO is an ontology for describing networked sensors and its output by introducing a minimal set of classes and relations centered around the notions of stimuli, sensor, and observations. It includes different operational, device related and quality of information attributes that are related to sensing devices, and it describes the operational range, battery and power and environmental ranges that are specified for sensor devices.

Upper Ontologies such as DOLCE can also play a role in extending other IoT ontologies. There are broader Device Ontologies which can leverage some of the Physics Domain Ontology available in DOLCE with its well organized, conceptbased vocabulary. DOLCE also has a pattern for situation ontologies.

Of course, sensors are only one small part of the picture. Ontologies for time, duration, and dates are needed in order to capture the distinction between snapshots of measurements and the dynamic behaviour of an embedded system. Ontologies for location are required for scenarios in which the smart objects on the network are widely distributed geographically.

Events are a key concept that play a critical role in many IoT applications. In some scenarios, events create context by connecting people, things, places, and time; approaches such as the Simple Event Ontology (SEM) can be used to annotate events in these contexts and support retrieval of information. However, there are many scenarios in which there is a need to compose events into larger activities and to link events together to recognize patterns of behaviour.

Finally, IoT systems are not all passive in many scenarios, smart objects are enabled to make decisions and act autonomously in particular contexts. Many existing event ontologies need to be extended to represent this notion of agency.

---

## How Ontologies are Used in IoT

There are several IoT applications that have utilized ontologies to various degrees. These applications include manufacturing, healthcare, and disaster management. Scenarios that include complex event processing require ontologies that have extensive axiomatizations in expressive logics such as firstorder logic. In particular, manufacturing processes have complex causal and temporal structures, and complex event processing requires reasoning over situations and events. Typical ontology use scenarios in ontology mapping and decision support are described below.

### Ontology Mapping

The wide array of sensors within an IoT application and the variety of data that they provide leads inexorably to the problem of integrating the ontologies that are associated with these sensors. A typical application requires the interconnection of algorithms and hardware for multiple existing networks (such as a medical network and a transportation network that provides traffic data). One approach is to select an existing ontology to bridge such networks, or to combine existing ontologies in various domains and use these ontologies to integrate systems [e.g., Quantities, Units, Dimensions; Semantic Sensor Networks; Foundation Model of Anatomy; Symptom Ontology; Human Disease Ontology]. Other approaches explicitly address the problem of mapping between ontologies. The simplest approaches manually map JSON entities to target ontologies. In the

Hyper/CAT approach (see <http://www.hypercat.io/standard.html>), servers provide catalogues – an array of URIs of resources, annotated with metadata to clients. In the most sophisticated approaches we find Inferencebased Mapping, in which the mappings between ontologies can be achieved using an inference engine (or AI theorem provers).

In many IoT applications, there are two fundamentally different approaches to interoperability. In the first approach, we find centralized processing of spatially distributed and heterogeneous sensor data (Semantics in the Cloud). Data is collected in different settings by various kinds of sensors/things/persons, and all sensor observations are sent to the cloud for semantic annotation and processing. The challenge is to describe the various sources correctly to allow semantic integration. In the second approach, there is local processing (Semantics at the Edge), in which local intelligent sensor networks perform inplace computing. The challenge here is in using ontologies to smartly aggregate, filter, process, access, and respond to sensor data.

## Standards Integration

MISSING

## Decision Support for IoT

Many IoT applications, ranging from complex event processing and situation awareness to manufacturing, use automated inference from ontologies to assist in the decision making and to implement smart objects that can automatically act and react to changing situations. The critical issues in the deployment of IoT focus on three questions:

1. What kinds of axiomatizations are required for IoT ontologies?
2. How are the axioms of an ontology used in IoT applications?
3. How can ontologybased solutions scale up to realistic IoT scenarios?

A commonplace maxim invoked by many Semantic Web practitioners is “A little semantics goes a long way.” The critical issue is to identify, for a given IoT application, exactly what ontological approach is adequate. If ontologies are being used to annotate IoT data, then lightweight taxonomies can have a major impact by enabling the interpretation of data by other software applications. Nevertheless, SPARQL and RDF models are not adequate for all tasks; while SPARQL is great for querying a knowledge base, it is less ideal for fetching objects, and it is cumbersome when working with dynamic data. Applications based on complex event processing require more expressive axiomatizations of events, states, and causality.

## Beyond Semantic Sensor Network Ontologies

MISSING

## Ontological Issues

MISSING

## Scalability

The number, volume and variety of sensor data, whether delivered in real time as data streams or processed as stored batches, results in Big Data challenges (e.g. heterogeneity challenging integration, interpolation and summarization, filtering, compression). Many Big Data issues are common to sensor networks, such as the explosion of standards and reliance on metadata vocabularies such as the idea of things within IoT like services, users, networks, concentrators/aggregators and devices called “resources.” In the face of these challenges we can ask whether lightweight sensor ontologies scale, and what are the realistic ontological commitments for big heterogeneous data.

One aspect that distinguishes IoT scenarios from other applications of ontologies is the role of physical constraints. A sensing/actuating task that requires the cooperation and coordination of thousands of devices (within an Internet of billions), might be impractical due to memory, processing, and energy constraints. The interplay between these constraints and the semantic content of the ontology remains to a large extent unexplored.

The challenge of scalability also arises in the design of ontologies. With the size and increasing complexity of IoT, extensible and modular approaches are useful, if not essential. Approaches for developing small, focused ontologies customized to the available sensors and sensor data might be necessary, but it is an open research question as to whether the combination and integration of a large number of such ontologies is feasible.

Scalability is influenced by the different application case studies that drive the need for more semantics in sensor networks, and these approaches can be contrasted in the following table:

Sensor data discovery and integration	In network data stream processing
``Offline": happens after the fact	``Online": happens when and where the data is collected
Somewhat centralized: only need to integrate data from different data collection servers	Completely decentralized: Each device is both sensor and data processor, with sensors making individual or collaborative decisions
Full datasets (with broad spatial and temporal scope) are available	Only small spatial and temporal window of data accessible
Can utilize full available computational power	Limited in processing power (sensor device limitations, including bandwidth and energy consumption)
Can employ complex ontologies	Limited to small tailored ontologies
Typical semantic problems: <ul style="list-style-type: none"> <li>• Integration problems arising from variety</li> <li>• Context of data and sensors</li> <li>• Provenance</li> </ul>	Typical semantic problems: <ul style="list-style-type: none"> <li>• Ontologies can be deployed on sensors</li> <li>• Integrating and maintaining ontologies across sensors.</li> </ul>

## Standards Integration

Ontology Summit 2009 explored ontologybased standards, and one of the key insights that arose from that work is that specifying an ontology for a standard enables more effective deployment of the standard and easier integration with other overlapping standards. There is also a symbiotic relationship between standards and ontologies the terminology within any standard provides the initial set of concepts which are axiomatized within an ontology, and the specification of the ontology provides rigorous, unambiguous semantics for the terminology of the standard.

What are the relevant or de facto standards involved in the adoption of ontologies for the Internet of Things? There have been several IOT Ontology success stories. The W3C Semantic Sensor Network Ontology (OWL 2) and the OGC Sensor Web Enablement project (including SensorML, a Transducer Model Language, a Sensor Observations Service, Sensor Planning Service) efforts were cited by speaker Henson (Bosch). The GraphOfThings project incorporates SPARQL and the Continuous Query Evaluation over Linked Stream (CQELS) tool. Intellego leverages OWL, RDF and the SSN Ontology.

A decadeold example that predated IoT's entry into common parlance was Project Drishti (Ran, Helal, & Moore, 2004). The investigators sought to integrate data streams from RFID tags, GPS and wireless networks to aid the visually impaired in common navigation tasks. There were numerous other integrations in the wearable and ubiquitous computing literature, even in science fiction.

Fast forward to the present and the number of data sources has multiplied. Big Data is competing with IoT for attention – and legitimately so, as noted in the 2014 Ontology Summit. This has created terrific momentum, especially for Big Data and the Apache stack which owns most of the developer mindshare about this paradigm shift. A convergence of open source projects, cloud computing and a steady march toward webenabled applications has facilitated big data, but has the same occurred for IoT? There does not seem to

It seems clear that there are many efforts underway, and that full coordination with standards or Standards Developing Organizations is not a prerequisite for building a workable system. Benefits from using ontologybased standards in IoT may be more evident as systems mature than at this early stage of IoT work simply because more things will be interconnected. A complex system requiring many different human and organizational roles, processing speed and volume might need an ontology as its associated sensor grid shifts beneath it.

## Challenges

Software Support We lack tools for a wide range of tasks, including for semantic annotation and ontology validation.

Furthermore, most applications still rely on manual methods for integration. There is also demand to create tools for ontology visualization and interoperability testing.

What ontologies are needed for supporting today's envisioned IoT applications? Much existing work for modeling IoT resources focuses primarily on sensors and sensor networks and is modeled by SSNO. Most of the existing IoT or sensorrelated ontologies represent IoT devices only partially (e.g. as sensing devices), so extensions will be required to include other entities and their relationship to actuator devices. A broader view of IoT resources including other important resources and devices such as actuators, IoT gateways, data aggregators and servers is needed. Work to develop ontologies for these is underway.

Beyond Semantic Sensor Network Ontologies How do we handle going beyond SSN with an Open Source Cloud solution for the Internet of Things (OpenIoT)? Challenges include sensor annotation, sensor mobility & efficient data harvesting and data quality.

What Kinds of Axioms are Needed? Is the priority work and opportunity for ontologies to be used to annotate IoT data, or to more fully represent and model sensors and data in order to analyse/understand it?

Semantic Annotation How can we provide an ontological base for generating semantic annotations of open source internetconnected objects? The challenge would be to obtain open sensor information in a standard encoding that is understandable by users and their software

Semantic Registry for IoT Entities, built on top of DUL and SSNO1. Besides the registration of IoT things, abstractions of technological heterogeneity are also required. Such abstract semantic heterogeneity leads to the need to use heterogeneous domain ontologies to semantically annotate data of IoT entities.

Ontology Evolution How can we characterize how ontologies change in order to address future IoT applications?

1 Some initial work along these lines can be found at <http://purl.org/IoT/iotontology.owl>

<http://ai-group.ds.unipi.gr/kotis/ontologies/IoTontology>

---

## Forecasts

Ontology Development There will be a number of efforts to enhance and extend IoT ontologies such as SSNO. More ambitious extensions of SSNO will support the extraction of knowledge from the raw sensor data, enabling the understanding of the ``big picture'' of what is happening by explicitly representing the interactions between complex processes and events that cannot be captured by a single signal alone.

Ontology Embedding The increased use of smart devices, storeandforward, embedded intelligence automated data fusion (perhaps especially for geospatial aspects) suggests that ontology embedding could become a design pattern. The pattern could be used in building intelligent IoT, but ontology embedding within sensor systems themselves is possible. Metadata for discovery and provenance from devices are possible starting points.

Automated Deployment of IoT Apps in Unknown Environments Approaches such as the Semantic Smart Gateway Framework will be extended to support full automation in terms of uncovering the semantics of IoT entities as well as aligning their semantics in cases of disagreement.

Exploitation of (Lazy) Developer Pain Points Known problem areas in IoT exist across many different types of sensors. These include security, privacy, signal noise, reliability, configuration management, infrastructure dependency and other known architectural nuisances. A standard solution in any of these areas could catch on because it would solve a welldefined problem that is tangential to an architect or sponsor's main system objectives.

Specialized Engines Reusable, highcomplexity solutions might take hold to implement mathematical solutions in certain spaces, such as Gruninger's work with PSL in ERP or Spencer Breiner's category theory.

Cloud Impact Because cloud engines such as Watson will provide complex building blocks for architects, the challenge may be taken up by small groups or even sole developers working in green field problem spaces.

Fun Hardware Syndrome Sometimes collateral innovations cooccur with fun hardware developments. The smart car, or low cost commercial unmanned vehicles could spur ontologyrich solutions. The reasons for such developments are connected both to standards and to the attitudes (plus and minus) about existing standards.

Integrated Development Environment Innovation Will IoT need its own integrated development environment? Test and development beds for IoT will likely require new combinations of devices, simulations, test data, standards, scalability exercises

and more.

---

## Recommendations

1. IoT ontologies need to deal with dynamic time varying data vs. the often static Semantic Web. In particular, more work is needed on the development of event ontologies for targeted domains, building from core ontologies.
  2. Use design patterns toward ontology virtualization: Given a set of ontology design patterns and their combination into microontologies, one can abstract the underlying axiomatization by: dynamically reconfiguring patterns in a plug and play style; bridging between different patterns as microtheories; providing ontological views and semantic shortcuts that suit particular provider, user, and use case needs by highlighting or hiding certain aspects of the underlying ontological model; and mapping between major modeling styles
- Integrating SSNO with other Web standards and ontologies is a nearterm focus for work. In particular, there is a need to support applications that combine SSNO with PROVO (for data provenance), CoAP (Constrained Application Protocol), and RDF Data Cube vocabulary. There are also many applications based on biomedical ontologies dealing with sensors in medical devices.
4. Ontology reuse is key. Of course, ontology reuse issues are not unique to IoT but there are some good ontologies such as SSN and PROV that provide some starting points for representing sensors, sensor also networks, observations, etc.
  5. Link your data and descriptions to other existing resources
  6. Semantics are only one part of the solution and often not the endproduct so the focus of the design should be on creating effective methods, tools and APIs to handle and process the semantics. Query methods, machine learning, reasoning and data analysis techniques and methods should be able to effectively use these semantics.
  7. A critical obstacle in the widespread adoption/application of ontologies to earth science and sensor systems is the lack of tools that address concrete use cases. Developers will need to focus on those tools and techniques that support the deployment of ontologies in IoT applications.
  8. Create an IoT equivalent to Google Search to identify the scope of available end points for different application domains.
  9. A more coordinated effort is required to compile IoT case studies which can serve as the basis for ontology reuse and the design of new ontologies. Key areas included Sensor integration, Smart Grid, and Smart Healthcare.
- 

## Terminology

**Internet of Things.** The Internet of Things (IoT) is a term that is being used to denote a network – typically the Internet of devices that constantly monitor the environment and can result in “intelligent actions.” These devices can range from simple sensors to complex systems such as automobiles and buildings. There are several views of IoT in vogue. For example, ITU (International Telecommunication Union) and IERC (IoT European Research Cluster) define IoT as “a global network infrastructure with selfconfiguring capabilities based on standard and interoperable communication protocols where physical and virtual things have identities, physical attributes and virtual personalities, use intelligent interfaces and are seamlessly integrated into the information network.” (See Internet of Things – From Research and Innovation to Market Deployment, Verma, O. and Friess, P. (editors), 2014, River Publishers, Aalborg, Denmark).

**CyberPhysical Systems.** Cyberphysical systems (CPSs) extend IoT by adding a control and decision making layer. Again, several views of CPSs exist. One commonly used definition is provided in <http://varma.ece.cmu.edu/summit/index.html>, which places an emphasis on embedded systems and the tight coupling between hardware and software. CPSs will play an increasingly important role in the next generation industrial systems.

**CyberPhysical Human Systems.** When humans take an active role in CPSs we have Cyberphysical Human Systems (CPHSs). These systems can be viewed as sociotechnical systems, with a symbiotic relationship between the human and the physical device.

**Cyberphysical Social Systems or Smart Networked Systems and Societies.** Social networks, such as Facebook and Twitter, primarily connect people to one another. These networks are playing very important roles in people's lives today, from how some of them behave and interact with one another, to change in human resources processes, how companies market and sell

products and services, developments in healthcare and smart (electrical) grid systems, and even roles in politics and democratic uprisings. Social networks have been used both to curtail and to propagate freedom of speech. When these networks are combined with CPSs, we have Smart Networked Systems and Societies (SNSS), which are also known as Cyberphysical Social Systems (CPSS) or Internet of Everything (IoE).

---

## Training Students to Extract Value from Big Data

Source: [https://www.nap.edu/login.php?record...ord\\_id%3D18981](https://www.nap.edu/login.php?record...ord_id%3D18981) (PDF)

---

### The National Academies Press

Source: <http://www.nap.edu/catalog/18981/training-students-to-extract-value-from-big-data>

#### Authors

Committee on Applied and Theoretical Statistics; [Board on Mathematical Sciences and Their Applications](#); [Division on Engineering and Physical Sciences](#); National Research Council

#### Description

As the availability of high-throughput data-collection technologies, such as information-sensing mobile devices, remote sensing, internet log records, and wireless sensor networks has grown, science, engineering, and business have rapidly transitioned from striving to develop information from scant data to a situation in which the challenge is now that the amount of information exceeds a human's ability to examine, let alone absorb, it. Data sets are increasingly complex, and this potentially increases the problems associated with such concerns as missing information and other quality concerns, data heterogeneity, and differing data formats.

The nation's ability to make use of data depends heavily on the availability of a workforce that is properly trained and ready to tackle high-need areas. Training students to be capable in exploiting big data requires experience with statistical analysis, machine learning, and computational infrastructure that permits the real problems associated with massive data to be revealed and, ultimately, addressed. Analysis of big data requires cross-disciplinary skills, including the ability to make modeling decisions while balancing trade-offs between optimization and approximation, all while being attentive to useful metrics and system robustness. To develop those skills in students, it is important to identify whom to teach, that is, the educational background, experience, and characteristics of a prospective data-science student; what to teach, that is, the technical and practical content that should be taught to the student; and how to teach, that is, the structure and organization of a data-science program.

Training Students to Extract Value from Big Data summarizes a workshop convened in April 2014 by the National Research Council's Committee on Applied and Theoretical Statistics to explore how best to train students to use big data. The workshop explored the need for training and curricula and coursework that should be included. One impetus for the workshop was the current fragmented view of what is meant by analysis of big data, data analytics, or data science. New graduate programs are introduced regularly, and they have their own notions of what is meant by those terms and, most important, of what students need to know to be proficient in data-intensive work. This report provides a variety of perspectives about those elements and about their integration into courses and curricula.

#### Topics

- [Math, Chemistry, and Physics — Math and Statistics](#)
- [Education — Math and Science Education](#)

#### Publication Info

66 pages | 7 x 10

Paperback

ISBN: 978-0-309-31437-4

#### Copyright Information

The National Academies Press (NAP) has partnered with Copyright Clearance Center's Rightslink service to offer you a variety of options for reusing NAP content. Through Rightslink, you may request permission to reprint NAP content in another publication, course pack, secure website, or other media. Rightslink allows you to instantly obtain permission, pay related fees, and print a license directly from the NAP website. The complete terms and conditions of your reuse license can be found in the license agreement that will be made available to you during the online order process. To request permission through Rightslink you are required to create an account by filling out a simple online form. The following list describes license reuses offered by the National Academies Press (NAP) through Rightslink:

- Republish text, tables, figures, or images in print
- Post on a secure Intranet/Extranet website
- Use in a PowerPoint Presentation
- Distribute via CD-ROM
- Photocopy

[Click here to obtain permission for the above reuses.](#) If you have questions or comments concerning the Rightslink service, please contact:

---

Rightslink Customer Care  
Tel (toll free): 877/622-5543  
Tel: 978/777-9929  
E-mail: [customercare@copyright.com](mailto:customercare@copyright.com)  
Web: <http://www.rightslink.com>

---

To request permission to distribute a PDF, please contact our Customer Service Department at 800-624-6242 for pricing.

To request permission to translate a book published by the National Academies Press or its imprint, the Joseph Henry Press, please [click here to view more information.](#)

---

## Cover Page

Training Students to Extract Value from Big Data: Summary of a Workshop

Maureen Mellody, Rapporteur

Committee on Applied and Theoretical Statistics

Board on Mathematical Sciences and Their Applications

Division on Engineering and Physical Sciences

TRAINING STUDENTS TO EXTRACT VALUE FROM  
**BIG DATA**

Summary of a Workshop

Maureen Mellody, *Rapporteur*

Committee on Applied and Theoretical Statistics

Board on Mathematical Sciences and Their Applications

Division on Engineering and Physical Sciences

NATIONAL RESEARCH COUNCIL  
OF THE NATIONAL ACADEMIES

---

## The National Academies Press

THE NATIONAL ACADEMIES PRESS 500 Fifth Street, NW Washington, DC 20001

NOTICE: The project that is the subject of this report was approved by the Governing Board of the National Research Council, whose members are drawn from the councils of the National Academy of Sciences, the National Academy of Engineering, and the Institute of Medicine.

This study was supported by Grant DMS-1332693 between the National Academy of Sciences and the National Science Foundation. Any opinions, findings, or conclusions expressed in this publication are those of the author and do not necessarily reflect the views of the organizations or agencies that provided support for the project.

International Standard Book Number-13: 978-0-309-31437-4

International Standard Book Number-10: 0-309-31437-2

This report is available in limited quantities from:

Board on Mathematical Sciences and Their Applications

500 Fifth Street NW

Washington, DC 20001

[bmsa@nas.edu](mailto:bmsa@nas.edu)

<http://www.nas.edu/bmsa>

Additional copies of this workshop summary are available for sale from the National Academies

Press, 500 Fifth Street, NW, Keck 360, Washington, DC 20001; (800) 624-6242 or (202) 334-3313;  
<http://www.nap.edu/>.

Copyright 2014 by the National Academy of Sciences. All rights reserved.

Printed in the United States of America

---

## Planning Committee on Training Students

### PLANNING COMMITTEE ON TRAINING STUDENTS TO EXTRACT VALUE FROM BIG DATA: A WORKSHOP

JOHN LAFFERTY, University of Chicago, Co-Chair

RAGHU RAMAKRISHNAN, Microsoft Corporation, Co-Chair

DEEPAK AGARWAL, LinkedIn Corporation

CORINNA CORTES, Google, Inc.

JEFF DOZIER, University of California, Santa Barbara

ANNA GILBERT, University of Michigan

PATRICK HANRAHAN, Stanford University

RAFAEL IRIZARRI, Harvard University

ROBERT KASS, Carnegie Mellon University

PRABHAKAR RAGHAVAN, Google, Inc.

NATHANIEL SCHENKER, Centers for Disease Control and Prevention

ION STOICA, University of California, Berkeley

#### Staff

NEAL GLASSMAN, Senior Program Officer

SCOTT T. WEIDMAN, Board Director

MICHELLE K. SCHWALBE, Program Officer

RODNEY N. HOWARD, Administrative Assistant

---

## Committee on Applied and Theoretical Statistics

### COMMITTEE ON APPLIED AND THEORETICAL STATISTICS

CONSTANTINE GATSONIS, Brown University, Chair

MONTSERRAT (MONTSE) FUENTES, North Carolina State University

ALFRED O. HERO III, University of Michigan

DAVID M. HIGDON, Los Alamos National Laboratory

IAIN JOHNSTONE, Stanford University

ROBERT KASS, Carnegie Mellon University

JOHN LAFFERTY, University of Chicago

XIHONG LIN, Harvard University

SHARON-LISE T. NORMAND, Harvard University

GIOVANNI PARMIGIANI, Harvard University

RAGHU RAMAKRISHNAN, Microsoft Corporation

ERNEST SEGLIE, Office of the Secretary of Defense (retired)

LANCE WALLER, Emory University

EUGENE WONG, University of California, Berkeley

#### Staff

MICHELLE K. SCHWALBE, Director

RODNEY N. HOWARD, Administrative Assistant

---

## Board of Mathematical Sciences and Their Applications

### BOARD ON MATHEMATICAL SCIENCES AND THEIR APPLICATIONS

DONALD SAARI, University of California, Irvine, Chair

DOUGLAS N. ARNOLD, University of Minnesota

GERALD G. BROWN, Naval Postgraduate School

L. ANTHONY COX, JR., Cox Associates, Inc.  
CONSTANTINE GATSONIS, Brown University  
MARK L. GREEN, University of California, Los Angeles  
DARRYL HENDRICKS, UBS Investment Bank  
BRYNA KRA, Northwestern University  
ANDREW W. LO, Massachusetts Institute of Technology  
DAVID MAIER, Portland State University  
WILLIAM A. MASSEY, Princeton University  
JUAN C. MESA, University of California, Merced  
JOHN W. MORGAN, Stony Brook University  
CLAUDIA NEUHAUSER, University of Minnesota  
FRED S. ROBERTS, Rutgers University  
CARL P. SIMON, University of Michigan  
KATEPALLI SREENIVASAN, New York University  
EVA TARDOS, Cornell University

#### Staff

SCOTT T. WEIDMAN, Board Director  
NEAL GLASSMAN, Senior Program Officer  
MICHELLE K. SCHWALBE, Program Officer  
RODNEY N. HOWARD, Administrative Assistant  
BETH DOLAN, Financial Associate

---

## Acknowledgment of Reviewers

This report has been reviewed in draft form by persons chosen for their diverse perspectives and technical expertise in accordance with procedures approved by the National Research Council's Report Review Committee. The purpose of this independent review is to provide candid and critical comments that will assist the institution in making its published report as sound as possible and to ensure that the report meets institutional standards of objectivity, evidence, and responsiveness to the study charge. The review comments and draft manuscript remain confidential to protect the integrity of the deliberative process. We thank the following individuals for their review of this report:

Michael Franklin, University of California, Berkeley,  
Johannes Gehrke, Cornell University,  
Claudia Perlich, Distillery, and  
Duncan Temple Lang, University of California, Davis.

Although the reviewers listed above have provided many constructive comments and suggestions, they were not asked to endorse the views presented at the workshop, nor did they see the final draft of the workshop summary before its release. The review of this workshop summary was overseen by Anthony Tyson, University of California, Davis. Appointed by the National Research Council, he was responsible for making certain that an independent examination of the summary was carried out in accordance with institutional procedures and that all review comments were carefully considered. Responsibility for the final content of this summary rests entirely with the author and the institution.

---

## 1 INTRODUCTION

Data sets—whether in science and engineering, economics, health care, public policy, or business—have been growing rapidly; the recent National Research Council (NRC) report [Frontiers in Massive Data Analysis](#) documented the rise of “big data,” as systems are routinely returning terabytes, petabytes, or more of information (National Research Council, 2013). Big data has become pervasive because of the availability of high-throughput data collection technologies, such as information-sensing mobile devices, remote sensing, radiofrequency identification readers, Internet log records, and wireless sensor networks. Science, engineering, and business have rapidly transitioned from the longstanding state of striving to develop information from scant data to a situation in which the challenge is now that the amount of information exceeds a human’s ability to examine, let alone absorb, it. Web companies—such as Yahoo, Google, and Amazon—commonly work with data sets that consist of billions of items, and they are likely to increase by an order of magnitude or more as the Internet of Things [1](#) matures. In other words, the size and scale of data, which can be overwhelming today, are only increasing. In addition, data sets

are increasingly complex, and this potentially increases the problems associated with such concerns as missing information and other quality concerns, data heterogeneity, and differing data formats.

Advances in technology have made it easier to assemble and access large amounts of data. Now, a key challenge is to develop the experts needed to draw reliable inferences from all that information. The nation's ability to make use of the data depends heavily on the availability of a workforce that is properly trained and ready to tackle these high-need areas. A report from McKinsey & Company (Manyika et al., 2011) has predicted shortfalls of 150,000 data analysts and 1.5 million managers who are knowledgeable about data and their relevance. It is becoming increasingly important to increase the pool of qualified scientists and engineers who can extract value from big data. Training students to be capable in exploiting big data requires experience with statistical analysis, machine learning, and computational infrastructure that permits the real problems associated with massive data to be revealed and, ultimately, addressed. The availability of repositories (of both data and software) and computational infrastructure will be necessary to train the next generation of data scientists. Analysis of big data requires cross-disciplinary skills, including the ability to make modeling decisions while balancing trade-offs between optimization and approximation, all while being attentive to useful metrics and system robustness. To develop those skills in students, it is important to identify whom to teach, that is, the educational background, experience, and characteristics of a prospective data science student; what to teach, that is, the technical and practical content that should be taught to the student; and how to teach, that is, the structure and organization of a data science program.

The topic of training students in big data is timely, as universities are already experimenting with courses and programs tailored to the needs of students who will work with big data. Eight university programs have been or will be launched in 2014 alone.<sup>2</sup> The workshop that is the subject of this report was designed to enable participants to learn and benefit from emerging insights while innovation in education is still ongoing.

## Workshop Overview

On April 11-12, 2014, the standing Committee on Applied and Theoretical Statistics (CATS) convened a workshop to discuss how best to train students to use big data. CATS is organized under the auspices of the NRC Board on Mathematical Sciences and Their Applications.

To conduct the workshop, a planning committee was first established to refine the topics, identify speakers, and plan the agenda. The workshop was held at the Keck Center of the National Academies in Washington, D.C., and was sponsored by the National Science Foundation (NSF). About 70 persons—including speakers, members of the parent committee and board, invited guests, and members of the public—participated in the 2-day workshop. The workshop was also webcast live, and at least 175 persons participated remotely.

A complete statement of task is shown in Box 1.1. The workshop explored the following topics:

- The need for training in big data.
- Curricula and coursework, including suggestions at different instructional levels and suggestions for a core curriculum.
- Examples of successful courses and curricula.
- Identification of the principles that should be delivered, including sharing of resources.

Although the title of the workshop was "Training Students to Extract Value from Big Data," the term big data is not precisely defined. CATS, which initiated the workshop, has tended to use the term massive data in the past, which implies data on a scale for which standard tools are not adequate. The terms data analytics and data science are also becoming common. They seem to be broader, with a focus on using data—maybe of unprecedented scale, but maybe not—in new ways to inform decision making. This workshop was not developed to explore any particular one of these definitions or to develop definitions. But one impetus for the workshop was the current fragmented view of what is meant by analysis of big data, data analytics, or data science. New graduate programs are introduced regularly, and they have their own notions of what is meant by those terms and, most important, of what students need to know to be proficient in data-intensive work. What are the core subjects in data science? By illustration, this workshop began to answer that question. It is clear that training in big data, data science, or data analytics requires a multidisciplinary foundation that includes at least computer science, machine learning, statistics, and mathematics, and that curricula should be developed with the active participation of at least these disciplines. The chapters of this summary provide a variety of perspectives about those elements and about their integration into courses and curricula.

### BOX 1.1 Statement of Task

An ad hoc committee will plan and conduct a public workshop on the subject of training undergraduate and graduate students to extract value from big data. The committee will develop the agenda, select and invite speakers and discussants, and

moderate the discussions. The presentations and discussions at the workshop will be designed to enable participants to share experience and perspectives on the following topics:

- What current knowledge and skills are needed by big data users in industry, government, and academia?
- What will students need to know to be successful using big data in the future (5-10 years out)?
- How could curriculum and training evolve to better prepare students for big data at the undergraduate and graduate levels?
- What options exist for providing the necessary interdisciplinary training within typical academic structures?
- What computational and data resources do colleges and universities need in order to provide useful training? What are some options for assembling that infrastructure?

Although the workshop summarized in this report aimed to span the major topics that students need to learn if they are to work successfully with big data, not everything could be covered. For example, tools that might supplant MapReduce, such as Spark, are likely to be important, as are advances in Deep Learning. Means by which humans can interact with and absorb huge amounts of information—such as visualization tools, iterative analysis, and human-in-the-loop systems—are critical. And such basic skills as data wrangling, cleaning, and integration will continue to be necessary for anyone working in data science. Educators who design courses and curricula must consider a wide array of skill requirements.

The present report has been prepared by the workshop rapporteur as a factual summary of what occurred at the workshop. The planning committee's role was limited to planning and convening the workshop. The views contained in the report are those of individual workshop participants and do not necessarily represent the views of all workshop participants, the planning committee, or the NRC.

## National Efforts in Big Data

Suzanne Iacono, National Science Foundation

Suzanne Iacono, of NSF, set the stage for the workshop by speaking about national efforts in big data, current challenges, and NSF's motivations for sponsoring the workshop. She explained that the workshop was an outgrowth of the national big data research and development (R&D) initiative. The federal government is interested in big data for three reasons:

- To stimulate commerce and the economy.
- To accelerate the pace of discovery and enable new activities.
- To address pressing national challenges in education, health care, and public safety.

Big data is of interest to the government now because of the confluence of technical, economic, and policy interests, according to Iacono. Advances in technology have led to a reduction in storage costs, so it is easier to retain data today. On the policy side, data are now considered to be assets, and government is pushing agencies to open data sets to the public. In other words, there has been a democratization of data use and tools.

Iacono described a recent book (Mayer-Schönberger and Cukier, 2012) that outlined three basic shifts in today's data:

- There are more data than ever before.
- Data are messy, and there must be an increased acceptance of imperfection.
- Correlations can help in making decisions.

She then described the national Big Data Research and Development Initiative in more detail. A 2010 report from the President's Council of Advisors on Science and Technology argued that the federal government was not investing sufficiently in big data research and development and that investment in this field would produce large returns. A working group in big data, under the interagency Networking and Information Technology Research and Development (NITRD) program and managed by the Office of Science and Technology Policy, was charged with establishing a framework for agency activity. The result was that in 2012, \$200 million was allocated for big data R&D throughout the NITRD agencies, including the Defense Advanced Research Projects Agency (DARPA), the Department of Energy (DOE) Office of Science, the National Institutes of Health (NIH), and NSF. Iacono showed the framework for moving forward with big data R&D, which included the following elements:

- Foundational research. Iacono stressed that this research is critical because data are increasing and becoming more heterogeneous.
- Cyberinfrastructure. New and adequate infrastructure is needed to manage and curate data and serve them to the larger research community.
- New approaches to workforce and education.

- New collaborations and outreach.

Iacono noted that policy envelops all four elements of the framework.

A 2013 White House memorandum directed executive branch agencies to develop plans to increase public access to the results of federally funded research, including access to publications and data, and plans are under way at the agency level to address this memorandum. Iacono noted that increased access to publications is not difficult, because existing publication-access methods in professional societies and some government agencies can be used as models. She also pointed out that NIH's PubMed [3](#) program may be a useful model in that it shares research papers. However, she noted that access to data will be much more difficult than access to publications because each discipline and community will have its own implementation plan and will treat data privacy, storage duration, and access differently.

Iacono described foundational R&D in more detail. She explained that NSF and NIH awarded 45 projects in big data in 2012 and 2013. About half were related to data collection and management and one-fourth to health and bioinformatics. The remaining awards were spread among social networks, physical sciences and engineering, algorithms, and cyberinfrastructure. Seventeen agencies are involved in the Big Data Senior Steering Group, and each is implementing programs of its own related to big data. For example, DARPA has implemented three new programs—Big Mechanism, Memex, and Big Data Capstone; the National Institute of Standards and Technology maintains a Big Data Working Group; DOE has an Extreme Scale Science initiative; and NSF and NIH each has a broad portfolio related to big data. Iacono stressed that big data is a national issue and that there is substantial interest now in industry and academe, so she believes that government should consider multistakeholder partnerships.

Iacono discussed three challenges related to big data:

- Technology. She emphasized that technology alone cannot solve big data problems, and she cited several recent popular books that discuss the folly of technological solutionism (Mayer-Schönberger and Cukier, 2012; Mele, 2013; Reese, 2013; Schmidt and Cohen, 2013; Surdak, 2014; Webb, 2013).
- Privacy. Iacono pointed out that many of our behaviors—including shopping, searching, and social interactions—are now tracked, and she noted that a White House 90-day review to examine the privacy implications of big data was under way. [4](#) In general, Iacono noted the importance of regulating data use, as opposed to data collection; balancing interests; and promoting data sharing.
- Education and workforce. As noted above, the 2011 report from the McKinsey & Company predicted large shortfalls of big data experts. Iacono noted that the Harvard Business Review labeled data science as “the sexiest job of the 21st century” (Davenport and Patil, 2012). The New York Times has recently hired a chief data scientist. The bottom line, Iacono explained, is that the talent pool in data science must be expanded to meet current and future needs.

Iacono pointed out that there are traditional ways to educate students through school curricula but that there are also other ways to learn. Such companies as DataKind and Pivotal are matching data scientists with data problems in the nonprofit community. Universities, such as the University of Chicago, as discussed by Rayid Ghani (see Chapter 2), are also working to connect data scientists to problems of social good. Iacono concluded by stressing the many opportunities and challenges in big data that lie ahead.

## Organization of This Report

The remaining chapters of this report summarize the workshop presentations and discussions. To assist the reader, each chapter begins with a short list of important statements made by speakers during the workshop session. Chapter 2 outlines the need for training. Chapter 3 discusses some of the principles of working with big data. Chapter 4 focuses on courses and curricula needed to support the use of big data. Chapter 5 discusses shared resources, and Chapter 6 summarizes the group discussion of lessons learned from the workshop. Finally, Appendix A lists the workshop participants, Appendix B shows the workshop agenda, and Appendix C defines acronyms used in this report.

---

## 2 THE NEED FOR TRAINING: EXPERIENCES AND CASE STUDIES

### Important Points Made by Individual Speakers

- Students often do not recognize that big data techniques can be used to solve problems that address societal good, such as those in education, health, and public policy; educational programs that foster relationships between data science and social problems have the potential to increase the number and types of students interested in data science. (Rayid Ghani)
- There may be a mismatch between some industry needs and related academic pursuits: current studies of

recommendation systems, such as off-line score prediction, do not always correlate well with important industry metrics, such as sales and user engagement. (Guy Lebanon)

- Academia does not have sufficient access to practical data scenarios in industry. (Guy Lebanon)

Big data is becoming pervasive in industry, government, and academe. The disciplines that are affected are as diverse as meteorology, Internet commerce, genomics, complex physics simulations, health informatics, and biologic and environmental research. The second session of the workshop focused on specific examples and case studies of real-world needs in big data. The session was co-chaired by John Lafferty (University of Chicago) and Raghu Ramakrishnan (Microsoft Corporation), the cochairs of the workshop's organizing committee. Presentations were made in Session 2 by Rayid Ghani (University of Chicago) and Guy Lebanon (Amazon Corporation).

## Training Students to Do Good with Big Data

Rayid Ghani, University of Chicago

Rayid Ghani explained that he has founded a summer program at the University of Chicago, known as the Data Science for Social Good Fellowship, to show students that they can apply their talents in data science to societal problems and in so doing affect many lives. He expressed his growing concern that the top technical students are disproportionately attracted to for-profit companies, such as Yahoo and Google, and posited that these students do not recognize that solutions to problems in education, health, and public policy also need data.

Ghani showed a promotional video for the University of Chicago summer program and described its applicant pool. Typically, half the applicants are computer science or machine learning students; one-fourth are students in social science, public policy, or economics; and one-fourth are students in statistics. Some 35 percent of the enrolled students are female (as Ghani pointed out, this is a larger proportion than is typical of a computer science graduate program). Many of the applicants are graduate students, and about 25 percent are undergraduate seniors. The program is competitive: in 2013, there were 550 applicants for 36 spots. Ghani hypothesized that the program would be appropriate for someone who had an affinity for mathematics and science but a core interest in helping others. Once in the program, students are matched with mentors, most of whom are computer scientists or economists with a strong background in industry.

He explained that the program is project-based, using real-world problems from government and nonprofit organizations. Each project includes an initial mapping of a societal problem to a technical problem and communication back to the agency or organization about what was learned. Ghani stressed that students need to have skills in communication and common sense in addition to technical expertise. The curriculum at the University of Chicago is built around tools, methods, and problem-solving skills. The program now consistently uses the Python language, and it also teaches database methods. Ghani emphasized the need to help students to learn new tools and techniques. He noted, for instance, that some of the students knew of regression only as a means of evaluating data whereas other tools may be more suitable for massive data.

Ghani described a sample project from the program. A school district in Arizona was experiencing undermatching—that is, students have the potential to go to college but do not, or students have the potential to go to a more competitive college than the one they ultimately select. The school district had collected several years of data. In a summer project, the University of Chicago program students built models to predict who would graduate from college, who would go to college, and who was not likely to apply. In response to the data analysis, the school district has begun a targeted career-counseling program to begin intervention.

## The Need for Training in Big Data: Experiences and Case Studies

Guy Lebanon, Amazon Corporation

Guy Lebanon began by stating that extracting meaning from big data requires skills of three kinds: computing and software engineering; machine learning, statistics, and optimization; and product sense and careful experimentation. He stressed that it is difficult to find people who have expertise and skills in all three and that competition for such people is fierce.

Lebanon then provided a case study in recommendation systems. He pointed out that recommendation systems (recommending movies, products, music, advertisements, and friends) are important for industry. He described a well-known method of making recommendations known as matrix completion. In this method, an incomplete user rating matrix is completed to make predictions. The matrix completion method favors low-rank (simple) completions. The best model is found by using a nonlinear optimization procedure in a high-dimensional space. The concept is not complex, but Lebanon indicated that its implementation can be difficult. Implementation requires knowledge of the three kinds referred to earlier. Specifically, Lebanon

noted the following challenges:

- Computing and software engineering: language skills (usually C++ or Java), data acquisition, data processing (including parallel and distributed computing), knowledge of software engineering practices (such as version control, code documentation, building tools, unit tests, and integration tests), efficiency, and communication among software services.
- Machine learning: nonlinear optimization and implementation (such as stochastic gradient descent), practical methods (such as momentum and step selection size), and common machine learning issues (such as overfitting).
- Product sense: an online evaluation process to measure business goals; model training; and decisions regarding history usage, product modification, and product omissions.

Lebanon described two problems that limit academic research in recommendation systems, both related to overlooking metrics that are important to industry. First, accuracy in academic, off-line score prediction does not correlate with important industry metrics, such as sales and increased user engagement. Second, academe does not have sufficient access to practical data scenarios from industry. Lebanon posited that academe cannot drive innovation in recommendation systems; research in recommendation systems does not always translate well to the real world, and prediction accuracy is incorrectly assumed to be equivalent to business goals.

He then described a challenge run by Netflix. In the early 2000s, Netflix held a competition to develop an improved recommendation system. It provided a data set of ratings that had been anonymized and offered a \$1 million prize to the top team. The competition created a boost in research, which saw a corresponding increase in research papers and overall interest. However, a group of researchers at the University of Texas, Austin, successfully deanonymized the Netflix data by joining them with other data. Netflix later withdrew the data set and is now facing a lawsuit. As a result of that experience, industry is increasingly wary about releasing any data for fear of inadvertently exposing private or proprietary data, but this makes it difficult for academe to conduct relevant and timely research.

Lebanon pointed out that the important result in a recommendation system is prediction of a user's reaction to a specific recommendation. For it to be successful, one needs to know the context in which the user acts—for instance, time and location information—but that context is not conveyed in an anonymized data set. In addition, methods that perform well on training and test data sets do not perform well in real environments when a user makes a single A/B comparison. [1](#) Lebanon proposed several new ideas to address those characteristics:

- Study the correlations between existing evaluation methods and increased user engagement in an A/B test.
- Develop new off-line evaluations to account for user context better.
- Develop efficient searches among the possibilities to maximize A/B test performance.

Few data sets are publicly available, according to Lebanon. Working with limited data, the research community may focus on minor improvements in incremental steps, not substantial improvements that are related to the additional contextual information that is available to the owners of the data, the companies. He pointed out that real-world information and context, such as user addresses and other profile information, could potentially be incorporated into a traditional recommendation system.

Lebanon concluded with a brief discussion of implicit ratings. In the real world, one often has implicit, binary-rating data, such as whether a purchase or an impression was made. Evaluating that type of binary-rating data requires a different set of tools and models, and scaling up from standard data sets to industry data sets remains challenging.

### 3 PRINCIPLES FOR WORKING WITH BIG DATA 13

Important Points Made by Individual Speakers

- MapReduce is an important programming method designed for easy parallel programming on commodity hardware. (Jeffrey Ullman)
- There is an expertise gap between domain scientists and data scientists: domain scientists do not know what is possible technically, and data scientists do not understand the domain. (Juliana Freire)
- A data scientist should have expertise in databases, machine learning and statistics, and visualization; it is challenging, and perhaps unrealistic, to find people who have expertise in all three. (Juliana Freire and other discussion participants)
- Data preparation is an important, time-consuming, and often overlooked step in data analysis, and too few people are trained in it. (Juliana Freire)

Through better understanding of the tools and techniques used to address big data, one can better understand the relevant

education and training needs. The third session of the workshop focused more specifically on how to work with big data. Presentations were made by Jeffrey Ullman (Stanford University), Alexander Gray (Skytree Corporation), Duncan Temple Lang (University of California, Davis), and Juliana Freire (New York University). The session was chaired by Brian Caffo (Johns Hopkins University).

## Teaching about MapReduce

Jeffrey Ullman, Stanford University

MapReduce (Dean and Ghemawat, 2004), explained Jeffrey Ullman, is a programming method designed for easy parallel programming on commodity hardware, and it eliminates the need for the user to implement the parallelism and to address recovery from failures. MapReduce uses a distributed file system that replicates chunks to protect against data loss, and it is architected so that hardware failures do not require that the job be restarted. Hadoop [1](#) is an open-source implementation of MapReduce, which is proprietary to Google.

MapReduce, Ullman said, consists of a map function and a reduce function. The map function converts a single element (such as a document, integer, or information record) into key-value pairs. The map tasks are executed in parallel; the code is sent to the data, and the task executes wherever chunks of input are. After the map function has been applied to all inputs, the key-value pairs are sorted by key. The reduce function takes a single key with its list of associated values and provides an output. Reduce tasks are also executed in parallel, and each key with its list of inputs is handled independently.

Ullman then described a data mining course being taught at Stanford University in which students are given access to Amazon Web Services, and many do choose to implement their algorithms by using Hadoop. The course uses real-world data from a variety of sources, including Twitter, Wikipedia, and other companies. Teams of three students propose projects, including the data set to use, the expected results, and how to evaluate their results. About a dozen teams are selected to participate in the course.

Ullman described a 2012 team project on drug interactions. The team used data from Stanford's medical school from which it extracted records for 3,000 drugs. It sought to identify drug interactions and examine each pair of drugs with a chisquared test, a statistical test to evaluate the likelihood that differences in data arise by chance. The team was able to identify 40 of the 80 known drug combinations that lead to an increased risk of heart attack. More important, it identified two previously unknown pairs on which there was very strong evidence of interaction. Ullman explained that the team recognized that to make the problem more tractable, it needed to address it with fewer keys and longer lists of values, and it combined the drugs into groups, thereby reducing the number of comparisons and correspondingly reducing the amount of network-use time needed. Ullman stated that this example illustrated how communication time can often be the bottleneck in MapReduce algorithms.

Ullman then spoke more broadly about the theory of MapReduce models. Such models require three elements:

- Reducer [2](#) size: the maximum number of inputs that a given reducer can have, which leads to an upper bound on the length of the value list.
- Replication rate: the average number of key-value pairs generated by a mapper on one input. This measures communication cost per input; it is common for the replication rate to measure the length of time needed to run the algorithm.
- Mapping schema: a description of how outputs for a problem are related to inputs or an assignment of inputs to sets of reducers. No reducer is assigned more inputs than the reducer size; and for every output, there is some reducer that receives all the inputs associated with it.

Ullman showed that replication rate is inversely proportional to reducer size; this forces a trade-off between the two variables and provides a bound on replication rate as a function of reducer size. Ullman pointed out that the inverse relationship makes sense: when more work is done by a single reducer, less parallelism is needed, and the communication cost becomes smaller.

## Big Data Machine Learning—Principles for Industry

Alexander Gray, Skytree Corporation

Alexander Gray began by briefly describing the first three phases of machine learning: artificial intelligence and pattern recognition (1950s-1970s), neural networks and data mining (1980s and 1990), and convergence of machine learning with statistics (middle 1990s to today). Gray considers that we are now seeing the beginning of a fourth phase, defined by big data with new scalable systems needed to support it.

Gray explained that almost every industry has big data and would be better served by understanding it. He noted a variety of situations in which machine learning is “mission-critical”; in general, this occurs when some extreme is needed, such as high volume, high speed, or extreme accuracy. Gray described a number of kinds of applications of big data, including science (the Search for Extra-Terrestrial Intelligence, the Sloan Digital Sky Survey, and the Large Hadron Collider), medicine (health-care cost reduction, predictive health, [3](#) and early detection), finance (improving derivative pricing, risk analysis, portfolio optimization, and algorithmic trading), and security (cybersecurity, crime prevention, and antiterrorism). In addition, Gray noted kinds of applications that he described as having lower stakes: recommendations, face tagging, dating matches, and online advertising. He posited that many companies would benefit from machine learning to compete and ultimately to survive.

Gray then asked how to maximize predictive accuracy and explained that overall prediction error decomposes into errors that result from the use of finite samples, the choice of model parameters (i.e., algorithmic accuracy), and the choice of models. He noted that one can increase computational speed by orders of magnitude by using smarter algorithms. In addition, speed is connected to accuracy in that speed allows the analyst more time to explore the parameter space. Gray then described weak and strong scaling, a high-performance computing concept that manages data either by using more machines (strong scaling) or by taking more time (weak scaling). With data sets that contain millions of items, parallelism can provide good scaling—for example, changing from one computer to five computers might lead to a 5-fold speed increase in calculation. Gray indicated that data sets that contain billions of items are not uncommon and said that his firm has worked with one client that had data sets that contained trillions of items. Gray noted that strong and weak scaling result in different errors.

In addressing algorithmic accuracy, Gray pointed out that stochastic methods are optimal but generally do not reach optimal results in a single iteration. That type of computation is useful for “quick and dirty” applications. In addressing model error, Gray emphasized the importance of understanding and using a variety of models, as the best model changes on the basis of the data set. He also indicated that the treatment of outliers can change the outcome of an analysis. And he pointed out the utility of visualizing data in a data-specific and domainspecific approach and indicated a need for improved exploratory data analysis and visualization tools. A workshop participant supported the use of visualization and emphasized the need to include the human in the loop; the user should be responsible for and involved in the visualization, not passive, and the visualization should enhance understanding of the data.

## Principles for the Data Science Process

Duncan Temple Lang, University of California, Davis

Duncan Temple Lang began by listing the core concepts of data science— items that will need to be taught: statistics and machine learning, computing and technologies, and domain knowledge of each problem. He stressed the importance of interpretation and reasoning—not only methods—in addressing data. Students who work in data science will have to have a broad set of skills—including knowledge of randomness and uncertainty, statistical methods, programming, and technology—and practical experience in them. Students tend to have had few computing and statistics classes on entering graduate school in a domain science.

Temple Lang then described the data analysis pipeline, outlining the steps in one example of a data analysis and exploration process:

1. Ask a general question.
2. Refine the question, identify data, and understand data and metadata. Temple Lang noted that the data used are usually not collected for the specific question at hand, so the original experiment and data set should be understood.
3. Access data. This is unrelated to the science but does require computational skill.
4. Transform to data structures.
5. Perform exploratory data analyses to understand the data and determine whether the results will scale. This is a critical step; Temple Lang noted that 80 percent of a data scientist’s time can be spent in cleaning and preparing the data.
6. Perform dimension reduction. Temple Lang stressed that it can be difficult or impossible to automate this step.
7. Perform modeling and estimation. Temple Lang noted that computer and machine learning scientists tend to focus more on predictive models than on modeling of physical behavior or characteristics.
8. Perform diagnostics. This helps to understand how well the model fits the data and identifies anomalies and aspects

for further study. This step has similarities to exploratory data analysis.

9. Quantify uncertainty. Temple Lang indicated that quantifying uncertainty with statistical techniques is important for understanding and interpreting models and results.

10. Convey results.

Temple Lang stressed that the data analysis process is highly interactive and iterative and requires the presence of a human in the loop. The next step in data processing is often not clear until the results of the current step are clear, and often something unexpected is uncovered. He also emphasized the importance of abstract skills and concepts and said that people need to be exposed to authentic data analyses, not only to the methods used. Data scientists also need to have a statistical understanding, and Temple Lang described the statistical concepts that should be taught to a student:

- Mapping the general question to a statistical framework.
- Understanding the scope of inference, sampling, biases, and limitations.
- Exploratory data analyses, including missing values, data quality, cleaning, matching, and fusing.
- Understanding randomness, variability, and uncertainty. Temple Lang noted that many students do not understand sampling variability.
- Conditional dependence and heterogeneity.
- Dimension reduction, variable selection, and sparsity.
- Spurious relationships and multiple testing.
- Parameter estimation versus “black box” prediction and classification.
- Diagnostics—residuals and comparing models.
- Quantifying the uncertainty of a model.
- Sampling structure and dependence for data reduction. Temple Lang noted that modeling of data becomes complicated when variables are not independent, identically distributed.
- Statistical accuracy versus computational complexity and efficiency.

Temple Lang then briefly discussed some of the practical aspects of computing, including the following:

- Accessing data.
- Manipulating raw data.
- Data structures and storage, including correlated data.
- Visualization at all stages (particularly in exploratory data analyses and conveying the results).
- Parallel computing, which can be challenging for a new student.
- Translating high-level descriptions to optimal programs.

During the discussion, Temple Lang proposed computing statistics on visualizations to examine data rigorously in a statistical and automated way. He explained that “scagnostics” (from scatter plot diagnostics) is a data analysis technique for graphically exploring the relationships among variables. A small set of statistical measures can characterize scatter plots, and exploratory data analysis can be conducted on the residuals. [4](#)

A workshop participant noted the difference between a data error and a data blunder. A blunder is a large, easily noticeable mistake. The participant gave the example of shipboard observations of cloud cover; blunders, in that case, occur when the location of the ship observation is given to be on land rather than at sea. Another blunder would be a case of a ship’s changing location too quickly. The participant speculated that such blunders could be generalized to detect problematic observations, although the tools would need to be scalable to be applied to large data sets.

## Principles for Working with Big Data

Juliana Freire, New York University

Juliana Freire began her presentation by discussing the tasks involved in addressing big data. She referred to a Computing Research Association (CRA) report [5](#) on the challenges posed by big data. CRA also documented the data analysis pipeline, which includes acquisition and recording; extraction, cleaning, and annotation; analysis and modeling; and interpretation. A simplified schematic of the pipeline is shown in Figure 3.1.

Freire posited that scaling for batch computation is not difficult—people have been working on this problem for several decades, and there is an infrastructure to support it. However, the human scalability is difficult; as the data size increases, it becomes more difficult for an analyst to explore the data space. The path from data to knowledge, she noted, is human-based

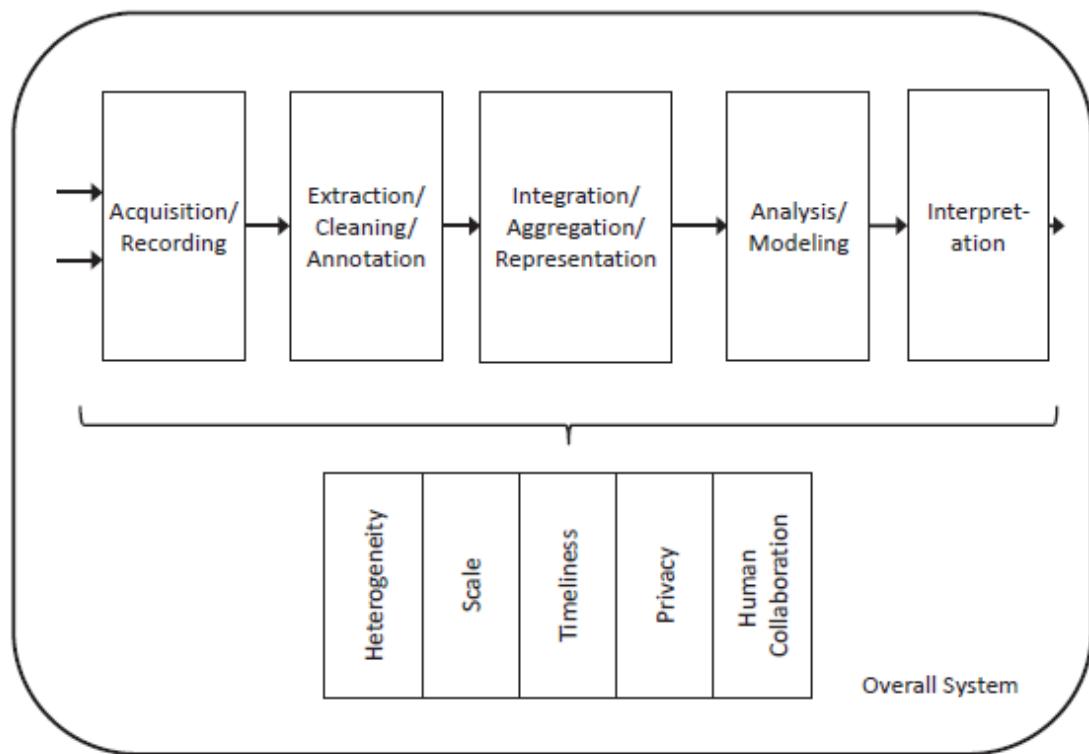
and has many complicated elements.

Freire explained that the CRA data analysis pipeline tasks can be classified into two categories: data preparation (which includes acquisition and recording; extraction, cleaning, and annotation; and integration, aggregation, and representation) and data analysis (which includes modeling and interpretation). Data science includes statistics, machine learning, data mining, and visualization, but Freire noted that in many institutions it is synonymous with machine learning, and less emphasis is placed on the other elements. She pointed out that data visualization has been growing in importance and that there is a corresponding need for additional training in it. Freire emphasized that the data pipeline is complex and that what is shown in Figure 3.1 is an oversimplification; for instance, the pipeline is not linear. She also stressed the importance of research provenance: provenance of the exploration process should be captured for transparency, reproducibility, and knowledge reuse. She noted that provenance management is not often taught.

Freire acknowledged that people underestimate the effort required in preparing data. Few people have the expertise to prepare data, but there is a high demand for data preparation. In contrast, there are many experts to conduct the analysis, but relatively little time is needed for this step. She stated that data preparation takes a long time, is idiosyncratic, and can limit analyses. She also noted that new data sets continually provide new challenges in big data, and many needs are not met by existing infrastructure.

**FIGURE 3.1 Simplified schematic of the big data analysis pipeline**

Major steps in the analysis of big data are shown in the flow at top. Below it are big data needs that make these steps challenging. SOURCE: Computing Community Consortium, February 2012.



**FIGURE 3.1 Simplified schematic of the big data analysis pipeline.** Major steps in the analysis of big data are shown in the flow at top. Below it are big data needs that make these steps challenging. SOURCE: Computing Community Consortium, February 2012.

Freire then provided an example of recent work in applying data science principles to New York City taxis. The raw data set consisted of 500,000 trips per day taken for more than 3 years, which yielded 150 GB of data. The data were not enormous, but they were complex and had spatial and temporal attributes. The data show an unusual degree of regularity; one can easily see temporal changes related to weekends and holidays. The goal was to allow city officials to explore the data visually. The work involved developing a spatiotemporal index that was based on an out-of-core k-dimensional tree (Ferreira et al., 2013) and a

new interactive map view.

Freire stated that domain scientists do not know what is possible to do with their data, and technologists do not know the domain, so there is an expertise gap. Freire quoted Alex Szalay (Faris et al., 2011), who described the ideal scientist as “π-shaped,” with deep knowledge in two fields and connections between them. Freire argued that although the data scientist is supposed to fill the expertise gap, in reality three people make up a “data scientist”: a database expert, a machine learning and statistics expert, and a visualization expert. She said that computer science and data management research have partly failed in that it has not been able to create usable tools for end users. Freire stated that the complexity of data science problems is often underestimated.

Freire was asked by a workshop participant how to prepare students in software while teaching them their domain science. She suggested adding a new course for students who do not have a computer science background. She noted that there were several boot-camp-style programs for Ph.D. science students but that their overall effectiveness is not known.

Participants also discussed the requirements for a data analyst, a topic discussed by Temple Lang during his presentation. One person posited that the single expert in databases, machine learning and statistics, and visualization that Freire described should also be knowledgeable in systems and tools. The database expertise should include computational environments, not just databases. Another participant described the data analyst as a “jazz player” rather than a “symphony player”—in other words, a data analyst should improvise and make decisions rapidly, which cannot be done if the analyst does not know the subject matter well.

Some participants discussed tools. One person noted that commercial tools (such as Spotfire [6](#) and Tableau [7](#)) exist in a polished form and work in a variety of applications. Others responded, however, that students need training on these tools, and that a single tool does not usually solve complex data problems. A participant noted that students cannot afford a subscription to Tableau and argued that the existing tools should be open-source; however, open-source tools may not always be well curated.

## 4 COURSES, CURRICULA, AND INTERDISCIPLINARY PROGRAMS

Important Points Made by Individual Speakers

- A residual effect of training students to work with data is that the training will empower them with a toolkit that they can use in multiple domains. (Joshua Bloom)
- Boot camps and other short courses appear to be successful in teaching data computing techniques to domain scientists and in addressing a need in the science community; however, outstanding questions remain about how to integrate these types of classes into a traditional educational curriculum. (Joshua Bloom)
- Educators should be careful to teach data science methods and principles and avoid teaching specific technologies without teaching the underlying concepts and theories. (Peter Fox)
- Massive online open courses (MOOCs) are one avenue for teaching data science techniques to a large population; thus far, data science MOOC participants tend to be computer science professionals, not students. (William Howe)

By the end of 2014, more than 30 major universities will have programs in data science. [1](#) Existing and emerging programs offer many opportunities for lessons learned and potential course and content models for universities to follow. The fourth workshop session focused on specific coursework, curricula, and interdisciplinary programs for teaching big data concepts. The session was chaired by James Frew (University of California, Santa Barbara). Presentations were made in this session by Joshua Bloom (University of California, Berkeley), Peter Fox (Rensselaer Polytechnic Institute), and William Howe (University of Washington).

### Computational Training and Data Literacy for Domain Scientists

Joshua Bloom, University of California, Berkeley

Joshua Bloom noted that the purpose of graduate school is to prepare students for a career in the forefront of science. A residual effect of training students to work with data is that the training will empower the students with a toolkit that they can use even if they leave a particular domain. He pointed out that the modern datadriven science toolkit is vast and that students are being asked to develop skills in both the domain science and the toolkit.

Bloom then described upcoming data challenges in his own domain of astronomy. The Large Synoptic Survey Telescope is expected to begin operations in 2020, and it will observe 800 million astronomical sources every 3 days. A large computational

framework is needed to support that amount of data, probably 20 TB per night. Other projects in radio astronomy have similar large-scale data production.

A goal in data science for time-domain astronomy in the presence of increasing data rates is to remove the human from the real-time data loop, explained Bloom—in other words, to develop a fully automated, state-of-the-art scientific stack to observe transient events. Often, the largest bottleneck is in dealing with raw data, but there are large-scale inference challenges further downstream.

Bloom pointed out that the University of California, Berkeley, has a long history of teaching parallel computing. The coursework is aimed at computer science, statistics, and mathematics students. Recently, “boot camps” that include two or three intensive training sessions have been initiated to teach students basic tools and common frameworks. The boot camp at Berkeley takes 3 full days. There are six to eight lectures per day, and hands-on programming sessions are interspersed. Bloom began teaching computing techniques to domain scientists, primarily physicalscience students. His first boot camp consisted of several all-day hands-on classes with nightly homework. The student needed to know a programming language before taking the boot camp. In 2010, the first year, 85 students participated. By 2013, the boot camp had grown to more than 250 students. Bloom uses live streaming and archiving of course material, and all materials used are open-source. The course has been widely used and repeated; for instance, NASA Goddard Space Flight Center used his materials to hold its own boot camp. Bloom noted, in response to a question, that instructors in his course walk around the room to assist students while they work. He posited that 90 percent of that interaction could be replaced with a well-organized chat among instructors and students; the course would probably take longer, and students would have to be self-directed.

Bloom explained that the boot camp is a prerequisite to Berkeley’s graduate-level follow-on seminar course in Python computing for science. The graduate seminar was the largest graduate course ever taught in Berkeley’s astronomy department; this indicated an unmet need for such a course at the graduate science level. Bloom said that the boot camps and seminars give rise to a set of education questions: Where do boot camps and seminars fit into a traditional domain-science curriculum? Are they too vocational or practical to be part of higher-education coursework? Who should teach them, and how should the instructors be credited? How can students become more (broadly) data literate before we teach them big data techniques? He emphasized that at the undergraduate level the community should be teaching “data literacy” before it teaches data proficiency. Some basic data-literacy ideas include the following:

- Statistical inference. Bloom noted that this is not necessarily big data; something as simple as fitting a straight line to data needs to be taught in depth.
- Versioning and reproducibility. Bloom noted that several federal agencies are likely to mandate a specific level of reproducibility in work that they fund.

Bloom suggested that there is a “novelty-squared” problem: what is novel in the domain science may not be novel in the data science methodology. He stressed the need to understand the forefront questions in various fields so that synergies can be found. For example, Berkeley has developed an ecosystem for domain and methodological scientists to talk and find ways to collaborate.

Bloom also noted that data science tends to be an inclusive environment that appeals to underrepresented groups. For instance, one-third of the students in the Python boot camps were women—a larger fraction than their representation in physical science graduate programs.

Bloom concluded by stating that domain science is increasingly dependent on methodologic competences. The role of higher education in training in data science is still to be determined. He stressed the need for data literacy before data proficiency and encouraged the creation of inclusive and collaborative environments to bridge domains and methodologies.

Bloom was asked what he seeks in a student. He responded that it depends on the project. He looked for evidence of prior research, even at the undergraduate level, as well as experience in programming languages and concepts. However, he noted that a top-quality domain scientist would always be desirable regardless of computational skills.

A participant commented that as much as 80 percent of a researcher’s time is spent in preparing the data. That is a large amount of time that could be spent on more fundamental understanding. Bloom responded that such companies and products as OpenRefine, [2](#) Data Wrangler, [3](#) and Trifacta [4](#) are working on data cleaning techniques. However, for any nontrivial question, it is difficult to systematize data preparation. He also suggested that a body of fundamental research should be accessible to practitioners. However, large-scale, human-generated data with interesting value do not typically flow to academe because of privacy and security concerns. He conjectured that the advent of the Internet of Things will allow greater data access, because those data will not be human data and therefore will have fewer privacy concerns.

## Data Science and Analytics Curriculum Development at Rensselaer (and the Tetherless World Constellation)

Peter Fox, Rensselaer Polytechnic Institute

Peter Fox began by describing the Tetherless World Constellation [5](#) at Rensselaer Polytechnic Institute (RPI). The research themes are divided loosely into three topics: future Web (including Web science, policy, and social issues), Xinformatics (including data frameworks and data science), and semantic foundations (including knowledge provenance and ontology engineering environments). Fox indicated that his primary focus is Xinformatics. He deliberately did not define X, saying that it can mean any number of things.

Fox explained that to teach data science, one must “pull apart” the ecosystem in which data science lives. Data, information, and knowledge are all related in a data science ecosystem; there is no linear pathway from data to information to knowledge. He explained that he teaches or is involved in classes on data science, Xinformatics, geographic information systems for the sciences, semantic eScience, data analytics, and semantic technologies. The students in those classes have varied backgrounds. Last semester, his data science class had 63 students (most of them graduate students), and Xinformatics had about 35 students. Fox structures his classes so that the first half of the semester focuses on individual work and gaining knowledge and skills. The second half focuses on team projects (with teams assigned by him) to demonstrate skills and share perspectives.

Fox explained that he teaches modern informatics and marries it with a method: the method that he teaches is iterative and is based on rapid prototyping applied to science problems. The framework for the iterative model is shown in Figure 4.1. Fox stressed that technology does not enter the method until well over halfway through the spiral; technology will change, so it is important to impart skills before adopting and leveraging technology.

Fox explained that a report was produced for NSF (Borgman et al., 2008) and that a diagram was developed that describes five types of mediation of increasing complexity (shown in Figure 4.2). The five generations of mediation were designed to apply to learning, but they hold true for science and research as well. Fox explained that, in contrast with most generational plots, all generations are present and active at once in both the learning and teaching environment and the science and research environment.

Fox explained that data analytics is a new course at RPI, and his desired prerequisites are not taught at the university; as a result, his class has no prerequisites. After teaching a variety of computer application languages simultaneously, Fox now uses the R and RStudio [6](#) environment exclusively. (Students preferred the simplicity of learning a single language.) The data analytics course builds from data to processing to reporting to analytics, both predictive and prescriptive. Fox explained that in the ideal scenario, value would be added as one progresses from one step to the next. Part of the challenge is to teach students to understand the value added and to teach them to identify when value is not being added. He emphasized the importance of understanding the value and application of data analysis, not just learning the individual steps. In response to a later question, Fox clarified that students work with self-selected, application-specific examples.

### FIGURE 4.1 Framework for modern informatics

A technology approach does not enter into the development spiral until over halfway through the process. SOURCE: Fox and McGuinness (2008), [http://tw.rpi.edu/media/latest/SemanticMethodologyPathwayPretty\\_v2.png](http://tw.rpi.edu/media/latest/SemanticMethodologyPathwayPretty_v2.png).

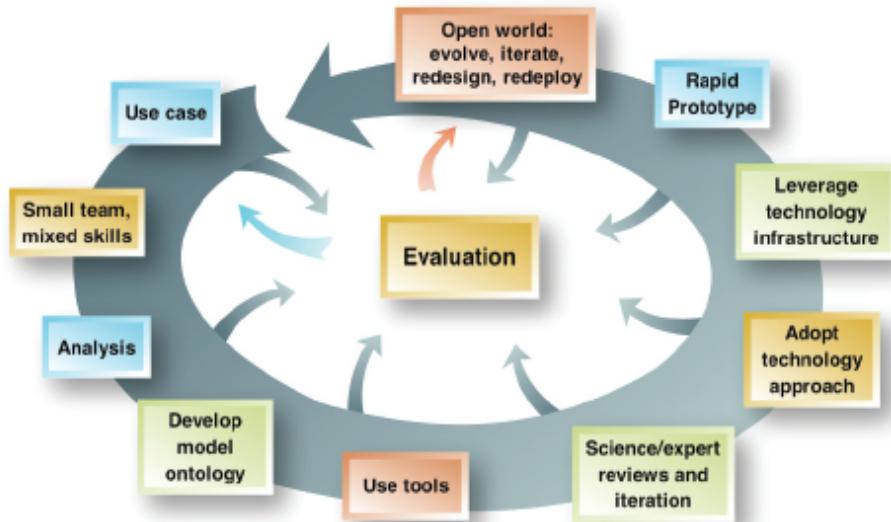
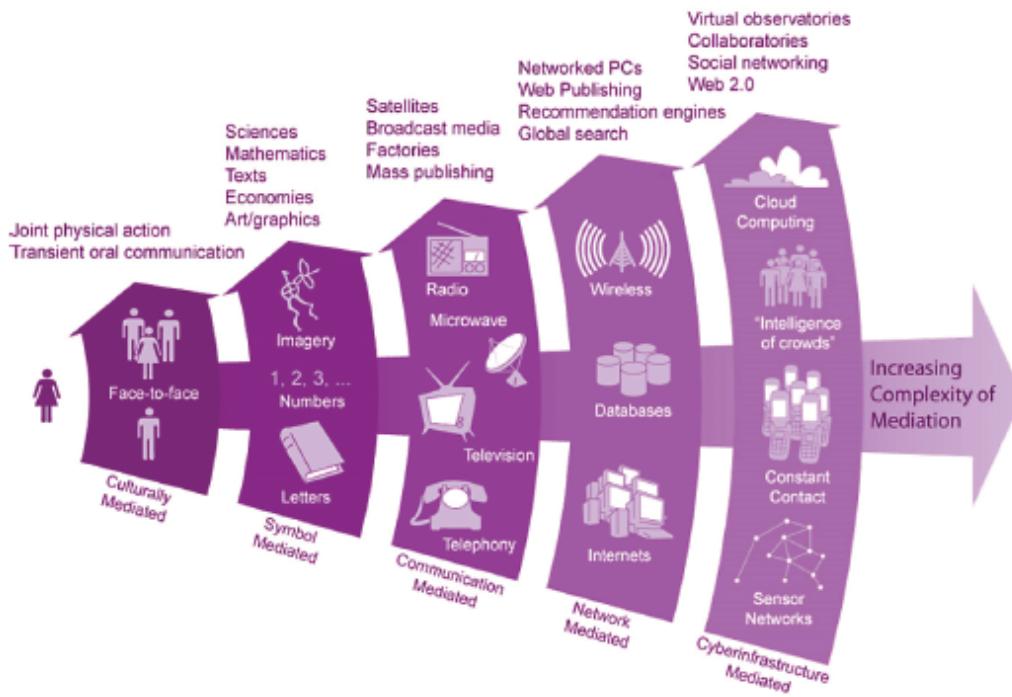


FIGURE 4.1 Framework for modern informatics. A technology approach does not enter into the development spiral until over halfway through the process. SOURCE: Fox and McGuinness (2008), [http://tw.rpi.edu/media/latest/SemanticMethodologyPathwayPretty\\_v2.png](http://tw.rpi.edu/media/latest/SemanticMethodologyPathwayPretty_v2.png).

FIGURE 4.2 Generations of mediation

applied to the learning and teaching environment and the science and research environment. SOURCE: Illustration by Roy Pea and Jillian C. Wallis, from Borgman et al. (2008).



**FIGURE 4.2** Generations of mediation, applied to the learning and teaching environment and the science and research environment. SOURCE: Illustration by Roy Pea and Jillian C. Wallis, from Borgman et al. (2008).

Fox then described information technology and Web-science coursework at RPI. RPI has an interdisciplinary program that consists of a B.S. with 20 concentrations, an M.S. with 10 concentrations, and a multidisciplinary Ph.D. offering. [Z](#) The program has four technical tracks—computer engineering, computer science, information systems, and Web science—with numerous concentrations in each track. Fox said that the M.S. was recently revised to include data analytics in the core curriculum and that the M.S. concentrations were updated. He noted in particular the addition of the Information Dominance concentration, which is designed to educate a select group of naval officers each year in skills needed to execute military cyberspace operations.

Fox talked about the Data Science Research Center [8](#) and the less formal Data Science Education Center at RPI. The centers are loosely organized; more than 45 faculty and staff are involved. RPI also maintains a data repository [9](#) for scientific data that result from on-campus research.

He listed some lessons learned after 5 years with these programs:

- Be interdisciplinary from the start; grow both technical and data skills simultaneously. Fox noted that teaching skills (such as how to manipulate data by using specific programming languages) can be difficult; skills need to be continually reinforced, and he cautioned that teaching skills may be perceived as training rather than education.
- Teach methods and principles, not technology.
- Make data science a skill in the same vein as laboratory skills.
- Collaboration is critical, especially in informatics.
- Teach foundations and theory.

Fox stated that access to data is progressing from provider-to-user to machine-to-user and finally to machine-to-machine; the burden of data access and usability shifts from the user to the provider. In the current research-funding paradigm, data are collected, data are analyzed by hand for several years, and the results are then published. Although that paradigm has served the research community well, Fox noted that it fails to reflect the change in responsibilities that is inherent in the new information era, in which the burden of access shifts from the user to the provider.

Fox concluded by positing that the terms data science and metadata will be obsolete in 10 years as researchers come to work with data as routinely as they use any other research tool.

Bloom noted that there was no mention of "big data" in Fox's presentation, only data. Fox stated that he does not distinguish big data from data. However, he acknowledged that, as a practical matter, size, heterogeneity, and structural representations will need to be parts of a student's course of study.

## Experience with a First Massive Online Open Course on Data Science

William Howe, University of Washington

William Howe stated that the University of Washington (UW) founded the eScience Institute in 2008 and that the institute is now engaged in a multi-institution partnership with the University of California, Berkeley, and New York University and is funded by the Gordon and Betty Moore Foundation and the Alfred P. Sloan Foundation to advance new data science techniques and technologies, foster collaboration, and create a cross-campus "data science environment." According to Howe, the strategy is to establish a "virtuous cycle" between the data science methodology researchers and the domain-science researchers in which innovation in one field will drive innovation in the other. The eScience Institute works to create and reinforce connections between the two sides, and six working groups act as bridges. One of the working groups is involved with education and training.

Howe explained that the education and training working group focuses on different ways to educate students and practitioners in data science. Through the working group, the UW eScience Institute has developed a data science certificate for working professionals, an interdisciplinary Ph.D. track in big data, new introductory courses, a planned data science master's degree, and a MOOC called "Introduction to Data Science." Howe focused in more detail on his experiences in developing and teaching the MOOC. Teaching a MOOC involves a large amount of work, he said, and the course is continuously developing. The goal of the MOOC is to organize a set of important topics spanning databases, statistics, machine learning, and visualization into a single introductory course. He provided some statistics on the data science MOOC:

- More than 110,000 students registered for the course. Howe noted that that is not a particularly relevant statistic, inasmuch as many people who register for MOOCs do not necessarily plan to participate.
- About 9,000 students completed the course assignments. Howe indicated that that is a typical level of attrition for a MOOC.
- About 7,000 students passed the course and earned the certificate.

He explained that the course had a discussion forum that worked well. Many comments were posted to it, and it was self-sustaining; Howe tried to answer questions posed there but found that questions were often answered first by other engaged students.

The syllabus that defined his 9-week MOOC consisted of the following elements:

- Background and scope of "data science."
- Data manipulation at scale.
- Analytics. Howe taught selected statistics concepts in 1 week and machine learning concepts in another week.
- Visualization.
- Graph and network analytics. Howe indicated this was a single, short module.

Howe explained that the selection of topics was motivated by a desire to develop the four dimensions of the course: tools versus abstractions (weighted toward abstractions), desktop versus cloud (weighted toward cloud), hackers versus analysts (balanced, although perhaps slightly in favor of hackers), and data structures and programming versus mathematics and statistics (weighted toward structures).

He conducted a demographic study of his MOOC participants and remarked that most of them were working professional software engineers, as has been reported for other MOOCs. He suggested that perhaps a MOOC could be used like a textbook, with instructors having their students watch some lectures and skip others, just as they do chapters of a book.

A teaching strategy that consists of both online and in-person components, he explained, has two possible approaches: offer the same course simultaneously online and in person or use the online component as a textbook and class time as an opportunity for practical application (juxtaposing the traditional roles of homework and classwork). There are examples of both teaching strategies, and it is unclear whether either will dominate. He also reiterated the importance of student-to-student learning that took place in his experience with a MOOC structure.

In the discussion period, a participant asked about the importance of understanding the foundations of programming and suggested that algorithms and data constructs should be taught to younger students, for example, in high school. (That last point generated some disagreement among participants. One suggested that even elementary school would be appropriate, but another was concerned that this might displace students from calculus and other critical engineering mathematics courses.) Howe replied that computer science enrollment is increasing at the undergraduate level by about 20 percent per year, and statistics departments are also seeing increased enrollment. Students understand that they need to understand the core concepts at the undergraduate level.

A workshop participant asked about other MOOC success stories in big data. Howe responded that he had only anecdotal evidence. Bloom concurred, noting that a graduate student who had participated in Berkeley's data science boot camp conducted large-scale parallel computing work that resulted in a seminal paper in his field (Petigura et al., 2014).

## 5 SHARED RESOURCES

### Important Points Made by Individual Speakers

- Synthetic knowledge bases for domain sciences, such as the PaleoDeepDive system at Stanford University, can be developed by using the automatic extraction of data from scientific journals. (Christopher Ré)
- Divide and recombine methods are powerful tools for analysts who conduct deep examinations of data, and such systems can be used by analysts without the need for complex programming or a profound understanding of the tools used. (Bill Cleveland)
- Yahoo Webscope is a reference library of large, scientifically useful, publicly available data sets for researchers to use. (Ron Brachman)
- Amazon Web Services (AWS) hosts large data sets in a variety of models (public, requestor-paid, and private or community) to foster large-scale data sharing. AWS also provides data computation tools, training programs, and grants. (Mark Ryland)

The fifth session of the workshop was chaired by Deepak Agarwal (LinkedIn Corporation). The session had four speakers: Christopher Ré (Stanford University), Bill Cleveland (Purdue University), Ron Brachman (Yahoo Labs), and Mark Ryland (Amazon Corporation).

### Can Knowledge Bases Help Accelerate Science?

Christopher Ré, Stanford University

Christopher Ré focused on a single topic related to data science: knowledge bases. He first discussed Stanford University's experience with knowledge bases. Ré explained that in general scientific discoveries are published and lead to the spread of ideas. With the advent of electronic books, the scientific idea base is more accessible than ever before. However, he cautioned, people are still limited by their eyes and brain power. In other words, the entire science knowledge base is accessible but not necessarily readable.

Ré noted that today's science problems require macroscopic knowledge and large amounts of data. Examples include health, particularly population health; financial markets; climate; and biodiversity. Ré used the latter as a specific example: broadly speaking, biodiversity research involves assembling information about Earth in various disciplines to make estimates of species extinction. He explained that this is "manually constructed" data—a researcher must input the data by examining and collating information from individual studies. Manually constructed databases are time-consuming to produce; with today's data sources, the construction exceeds the time frame of the typical research grant. Ré posited that the use of samplebased data and their synthesis constitute the only way to address many important questions in some fields. A system that synthesizes sample-based data could "read" journal articles and automatically extract the relevant data from them. He stated that "reading" machines may be coming in the popular domain (from such Web companies as IBM, Google, Bing, and Amazon). The concept of these machines could be extended to work in a specific scientific domain. That would require highquality reading—reading of a higher quality than is needed in a popular-domain application—in that mistakes can be more harmful in a scientific database.

Ré described a system that he has developed, PaleoDeepDive, [1](#) a collaborative effort with geologist Shanan Peters, also of Stanford University. The goal of PaleoDeepDive is to build a higher-coverage fossil record by extracting paleobiologic facts from research papers. The system considers every character, word, or fragment of speech from a research paper to be a variable and then conducts statistical inference on billions of variables defined from the research papers to develop relationships between biologic and geologic research. PaleoDeepDive has been in operation for about 6 months, and preliminary results of occurrence

relations extracted by PaleoDeepDive show a precision of around 93 percent; Ré indicated that this is a very high-quality score.

Ré then stated the challenges for domain scientists related to synthetic knowledge bases:

- Students are not trained to ask questions of synthetic data sets. Ré noted that this may be changing; the University of Chicago, for instance, includes such training in its core curriculum. Stanford has an Earth-science class on how to use PaleoDeepDive.
- Students lack skills in computer science and data management. Ré indicated that this also may be changing; 90 percent of Stanford students now take at least one computer science class.
- Some people are skeptical of human-generated synthetics. Ré suggested that this is also changing as statistical methods take stronger hold.

Ré noted challenges for computer scientists related to synthetic knowledge bases:

- Finding the right level of abstraction. Ré posited that approaches to many interesting questions would benefit from the use of synthetic knowledge bases. However, PaleoDeepDive is not necessarily scalable or applicable to other disciplines.
- Identifying features of interest. Computer scientists, Ré noted, focus on algorithms rather than on features. However, synthetic knowledge bases are feature based and require a priori knowledge of what is sought from the data set.

A participant noted that noise, including misspelled words and words that have multiple meanings, is a standard problem for optical character recognition (OCR) systems. Ré acknowledged that OCR can be challenging and even state-of-the-art OCR systems make many errors. PaleoDeepDive uses statistical inference and has a package to improve OCR by federating open-source material together and using probabilistic inputs. Ré indicated that Stanford would be releasing tools to assist with OCR.

## Divide and Recombine for Large, Complex Data

Bill Cleveland, Purdue University

Bill Cleveland explained the goals of divide and recombine for big data:

- Methods and environments that do not require reductions in dimensionality should be used for analyzing data at the finest level of granularity possible. The data analysis could include visualization.
- At the front end, analysts can use a language for data analysis to tailor the data and make the system efficient.
- At the back end, a distributed database is accessible and usable without the need for the analyst to engage in the details of computation.
- Within the computing environment, there is access to the many methods of machine learning and visualization.
- Software packages enable communication between the front and back ends.
- The system can be used continuously to analyze large, complex data sets, generate new ideas, and serve as a test bed.

Cleveland then described the divide and recombine method. He explained that a division method is used first to divide the data into subsets. The subsets are then treated with one of two categories of analytic methods:

- Number-category methods. Analytic methods are applied to each of the subsets with no communication among the computations. The output from this method is numeric or categoric.
- Visualization. The data are organized into images. The output from this method is plots. It is not feasible to examine all the plots, so the images are sampled. That can be done rigorously; sampling plans can be developed by computing variables with one value per subset.

Cleveland described several specific methods of division. In the first, conditioning-variable division, the researcher divides the data on the basis of subject matter regardless of the size of the subsets. That is a pragmatic approach that has been widely used in statistics, machine learning, and visualization. In a second type of division, replicate division, observations are exchangeable, and no conditioning variables are used. The division is done statistically rather than by subject matter. Cleveland stated that the statistical division and recombination methods have an immense effect on the accuracy of the divide and recombine result. The statistical accuracy is typically less than that with other direct methods. However, Cleveland noted that this is a small price to pay for the simplicity in computation; the statistical computation touches subsets no more than once. Cleveland clarified that the process is not MapReduce; statistical methods in divide and recombine reveal the best way to separate the data into subsets and put them back together.

Cleveland explained that the divide and recombine method uses R for the front end, which makes programming efficient. R saves the analyst time, although it is slower than other options. It has a large support and user community, and statistical packages are readily available. On the back end, Hadoop is used to enable parallel computing. The analyst specifies, in R, the

code to do the division computation with a specified structure. Analytic methods are applied to each subset or each sample. The recombination method is also specified by the analyst. Cleveland explained that Hadoop schedules the microprocessors effectively. Computation is done by the mappers, each with an assigned core for each subset. The same is true for the reducers; reducers carry out the recombination. The scheduling possibilities are complex, Cleveland said. He also noted that this technique is quite different from the high-performance computing systems that are prevalent today. In a high-performance computing application, time is reserved for batch processing; this works well for simulations (in which the sequence of steps is known ahead of time and is independent of the data), but it is not well suited to sustained analyses of big data (in which the process is iterative and adaptive and depends on the data).

Cleveland described three divide and recombine software components between the front and back ends. They enable communication between R and Hadoop to make the programming easy and insulate the analyst from the details of Hadoop. They are all open-source.

- R and Hadoop Integrated Programming Environment (RHIPE [2](#)). This is an R package available on GitHub. [3](#) Cleveland noted that RHIPE can be too strenuous for some operating systems.
- Datadr. [4](#) Datadr is a simple interface for division, recombination, and other data operations, and it comes with a generic MapReduce interface.
- Trelliscope. [5](#) This is a trellis display visualization framework that manages layout and specifications; it extends the trellis display to large, complex data.

Cleveland explained that divide and recombine methods are best suited to analysts who are conducting deep data examinations. Because R is the front end, R users are the primary audience. Cleveland emphasized that the complexity of the data set is more critical to the computations than the overall size; however, size and complexity are often correlated.

In response to a question from the audience, Cleveland stated that training students in these methods, even students who are not very familiar with computer science and statistics, is not difficult. He said that the programming is not complex; however, analyzing the data can be complex, and that tends to be the biggest challenge.

## Yahoo's Webscope Data Sharing Program

Ron Brachman, Yahoo Labs

Ron Brachman prefaced his presentation by reminding the audience of a 2006 incident in which AOL released a large data set with 20 million search queries for public access and research. Unfortunately, personally identifiable information was present in many of the searches, and this allowed the identification of individuals and their Web activity. Brachman said that in at least one case, an outside party was able to identify a specific individual by cross-referencing the search records with externally available data. AOL withdrew the data set, but the incident caused shockwaves throughout the Internet industry. Yahoo was interested in creating data sets for academics around the time of the AOL incident, and the AOL experience caused a slow start for Yahoo. Yahoo persisted, however, working on important measures to ensure privacy, and has developed the Webscope [6](#) data sharing program. Webscope is a reference library of interesting and scientifically useful data sets. It requires a license agreement to use the data; the agreement is not burdensome, but it includes terms whereby the data user agrees not to attempt to reverse-engineer the data to identify individuals.

Brachman said that Yahoo has just released its 50th Webscope data set. Data from Webscope have been downloaded more than 6,000 times. Webscope has a variety of data categories available, including the following:

- Language and content. These can be used to research information-retrieval and natural-language processing algorithms and include information from Yahoo Answers. (This category makes up 42 percent of the data in Webscope.)
- Graph and social data. These can be used to research matrix, graph, clustering, and machine learning algorithms and include information from Yahoo Instant Messenger (16 percent).
- Ratings, recommendation, and classification data. These can be used to research collaborative filtering, recommender systems, and machine learning algorithms and include information on music, movies, shopping, and Yelp (20 percent).
- Advertising data. These can be used to research behavior and incentives in auctions and markets (6 percent).
- Competition data (6 percent).
- Computational-system data. These can be used to analyze the behavior and performance of different types of computer systems architectures, such as distributed systems and networks and include data from the Yahoo Sherpa database system (6 percent).
- Image data. These can be used to analyze images and annotations and are useful for image-processing research (less than 4 percent).

Brachman explained that in many cases there is a simple click-through agreement for accessing the data, and they can be downloaded over the Internet. However, downloads are becoming impractical as database size increases. Yahoo had been asking for hard drives to be sent through the mail; now, however, it is hosting some of its databases on AWS.

In response to questions, Brachman explained that each data set is accompanied by a file explaining the content and its format. He also indicated that the data provided by Webscope are often older (around a year or two old), and this is one of the reasons that Yahoo is comfortable with its use for academic research purposes.

Brachman was asked whether any models fit between the two extremes of Webscope (with contracts and nondisclosure agreements) and open-source. He said that the two extremes are both successful models and that the middle ground between them should be explored. One option is to use a trusted third party to hold the data, as is the case with the University of Pennsylvania's Linguistic Data Consortium data. [7](#)

## Resource Sharing

### Mark Ryland, Amazon Corporation

Mark Ryland explained that resource sharing can mean two things: technology capabilities to allow sharing (such as cloud resources) and economic and cost sharing, that is, how to do things less expensively by sharing. AWS is a system that does both. AWS is a cloud computing platform that consists of remote computing storage and services. It holds a large array of data sets with three types of product. The first is public, freely available data sets. These data sets consist of freely available data of broad interest to the community and include Yahoo Webscope data, Common Crawl data gathered by the open-source community (240 TB), Earthscience satellite data (40 TB of data from NASA), 1000 Genomes data (350 TB of data from NIH), and many more. Ryland stated that before the genome data were publicly stored in the cloud, fewer than 20 researchers worked with those data sets. Now, more than 200 are working with the genome data because of the improved access.

A second type of AWS data product is requestor-paid data. This is a form of cost-sharing in which data access is charged to the user account but data storage is charged to the data owner's account. It is fairly popular but perhaps not as successful as AWS would like it to be, and AWS is looking to broaden the program.

The third type of AWS data product is community and private. AWS may not know what data are shared in this model. The data owner controls data access. Ryland explained that AWS provides identity control and authentication features, including Web Identity Federation. He also described a science-oriented data service ([Globus 8](#)), which provides cloud-based services to conduct periodic or episodic data transfers. He explained that an ecosystem is developing around data sharing.

Sharing is also taking place in computation. Ryland noted that people are developing Amazon Machine Images with tools and data "prebaked" into them, and he provided several examples, including Neuroimaging Tools and Resources Clearinghouse and scientific Linux tools. Ryland indicated that there are many big data tools, some of which are commercial and some of which are open-source. Commercial tools can be cost-effective when accessed via AWS in that a user can access a desired tool via AWS and pay only for the time used. Ryland also pointed out that AWS is not limited to single compute nodes and includes cluster management, cloud formation, and cross-cloud capacities. AWS also uses spot pricing, which allows people to bid on excess capacity for computational resources. That gives users access to computing resources cheaply, but the resource is not reliable; if someone else bids more, then the capacity can be taken away and redistributed. Ryland cautioned that projects must be batch-oriented and assess their own progress. For instance, MapReduce is designed so that computational nodes can appear and disappear.

Ryland explained that AWS offers a number of other managed services and provides higher-level application program interfaces. These include Kinesis [9](#) (for massive-scale data streaming), Data Pipeline [10](#) (managed datacentric workflows), and RedShift [11](#) (data warehouse).

AWS has a grants program to benefit students, teachers, and researchers, Ryland said. It is eager to participate in the data science community to build an education base and the resulting benefits. Ryland reported that AWS funds a high percentage of student grants, a fair number of teaching grants, but few research grants. Some of the research grants are high in value, however. In addition to grants, AWS provides spot pricing, volume discounting, and institutional cooperative pricing. In the latter, members of the cooperative can receive shared pricing; AWS intends to increase its cooperative pricing program.

Ryland explained that AWS provides education and training in the form of online training videos, papers, and hands-on, self-paced laboratories. AWS recently launched a fee-based training course. Ryland indicated that AWS is interested in working with the community to be more collaborative and to aggregate opensource materials and curricula. In response to a question, Ryland clarified that the instruction provided by Amazon is on how to use Amazon's tools (such as RedShift), and the instruction is

product-oriented, although the concepts are somewhat general. He said that the instruction is not intended to be revenuegenerating, and AWS would be happy to collaborate with the community on the most appropriate coursework.

A workshop participant posited that advanced tools, such as the AWS tools, enable students to use systems to do large-scale computation without fully understanding how it works. Ryland responded that this is a pattern in computer science: a new level of abstraction develops, and a compiled tool is developed. The cloud is an example of that. Ryland posited that precompiled tools should be able to cover 80 percent or more of the use cases although some researchers will need more profound access to the data.

## 6 WORKSHOP LESSONS

Robert Kass (Carnegie Mellon University) led a final panel discussion session at the end of the workshop. Panelists included James Frew (University of California, Santa Barbara), Deepak Agarwal (LinkedIn Corporation), Claudia Perlich (Dstillery), Raghu Ramakrishnan (Microsoft Corporation), and John Lafferty (University of Chicago). Panelists and participants were invited to add their comments to the workshop; final comments tended to focus in four categories: types of students, organizational structures, course content, and lessons learned from other disciplines.

### Whom to Teach: Types of Students to Target in Teaching Big Data

Robert Kass opened the discussion session by noting that the workshop had shown that there are many types of potential students and that each type would have different training challenges. One participant suggested that business managers need to understand the potential and realities of big data better to improve the quality of communication. Another pointed out that older students may be attracted to big data instruction to pick up missing skill sets. And another suggested pushing instruction into the high-school level. Several participants posited that the background of the student, more than the age or level, is the critical element. For instance, does the student have a background in computer science or statistics? Workshop participants frequently mentioned three main subjects related to big data: computation, statistics, and visualization. The student's background knowledge in each of the three will have the greatest effect on the student's learning.

### How to Teach: The Structure of Teaching Big Data

Numerous participants discussed the types of educational offerings, including massive online open courses (MOOCs), certificate programs, degree-granting programs, boot camps, and individual courses. Participants noted that certificate programs would typically involve a relatively small investment in a student's time, unlike a degree-granting program. One participant proposed a structure consisting of an introductory data science course and three or four additional courses in the three domains (computation, statistics, and visualization). Someone noted that the University of California, Santa Barbara, has similar "emphasis" programs in information technology and technology management. These are sought after because students wish to demonstrate their breadth of understanding. In the case of data science, however, students may wish to use data science to further their domain science. As a result, the certificate model in data science may not be in high demand, inasmuch as students may see value in learning the skills of data science but not in receiving the official recognition of a certificate.

A participant reiterated Joshua Bloom's suggestion made during his presentation to separate data literacy from data fluency. Data fluency would require several years of dedicated study in computing, statistics, visualization, and machine learning. A student may find that difficult to accomplish while obtaining a domain- science degree. Data literacy, in contrast, may be beneficial to many science students and less difficult to obtain. A participant proposed an undergraduate-level introductory data science course focused on basic education and appreciation to promote data literacy.

Workshop participants discussed the importance of coordinating the teaching of data science across multiple disciplines in a university. For example, a participant pointed out that Carnegie Mellon University has multiple master's degree offerings (as many as nine) around the university that are related to data science. Each relevant discipline, such as computer science and statistics, offers a master's degree. The administrative structure is probably stovepiped, and it may be difficult to develop multidisciplinary projects. Another participant argued that an inherently interdisciplinary field of study is not well suited to a degree crafted within a single department and proposed initiating task forces across departments to develop a degree program jointly. And another proposed examining the Carnegie Mellon University data science master's degrees for common topics taught; those topics probably are the proper subset of what constitutes data science.

A workshop participant noted that most institutions do not have nine competing master's programs; instead, most are struggling to develop one. Without collective agreement in the community about the content of a data science program of study, he cautioned that there may be competing programs in each school instead of a single comprehensive program. The

participant stressed the need to understand the core requirements of data science and how big data fits into data science.

Someone noted the importance of having building blocks—such as MOOCs, individual courses, and course sequences—to offer students who wish to focus on data science. Another participant pointed out that MOOCs and boot camps are opposites: MOOCs are large and virtual, whereas boot camps are intimate and hands-on. Both have value as nontraditional credentials.

Guy Lebanon stated that industry finds the end result of data science programs to be inconsistent because they are based in different departments that have different emphases. As a result, industry is uncertain about what a graduate might know. It may be useful to develop a consistent set of standards that can be used in many institutions.

Ramakrishnan stated that “off-the-shelf ” courses in existing programs cannot be stitched together to make a data science curriculum. He suggested creating a wide array of possible prerequisites; otherwise, students will not be able to complete the course sequences that they need.

## What to Teach: Content in Teaching Big Data

The discussion began with a participant noting that it would be impossible to lay out specific topics for agreement. Instead, he proposed focusing on the desired outcomes of training students. Another participant agreed that the fields of study are well known (and typically include databases, statistics and machine learning, and visualization), but said that the specific key components of each field that are needed to form a curriculum are unknown.

Several participants noted the importance of team projects for teaching, especially the creation of teams of students who have different backgrounds (such as a domain scientist and a computer scientist). Team projects foster creativity and encourage new thinking about data problems. Several participants stressed the importance of using real-world data, complete with errors, missing data, and outliers. To some extent, data science is a craft more than a science, so training benefits from the incorporation of real-world projects.

A participant stated that an American Statistical Association committee had been formed to propose a data science program model for a statistical data science program; it would probably include optimization and algorithms, distributed systems, and programming. However, other participants pointed out that that initiative did not include computer science experts in its curriculum development and that that would alter the emphases.

One participant proposed including data security and data ethics in a data science curriculum.

Several participants discussed how teaching data science might differ from teaching big data. One noted that data science does not change its principles when data move into the big data regime, although the approach to each individual step may differ slightly. Temple Lang said that with large data sets, it is easy to get mired in detail, and it becomes even more important to reason through how to solve a problem.

Ramakrishnan recommended including algorithms and analysis in computer science. He noted that although grounding instruction in a specific tool (such as R, SAS, or SQL) teaches practical skills, teaching a tool can compete with teaching of the underlying principles. He endorsed the idea of adding a project element to data science study.

## Parallels in Other Disciplines

Two examples in other domains that were discussed by participants could provide lessons learned to the data science community.

Computational science. A participant noted that computational science was an emerging field 25 years ago. Interdisciplinary academic programs seemed to serve the community best although that model did not fit every university. The participant discussed specifically how the University of Maryland structured its computational-science instruction, which consisted of core coursework and degrees managed through the domain departments. The core courses were co-listed in numerous departments. That model does not require new hiring of faculty or any major restructuring.

Environmental science. Participants discussed an educational model used in environmental science. An interdisciplinary master's-level program was developed so that students could obtain a master's degree in a related science (such as geography, chemistry, or biology). The program involved core courses, research projects, team teaching, and creative use of the academic calendar to provide students with many avenues to an environmental-science degree.

## Footnotes

1

The Internet of Things is the network of uniquely identifiable physical objects embedded throughout a network structure, such as home appliances that can communicate for purposes of adjusting their settings, ordering replacement parts, and so on.

2

See the Master's in Data Science website at <http://www.mastersindatascience.org/> for more information, accessed June 5, 2014.

3

See the National Center for Biotechnology Information's PubMed database at <http://www.ncbi.nlm.nih.gov/pubmed> (accessed May 25, 2014) for more information.

4

That study has since been completed and can be found at Executive Office of the President, Big Data: Seizing Opportunities, Preserving Values, May 2014,

[http://www.whitehouse.gov/sites/default/files/omb/data\\_privacy\\_report\\_may\\_1\\_2014.pdf](http://www.whitehouse.gov/sites/default/files/omb/data_privacy_report_may_1_2014.pdf).

1

In A/B testing, more formally known as two-sample hypothesis testing, two variants are presented to a user, and the user determines a winner.

1

See Apache Software Foundation, "Apache Hadoop," <http://hadoop.apache.org/> (accessed May 14, 2014), for more information.

2

A reducer is a function that typically maps a larger set of values to a smaller set of values.

3

The goal of "predictive health" is to predict the probability of future diseases to identify useful proactive lifestyle modifications and surveillance.

4

More information about scagnostics can be found in Wilkinson et al. (2005, 2006).

5

"Challenges and Opportunities with Big Data—A Community White Paper Developed by Leading Researchers Across the United States," [http://www.cra.org/ccc/files/docs/information\\_technology/big\\_data\\_challenges\\_and\\_opportunities.pdf](http://www.cra.org/ccc/files/docs/information_technology/big_data_challenges_and_opportunities.pdf), accessed May 19, 2014.

6

See TIBCO Software, Inc., "Spotfire," <http://spotfire.tibco.com/>, accessed June 9, 2014, for more information.

7

See the Tableau Software website at <http://www.tableausoftware.com/> (accessed June 9, 2014) for more information.

1

See the Master's in Data Science website at <http://www.mastersindatascience.org/> (accessed June 5, 2014) for more information.

2

See the OpenRefine website at <http://openrefine.org/> (accessed June 9, 2014) for more information.

3

See Stanford Visualization Group, "Data Wrangler alpha," <http://vis.stanford.edu/wrangler/>, accessed June 9, 2014, for more information.

4

See the Trifacta website at <http://www.trifacta.com/> (accessed June 9, 2014) for more information.

5

See Rensselaer Polytechnic Institute (RPI), "Tetherless World Constellation," <http://tw.rpi.edu>, accessed May 22, 2014, for more information.

6

RStudio is an open source, professional user interface for R. See the RStudio website at <http://www.rstudio.com/>, accessed September 20, 2014.

7

See the RPI Information Technology and Web Science website at <http://itws.rpi.edu> (May 22, 2014) for more information.

8

See the RPI Data Science Research Center website at <http://dsrc.rpi.edu/> (May 22, 2014) for more information.

9

See RPI, "Rensselaer Data Services," <http://data.rpi.edu>, accessed May 22, 2014, for more information.

1

See the DeepDive website at <http://deepdive.stanford.edu/> (accessed June 9, 2014) for more information.

2

See Purdue University, Department of Statistics, "Divide and Recombine (D&R) with RHipe," <http://www.datadr.org/>, accessed June 9, 2014, for more information.

3

See GitHub, Inc., "R and Hadoop Integrated Programming Environment," <http://github.com/saptarshiguha/RHipe/>, accessed June 9, 2014, for more information.

4

See Tessera, "datadr: Divide and Recombine in R," <http://hafen.github.io/datadr/>, accessed June 9, 2014, for more information.

5

See Tessera, "Trelliscope: Detailed Vis of Large Complex Data in R," <http://hafen.github.io/trelliscope/>, accessed June 9, 2014, for more information.

6

See Yahoo! Labs, "Webscope," <http://webscope.sandbox.yahoo.com>, accessed May 20, 2014, for more information.

7

See University of Pennsylvania, Linguistic Data Consortium, "LDC Catalog," <http://catalog.ldc.upenn.edu/>, accessed May 14, 2014, for more information.

8

See the Computation Institute, University of Chicago, and Argonne National Laboratory "Globus" website at <https://www.globus.org/> (accessed May 14, 2014) for more information.

9

See Amazon Web Services, "Amazon Kinesis," <http://aws.amazon.com/kinesis/>, accessed June 9, 2014, for more information.

10

See Amazon Web Services, "AWS Data Pipeline," <http://aws.amazon.com/datapipeline/>, accessed June 9, 2014, for more information.

11

See Amazon Web Services, "Amazon Redshift," <http://aws.amazon.com/redshift/>, accessed June 9, 2014, for more information.

## REFERENCES

- Borgman, C., H. Abelson, L. Dirks, R. Johnson, K.R. Koedinger, M.C. Linn, C.A. Lynch, D.G. Oblinger, R.D. Pea, K. Salen, M.S. Smith, and A. Szalay. 2008. Fostering Learning in the Networked World: The Cyberlearning Opportunity and Challenge. Report of the National Science Foundation Task Force on Cyberlearning. National Science Foundation, Washington, D.C.
- Davenport, T.H., and D.J. Patil. 2012. Data scientist: The sexiest job of the 21st century. Harvard Business Review 90(10):70-76.
- Dean, J., and S. Ghemawat. 2004. MapReduce: Simplified data processing on large clusters. Proceedings of the Sixth Symposium on Operating Systems Design and Implementation. <https://www.usenix.org/legacy/public...s/osdi04/tech/>.
- Faris, J., E. Kolker, A. Szalay, L. Bradlow, E. Deelman, W. Feng, J. Qiu, D. Russell, E. Stewart, and E. Kolker. 2011. Communication and data-intensive science in the beginning of the 21st century. OMICS: A Journal of Integrative Biology 15(4):213-215.
- Fox, P., and D.L. McGuinness. 2008. "TWC Semantic Web Methodology." [http://tw.rpi.edu/web/doc/TWC\\_SemanticWebMethodology](http://tw.rpi.edu/web/doc/TWC_SemanticWebMethodology).
- Ferreira, N., J. Poco, H.T. Vo, J. Freire, and C.T. Silva. 2013. Visual exploration of big spatio-temporal urban data: A study of New York City taxi trips. IEEE Transactions on Visualization and Computer Graphics 19(12):2149-2158.
- Manyika, J., M. Chu, B. Brown, J. Bughin, R. Dobbs, C. Roxburgh, and A. Hung Byers. 2011. Big Data: The Next Frontier for Innovation, Competition, and Productivity. McKinsey and Company, Washington, D.C.
- Mayer-Schönberger, V., and K. Cukier. 2012. Big Data: A Revolution That Transforms How We Work, Live, and Think. Houghton Mifflin Harcourt, Boston, Mass.
- Mele, N. 2013. The End of Big: How the Internet Makes David the New Goliath. St. Martin's Press, New York.
- National Research Council. 2013. Frontiers in Massive Data Analysis. The National Academies Press, Washington, D.C.
- Petigura, E.A., A.W. Howard, and G.W. Marcy. 2014. Prevalence of Earth-like planets orbiting Sunlike stars. Proceedings of the National Academy of Sciences 110(48):19273.
- President's Council of Advisors on Science and Technology. 2010. Federally Funded Research and Development in Networking and Information Technology. Executive Office of the President, Washington, D.C.
- Reese, B. 2013. Infinite Progress: How the Internet and Technology Will End Ignorance, Disease, Poverty, Hunger, and War. Greenleaf Book Group Press, Austin, Texas.
- Schmidt, E., and J. Cohen. 2013. The New Digital Age: Reshaping the Future of People, Nations and Business. Knopf Doubleday, New York.

Surdak, C. 2014. Data Crush: How the Information Tidal Wave Is Driving New Business Opportunities. AMACOM Books, Saranac Lake, N.Y.

Webb, A. 2013. Data, A Love Story: How I Gamed Online Dating to Meet My Match. Dutton, New York.

Wilkinson, L., A. Anand, and R. Grossman. 2005. Graph-theoretic scagnostics. Pp. 157-164 in IEEE Symposium on Information Visualization. doi:10.1109/INFVIS.2005.1532142.

Wilkinson, L., A. Anand, and R. Grossman. 2006. High-dimensional visual analytics: Interactive exploration guided by pairwise views of point distributions. IEEE Transactions on Visualization and Computer Graphics 12(6):1363-1372.

---

## APPENDIXES

### A Registered Workshop Participants

Agarwal, Deepak – LinkedIn Corporation

Albrecht, Jochen – Hunter College, City University of New York (CUNY)

Asabi, Faisal – Student / No affiliation known

Bailer, John – Miami University

Begg, Melissa – Columbia University

Bloom, Jane – International Catholic Migration Commission

Bloom, Joshua – University of California, Berkeley

Brachman, Ron – Yahoo Labs

Bradley, Shenae – National Research Council (NRC)

Bruce, Peter – Statistics, Inc.

Buechler, Steven – University of Notre Dame

Caffo, Brian – Johns Hopkins University

Christman, Zachary – Rowan University

Cleveland, Bill – Purdue University

Costello, Donald – University of Nebraska

Curry, James – National Science Foundation

Dell, Robert – Naval Postgraduate School

Dent, Gelonia –Medgar Evers College, CUNY

Desaraju, Kruthika – George Washington University

Dobbins, Janet – Statistics, Inc.

Donovan, Nancy – Government Accountability Office

Dozier, Jeff – University of California, Santa Barbara

Dreves, Harrison – NRC

Dutcher, Jennifer – University of California, Berkeley

Eisenberg, Jon – NRC

Eisner, Ken – Amazon Corporation

Fattah, Hind – Chipotle

Feng, Tingting – University of Virginia

Fox, Peter – Rensselaer Polytechnic Institute

Freire, Juliana – New York University

Freiser, Joel – John Jay College of Criminal Justice

Frew, James – University of California, Santa Barbara

Fricker, Ron – Naval Postgraduate School

Gatsonis, Constantine – Brown University

Ghani, Rayid – University of Chicago

Ghosh, Sujit – National Science Foundation

Glassman, Neal – NRC

Gray, Alexander – Skytree Corporation

Haque, Ubydul – Johns Hopkins University

Howard, Rodney – NRC

Howe, William – University of Washington

Hughes, Gary – Statistics, Inc.  
Huo, Xiaoming – Georgia Tech, National Science Foundation  
Iacono, Suzanne – National Science Foundation  
Kafadar, Karen – Indiana University  
Kass, Robert – Carnegie Mellon University  
Khaloua, Asmaa – Macy  
Kong, Jeongbae – Enanum, Inc.  
Lafferty, John – University of Chicago  
Lesser, Virginia – Oregon State University  
Lebanon, Guy – Amazon Corporation  
Levermore, David – University of Maryland  
Liu, Shiyong – Southwestern University of Finance and Economics  
Mandl, Kenneth – Harvard Medical School Boston Children's Hospital  
Marcus, Stephen – National Institute of General Medical Sciences, National Institutes of Health (NIH)  
Martinez, Waldyn – Miami University  
Mellody, Maureen – NRC  
Neerchal, Nagaraj – University of Maryland, Baltimore County  
Orwig, Jessica – American Physical Society  
Pack, Quinn – Mayo Clinic  
Parmigiani, Giovanni – Dana Farber Cancer Institute  
Pearl, Jennifer – National Science Foundation  
Pearsall, Hamil – Temple University  
Perlich, Claudia – Dstillery  
Rai, Saatvika – University of Kansas  
Ralston, Bruce – University of Tennessee  
Ramakrishnan, Raghu – Microsoft Corporation  
Ranakrishnan, Raghunath – University of Texas, Austin  
Ravichandran, Veerasamy – NIH  
Ré, Christopher – Stanford University  
Ryland, Mark – Amazon Corporation  
Schwalbe, Michelle – NRC  
Schou, Sue – Idaho State University  
Shams, Khawaja – Amazon Corporation  
Sharman, Raj –University at Buffalo, State University of New York (SUNY)  
Shekhar, Shashi – University of Minnesota  
Shipp, Stephanie – VA Bioinformatics Institute at Virginia Tech University  
Shneiderman, Ben – University of Maryland  
Spencer Huang, ChiangChing – University of Wisconsin, Milwaukee  
Spengler, Sylvia – National Science Foundation  
Srinivasarao, Geetha – Information Technology Specialist, Department of Health and Human Services  
Szewczyk, Bill – National Security Agency  
Tannouri, Ahlam – Morgan State University  
Tannouri, Charles – Department of Homeland Security  
Tannouri, Sam – Morgan State University  
Temple Lang, Duncan – University of California, Davis  
Torrens, Paul – University of Maryland, College Park  
Ullman, Jeffrey – Stanford University  
Vargas, Juan – Georgia Southern University  
Wachowicz, Monica – University of New Brunswick, Fredericton  
Wang, Rong – Illinois Institute of Technology  
Wang, Youfa – University at Buffalo, SUNY  
Wee, Brian – National Ecological Observatory Network (NEON), Inc.  
Weese, Maria – MIA  
Weidman, Scott – NRC  
Weiner, Angelica – Amazon Corporation

Wynn, Sarah – NRC Christine Mirzayan Science and Technology Policy Graduate Fellow  
Xiao, Ningchuan – Ohio State University  
Xue, Hong – University at Buffalo, SUNY  
Yang, Ruixin – George Mason University  
Zhang, Guoping – Morgan State University  
Zhao, Fen – National Science Foundation

## B Workshop Agenda

APRIL 11, 2014

8:30 a.m. Opening Remarks

Suzanne Iacono, Deputy Assistant Director, Directorate for Computer and Information Science and Engineering, National Science Foundation

8:40 The Need for Training: Experiences and Case Studies

Co-Chairs: Raghu Ramakrishnan, Microsoft Corporation and John Lafferty, University of Chicago

Speakers: Rayid Ghani, University of Chicago Guy Lebanon, Amazon Corporation

10:15 Principles for Working with Big Data

Chair: Brian Caffo, Johns Hopkins University

Speakers: Jeffrey Ullman, Stanford University

Alexander Gray, Skytree Corporation

Duncan Temple Lang, University of California, Davis

Juliana Freire, New York University

12:45 p.m. Lunch

1:45 Courses, Curricula, and Interdisciplinary Programs

Chair: James Frew, University of California, Santa Barbara

Speakers: William Howe, University of Washington

Peter Fox, Rensselaer Polytechnic Institute

Joshua Bloom, University of California, Berkeley

4:30 Q&A/Discussion

APRIL 12, 2014

8:30 a.m. Shared Resources

Chair: Deepak Agarwal, LinkedIn Corporation

Speakers: Christopher Ré, Stanford University

Bill Cleveland, Purdue University

Ron Brachman, Yahoo Labs

Mark Ryland, Amazon Corporation

11:15 Panel Discussion: Workshop Lessons

Chair: Robert Kass, Carnegie Mellon University

Panel Members: James Frew, University of California, Santa Barbara

Deepak Agarwal, LinkedIn Corporation

Claudia Perlich, Dstillery

Raghu Ramakrishnan, Microsoft Corporation

John Lafferty, University of Chicago

1:00 p.m. Workshop Adjourns

## C Acronyms

AOL America OnLine  
 AWS Amazon Web Services  
 BMSA Board on Mathematical Sciences and Their Applications  
 CATS Committee on Applied and Theoretical Statistics  
 CRA Computing Research Association  
 DARPA Defense Advanced Research Projects Agency  
 DOE Department of Energy  
 MOOC massive online open course  
 NASA National Aeronautics and Space Administration  
 NIH National Institutes of Health  
 NITRD Networking and Information Technology Research and Development  
 NRC National Research Council  
 NSF National Science Foundation  
 OCR optical character recognition  
 RHipe R and Hadoop Integrated Programming Environment

---

## Getting Data Right: Tackling The Challenges of Big Data Volume and Variety

[PDF](#)

Edited by Shannon Cutt

This Preview Edition of Getting Data Right, Chapter 4, is a work in progress. The final book is currently scheduled for release in fall 2015.

---

## CHAPTER 1 The Solution: Data Curation at Scale

—Michael Stonebraker, Ph.D.

Integrating data sources isn't a new challenge. But the challenge has intensified in both importance and difficulty, as the volume and variety of usable data - and enterprises' ambitious plans for analyzing and applying it - have increased. As a result, trying to meet today's data integration demands with yesterday's data integration approaches is impractical.

In this chapter, we look at the three generations of data integration products and how they have evolved. We look at new third generation products that deliver a vital missing layer in the data integration "stack": data curation at scale. Finally, we look at five key tenets of an effective data curation at scale system.

### Three Generations of Data Integration Systems

Data integration systems emerged to enable business analysts to access converged data sets directly for analyses and applications.

First-generation data integration systems - data warehouses - arrived on the scene in the 1990s. Led by the major retailers, customer-facing data (e.g., item sales, products, customers) were assembled in a data store and used by retail buyers to make better purchase decisions. For example, pet rocks might be out of favor while Barbie dolls might be "in." With this intelligence, retailers could discount the pet rocks and tie up the Barbie doll factory with a big order. Data warehouses typically paid for themselves within a year through better buying decisions.

First-generation data integration systems were termed ETL (Extract, Transform and Load) products. They were used to assemble the data from various sources (usually fewer than 20) into the warehouse. But enterprises underestimated the "T" part of the process - specifically, the cost of data curation (mostly, data cleaning) required to get heterogeneous data into the proper format for querying and analysis. Hence, the typical data warehouse project was usually substantially over-budget and late because of the difficulty of data integration inherent in these early systems.

This led to a second generation of ETL systems, whereby the major ETL products were extended with data cleaning modules, additional adapters to ingest other kinds of data, and data cleaning tools. In effect, the ETL tools were extended to become data

curation tools.

Data curation involves:

- ingesting data sources
- cleaning errors from the data (-99 often means null)
- transforming attributes into other ones (for example, Euros to dollars)
- performing schema integration to connect disparate data sources
- performing entity consolidation to remove duplicates

In general, data curation systems followed the architecture of earlier first-generation systems: toolkits oriented toward professional programmers. In other words, they were programmer productivity tools.

Second-generation data curation tools have two substantial weaknesses:

**Scalability**

They have several thousand data sources, everything from company budgets in the CFO's spreadsheets to peripheral operational systems. There is "business intelligence gold" in the long tail, and enterprises wish to capture it - for example, for crossselling of enterprise products. Furthermore, the rise of public data on the web leads business analysts to want to curate additional data sources. Anything from weather data to customs records to real estate transactions to political campaign contributions are readily available. However, in order to capture longtail enterprise data as well as public data, curation tools must be able to deal with hundreds to thousands of data sources rather than the typical few tens of data sources.

**Architecture**

Second-generation tools typically are designed for central IT departments. A professional programmer does not know the answers to many of the data curation questions that arise. For example, are "rubber gloves" the same thing as "latex hand protectors?" Is an "ICU50" the same kind of object as an "ICU?" Only business people in line-of-business organizations can answer these kinds of questions. However, business people are usually not in the same organization as the programmers running data curation projects. As such, second-generation systems are not architected to take advantage of the humans best able to provide curation help.

These weaknesses led to a third generation of data curation products, which we term scalable data curation systems. Any data curation system should be capable of performing the five tasks noted above. However, first- and second-generation ETL products will only scale to a small number of data sources, because of the amount of human intervention required.

To scale to hundreds or even thousands of data sources, a new approach is needed - one that:

1. Uses statistics and machine learning to make automatic decisions wherever possible.
2. Asks a human expert for help only when necessary.

Instead of an architecture with a human controlling the process with computer assistance, move to an architecture with the computer running an automatic process, asking a human for help only when required. And ask the right human: the data creator or owner (a business expert) not the data wrangler (a programmer).

Obviously, enterprises differ in the required accuracy of curation, so third-generation systems must allow an enterprise to make tradeoffs between accuracy and the amount of human involvement. In addition, third-generation systems must contain a crowdsourcing component that makes it efficient for business experts to assist with curation decisions. Unlike Amazon's Mechanical Turk, however, a data-curation crowdsourcing model must be able to accommodate a hierarchy of experts inside an enterprise as well as various kinds of expertise. Therefore, we call this component an expert sourcing system to distinguish it from the more primitive crowdsourcing systems.

In short: a third-generation data curation product is an automated system with an expert sourcing component. Tamr is an early example of this third generation of systems.

Third-generation systems can co-exist with currently-in-place second-generation systems, which can curate the first tens of data sources to generate a composite result that in turn can be curated with the "long tail" by third-generation systems.

Table 1-1. Evolution of Three Generations of Data Integration Systems

	First Generation	Second Generation	
--	------------------	-------------------	--

Evolution/Generation	1990s	2000s	Third Generation 2010s
Approach	ETL	ETL+>Data Curation	Scalable Data Curation
Target Data Environment(s)	Data Warehouse	Data Warehouses or Data Marts	Data Lakes & Self-Service Data Analytics
Users	IT/Programmers	IT/Programmers	Data Scientists, Data Stewards, Data Owners, Business Analysts
Integration Philosophy	Top-down/rules-based/IT-driven	Top-down/rules-based/IT-driven	Bottom-up/demand-based/business-driven
Architecture	Programmer productivity tools (task automation)	Programming productivity tools (task automation with machine assistance)	Machine-driven, human-guided process
Scalability (# of data sources)	10s	10s to 100s	100s to 1000s+

To summarize: ETL systems arose to deal with the transformation challenges in early data warehouses. They evolved into second generation data curation systems with an expanded scope of offerings. Third-generation data curation systems, which have a very different architecture, were created to address the enterprise's need for data source scalability.

## Five Tenets for Success

Third-generation scalable data curation systems provide the architecture, automated workflow, interfaces and APIs for data curation at scale. Beyond this basic foundation, however, are five tenets that are desirable in any third-generation system.

### Tenet 1: Data curation is never done

Business analysts and data scientists have an insatiable appetite for more data. This was brought home to me about a decade ago during a visit to a beer company in Milwaukee. They had a fairly standard data warehouse of sales of beer by distributor, time period, brand and so on. I visited during a year when El Niño was forecast to disrupt winter weather in the US. Specifically, it was forecast to be wetter than normal on the West Coast and warmer than normal in New England. I asked the business analysts: "Are beer sales correlated with either temperature or precipitation?" They replied, "We don't know, but that is a question we would like to ask." However temperature and precipitation were not in the data warehouse, so asking was not an option.

The demand from warehouse users to correlate more and more data elements for business value leads to additional data curation tasks. Moreover, whenever a company makes an acquisition, it creates a data curation problem (digesting the acquired's data). Lastly, the treasure trove of public data on the web (such as temperature and precipitation data) is largely untapped, leading to more curation challenges.

Even without new data sources, the collection of existing data sources is rarely static. Hence, inserts and deletes to these sources generates a pipeline of incremental updates to a data curation system. Between the requirements of new data sources and updates to existing ones, it is obvious that data curation is never done, ensuring that any project in this area will effectively continue indefinitely. Realize this and plan accordingly.

One obvious consequence of this tenet concerns consultants. If you hire an outside service to perform data curation for you, then you will have to rehire them for each additional task. This will give the consultant a guided tour through your wallet over time. In my opinion, you are much better off developing in-house curation competence over time.

### Tenet 2: A PhD in AI can't be a requirement for success

Any third-generation system will use statistics and machine learning to make automatic or semi-automatic curation decisions. Inevitably, it will use sophisticated techniques such as T-tests, regression, predictive modeling, data clustering, and classification. Many of these techniques will entail training data to set internal parameters. Several will also generate recall and/or precision estimates.

These are all techniques understood by data scientists. However, there will be a shortage of such people for the foreseeable future, until colleges and universities produce substantially more than at present. Also, it is not obvious that one can "retread" a

business analyst into a data scientist. A business analyst only needs to understand the output of SQL aggregates; in contrast, a data scientist is typically knowledgeable in statistics and various modeling techniques.

As a result, most enterprises will be lacking in data science expertise. Therefore, any third-generation data curation product must use these techniques internally, but not expose them in the user interface. Mere mortals must be able to use scalable data curation products.

### Tenet 3: Fully automatic data curation is not likely to be successful

Some data curation products expect to run fully automatically. In other words, they translate input data sets into output without human intervention. Fully automatic operation is very unlikely to be successful in an enterprise for a variety of reasons. First, there are curation decisions that simply cannot be made automatically. For example, consider two records; one stating that restaurant X is at location Y while the second states that restaurant Z is at location Y. This could be a case where one restaurant went out of business and got replaced by a second one or it could be a food court. There is no good way to know the answer to this question without human guidance.

Second, there are cases where data curation must have high reliability. Certainly, consolidating medical records should not create errors. In such cases, one wants a human to check all (or maybe just some) of the automatic decisions. Third, there are situations where specialized knowledge is required for data curation. For example, in a genomics application one might have two terms: ICU50 and ICE50. An automatic system might suggest that these are the same thing, since the lexical distance between the terms is low. However, only a human genomics specialist can decide this question.

For these reasons, any third-generation data curation system must be able to ask a human expert - the right human expert - when it is unsure of the answer. Therefore, one must have multiple domains in which a human can be an expert. Within a single domain, humans have a variable amount of expertise, from a novice level to enterprise expert. Lastly, one must avoid overloading the humans that it is scheduling. Therefore, when considering a third generation data curation system, look for an embedded expert system with levels of expertise, load balancing and multiple expert domains.

### Tenet 4: Data curation must fit into the enterprise ecosystem

Every enterprise has a computing infrastructure in place. This includes a collection of DBMSs storing enterprise data, a collection of application servers and networking systems, and a set of installed tools and applications. Any new data curation system must fit into this existing infrastructure. For example, it must be able to extract from corporate databases, use legacy data cleaning tools, and export data to legacy data systems. Hence, an open environment is required whereby callouts are available to existing systems. In addition, adapters to common input and export formats is a requirement. Do not use a curation system that is a closed "black box."

### Tenet 5: A scheme for "finding" data sources must be present

A typical question to ask CIOs is "How many operational data systems do you have?". In all likelihood, they do not know. The enterprise is a sea of such data systems connected by a hodgepodge set of connectors. Moreover, there are all sorts of personal datasets, spreadsheets and databases, as well as datasets imported from public web-oriented sources. Clearly, CIOs should have a mechanism for identifying data resources that they wish to have curated. Such a system must contain a data source catalog with information on a CIO's data resources, as well as a query system for accessing this catalog. Lastly, an "enterprise crawler" is required to search a corporate internet to locate relevant data sources. Collectively, this represents a schema for "finding" enterprise data sources.

Collectively, these five tenets indicate the characteristics of a good third generation data curation system. If you are in the market for such a product, then look for systems with these characteristics.

## NEXT

### Table of Contents

1. [Story](#)
  1. [Data Science for Cyber Physical Systems-Internet of Things](#)
  2. [Training Students to Extract Value from Big Data](#)

1. [Slide 1 An Internet of Things: People, Processes, and Products in the Spotfire Cloud Library](#)
2. [Slide 2 Outline](#)
3. [Slide 3 Semantic Community: Spotfire Cloud Library MindTouch](#)
4. [Slide 4 Semantic Community: Spotfire Cloud Library Recent Analyses](#)
5. [Slide 5 Semantic Community: Spotfire Cloud Library Browse Library](#)
6. [Slide 6 Semantic Community: Spotfire Cloud Library Shared Folders](#)
7. [Slide 7 Semantic Community: Spotfire Web Player](#)
8. [Slide 8 Semantic Community: Spotfire Web Player Edit](#)
9. [Slide 9 Semantic Community: Spotfire Web Player Menu](#)
10. [Slide 10 Semantic Community: Spotfire Web Player Menu and New Visualizations](#)
11. [Slide 11 TIBCO Spotfire: Cloud User's Guide](#)
12. [Slide 12 TIBCO: Spotfire Cloud User's Guide Getting Started](#)
13. [Slide 13 TIBCO: Spotfire Cloud User's Guide Setting Up Analyses](#)
14. [Slide 14 TIBCO: Spotfire Cloud User's Guide Data Preparation in Microsoft Excel](#)
15. [Slide 15 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 1](#)
16. [Slide 16 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 2](#)
17. [Slide 17 TIBCO: Spotfire Cloud User's Guide Creating a Visualization 3](#)
18. [Slide 18 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Markers](#)
19. [Slide 19 TIBCO: Spotfire Cloud User's Guide Creating a Map Chart with Colored Regions](#)
20. [Slide 20 TIBCO: Spotfire Cloud User's Guide Exporting Your Analysis](#)
21. [Slide 21 TIBCO: Spotfire Cloud User's Guide Sharing Your Analysis](#)
22. [Slide 22 TIBCO Spotfire: Education](#)
3. [Summit Highlights](#)
  1. [Slide 1 Opening Remarks 1](#)
  2. [Slide 2 Opening Remarks 2](#)
  3. [Slide 3 Dr. Shoumen Data 1](#)
  4. [Slide 4 Dr. Shoumen Data 2](#)
  5. [Slide 5 Dr. Shoumen Data 3](#)
  6. [Slide 6 Dr. Shoumen Data 4](#)
  7. [Slide 7 Track A 1](#)
  8. [Slide 8 Track A 2](#)
  9. [Slide 9 Track B 1](#)
  10. [Slide 10 Track B 2](#)
  11. [Slide 11 Keith Mazullo 1](#)
  12. [Slide 12 Keith Mazullo 2](#)
  13. [Slide 13 Keith Mazullo 3](#)
  14. [Slide 14 Keith Mazullo 4](#)
  15. [Slide 15 Track C 1](#)
  16. [Slide 16 Track C 2](#)
  17. [Slide 17 Track D 1](#)
  18. [Slide 18 Track D 2](#)
  19. [Slide 19 Dr. Bradford Hess 1](#)
  20. [Slide 20 Dr. Bradford Hess 2](#)
  21. [Slide 21 Dr. Harry Foxwell 1](#)
  22. [Slide 22 Dr. Harry Foxwell 2](#)
  23. [Slide 23 Eric Simmon 1](#)
  24. [Slide 24 Eric Simmon 2](#)
  25. [Slide 25 William Miller 1](#)
  26. [Slide 26 William Miller 2](#)
2. [Slides](#)
  1. [Slide 1 Data Science for Big Data](#)
  2. [Slide 2 Overview](#)
  3. [Slide 3 The Profit and Data Enterprises](#)
  4. [Slide 4 Federal Big Data Working Group Meetup](#)
  5. [Slide 5 Silicon Valley to Washington](#)

6. [Slide 6 First White House Data Chief Discusses His Top Priorities](#)
  7. [Slide 7 Precision Medicine and Natural Medicine](#)
  8. [Slide 8 Tech Meetup at White House](#)
  9. [Slide 9 USDA Data Science MOOC](#)
  10. [Slide 10 Upcoming Meetups](#)
  11. [Slide 11 Summary 1](#)
  12. [Slide 12 Summary 2](#)
3. [Slides](#)
1. [Slide 1 Semantic Data Discovery: Proof of Concept for DHS](#)
  2. [Slide 2 Information Sharing at DHS](#)
  3. [Slide 3 NIEM as Big Data in a Network with Data Science](#)
  4. [Slide 4 NIEM 3.0 Alpha 2 Release and Thetus Savanna Review](#)
  5. [Slide 5 NIEM and UCore 2.0 Semantic Layer for Information Sharing](#)
  6. [Slide 6 Global Terrorism Database Experience](#)
  7. [Slide 7 A Quint for Cross Information Sharing and Integration in the Intelligence Community](#)
  8. [Slide 8 Dynamic Case Management Pilot for Healthcare.gov](#)
  9. [Slide 9 Proof of Concept Steps](#)
  10. [Slide 10 Semantic Community](#)
4. [Spotfire Dashboard](#)
5. [Research Notes](#)
6. [Ontology Summit 2015 Agenda](#)
  1. [Monday, April 13](#)
  2. [Tuesday, April 14](#)

7. [Ontology Summit 2015 Background](#)
  1. [Prepared Presentation Materials](#)
  2. [Audio Recordings](#)
  3. [Additional Resources](#)
  4. [Abstract](#)

8. [Ontology Summit 2015 Communique](#)
  1. [Introduction](#)
  2. [The Case for IoT Ontologies](#)
  3. [How Ontologies are Used in IoT](#)
    1. [Ontology Mapping](#)
    2. [Standards Integration](#)
    3. [Decision Support for IoT](#)
  4. [Beyond Semantic Sensor Network Ontologies](#)
  5. [Ontological Issues](#)
    1. [Scalability](#)
    2. [Standards Integration](#)
  6. [Challenges](#)
  7. [Forecasts](#)
  8. [Recommendations](#)
  9. [Terminology](#)

9. [Training Students to Extract Value from Big Data](#)
  1. [The National Academies Press](#)
    1. [Authors](#)
    2. [Description](#)
    3. [Topics](#)
    4. [Publication Info](#)
    5. [Copyright Information](#)
  2. [Cover Page](#)
  3. [The National Academies Press](#)
  4. [Planning Committee on Training Students](#)
  5. [Committee on Applied and Theoretical Statistics](#)
  6. [Board of Mathematical Sciences and Their Applications](#)

7. [Acknowledgment of Reviewers](#)

8. [\*\*1 INTRODUCTION\*\*](#)

1. [Workshop Overview](#)

1. [BOX 1.1 Statement of Task](#)

2. [National Efforts in Big Data](#)

1. [Suzanne Iacono, National Science Foundation](#)

3. [Organization of This Report](#)

9. [\*\*2 THE NEED FOR TRAINING: EXPERIENCES AND CASE STUDIES\*\*](#)

1. [Training Students to Do Good with Big Data](#)

1. [Rayid Ghani, University of Chicago](#)

2. [The Need for Training in Big Data: Experiences and Case Studies](#)

1. [Guy Lebanon, Amazon Corporation](#)

10. [\*\*3 PRINCIPLES FOR WORKING WITH BIG DATA\*\*](#)

1. [Teaching about MapReduce](#)

1. [Jeffrey Ullman, Stanford University](#)

2. [Big Data Machine Learning—Principles for Industry](#)

1. [Alexander Gray, Skytree Corporation](#)

3. [Principles for the Data Science Process](#)

1. [Duncan Temple Lang, University of California, Davis](#)

4. [Principles for Working with Big Data](#)

1. [Juliana Freire, New York University](#)

2. [FIGURE 3.1 Simplified schematic of the big data analysis pipeline](#)

11. [\*\*4 COURSES, CURRICULA, AND INTERDISCIPLINARY PROGRAMS\*\*](#)

1. [Computational Training and Data Literacy for Domain Scientists](#)

1. [Joshua Bloom, University of California, Berkeley](#)

2. [Data Science and Analytics Curriculum Development at Rensselaer \(and the Tetherless World Constellation\)](#)

1. [Peter Fox, Rensselaer Polytechnic Institute](#)

2. [FIGURE 4.1 Framework for modern informatics](#)

3. [FIGURE 4.2 Generations of mediation](#)

3. [Experience with a First Massive Online Open Course on Data Science](#)

1. [William Howe, University of Washington](#)

12. [\*\*5 SHARED RESOURCES\*\*](#)

1. [Can Knowledge Bases Help Accelerate Science?](#)

1. [Christopher Ré, Stanford University](#)

2. [Divide and Recombine for Large, Complex Data](#)

1. [Bill Cleveland, Purdue University](#)

3. [Yahoo's Webscope Data Sharing Program](#)

1. [Ron Brachman, Yahoo Labs](#)

4. [Resource Sharing](#)

1. [Mark Ryland, Amazon Corporation](#)

13. [\*\*6 WORKSHOP LESSONS\*\*](#)

1. [Whom to Teach: Types of Students to Target in Teaching Big Data](#)

2. [How to Teach: The Structure of Teaching Big Data](#)

3. [What to Teach: Content in Teaching Big Data](#)

4. [Parallels in Other Disciplines](#)

14. [\*\*Footnotes\*\*](#)

1. [1](#)

2. [2](#)

3. [3](#)

4. [4](#)

5. [1](#)

6. [1](#)

7. [2](#)

8. [3](#)

9. [4](#)

10. [5](#)

11. [6](#)

12. [7](#)

13. [1](#)

14. [2](#)

15. [3](#)

16. [4](#)

17. [5](#)

18. [6](#)

19. [7](#)

20. [8](#)

21. [9](#)

22. [1](#)

23. [2](#)

24. [3](#)

25. [4](#)

26. [5](#)

27. [6](#)

28. [7](#)

29. [8](#)

30. [9](#)

31. [10](#)

32. [11](#)

## 15. REFERENCES

## 16. APPENDIXES

1. [A Registered Workshop Participants](#)

2. [B Workshop Agenda](#)

1. [APRIL 11, 2014](#)

2. [APRIL 12, 2014](#)

3. [C Acronyms](#)

## 10. Getting Data Right: Tackling The Challenges of Big Data Volume and Variety

1. [CHAPTER 1 The Solution: Data Curation at Scale](#)

1. [Three Generations of Data Integration Systems](#)

1. [Table 1-1. Evolution of Three Generations of Data Integration Systems](#)

2. [Five Tenets for Success](#)

1. [Tenet 1: Data curation is never done](#)

2. [Tenet 2: A PhD in AI can't be a requirement for success](#)

3. [Tenet 3: Fully automatic data curation is not likely to be successful](#)

4. [Tenet 4: Data curation must fit into the enterprise ecosystem](#)

5. [Tenet 5: A scheme for "finding" data sources must be present](#)

## 11. NEXT