

X-Informatics Introduction: What is Big Data, Data Analytics and X-Informatics? Part I

July 6 2013

Geoffrey Fox

gcf@indiana.edu

<http://www.infomall.org/X-InformaticsSpring2013/index.html>

Associate Dean for Research, School of Informatics and
Computing

Indiana University Bloomington
2013

X-Informatics and its Motto

Some Trends

- ➊ **The Data Deluge** is clear trend from Commercial (Amazon, e-commerce) , Community (Facebook, Search) and Scientific applications
- ➋ **Light weight clients** from smartphones, tablets to sensors
- ➌ **Multicore** reawakening parallel computing
- ➍ **Exascale initiatives** will continue drive to high end with a simulation orientation
- ➎ **Clouds with cheaper, greener, easier to use IT for (some) applications**
- ➏ **New jobs** associated with new curricula
 - ➐ **Clouds as a distributed system** (classic CS courses)
 - ➑ **Data Analytics** (Important theme in academia and industry)
 - ➒ **Social Media**

Some Terms

- **Data:** the raw bits and bytes produced by instruments, web , e-mail, social media
- **Information:** The cleaned up data without deep processing applied to it
- **Knowledge/wisdom/decisions** comes from sophisticated analysis of Information
- **Data Analytics** is the process of converting data to Information and Knowledge and then decisions or policy
- **Data Science** describes the whole process
- **X-Informatics** is use of Data Science to produce wisdom in field X

Big Data Ecosystem in One Sentence

Use **Clouds** running **Data Analytics Collaboratively**
processing **Big Data** to solve problems in
X-Informatics (or e-X)

X = Astronomy, Biology, Biomedicine, Business, Chemistry, Climate,
Crisis, Earth Science, Energy, Environment, Finance, Health,
Intelligence, Lifestyle, Marketing, Medicine, Pathology, Policy, Radar,
Security, Sensor, Social, Sustainability, Wealth and Wellness with
more fields (physics) defined implicitly
Spans Industry and Science (research)

Education: **Data Science** see recent New York Times articles
<http://datascience101.wordpress.com/2013/04/13/new-york-times-data-science-articles/>

X-Informatics already Defined

- Biomedical, Medical, Bio, Chem(istry), Health, Pathology, Wellness Informatics
- Life Style Informatics (from IT for Facebook to IT for life or Health – that's better Life Style Informatics))
- Astro(nomy), Climate, Earth Science, Energy, Environment, Radar, Sensor, Sustainability Informatics
 - Physics Informatics ought to exist but doesn't
- Social Informatics in our school
- Business, Wealth, Financial, Marketing Informatics
- Security (also in School), Crisis, Intelligence Informatics
- Policy Informatics (many X-Informatics impact policies)



Climate Informatics
network

How Wealth Informatics can help
with your financial freedom?



Xinformatics

xinfor
XIU TOU

Biomedical Informatics
Computer Applications in Health Care
and Biomedicine

AstroInformatics2012

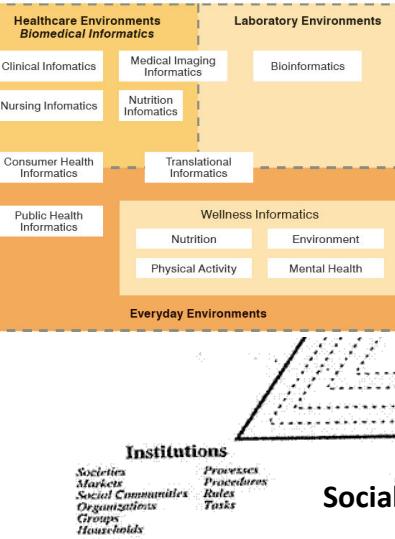
Redmond, WA, September 10 - 14, 2012

RICHARD E. NEAPOLITAN • XIA JIANG

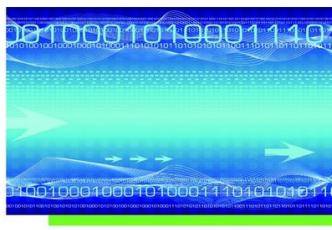
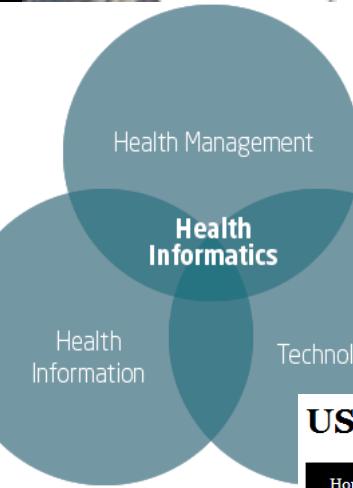
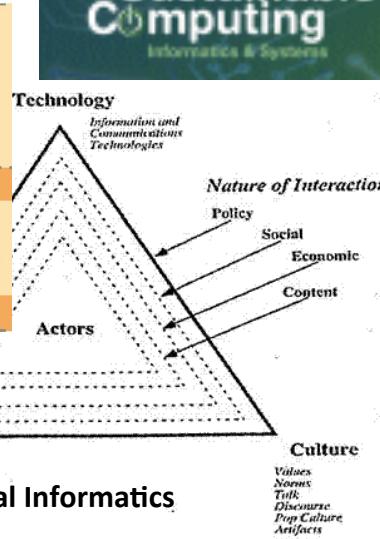
PROBABILISTIC
METHODS
FOR FINANCIAL AND
MARKETING
INFORMATICS



Sustainable
Computing
Informatics & Systems



Social Informatics



Noelia Penelope Greer (Ed.)
Business Informatics
Information technology, Management,



ASU School of Public Affairs
ARIZONA STATE UNIVERSITY

USC Center For Energy Informatics

Home Research Publications Smart Grids

GEO Informatics
Knowledge for Surveying, Mapping & GIS Professionals

About the Center

Welcome to the Center For Energy Informatics (CEI) at USC, an Organized Research Unit (ORU) housed in the [Viterbi School of Engineering](#). Energy Informatics is the application of information technologies to energy systems.

Lifestyle Informatics

Applications of Life Sciences in the field of healthcare
How is the training classified
Occupation Professions
Further study
Student at the University
Watch the movie
Studying abroad

Admission and registration
VU Honours Programme

BACHELOR-VOORLEIDINGSDAG
ZATERDAG 3 NOVEMBER

LOOP EEN DAG MEE MET EEN STUDENT

ENVIRONMENTAL INFORMATICS

Combine body, mind, and soul for a healthier, more active life through training

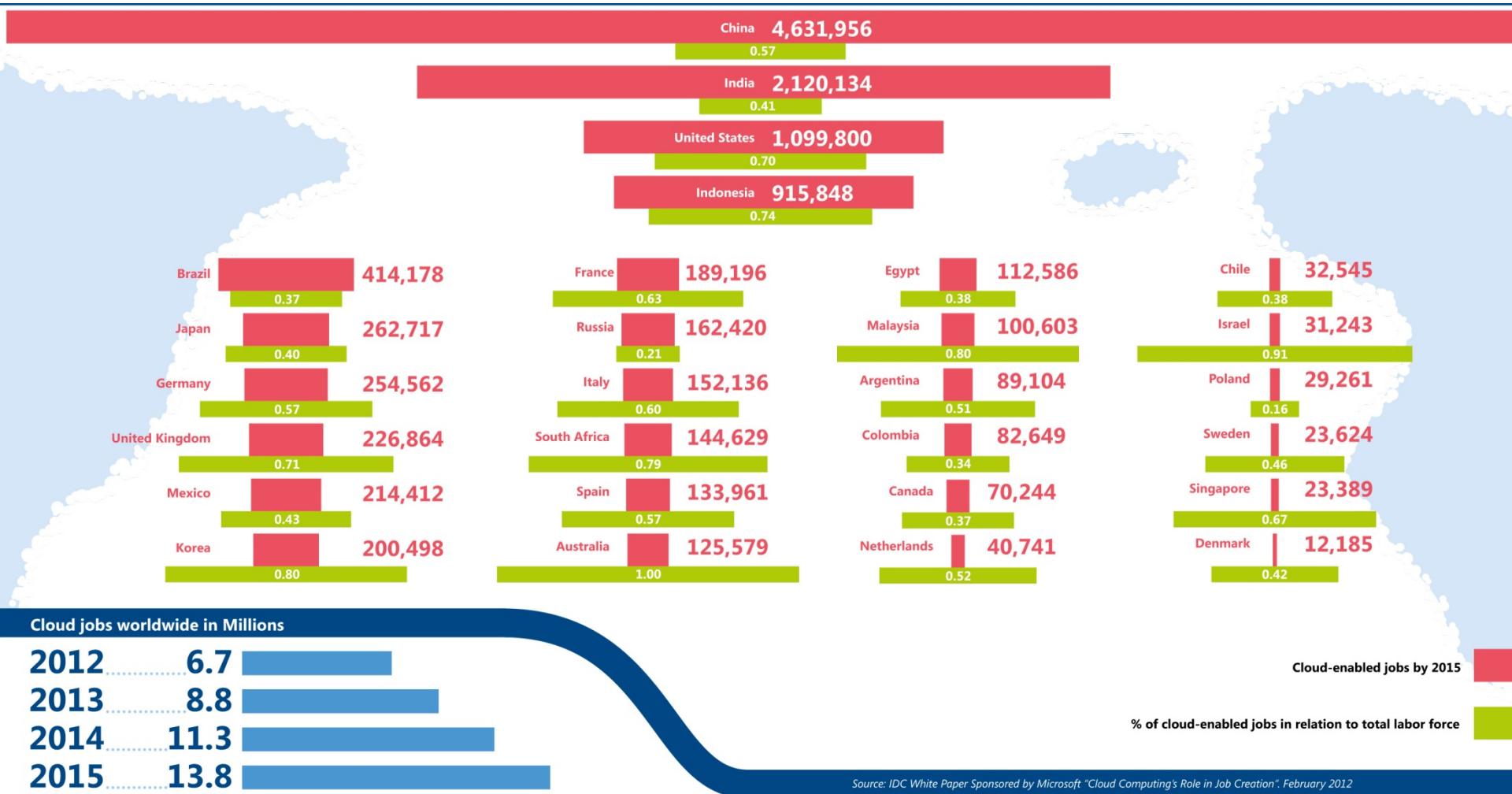
Lifestyle Informatics: Let people live longer and better

The study Lifestyle Informatics is about studying the impact of lifestyle on health. This bachelor including applied psychology, ergonomics, and informatics. It provides knowledge about language and information processing, how people live their lives, and how they can live better.

Lifestyle Informatics

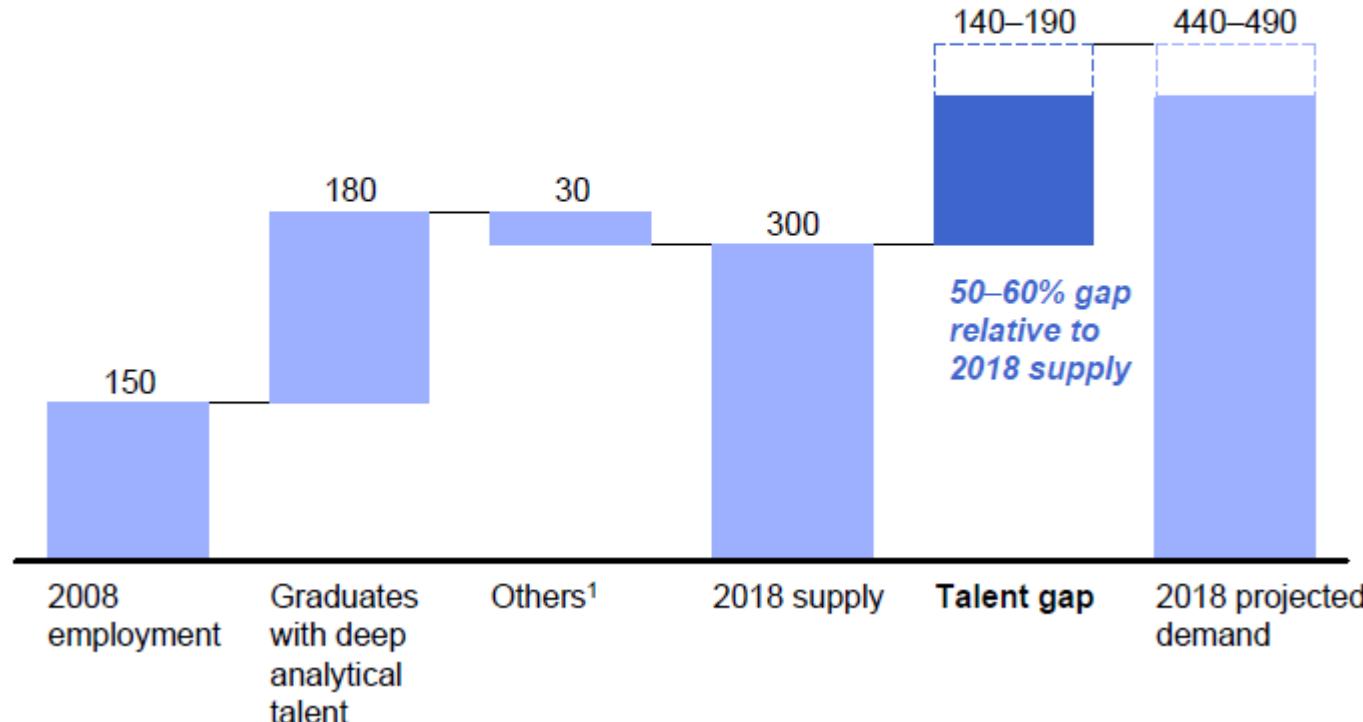
Jobs

Jobs v. Countries



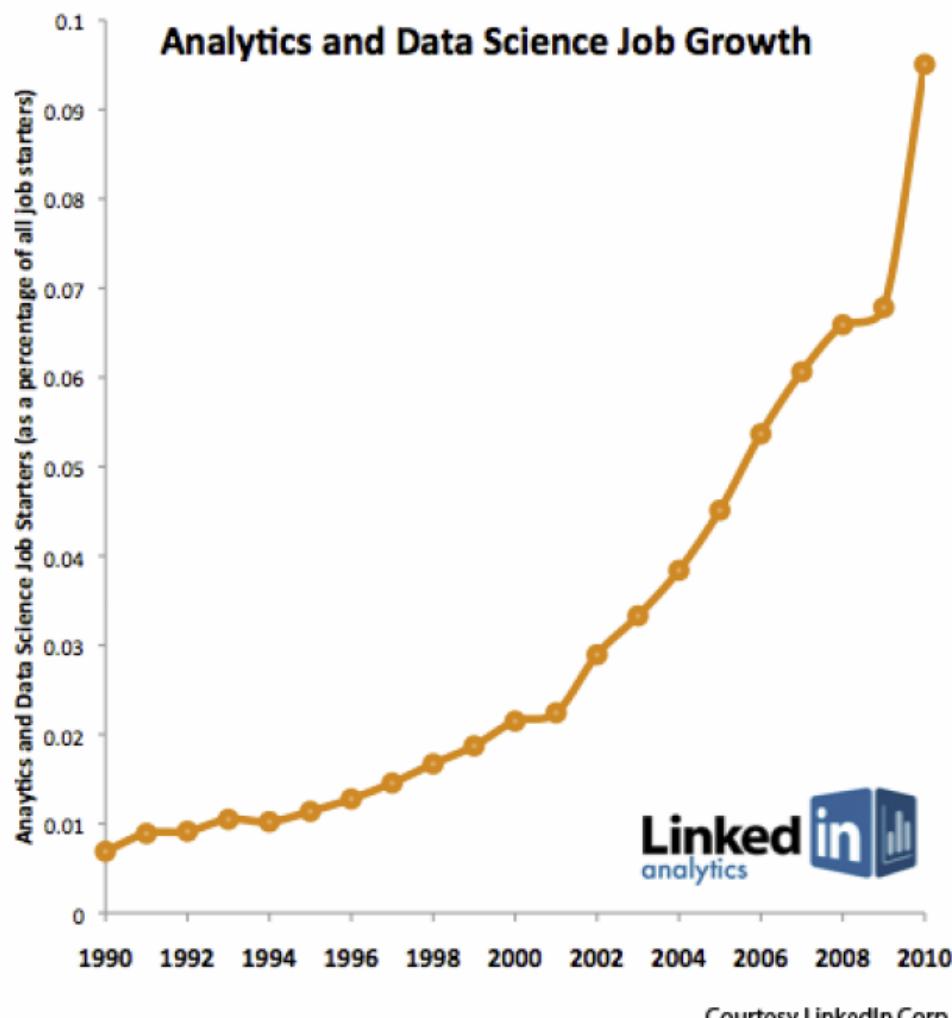
<http://www.microsoft.com/en-us/news/features/2012/mar12/03-05CloudComputingJobs.aspx>

McKinsey Institute on Big Data Jobs



- There will be a shortage of talent necessary for organizations to take advantage of big data. By 2018, the United States alone could face a shortage of 140,000 to 190,000 people with deep analytical skills as well as 1.5 million managers and analysts with the know-how to use the analysis of big data to make effective decisions.
- This course aimed at 1.5 million jobs. Computer Science covers the 140,000 to 190,000

The Rise of Data Scientists and Analysts



Tom Davenport Harvard Business School http://fisheritcenter.haas.berkeley.edu/Big_Data/index.html Nov 2012

Data Deluge General Structure

Some Data sizes

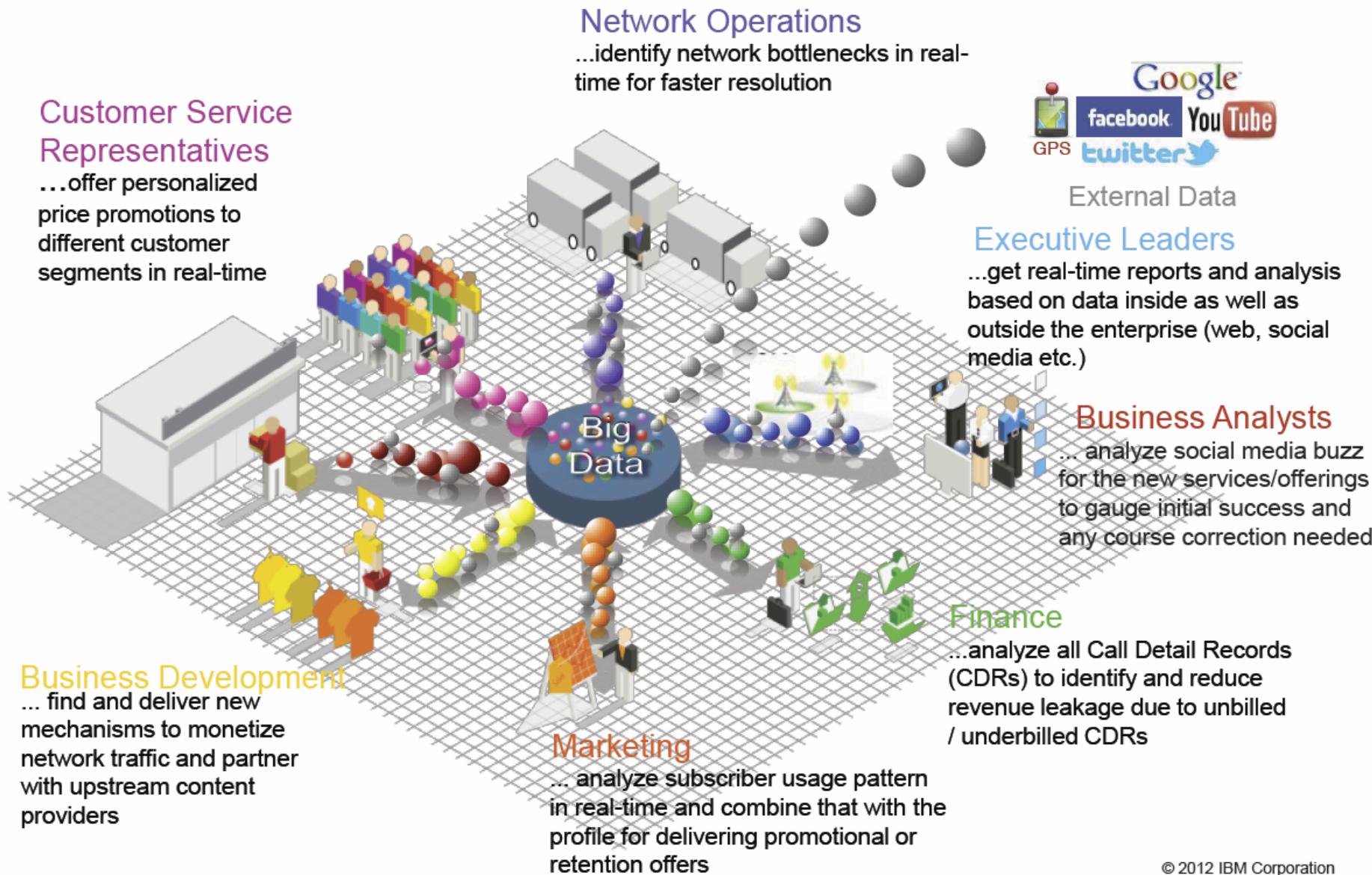
- ➊ ~ $40 \cdot 10^9$ **Web pages** at ~300 kilobytes each = 10 Petabytes
- ➋ **Youtube** 48 hours video uploaded per minute;
 - ➌ in 2 months in 2010, uploaded more than total NBC ABC CBS
 - ➍ ~2.5 petabytes per year uploaded?
- ➎ **LHC** 15 petabytes per year
- ➏ **Radiology** 69 petabytes per year
- ➐ **Square Kilometer Array Telescope** will be 100 terabits/second
- ➑ **Earth Observation** becoming ~4 petabytes per year
- ➒ **Earthquake Science** – few terabytes **total** today
- ➓ **PolarGrid** – 100's terabytes/year
- ➔ **Exascale simulation** data dumps – terabytes/second

Some Use Cases for Big Data

- Social media analytics – "People You May Know" at LinkedIn
- Voice analytics – Call center triage
- Text analytics – Voice of customer, sentiment analysis, warranty analysis
- Video analytics – Intelligence, policing, retail applications
- Genome data – what genetic profiles are associated with certain cancers?



The new era of analytics delivers value across the enterprise



The Rise of Big Data

What is it?

- Data that's too big (petabytes), too unstructured (not in rows and columns), or too diverse (mashups) to be stored and analyzed by conventional means (also relative)

Where does it come from?

- Internet/social media
- Genomic analysis
- Voice and video
- Sensors everywhere

What is to be done with it?

- Structure, classify, and count it
- Then analyze it (just as you would small data)



What's Different About Big Data?

The need for continuous flows of data, not stocks

- Stocks may be useful to develop models, but big data eventually requires a continuous process of analysis on moving data

Data scientists, not analysts

- IT “hacking” abilities in addition to the usual analyst attributes
- Scientific and exploration focus
- Closer to the product or process

New ways of deciding and acting on it

- It just keeps on coming, so have to establish ongoing processes to manage or decide on it



What's Different About Big Data? (cont.)

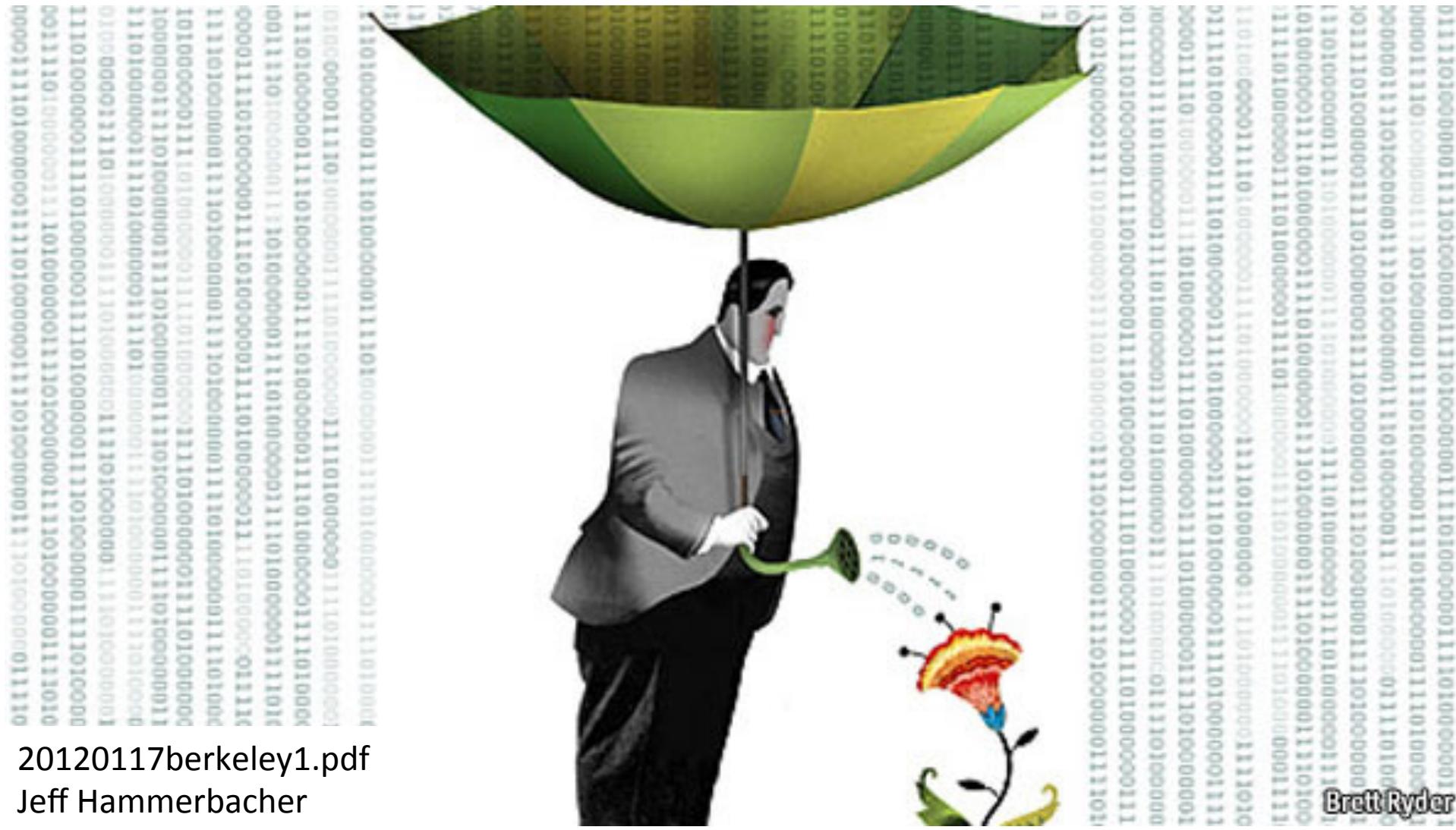
New technologies
to manage it

- Filtering, structuring, and classification tools – MapReduce, Hadoop, etc.
- Content analytics tools--NLP
- Data redundancy management
- Cloud analytics
- Machine learning
- Open source everything, including R (it's capable, it's free, and that's what everybody coming out of school wants to use)



The data deluge: The Economist Feb 25 2010 <http://www.economist.com/node/15579717>

According to one estimate, mankind created 150 exabytes (billion gigabytes) of data in 2005. This year(2010), it will create 1,200 exabytes. Merely keeping up with this flood, and storing the bits that might be useful, is difficult enough. Analysing it, to spot patterns and extract useful information, is harder still. Even so, the data deluge is already starting to transform business, government, science and everyday life



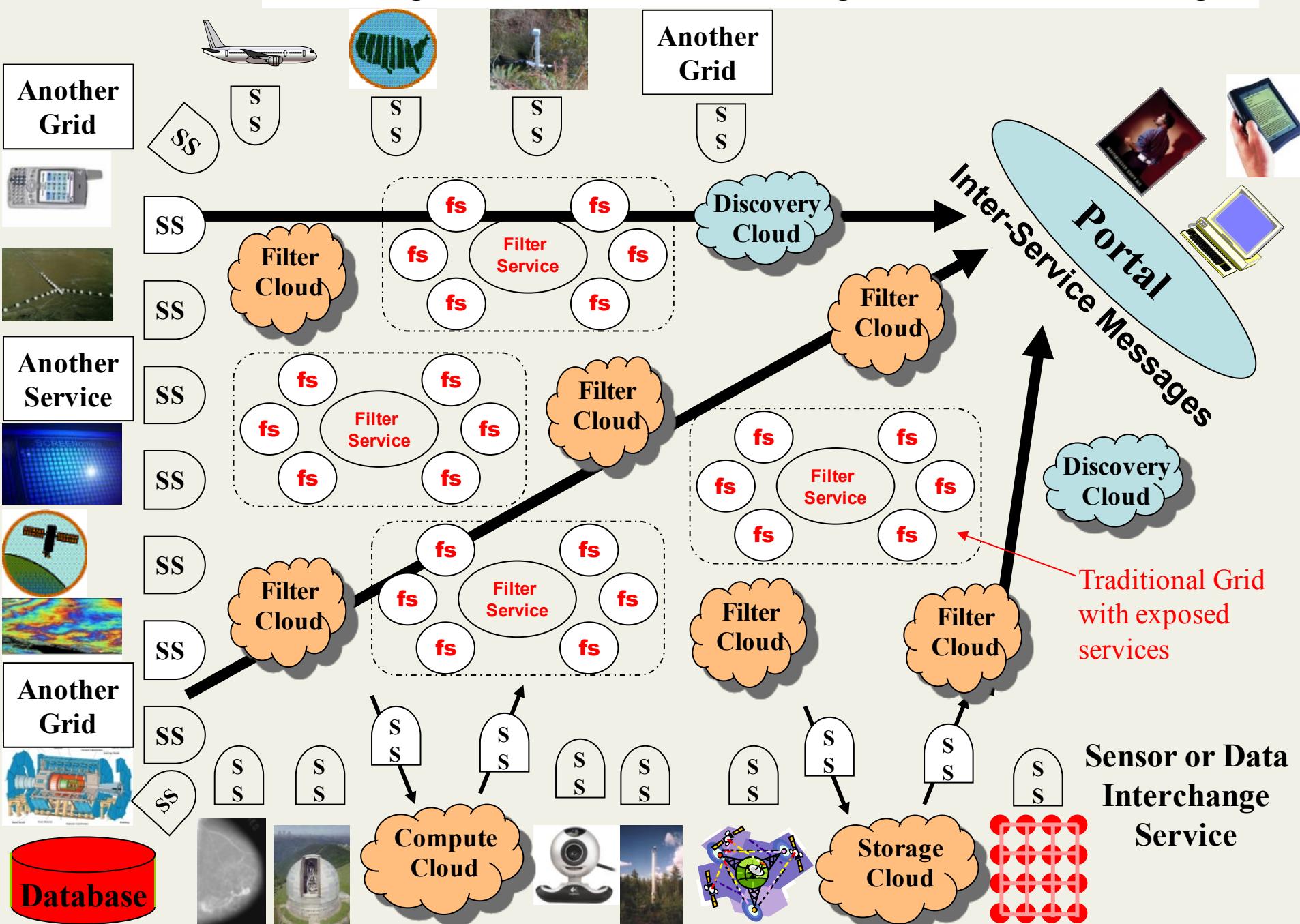
Data Science Process

DIKW Process

- **Data** becomes
- **Information** becomes
- **Knowledge** becomes
- **Wisdom or Decisions**
 - Community acceptance of results or approach important here
 - Volume of bits&bytes decreases as we proceed down DIKW pipeline

Raw Data →

Data Deluge is also Information/Knowledge/Wisdom/Decision Deluge? s



Example of Google Maps/Navigation

- Data comes from traditional maps (US Geological Survey), Satellites (overlays) and street cams
- Information is presented by basic Google Maps web page
- Knowledge is a particular optimized route
- Decisions (wisdom) comes from deciding to drive a particular route

Data Deluge Internet

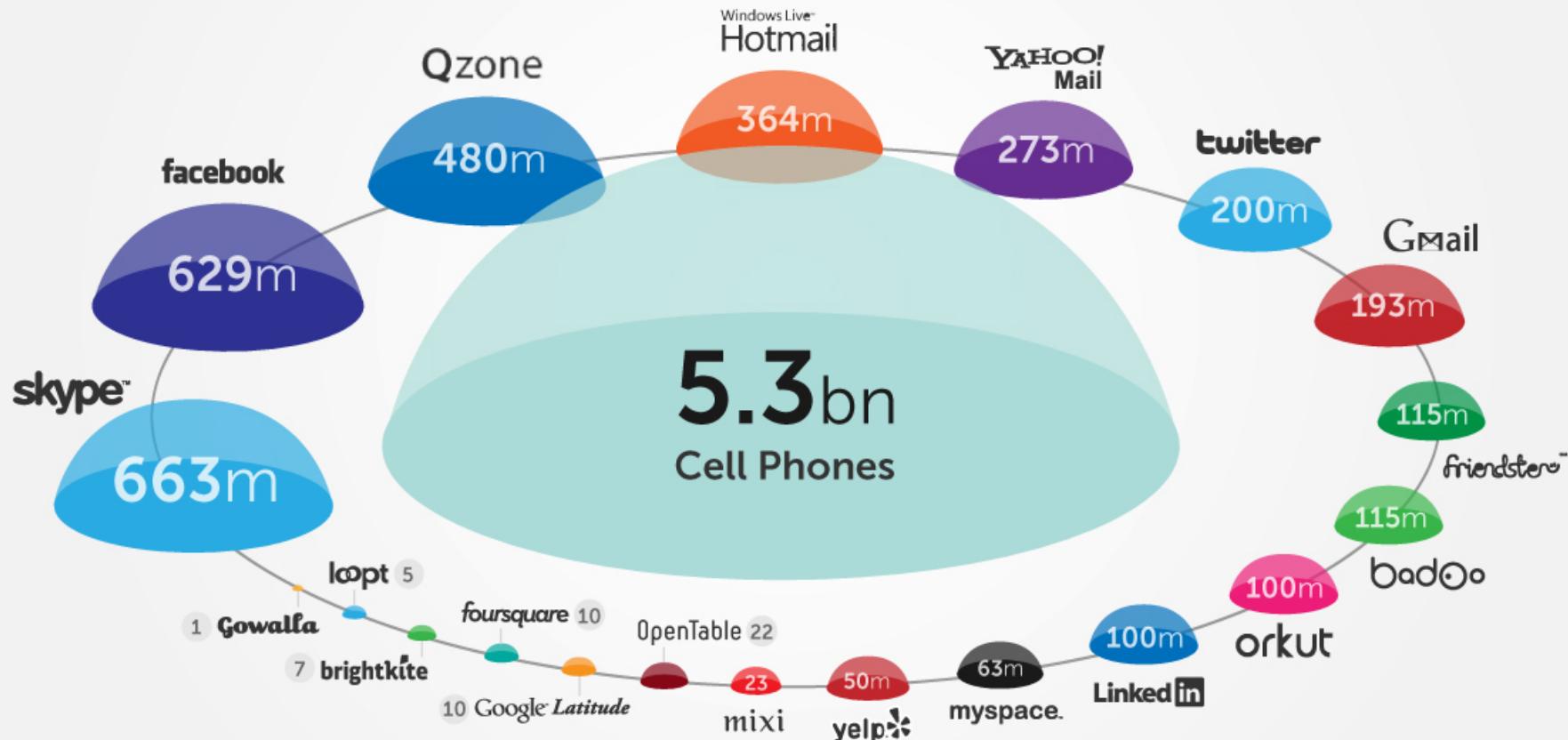


Where is big data coming from?



the geosocial universe

Brought to you by JESS3



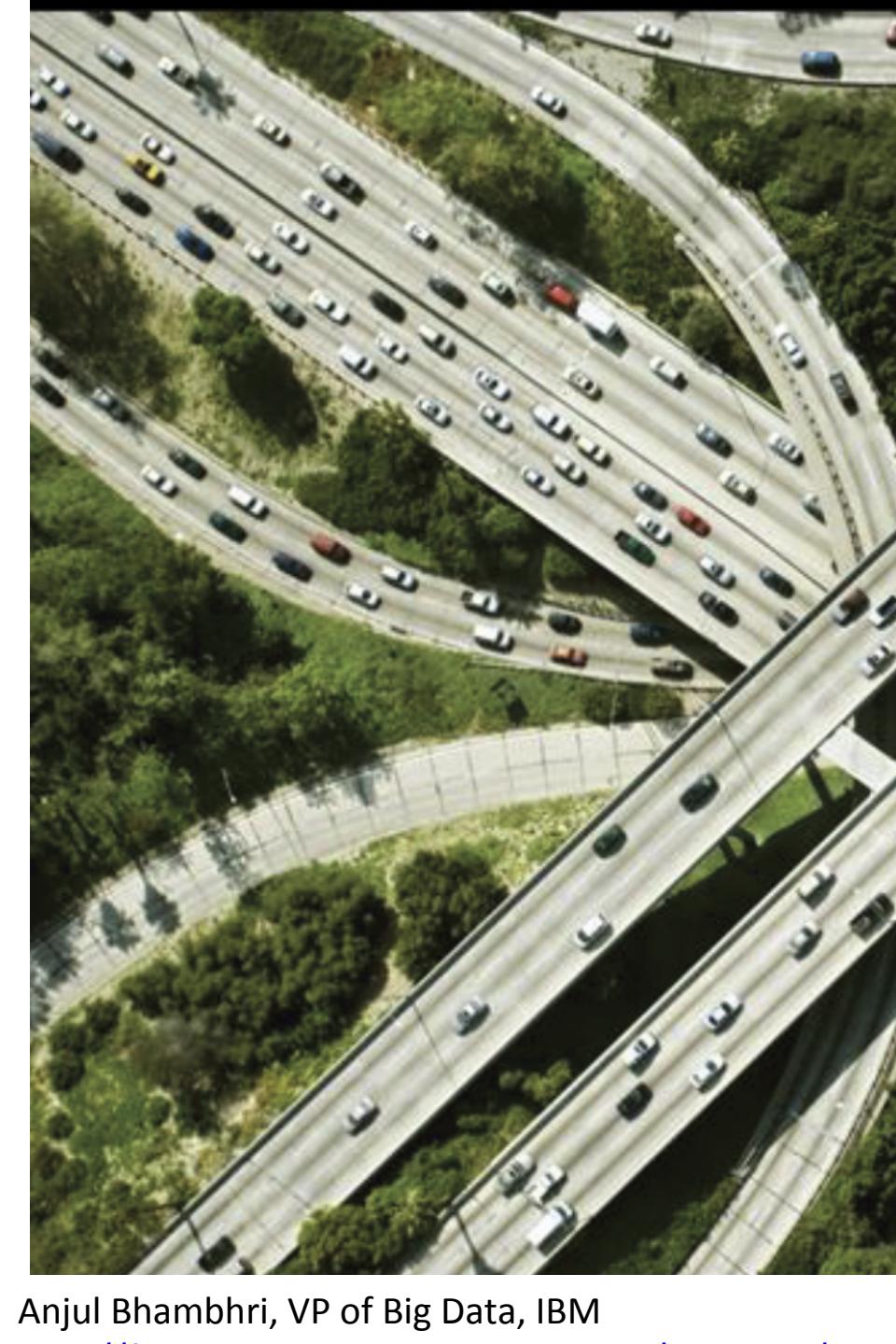
Data Deluge Business



Vestas optimizes capital investments based on **2.5 Petabytes** of information.

- Model the weather to optimize placement of turbines, maximizing power generation and longevity.
- Reduce time required to identify placement of turbine from weeks to hours.
- Incorporate 2.5 PB of structured and semi-structured information flows. Data volume expected to grow to 6 PB.

Vestas



Dublin City Centre Increases Bus Transportation Performance

Capabilities Utilized:

Stream Computing

- Public transportation awareness solution improves on-time performance and provides real-time bus arrival info to riders
- Continuously analyzes bus location data to infer traffic conditions and predict arrivals
- Collects, processes, and visualizes location data of all bus vehicles
- Automatically generates transportation routes and stop locations

Results:

- Monitoring 600 buses across 150 routes
- Analyzing 50 bus locations per second
- Anticipated to Increase bus ridership





Asian telco reduces billing costs and improves customer satisfaction.

Capabilities:

Stream Computing
Analytic Accelerators

Real-time mediation and analysis of
6B CDRs per day

Data processing time reduced from
12 hrs to 1 sec

Hardware cost reduced to 1/8th

Proactively address issues
(e.g. dropped calls) impacting customer
satisfaction.

CDR = Call Data Record

Anjul Bhambhani, VP of Big Data, IBM

Forces Shaping the Future

GE is a company that builds the machines that make the world work and has access to and deep understanding of the information that can make them work better

1. Internet

Hyper-connectivity: a living network of machines data and people

Internet of things: more devices tap into the Internet than people on Earth to use them

Internet of Things

2. Intelligent Machines

Increasing system intelligence through embedded software

Rise of machines: networked devices overtook the global population in 2011

3. Big Data

Democratization of data

Data overload: 2.5 quintillion bytes of data created every day

4. Analytics

Generating data-driven insights

Enhancing asset performance by detecting & predicting forecasts

Algorithms on installed base

10011010
10101010
10101010
11000101



imagination at work

Scale of Industrial Internet

Social media versus electric generating power source

2012 Twitter Usage

Gas Turbine Compressor Blade Monitoring potential*

vs.



80 Gigabytes per day

enabling social connections



588 Gigabytes per day

enabling capital asset productivity

Data volume potential is 7x greater from a gas turbine than current Twitter usage



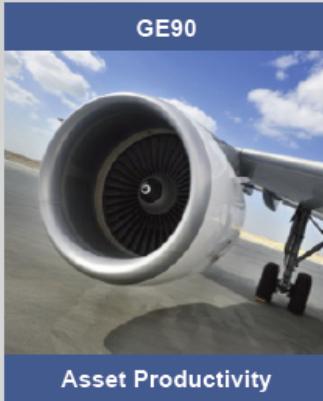
imagination at work

Value of Data & Analytics

Monitor fleet of ~25,000* engines ... 3.6MM flight records/month



B777



GE90



Prognostics

- ✓ Dispatch reliability
- ✓ Preventive maintenance
- ✓ Asset utilization

Prevent failures = customer efficiency

DATA

90,000 flight records analyzed
~200 parameters per flight record
~18MM parameters per month

System & Optimization

- ✓ Time & space management
- ✓ Fuel efficiency
- ✓ Airspace capacity

Drives strong alignment with customers

Creates productivity in long-term service agreements

Value-added services fuels growth

MM = Million

Asset Productivity

- ✓ Enhanced service offerings
- ✓ Airline cost structure
- ✓ Fuel performance

Streamline operations = increased airline productivity

Integrated systems = value-added services



imagination at work

Trend: Analytical Software Solutions

- Integrated software solutions for business problems

Problem	Industry
Fraud detection	Banking
Credit and operational risk	Banking
Credit scoring	Financial services
Warranty analysis	Manufacturing
Customer retention	Telecommunications
Markdown optimization	Retail

- What are the opportunities for statisticians?

every
49 minutes
a Ford Mustang
is sold



every
5 seconds
a cell phone
is sold



every
6 seconds
a pair of shoes
is sold



108+ million

Active buyers and sellers
worldwide

250+ million

Queries every day to the eBay
search engine

350+ million

Live global listings

20+ petabytes

Of data in our Hadoop and Teradata
clusters

2 billion

Page views each day

75 billion

Database calls each day

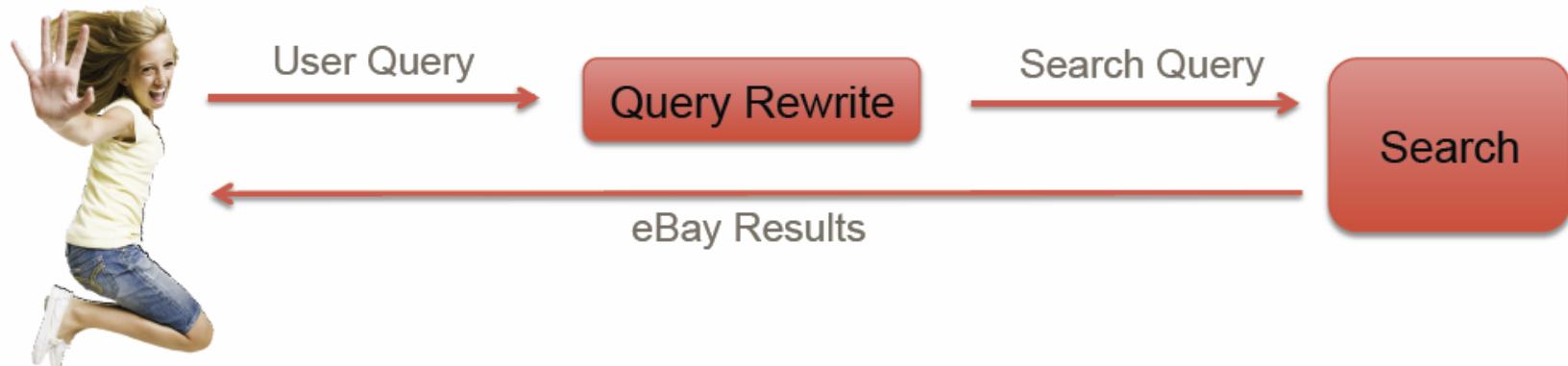


BIG DATA IS TRANSFORMATIONAL

- Big data informs on:
 - *Patterns*
 - Anomalies and outliers
 - Generalizations
 - *Predictions*
 - *Relative performance*
 - An holistic customer picture
- Vast array of applications at eBay:
 - *Product development, A/B testing, system performance, fraud and risk detection, purchase prediction, customer support, buyer demand, seller intelligence, financial performance, ...*

PATTERNS: QUERY REWRITES

- In 2010, our search engine was very literal: it matched exactly what you typed
- We're on a journey to make it more intuitive, so it does a great job of understanding user intent and finding all of the relevant results
- Idea: Mine our extreme data, look for patterns, and use these to map words in user queries to **synonyms** and **structured data** associated with items for sale at eBay



PATTERNS: QUERY REWRITES ...



HOW DO BUYERS PURCHASE THE PILZLAMPE?

- It turns out, they try one (or more) of a few things:
 - Type **pilzlampe**, and purchase
 - Type **pilzlampe**, ... , **pilz lampe**, and purchase
 - Type **pilzlampe**, ... , **pilzlampen**, and purchase
 - Type **pilz lampen**, ... , **pilzlampe**, and purchase
 - ...

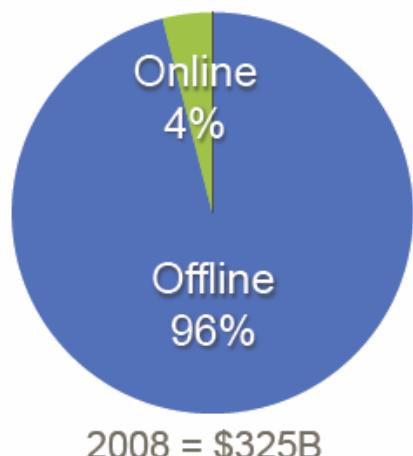
PATTERNS: QUERY REWRITES ...

- From our big data mining:
 - We automatically discover that *pilz lampe* and *pilzlampe* are the same
 - We also discover that *pilz* and *pilze* are the same, and *lampe* and *lampen* are the same
- From these patterns, we rewrite the user's query *pilzlampe* as:
*pilzlampe OR "pilz lampe" OR "pilz lampen" OR pilzlampen
OR "pilze lampe" OR pilzelampe OR "pilze lampen" OR
pilzelampen*

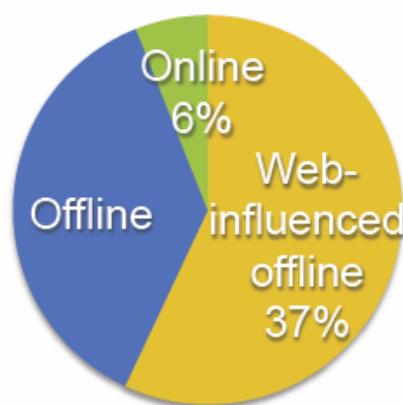
IT'S JUST COMMERCE

There's no longer online and offline

Yesterday



Today



Tomorrow



■ Offline ■ Online ■ Web-influenced Offline

Source: Forrester, Euromonitor and
Economist Intelligence Unit

Source: Forrester

Source: Economist Intelligence Unit