



What's Next for Big Data?

Lessons from a Decade+ Experiment in Big Data

A Discussion with IEEE BDIW

David Belanger PhD

Senior Research Fellow – Stevens Institute of Technology

dbelange@stevens.edu

Introduction

What, How, Why

Big Data 1880

Census

Population: 50,189,209
Size: Low Gigabytes

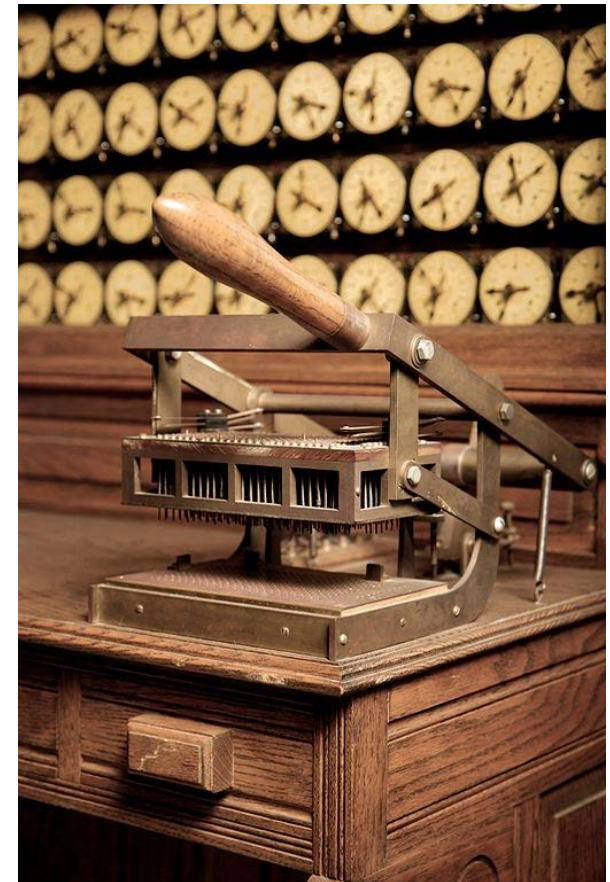
Hollerith Tabulating Machine

Page No. 1
Superior's Dist. No. 1
Enumeration Dist. No. 1

SCHEDULE 1—Inhabitants in Appler in the County of Delaware, State of Wisconsin
enumerated by me on the twentieth day of June, 1880.

1880

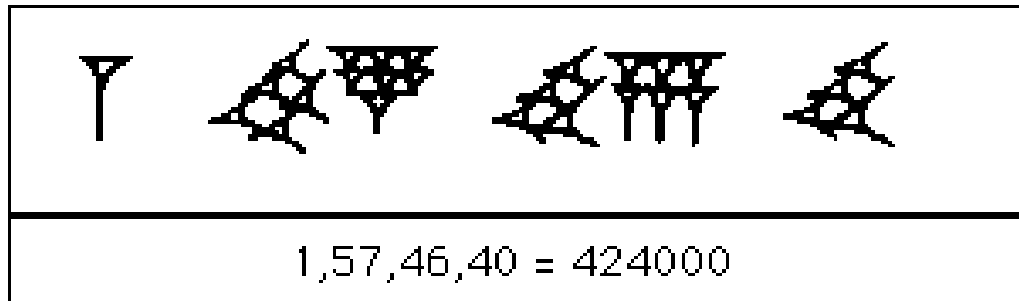
No.	Name	Age	Sex	Color	Mar.	Prof.	Ind.	Ag.	Dom.	Foreign	Imm.	Other
1	John H. Jones	25	M	W	M							
2	John H. Jones	25	M	W	M							
3	John H. Jones	25	M	W	M							
4	John H. Jones	25	M	W	M							
5	John H. Jones	25	M	W	M							
6	John H. Jones	25	M	W	M							
7	John H. Jones	25	M	W	M							
8	John H. Jones	25	M	W	M							
9	John H. Jones	25	M	W	M							
10	John H. Jones	25	M	W	M							
11	John H. Jones	25	M	W	M							
12	John H. Jones	25	M	W	M							
13	John H. Jones	25	M	W	M							
14	John H. Jones	25	M	W	M							
15	John H. Jones	25	M	W	M							
16	John H. Jones	25	M	W	M							
17	John H. Jones	25	M	W	M							
18	John H. Jones	25	M	W	M							
19	John H. Jones	25	M	W	M							
20	John H. Jones	25	M	W	M							
21	John H. Jones	25	M	W	M							
22	John H. Jones	25	M	W	M							
23	John H. Jones	25	M	W	M							
24	John H. Jones	25	M	W	M							
25	John H. Jones	25	M	W	M							
26	John H. Jones	25	M	W	M							
27	John H. Jones	25	M	W	M							
28	John H. Jones	25	M	W	M							
29	John H. Jones	25	M	W	M							
30	John H. Jones	25	M	W	M							
31	John H. Jones	25	M	W	M							
32	John H. Jones	25	M	W	M							
33	John H. Jones	25	M	W	M							
34	John H. Jones	25	M	W	M							
35	John H. Jones	25	M	W	M							
36	John H. Jones	25	M	W	M							
37	John H. Jones	25	M	W	M							
38	John H. Jones	25	M	W	M							
39	John H. Jones	25	M	W	M							
40	John H. Jones	25	M	W	M							
41	John H. Jones	25	M	W	M							
42	John H. Jones	25	M	W	M							
43	John H. Jones	25	M	W	M							
44	John H. Jones	25	M	W	M							
45	John H. Jones	25	M	W	M							
46	John H. Jones	25	M	W	M							
47	John H. Jones	25	M	W	M							
48	John H. Jones	25	M	W	M							
49	John H. Jones	25	M	W	M							
50	John H. Jones	25	M	W	M							
51	John H. Jones	25	M	W	M							
52	John H. Jones	25	M	W	M							
53	John H. Jones	25	M	W	M							
54	John H. Jones	25	M	W	M							
55	John H. Jones	25	M	W	M							
56	John H. Jones	25	M	W	M							
57	John H. Jones	25	M	W	M							
58	John H. Jones	25	M	W	M							
59	John H. Jones	25	M	W	M							
60	John H. Jones	25	M	W	M							
61	John H. Jones	25	M	W	M							
62	John H. Jones	25	M	W	M							
63	John H. Jones	25	M	W	M							
64	John H. Jones	25	M	W	M							
65	John H. Jones	25	M	W	M							
66	John H. Jones	25	M	W	M							
67	John H. Jones	25	M	W	M							
68	John H. Jones	25	M	W	M							
69	John H. Jones	25	M	W	M							
70	John H. Jones	25	M	W	M							
71	John H. Jones	25	M	W	M							
72	John H. Jones	25	M	W	M							
73	John H. Jones	25	M	W	M							
74	John H. Jones	25	M	W	M							
75	John H. Jones	25	M	W	M							
76	John H. Jones	25	M	W	M							
77	John H. Jones	25	M	W	M							
78	John H. Jones	25	M	W	M							
79	John H. Jones	25	M	W	M							
80	John H. Jones	25	M	W	M							
81	John H. Jones	25	M	W	M							
82	John H. Jones	25	M	W	M							
83	John H. Jones	25	M	W	M							
84	John H. Jones	25	M	W	M							
85	John H. Jones	25	M	W	M							
86	John H. Jones	25	M	W	M							
87	John H. Jones	25	M	W	M							
88	John H. Jones	25	M	W	M							
89	John H. Jones	25	M	W	M							
90	John H. Jones	25	M	W	M							
91	John H. Jones	25	M	W	M							
92	John H. Jones	25	M	W	M							
93	John H. Jones	25	M	W	M							
94	John H. Jones	25	M	W	M							
95	John H. Jones	25	M	W	M							
96	John H. Jones	25	M	W	M							
97	John H. Jones	25	M	W	M							
98	John H. Jones	25	M	W	M							
99	John H. Jones	25	M	W	M							
100	John H. Jones	25	M	W	M							



Big Data 2000 BC

Base 60 Positional Arithmetic

1,57,46,40 in Babylonian numerals



Source: http://www-history.mcs.st-and.ac.uk/HistTopics/Babylonian_numerals.html

ComSoc – Big Data Contributions

IEEE
NETWORK

This month's complimentary access to a select article is from current issue of IEEE Network

"Big Data: Transforming the Design Philosophy of Future Internet"

- ✓ Communications
- ✓ Conferences
- ✓ Publications
- ✓ Training

EXPLORE THE WORLD OF BIG DATA DURING AT IEEE/IFIP IM'15 CONFERENCE IN OTTAWA, CANADA

Like 5 Tweet 2 +1 1 Share 4



IEEE COMSOC **TRAINING**™ CONTINUING EDUCATION FOR COMMUNICATION PROFESSIONALS

HOME CALENDAR INSTRUCTORS CUSTOM TRAINING IEEE CEUS / PDHs

THE BIG PICTURE FOR SMALL CELLS

Create a Small Cell Deployment Strategy

Instructor: [Jonathan Levine](#)

Wednesday, October 15, 2014 - 9:00am to 4:30pm EDT

Online via WebEx

IEEE Communications MAGAZINE

CURRENT ISSUE ▼ CALL FOR PAPERS GUEST EDITOR KIT REVIEWER GUIDELINES

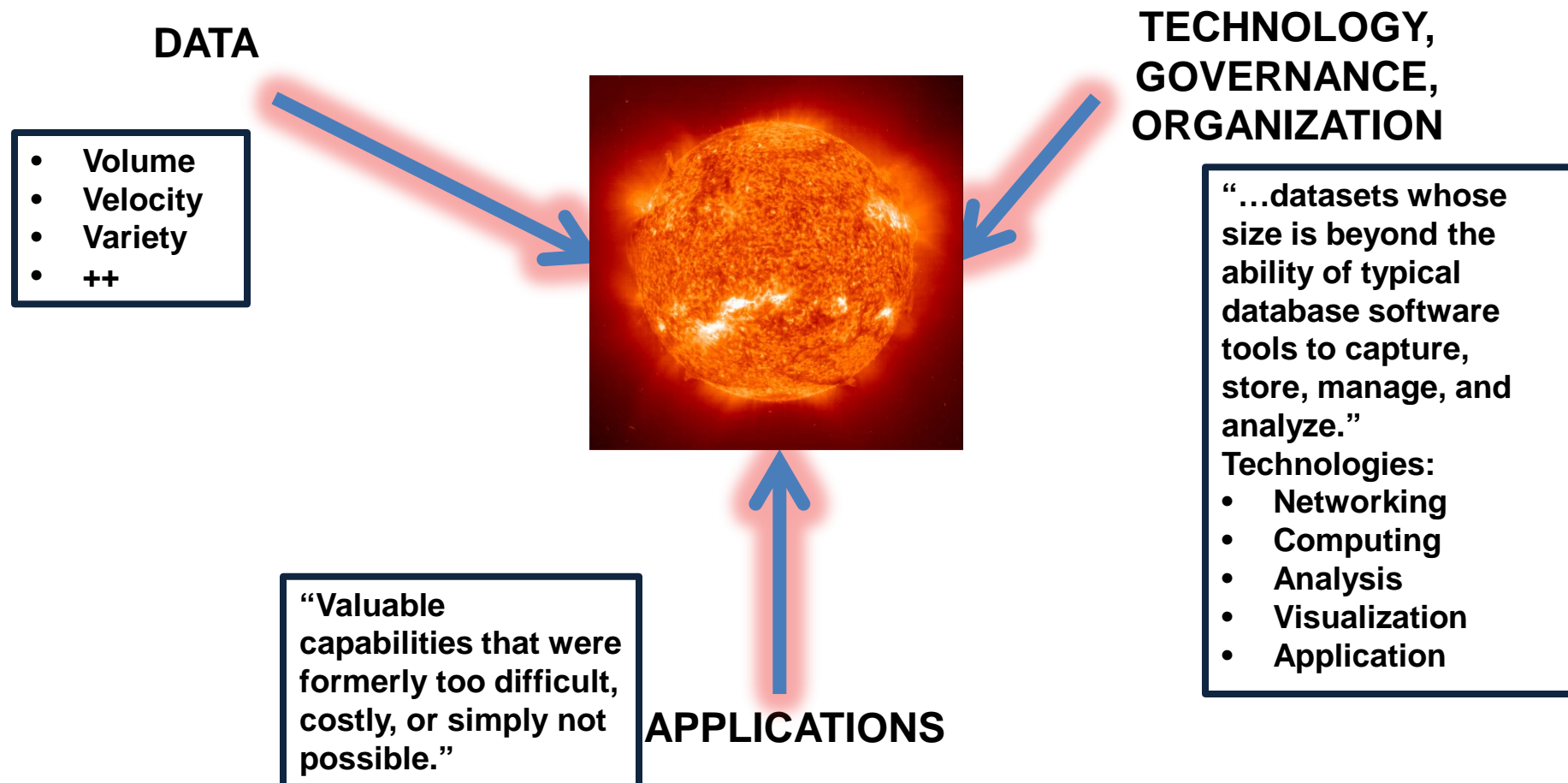
AUTHOR GUIDELINES ▼ EDITORIAL BOARD ▼ THE PRESIDENT'S PAGE ADVERTISING

HOME » THE PRESIDENT'S PAGE - AUGUST 2014

THE PRESIDENT'S PAGE - AUGUST 2014

Every year a full day in June is spent with Society Presidents or their delegates, discussing promising areas that may be worth exploring. The June meeting is open to all the operating units of IEEE. We typically have Educational Activities, Member and Geographic Areas, Marketing, Corporate, the Standards Association, as well as other unique areas, attending this open meeting. As an example, this June the area of Leveraging Big Data was identified. This is an area where a number of Societies are already involved, and what FDC will be doing, together with interested Societies, is considering how the theme of Leveraging Big Data can be addressed in a transversal way, including how to leverage data being created and harvested by IEEE and its Societies. The Future Directions Committee was charged by the NIC in 2013 and 2014 to develop a stand-alone Big Data initiative, as the Cloud Computing Initiative moves out of the FDC and into a sustainable environment within the IEEE operational units. Typically the FDC first begins incubating a new initiative within Technical Activities, but in the case of Big Data, the incubation has taken place within the Cloud Computing Initiative.

Big Data: What?, How?, Why?



Big Data Directions

Networked & Programmable World

Pervasive Monitoring/Control

Internet of Things

OPEN DATA

COSM Xively

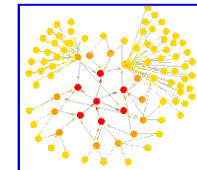
Immersive, Augmented Reality Interfaces



Diversification of Applications & Users



Next Gen Analytics, Prediction



Multiple RT Streams / Integration

Anywhere, AnyDevice

DGB 5/2013



Source	Target	Amount	Category
Source 1	Target 1	100	Category 1
Source 2	Target 2	200	Category 2
Source 3	Target 3	300	Category 3
Source 4	Target 4	400	Category 4
Source 5	Target 5	500	Category 5
Source 6	Target 6	600	Category 6
Source 7	Target 7	700	Category 7
Source 8	Target 8	800	Category 8
Source 9	Target 9	900	Category 9
Source 10	Target 10	1000	Category 10



Data Stream Mining

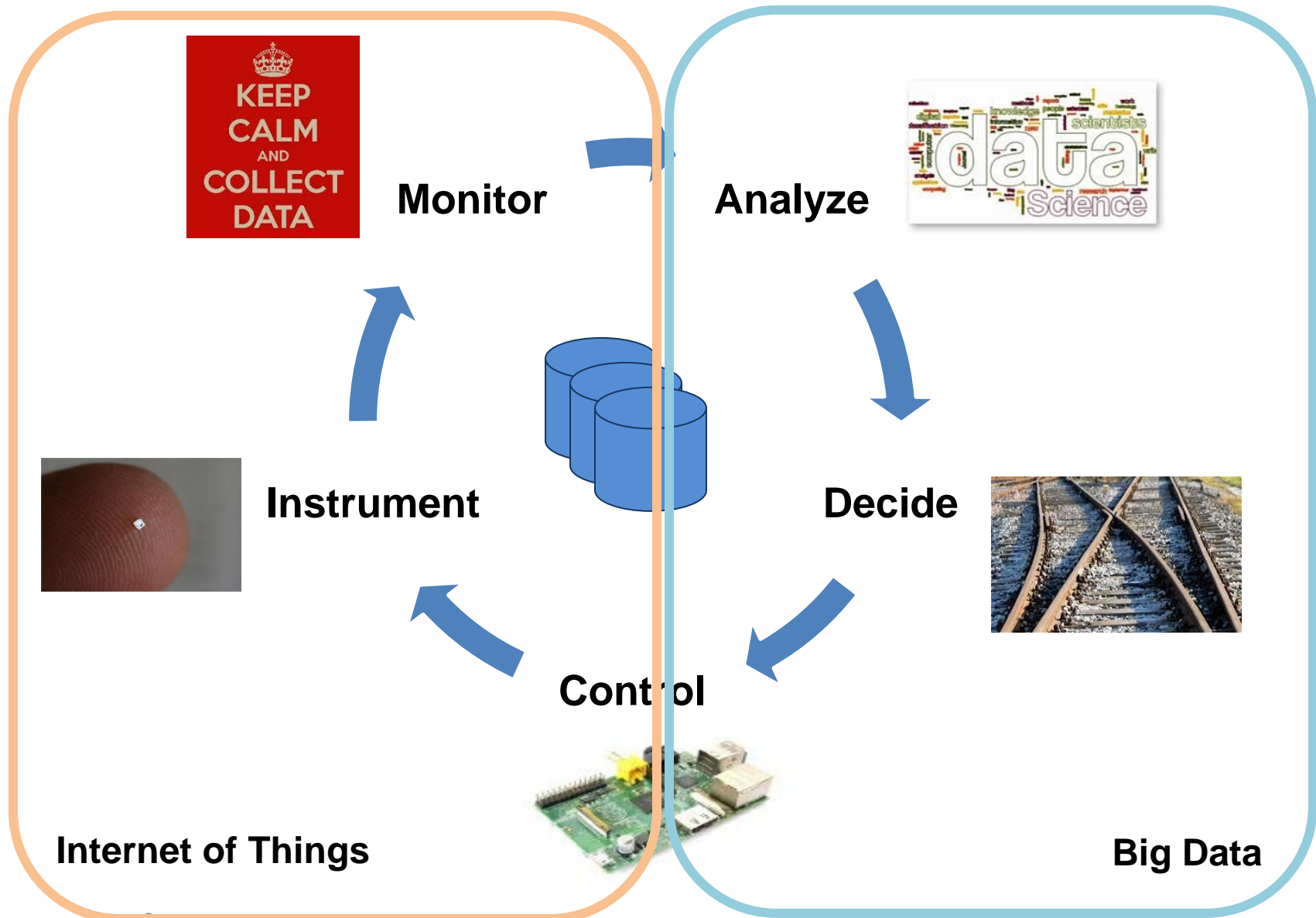


Speech/Text Mining



DATA

Networked & Programmable World: IoT + Big Data



Some Things That Make a Difference

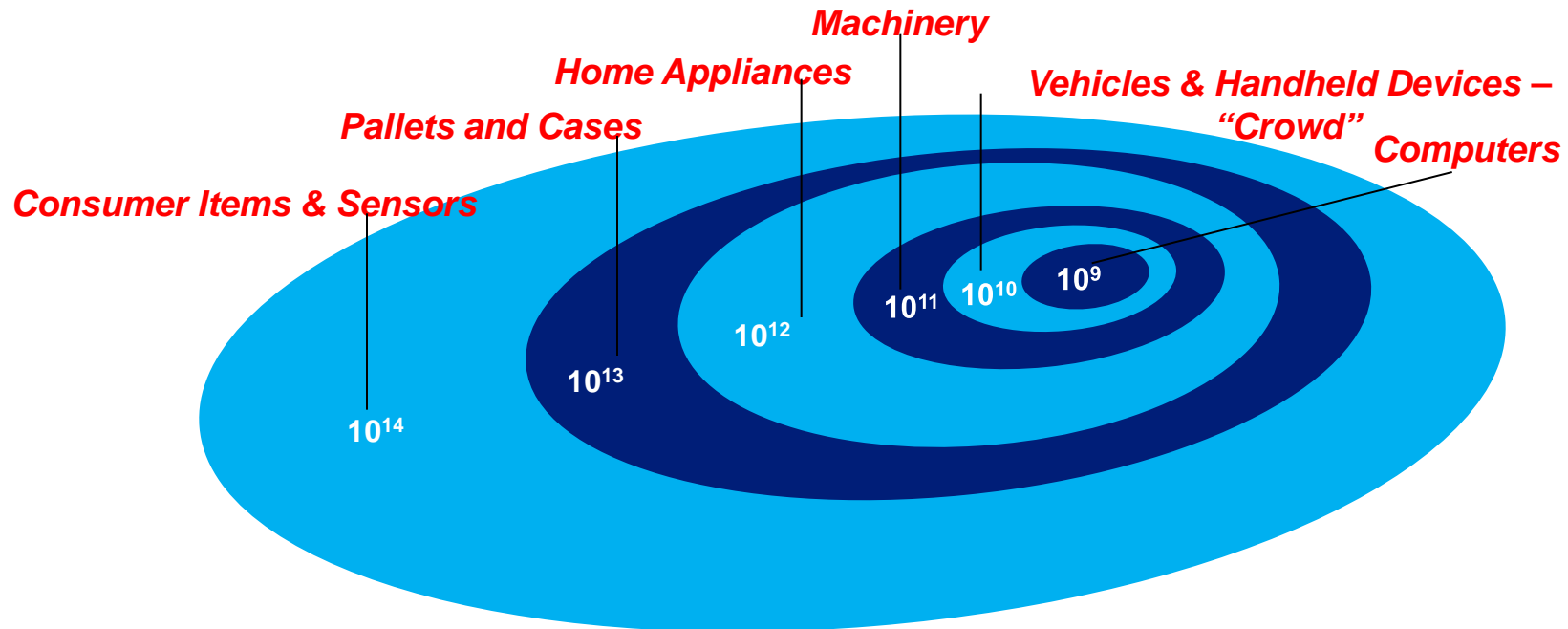
DATA

Change	Classical	Big Data	Example
Granularity	Transactional or Aggregate	Elementary, Personalized	Web Transaction vs Click Stream
Signal Strength	Strong	Σ Weak	Google Flu
Latency	Transactional or Long	Streams and Real time	Location
Location	ZIP Code, Area Code, Nxx,	GPS, Lat/Long	Location Based Systems
Structure	Relational	Structured, SemiStructured, Unstructured, Graph	Text/Speech Analysis, (e.g. Twitter), Social Networks
Sources of Data	Operational	Operational + Crowd + Sensors + IoT + Open	Dallas Museum of Art, data.gov - 156,584 datasets
Data Integration @ Scale	Hard, Join, Structured	Structured & Unstructured, Large Scale, Still not easy	

The next “crowd” - Objects

Devices That Can Be Networked & IP Addressable

How can we best exploit the billions of devices, many mobile, intelligent, and video enabled as computing, sensing, and communications platforms?



Invisible Computing - Will Far outnumber current IT Devices and People

- Home Appliances
- Consumer Items
- Pallets and Cases
- Machinery
- Vehicles and Handheld Devices
- Sensors – Machine to Machine Internet

Integration

- Communications
- Entertainment

Technology

Some Things That Make a Difference Technology


Change	Classical	Big Data	Example
Computing Platforms	Large Symmetric Multiprocessors, expensive	Parallelism using commodity hardware	Cost reduced by x10, Cloud
Software Platforms	RDBMS, Analytic Software, Viz Software	Oriented to massive parallelism, Often Open Source	Map/Reduce, Hadoop, Storm,...
Data Base Systems	RDBMS, Transactional	Often Column Oriented, Availability Oriented, Document	Hbase, MongoDB, Cassandra,
Visualization	Static, Aggregate, Dashboards	Interactive, Drill Down	Swift
Streams vs Warehouses	Warehouse: Size: N View: Query Latency: High	Streams: Size: ∞ View: Window Latency: Low	Fraud, High Frequency Trading, Targeted Marketing
Networks	Access speed limited,	Access: Variety and Bandwidth; Core: Bandwidth & Low Latency	Ubiquitous.
Standards	Common	Evolving, Needed	Data Query Lang.

Dimensions of Technology

DATA

Structured – SemiStructured - Unstructured

Networking: Efficient, Reliable, Secure Data Transport.



Computing: Storage and Processing Architectures that operate at scale, and in Real Time



Analysis: Industry leading information mining and analysis technology.



Visualization: The most effective way to deliver information & alerts to decision makers

APPLICATIONS

Fraud/Security – Marketing – CFO – Hosting – Mobility - ...

What Got Us Here: Technology Startups/Open Source Landscape

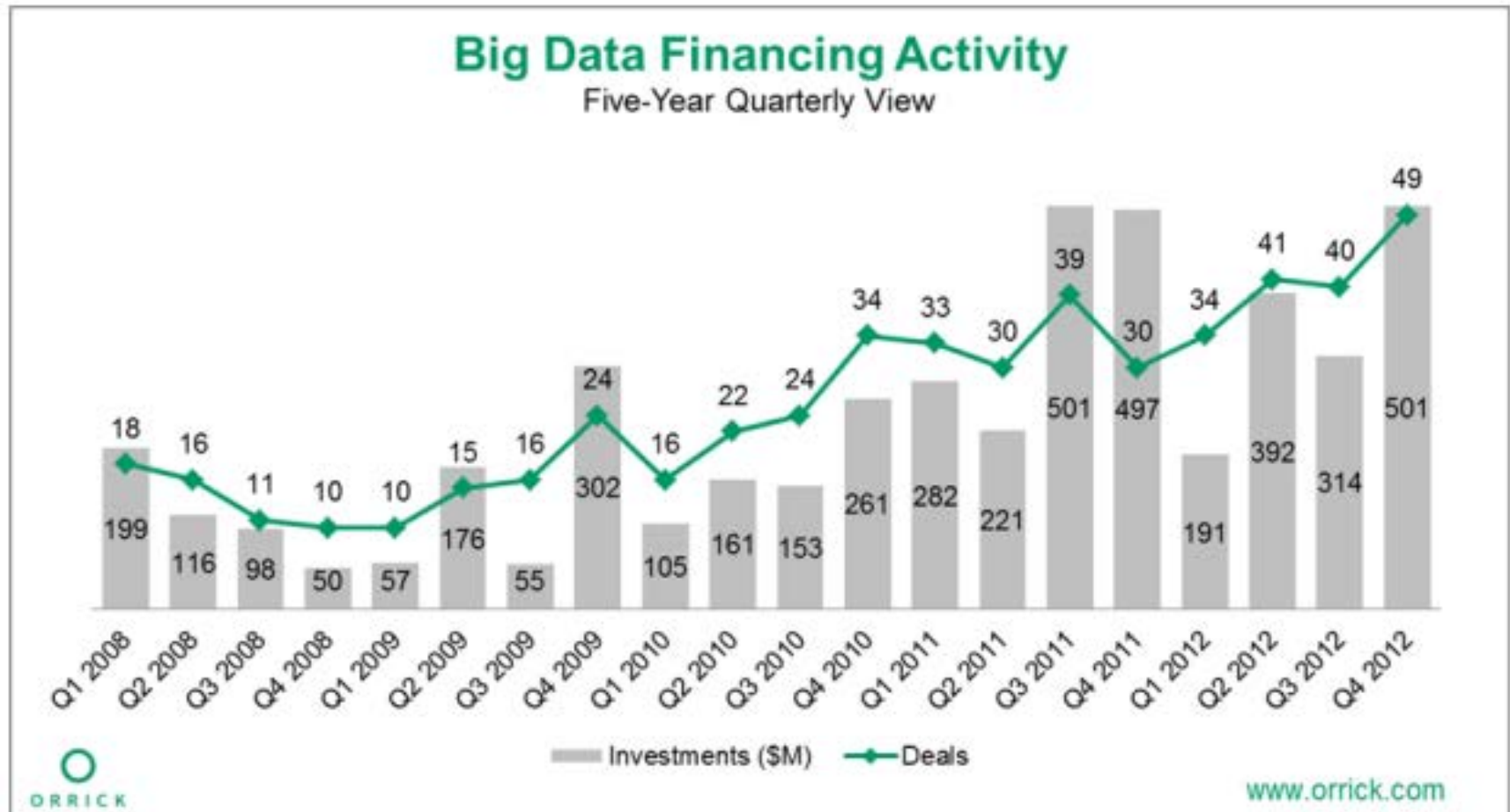
The BigData-Startups Open Source Landscape 2.0



Created by: www.bigdata-startups.com

Investment in Big Data

Investment in Big Data Companies in the Big Data industry have received [\\$4.9 billion in funding](#) since 2008, with the bulk of financing in 2011 and 2012.



Source: <http://www.kdnuggets.com/2013/03/top-big-data-vc-investors.html>

DGB 5/2013

At the end of the day, Big Data is about taking advantage of scale rather than letting scale take advantage of you.

<http://www.google.com/trends/explore#q=big%20data&cmpt=q>

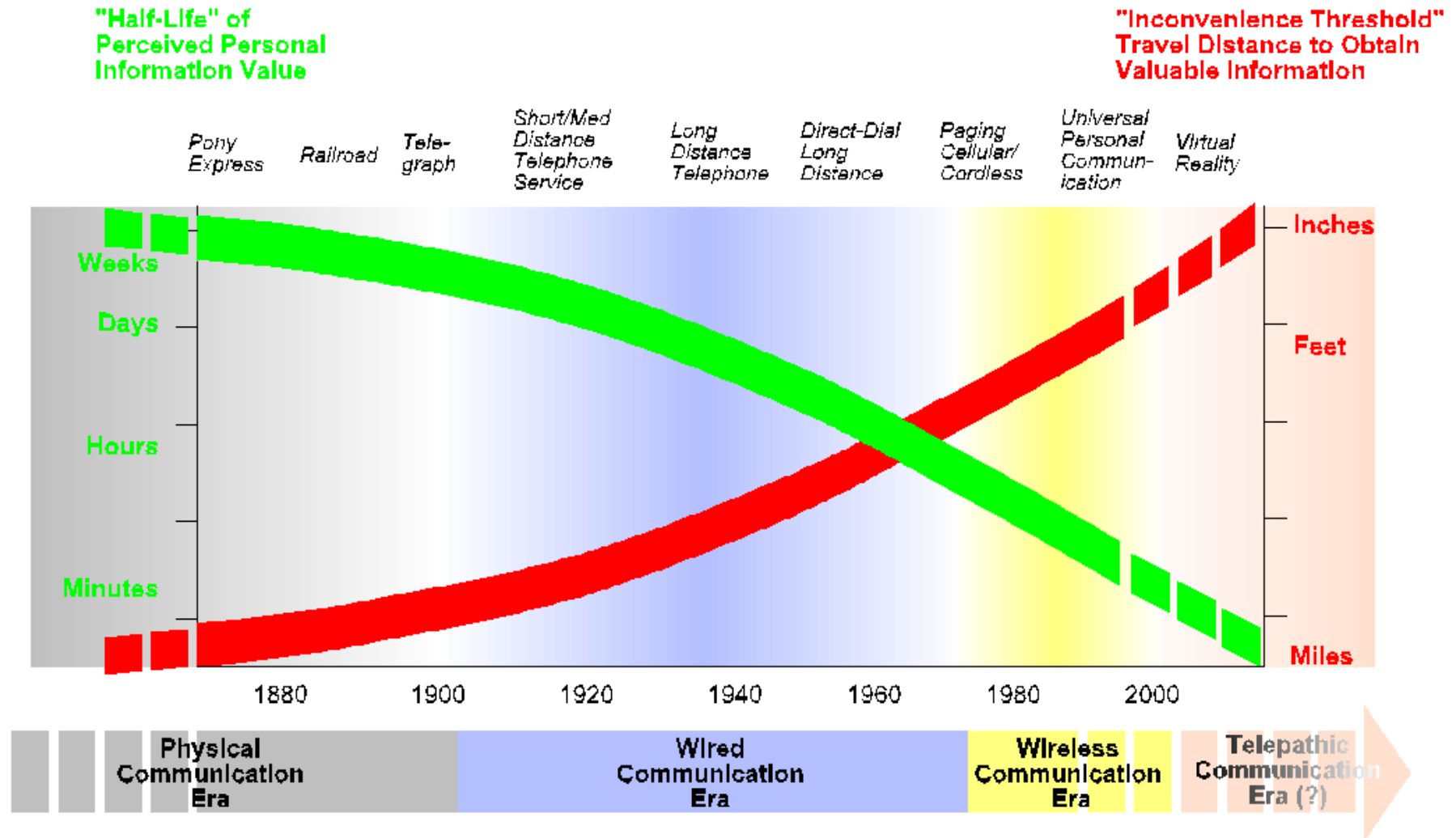
<http://www.google.com/trends/explore#q=data%20scientist&cmpt=q>

Applications

Making Scale Your Friend

Mobility Trends:

“Half Life” and “Inconvenience Threshold”



Each communication technology advance has shortened the useful life of information and increased the need to obtain new information more rapidly regardless of the situation or location...

Convergence of Communication and Entertainment

Can I find what I want and watch, when I want?

•Base Issue

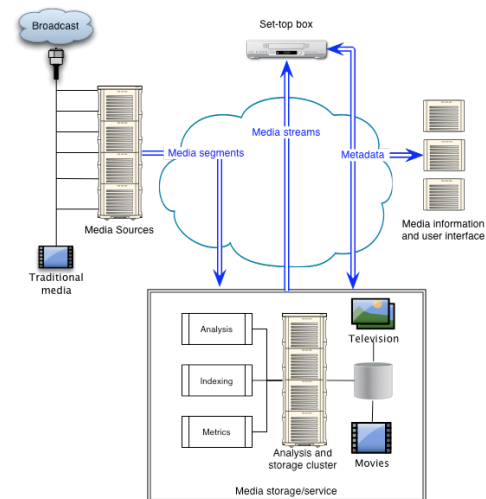
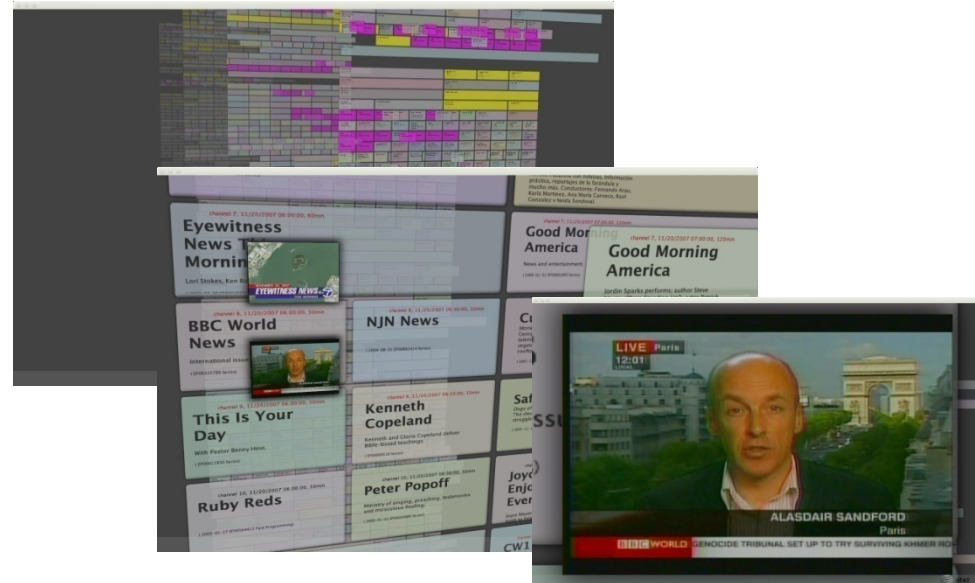
- How do we differentiate our media services in a world of commodity media services?

•Vision

- Access to all media, from all time, from anywhere

•Challenges

- User interface, especially search
- Indexing
- Storage
- Distribution



Some Things That Make a Difference Applications

Change	Classical	Big Data	Example
Granularity	Transactional or Aggregate	Atomic, Personalized	Individual's location at a specific time
Rare Events	Sample Based	Population, Signatures, Classification	Fraud, Security, Social Location
Broad Search	Indexed Retrieval, Flat Files	Structured and Unstructured data	Search Engines
Relationships	Customized	Graphs and (social) Networks	Recommender systems, Col
Analytics at Scale	Sample	Population	Call Detail, Call Center
Personalized Traits	Survey based	Behavioral	Signatures, Targeted Marketing
Process Control	Silo'ed	End to End, Feedback Control	Virtual Integration
Recommendations, Learning	Rules	ML, Collaborative, Relationship based	Translation, Recommenders
Transparency	Within Silo	Across Processes	Provisioning
Prediction	Limited, Aggregate	Real time, Specific	Slipper
Text/Speech/Video/Image	Limited, often special purpose	Text, Speech Now; Others soon	Many uses of Twitter, Contracts.
Control	Customized, Rule Based	Learning Based	Driverless Cars

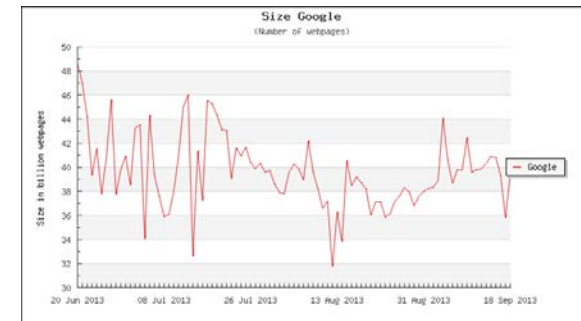
Canonical Examples of Big Data (1)

Search

Volume

```
<!DOCTYPE HTML
PUBLIC "-//W3C//DTD
HTML 4.01//EN"
"http://www.w3.org/TR/html
4/strict.dtd"> <HTML>
<HEAD> <TITLE>My first
HTML document</TITLE>
</HEAD> <BODY>
<P>Hello world! </BODY>
</HTML>
```

Variety



Canonical Examples of Big Data (2)

Fraud

Volume

Velocity

Initiating Phone Number	Terminating Phone Number	Time of Day	Day of Week	Duration	Bus/Res	Lat	Long		

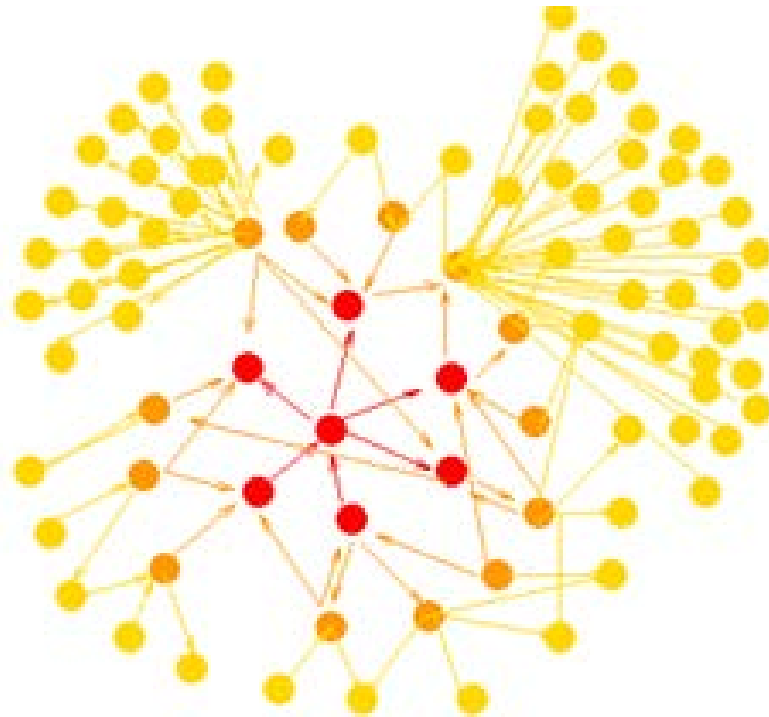
Canonical Examples of Big Data (3)

Marketing (Targeted/Viral/Recommenders)

Volume

Velocity

Variety



Canonical Examples of Big Data (4)

Call Center

Velocity

Variety



Canonical Examples of Big Data (5)

Machine Translation

Volume

Variety

Watson, come here, I want you

translate.google.com

Native

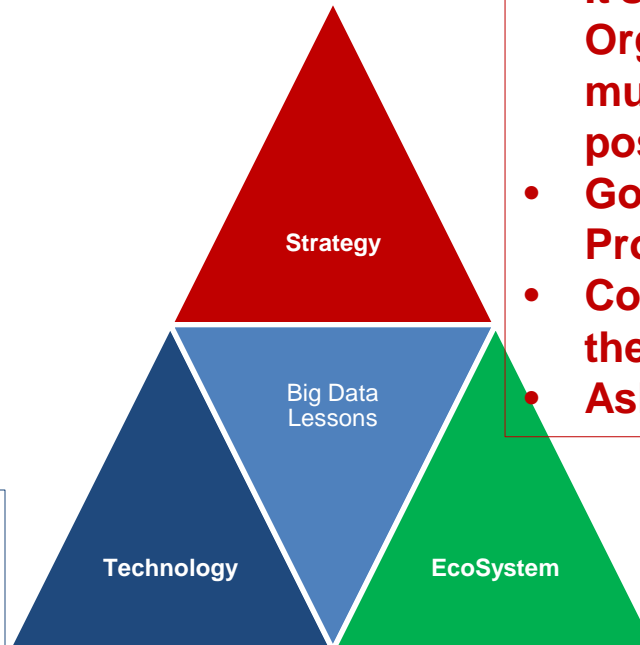
沃森，來到這裡，我要你
沃森，過來，我需要你。

Lessons from Experience

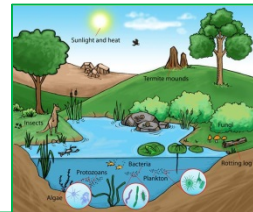
Lessons from Experience



- **Multidimensional Expertise is Essential**
-> Organizational Change
- The Nature of Analytics has Changed
- There are tradeoffs:
e.g. Speed vs. Accuracy
- Think Steams of Data

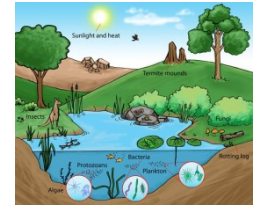


- **It's All About the Data:**
Organic & Inorganic, As much & as fresh as possible
- **Goal: Fundamental Process & Product Change**
- **Concentrate on Framing the Right Questions.**
- **Ask Big Questions**



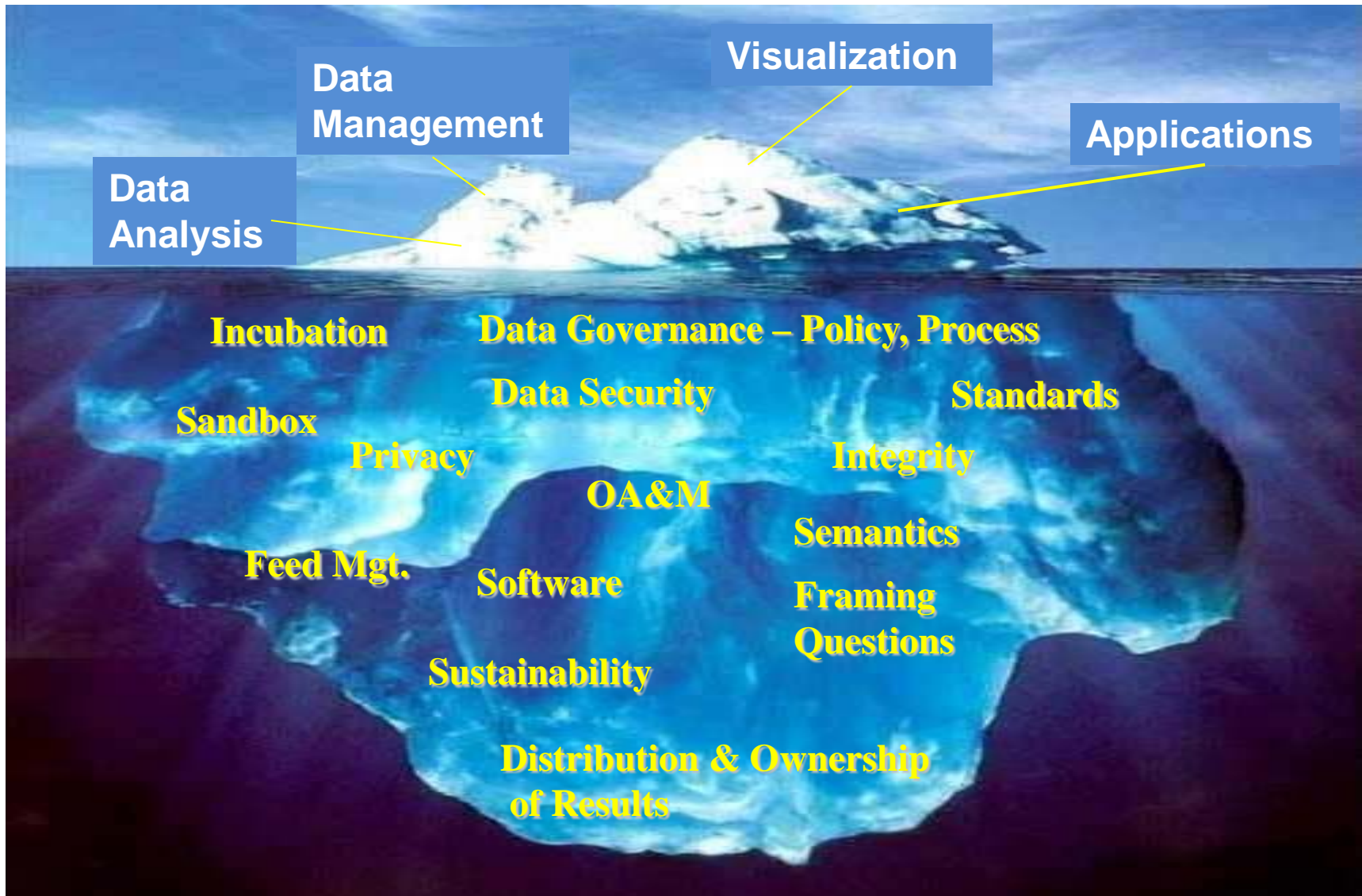
- **Do Not Ignore the Broader World:**
Data, Technology, Practices
- **We are entering a second generation, which will include many more players**
- **Customer Focus: Make Heroes of Your Partners**
- **RESEARCH – INCUBATION – PRODUCTION**

A Few More Lessons Learned



- **Technology (Leverage Aggressively):**
 - **Multidimensional technical expertise** is essential:
 - Network – Computing – Data Analysis – Visualization - Domains
 - **Data Feed management** can scale linearly. That is really bad.
 - Interactive visualization can be the most effective tool for both the Understanding and Integration of data.
 - Tradeoffs: **Speed vs. Accuracy : Rules vs. Differences : Weak & Strong Signals**
- **Strategy:**
 - It's all about the Data
 - Collect all the data you can.
 - Collect, and use, the freshest data available.
 - Do not ignore unstructured data.
 - Occam's Razor will win over time, though it may be a while.
- **Ecosystem (Leverage Expertise Wherever You Can):**
 - In many areas, **Standards** have not become firm. Look for them where possible.
 - Open Data and Self Generated Data are competitive necessities

A More Complete Picture



Meta Challenges

Then	Now
Significant systems containing sensitive data are not easily accessed. Complex semantics and poor integrity often exist, but impact is hidden because data is relatively closed. Integration, outside of joins, uncommon at scale.	Protection of SPI data a constant problem. Transparency of use, integrity of data a concern. Open data provides much more opportunity for interesting new apps from integration, and semantic confusion. Integration complex.

Security
Privacy
Integrity
Semantics
Integration

Data Quality and Integrity

Then	Now
Much of the burden of quality and integrity lies in the fact that the ACID properties and input rules of transactional systems are strictly enforced. There is a very rich technical ecosystem that has been built around integrity and is made available in most mature data management systems.	In many systems, at the volume, velocity, latency, and complexity expected, the levels of correctness required of transactional systems are neither possible nor necessary. Analytic techniques must take these changes into account. Given the very diverse nature of potential data sources, and consequent reduction in control over the data, this becomes a very challenging problem.

Organizing for Innovation - Maturity

Then	Now
Classical research or exploratory development teams create new products, often in large teams with significant timelines. Careful attention is paid to decision gates to prevent runaway costs. Due to costs, decisions often top down.	Small, elite teams create prototypes of potential products quickly, trial the prototypes, and, when successful, present for funding to go to market. Go to trial, very quick. Classical research provides technology base to prototyping teams, and partners with them.

Incubation and Production
Data OA&M
Data Governance

Data Governance

Then	Now
Data governance is typically the responsibility of the IT department, who, in cooperation with various stakeholders, develops policy, along with the processes to implement the policy.	Driven by the broader use models of big data, and the significance of non-IT stakeholders, policy should be done at a higher level. Protection of SPI data a constant problem. Transparency of use, and integrity of data are concerns. In general, the level of risk associated with the creation, curation, management, and use of data is considerably higher, hence governance is an even more critical task. The state of the art in this area is not yet mature.

Organization

Production



Incubation



R&D



Source: <http://crocodilian.com/cnhc/potm-apr02.html>
<http://www.frauleindi.com/HiltonHeadNature/Wildlife/wildlife.htm>
<http://www.animaltrial.com/alligator.html>

Some Thoughts

- Use of existing database, and expansion of data sources.
- Standards informed by breadth of expertise
- Leadership in second generation breadth and depth
- Integration with other emerging technologies such as IoT, Cloud, etc.

Thank You!

