

Big Data is all about data that we don't have

David Taniar

Monash University, Australia

EGC'2015, Luxembourg, 28 January 2015











🔀 Jakarta to Amsterdam - Go 🗙 🔪





k - dtaniar@gmail.com ×

Amsterdam to Luxembourn × 🕔

tps://www.google.nl/maps/dir/Amsterdam/Luxembourg+City,+Luxembourg/@51.0498988,3.8252892,7z/data=!4m14!4m13!1m5!1m1!1s0x47c63fb5949a77 Leeuwarden Groningen Lincoln Peak District Leer Bremer Oldenburgo dtaniar@gmai 50 Q × Sneek / Drachten Assen Den Helder Delmenhorst 1 h.from €183 Vetherlands—Luxembourg City. Emmen E45 Leicester Norwich Hoorn Cell 5 Toll Peterborough Alkmaaro Lelystad Zwolle Lingen E22 Lowestoft GLAND Hanover 31 Corby Amsterdam O Amsterdam Rheine Osnabrück Minden Deventer Northampton Cambridge Schiphol Airport Enschede Netherlands Hildesheimo Ipswich Bielefeld Si Arnhem Münster O Detmold Rotterdam M40 Colchester Pelixstowe am Luton Gütersloho 33 Nijmegens Oxford Chelmsford Dordrechto OClacton-on-Sea Paderborn E35 Hamm 0 Breda vindon Southend-on-Sea Göttinger London 0 57 Essen Reading oDortmund Vlissingen Tilburg M4 0 Gillingham Margate 0 M25 Duisburgo Eindhoven Kassel Basingstoke oWoking Düsseldorf Hagen Antwerp Maidstone Deal Bruges E34 Bad 0 oSolingen Wildungen Crawley Ostend Ghent 0 Ger Ashford Dover Dunkirk Mechelen E313 1 h Diksmuide Cologne Calais Hasselt from €183 Southampton O E40 Siegen E40 0 Marburg Alsfeld Hastings Brussels Brighton Roeselare Maastrichto 0 Bonn emouth Portsmouth Boulogne-sur-Mer Roubaix Aachen Liège E411 Giessen Fulda Béthune Lille oTournai 61 Verviers Belgium Neuwied Wetzlar Berck Arras ODouai Mons 6 5 Koblenz Cambrai lish Channel Frankfurt Abbeville A1 A26 66 0 Schw Wiesbadeno E44 Dieppe Aschaffenburg Hirson Amiens E42 Mainz Saint-Quentin Bad oDarmstadt Würzbur Charleville-Mézières Fécamp Trier A29 Kreuznach Worms Cherbourg A1 Luxembourg City O Luxembourg Laon 63 Le Havre Bolbec Mannheim Rouen Beauvais Compiègne Merzia 0 Esch-sur-Alzette Elbeuf Neustadto Heidelberg Reims Soissons Saarbrücken A13 Bayeuxo Caen Thionville A16 Lisieux Louviers Heilbronn E50 Metz Creutzwald Sarreguemines 65 Saint-Lô O. Évreux Cergy Karlsruhe 81 Bernay Châlons-en-Champagne Pont-à-Mousson Paris Pforzhelmo Vire Vire Google Granville 0 Flers Versailleso Argentan Haquenauo Aas Vitry-le-Francois Bar-le-Duc Nancy Antony 20 O - Feeling 0 Map data @2015 GeoBasis-DE/BKG (@2009), Google Privacy maps.google.com Report a problem 50 km L Lite mode Terms

ibox - dtaniar@gmail.com × / 🏹 Google Maps

https://www.google.nl/maps/@14.7914948,88.7970029,3z?hl=en

×





Big Data is all about data that we don't have

David Taniar

Monash University, Australia

EGC'2015, Luxembourg, 28 January 2015



Big Data is all about data that we don't have

David Taniar

Monash University, Australia

EGC'2015, Luxembourg, 28 January 2015



Call For Papers of conferences, workshops, journals in the period of 6 months (July-Dec last year)

Source: DBWorld



Call For Papers of conferences, workshops, journals in the period of 6 months (July-Dec last year)



Source: DBWorld

International Symposium on Big Data Computing (BDC 2014)

In conjunction with:

7th IEEE/ACM International Conference on Utility and Cloud Computing (UCC 2014) Sponsored by: IEEE Computer Society and ACM (Association for Computing Machinery), USA

http://www.cloudbus.org/bdc2014



The 6th International Workshop on Workflow Management in Cloud and Big Data 3 - 5 Dec. 2014, Sydney, Australia

in conjunction with the International Conference on Big Data and Cloud Computing (BDCloud 2014)

2014 IEEE Internatio

Supported by IEEE TCSC Technical Area on Workflow Management in Scalable Computing Environments



International Symposium on Big Data Computing (BDC 2014)

In conjunction with:

7th IEEE/ACM International Conference on Utility and Cloud Computing (UCC 2014) Sponsored by: IEEE Computer Society and ACM (Association for Computing Machinery), USA

http://www.cloudbus.org/bdc2014

CIKM 2014 Workshop on Interactive Mining for Big Data (ImBig)

Special Session on "Warehousing and Intelligent Analysis Of Complex Network Big Data (WIBIG 2014)" @DSAA'2014

1st International Workshop on Big Data Discovery & Curation

Co-located with KDD 2014

Sunday August 24, 2014, New York, USA



IEEE Transactions on Services Computing

Special Issue on Big Data Analytics, Infrastructure, and Applications

Pervasive and Mobile Computing

Special Issue on Big Data Analytics for Smarter Health Care - Call for Papers -

CFP: ACM TKDD Special Issue: Connected Health at Big Data Era

The availability of big data and the emergence of wearable computing, network science, and computational landscape of how we decipher our lives, our social interactions, and our day-to-day activities. This well transforming healthcare from reactive and hospital-centered, to preventive, proactive, evidence-based, pe ailment recovery.

Call for Chapters: Managing Big Data in Cloud Computing Environments

2014 International Conference on Cloud Computing and Big Data

November 12-14, 2014 Wuhan, China

The Third International Conference on Big Data Analytics

December 20 to 23, 2014 | Jawaharlal Nehru University, Delhi.

SIMBig'2014 - 1st Symposium on Information Management and Big Data 8-10 September 2014 - Cusco, PERU

Paper/Demos Submission Deadline: extended till July 11, 2014

http://www.lirmm.fr/simbig2014/ simbig2014@lirmm.fr

The International Conference on Data Mining, Internet Computing, and Big Data (BigData2014) Asia Pacific University of Technology and Innevation (APU) November 17–19, 2014 (Kuala:Lumpur, Malaysia





Big Data in Motion and Big Data at Rest (BD-MR)

In conjunction with the IEEE International Conference on Big Data, October 27, 2014, Washington DC, USA



3rd ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data (BigSpatial-2014), Nov 4, 2014, Dallas, TX, USA.





Big Data in Motion and Big Data at Rest (BD-MR)

In conjunction with the IEEE International Conference on Big Data, October 27, 2014, Washington DC, USA



3rd ACM SIGSPATIAL International Workshop on Analytics for Big Geospatial Data (BigSpatial-2014), Nov 4, 2014, Dallas, TX, USA.

1st International Workshop on Privacy and Security of Big Data

(PSBD 2014)

in conjunction with

The 23rd ACM International Conference on Information and Knowledge Management (CIKM 2014)

November 7, 2014, Shanghai, China

The 3rd International Symposium on Privacy and Security in Cloud and Big Data (PriSec 2014)

3-5 December 2014, Sydney, Australia

http://www.swinflow.org/confs/prisec2014

co-located with The 4th IEEE International Conference on Big Data and Cloud Computing (BdCloud2014)

Sponsored by IEEE $\underline{\mathsf{TCSC}}$ Technical Area on Privacy and Security in Cloud and Big Data





Overview

IEEE BigData 2014

(CASK-14): Collaborative methodologies to Accelerate Scientific Knowledge discovery in big Data

International Workshop on Collaborative Big Data (C-Big 2014)

October 22, 2014 Miami, Florida, US In conjunction with the IEEE International Conference on Collaborative Computing: Networking, Applications and Worksharing (CollaborateCom 2014)

Workshop on Semantics for Big Data on the Internet of Things (SemBIoT 2014)

BigData in Bioinformatics and Healthcare Informatics

10110 01ACG TTCGA

Washington D.C., Oct 27th, 2014 in conjuction with The IEEE International Conference on BigData



BigData in Bioinformatics and Healthcare Informatics

Washington D.C., Oct 27th, 2014 in conjuction with The IEEE International Conference on BigData



→ C Attps://mail.goog	gle.com/mail/u/0/#inbox/1482d9635f16e3a0						
Google	Click here to enable desktop patifications for Gmail Learn more Hide						
Gmail -	← E C Î More → More → 2 of 136 < > ♥ →						
COMPOSE	Yamaha & Kawai Clearance - www.australianpianowarehouse.com.au - Warehouse Prices Save \$1000's Uprights Grands I Why this ad?						
Inbox (1) Sent Mail Drafts (7)	IEEE Workshop – Big Data in Computational Epidemiology 🖶 🖻 People (3)						
All Mail Spam Trash CFP-PCMember-Revi	Sandeep Gupta <sandeep@vbi.vt.edu> to dbworld •</sandeep@vbi.vt.edu>						
Coursera ICCSA Indo-DM Nihon-go Research-Collaboration Research-General	Computational epidemiology aims to understand the spread of diseases and efficient strategies to mitigate their outbreak. It studies dynamics in socio-technical systems, where disease spread co-evolves with public health interventions as well as individual behavior. It has evolved from ODE models to networked models which apply agent based modeling and simulation methodologies. Computation of such high resolution models involves processing data sets that are massive, disparate, heterogeneous, evolving (at an ever increasing rate), and potentially unstructured and of various quality.						
Research-Students Teaching Travel-others More	The workshop brings together researchers from epidemiology, data science, computational science, and health IT domains to tap the potential of emerging technologies in data intensive computations and analytical processing to advance the state of art in computational epidemiology. The central theme of how to manage, integrate, analyze, and visualize vast array of datasets has wider applications in the bio- and physical- simulation and informatics based sciences such as immunology, high energy physics, and, medical informatics.						

The workshop welcomes original research related to computational models and methodologies developed for handling big data and their application to epidemiology.





How many CFPs in Big Data are shown in the previous slides?





26

How many CFPs in Big Data are shown in the previous slides?







Big Data - a new buzzword The Best-Run Businesses Run SAP

					+1-800-872-1
	gn Out Help Sel	Big Data		Big Data Solutions	>
		#Bigdata		Customer Testimonials	
	ORE SUPPOI			Community Experts	3
		SAP.com Solut	ons Big Da	Ita	sole 🔻
Cloud Platforr	n	Make I	Big Da	ta real: Real t	ime, 🕞
	port Custo	real wo	orid, re	eal results	
	igQuery C	Transform the w Transact, analyz real-time platforr Big Data apps a	ay you do busir e, and act on m n. Create new t nd analytics for	ness with Big Data offerings from S nassive volumes of data – instantly pusiness models and revenue stre- the real world. And get the suppor	AP. – with our ams with t you need

to drive real results - with our Big Data services.

Turn Big Data into big progress >

United States

Share 💟 f in 💱 🖂

Newsletter



Big Data – a new buzzword





Big Data – a new buzzword

SAMSUNG BIG DATA IMPLEMENTATION

"SCENARIO ANALYSIS OF DECISION SUPPORT SYSTEM WITH MICROSOFT WINDOWS SERVER 2012 OS & SQL SERVER 2012 AND SAMSUNG 20nm-class DRAM & SSD ON DELL POWEREDGE R910 SERVER "

LEARN MORE>

Data growth is forcing organizations to seek better performance coupled with increased power savings. Samsung Semiconductor and Microsoft performed a Proof of Concept (using four predefined queries) to show how well Samsung Green Memory meets this challenge.

LEARN MORE>

Intel IT Center

Big Data Intelligence Begins with Intel

Learn how big data analytics delivers fre

	Services / IT Consulting Services									
(I)	Contact u	s 🗸	+ Share							
PRODUCTS	SOLUTIONS	SERVICES	CUSTOMERS	PARTNERS						
Applications	Infrastructure	Industries	IT Strategies							
Home Solutions IT Strategies Big Data										
BIG DATA STRATEGIES AND SOLUTIONS										

Capitalize on Big Data Today and Tomorrow



- Use the right infrastructure to gain value from big data.
- Discover, search and analyze all of your data.
- Start building tomorrow's big data framework now.

Big Data – a new buzzword



Australian Government

Department of Finance and Deregulation Australian Government Information Management Office

HOME BLOG POLICY, GUIDES & PROCUREMENT COLLABORATION, SERVICES &

Home » Blog » For public consultation: Big Data Strategy (draft)

For public consultation: Big Data Strategy (draft)

Category: AGCIO

Tags: AGCIO, Big Data, Consultation, ICT Strategy, NDES

Big Data

WHY GARTNER ANALYSTS

Gartner.

What information, if you had it, would change the way you run your business?

VIEW FREE RESEARCH Big Data Strategy: IT Essentials

Big data – information of extreme size, diversity and complexity – is everywhere. This disruptive phenomenon is destined to help organizations drive innovation by gaining new and faster insight into their customers. So, what are the business opportunities? And what will they cost?



DID YOU KNOW?

RESEARCH EVENTS CONSULTING ABOUT

15% of enterprises will adapt their information technology abilities for extreme data, socially mediated content and new connected devices.



- 1. Misconceptions
- 2. Motivating Examples
- 3. Challenges



"**Everyday**, 2.5 **quintillion** bytes of data are created and 90% of the data in the world today was created within the past two years".

IBM Corporation

 $10^{6} = \text{million (megabytes)}$ $10^{9} = \text{billion (gigabytes)}$ $10^{12} = \text{trillion (terabytes)}$ $10^{15} = \text{quadrillion (petabytes)}$ $10^{18} = \textbf{quintillion (exabytes)}$



"**Everyday**, 2.5 **quintillion** bytes of data are created and 90% of the data in the world today was created within the past two years".

IBM Corporation

"Worldwide information is more than **doubling every two years**, with **1.8 zettabytes** or 1.8 trillion gigabytes projected to be created and replicated this year alone".

ZDNet news

 10^{15} = quadrillion (petabytes) 10^{18} = **quintillion** (exabytes) 10^{21} = sextillion (**zettabytes**)



Data comes from everywhere:

Post to social media sites



"As of September 2012, Facebook has more than **one billion** active users"

Facebook wikipedia



"Twitter has over 500 million users in 2012, generating over **340 million tweets** and handling over 1.6 billion search queries per day".

Twitter wikipedia


Data comes from everywhere:

- Post to social media sites
- Digital pictures and videos posted online



"There has been more video uploaded to YouTube in the last 2 months than if ABC, NBC, and CBS had been airing content 24/7/365 continuously since 1948".

Gartner Research



Facebook handles **50 billion photos** from its use base.



Data comes from everywhere:

- Post to social media sites
- Digital pictures and videos posted online
- Transaction record of online purchases

"Walmart handles more than **1 million customer transactions every hour**, which is imported into database estimated to contain more than 2.5 petabytes of data".

Save money. Live better HANAGER TRENE BROWN (360) 532 - 7595 ST# 2037 OP# 00003048 TE# 18 TR# 05704 1.68 0 HALF HALF .68 0 HALF GU BRUN RICE 007874212222 0 76 GU LONG RICE 0078 GU 2 RF MLK 007874235183 3.08 AMEN-BEEF12 004178900232 1.82 URTOTAL 10.30 10.30 10 00 SHAPPING CARD TEND CASH TEND 0.50 0.20 CHANGE DUE 10.00 SHOP.CARD REDEMPTION ACCOUNT 605214515600 APPR. CODE = 037453REF #0571931 Beg Bal Tran Ant End Bal 10.00 0.00 10.00 21:20:50



Washington E-Cycles: Free Recycling For Computers, Monitors, and TV's www.ecyclewashington.org 1-800-RECYCLE

THANK YOU FOR SHOPPING WITH US 10/14/10 21:20:58

The Economist



Data comes from everywhere:

- Post to social media sites
- Digital pictures and videos posted online
- Transaction record of online purchases



"Amazon's websites averaged **94 million unique visitors** per month, compared to 22 million visitors to Target's website and 17 million visitors to Best Buy's online presence".

Statista, Germany

+ How big is "BIG"? - The Big Misconception

Data comes from everywhere:

- Post to social media sites
- Digital pictures and videos posted online
- Transaction record of online purchases
- Mobile phone use and GPS signals



"There are **6 billion** mobile-phone subscriptions worldwide, and there are between 1 and 2 billion people accessing the Internet".

The Economist

+ How big is "BIG"? - The Big Misconception

Data comes from everywhere:

- Post to social media sites
- Digital pictures and videos posted online
- Transaction record of online purchases
- Mobile phone use and GPS signals
- Various sensor and satellite data



"The Sloan Digital Sky Survey produced **13 petabytes** of data in one year".

+ But... is it **Big Data**?

Are all these Big Data?

- Post to social media sites
- Digital pictures and videos posted online
- Transaction record of online purchases
- Mobile phone use and GPS signals
- Various sensor and satellite data



Misconception #1: Big Data is Big, but it is not about big data volume.



Misconception #1: Big Data is Big, but it is not about big data volume.



If Big Data is not about Big Data Volume,

Then... Is Big Data about data store?

Big Data certainly needs big data storage.
But... Is Big Data about data storage?



Capabilities



Big Data: end to end

"Big Data" refers to a collection of tools, techniques and technologies for working with data productively, at any scale. The tools to support data collection, computation along with collaboration and sharing are all available in a couple of clicks, with Amazon Web Services.

Learn more about Elastic MapReduce »

Another Misconception

- Most data are in files:
 - Web is file centric
 - Email is file centric
 - Science is file centric
 - ...

Big data is not in databases

- DBMS delivers better performance than other systems
- So perhaps, a new database system, like **NoSQL**, is the answer

Another Misconception

NOSQL, the new DBMS?

- NoSQL = no-SQL or Not Only SQL. Non-relational data storage systems.
- NoSQL movement was probably inspired by Google's Big Table or even Amazon's S3.
- Explosion of social media sites (Facebook, Twitter) with large data needs
- NoSQL feature: shared-nothing horizontal scaling replicating and partitioning data over many servers.

Definition: "Next generation databases mostly addressing some of the points: being **non-relational**, **distributed**, **open-source**, and **horizontal scalable**. The original intention has been modern web-scale databases. The characteristics are: eventually consistent/BASE (not ACID), a huge data amount, and more."

www.nosql-database.org

Another Misconception

- Varieties of NoSQL:
 - Column-stores: each storage block contains data from only one column
 - Document stores: stores documents made up of tagged elements
 - Key-value stores: hash table of keys

Facebook's Cassandra	MongoDB	Redis
Google's Big Table	CouchDB	Riak
Yahoo's PNUTS		Tokyo Cabinet
Apache's Hbase		Scalaris
Hypertable		Memcached, Membrain, Membase
HadoopDB		

Another Misconception

- Leading users of NoSQL
 - Social networking sites (Twitter, Facebook, LinkedIn)
- Most corporate IT companies are not social networking industries.

Misconception #2: Big Data is not about Big Data Store, nor a new wave of DBMS, like NoSQL, although Big Data comes in a variety of data formats (e.g. social media formats)

- Although big companies have now been offering data storage for Big Data; and
- DBMS vendors are now selling data management systems for Big Data.

The Last Misconception

If Big Data were about big data volume, then we would need parallel processing to process Big Data.

But...

- Parallel Databases: started in mid-late 80s.
- The trend was to build specialized database machines (e.g. Gamma, Bubba).
- Various parallel machine architectures were built:
 - Shared-memory architecture
 - Shared-nothing architecture
 - Shared-something architecture
- Due to the extensive research in parallel databases, now most commercial DBMSs (Oracle) have parallelization capabilities.

The Last Misconception

- Parallelization through data partitioning
- Hence, parallel scans, yield I/O parallelism



Search U

- Parallel Databases
- Then...
 - Cluster Computing
 - Grid Computing
 - Cloud Computing

The Last Misconception

- Parallel Databases
- Then...
 - Cluster Computing
 - Grid Computing
 - Cloud Computing
- And Now...
 - MapReduce
 - Hadoop



Big Data is often associated with MapReduce/Hadoop. But, is Big Data about MapReduce/Hadoop?



My Account / Console 🔻 English 🔻

AWS Products & Solutions -

AWS Product Information -

Sign Up

Q

Developers -Support -

Big Data on AWS

Drive innovation through data, with scalable services for data collection, storage, integration, analytics and collaboration.

Get started for free »

Learn about the AWS Free Tier

"Amazon Elastic MapReduce enables us to focus on our Hadoopbased analysis without worrying about the underlying infrastructure."

- Jason Davis, Director of Search & Personalization



Read the Story »











- Parallel computing
 - Constructing high performance parallel computers using a large number of (low-end) commodity processors.
 - Commodity machines (cheap, but unreliable).
 - Commodity network.
 - Scalable (1000's of machines, 10,000's of disks)



- Parallel programming
 - What happens with parallel programming that has existed for many decades (e.g. MPI)?
 - A new parallel programming paradigm: MapReduce
- MapReduce: a simple data-parallel programming model designed for scalability and fault-tolerance.
- Pioneered by Google
 - Processes 20 Petabytes of data per day
- Popularized by open-source Hadoop project
 - Used at Yahoo!, Facebook, Amazon, ...

- Cheap nodes fail, especially if you have many of them
 - Mean time between failures for 1 node = 3 years
 - Mean time between failures for 1000 nodes = 1 day
 - **Solution**: Build fault-tolerance into system
- Commodity network = low bandwidth
 - Solution: Push computation to the data
- Programming distributed systems is hard
 - Solution: Data-parallel programming model users write "map" and "reduce" functions, system distributes work and handles faults.





- Single *master* controls job execution on multiple *slaves*
- Mappers preferentially placed on same node or same rack as their input block
 - Minimizes network usage
- Mappers save outputs to local disk before serving them to reducers
 - Allows recovery if a reducer crashes
- Allows having more reducers than nodesMappers save outputs to local disk before serving them to reducers
 - Allows recovery if a reducer crashes
 - Allows having more reducers than nodes

+ MapReduce Fault Tolerance



- 1. If a task crashes, Retry on another node
 - OK for a map because it has no dependencies
 - OK for reduce because map outputs are on disk

2. If a node crashes:

- Re-launch its current tasks on other nodes
- Re-run any maps the node previously ran (Necessary because their output files were lost along with the crashed node)

3. If a **task is going slowly** (straggler):

- Launch second copy of task on another node
- Take the output of whichever copy finishes first, and kill the other

Surprisingly important in large clusters

- Stragglers occur frequently due to failing hardware, software bugs, misconfiguration, etc
- Single straggler may noticeably slow down a job

The Last Misconception

- Parallel Databases
- Then...
 - Cluster Computing
 - Grid Computing
 - Cloud Computing
- And Now...
 - MapReduce
 - Hadoop



Misconception #3: Parallel processing certainly helps. But Big Data processing is not about MapReduce/Hadoop or any parallel technologies.

The Big Data Misconceptions

- 1. It is not about Big Data Volume, although Big Data is big;
- 2. It is not about new Data Stores, although Big Data needs big data stores as data comes in a rich data formats; and
- 3. It is not about new trends in Parallel Processing, although parallel processing is needed to process Big Data.

So... what is then Big Data?



- 1. Misconceptions
- 2. Motivating Examples
- 3. Challenges



- 1. A Clothing Store
- 2. An Airline
- 3. Health

- Sales of a clothing shop decreasing
- They analyzed their customers data and purchases to target certain customers
- Even after launching promotions, the sales was not improved



- Sales of a **clothing shop** decreasing
- They analyzed their customers data and purchases to target certain customers
- Even after launching promotions, the sales was not improved
- They also used geofencing
- And still unsuccessful
- Puzzled ???!!!!###%%%%





- Their own data did not and could not solve their sales problems.
- So, they went out and searched for outside data in this case social media tweets.
- They found from tweets from teenagers that the clothes currently being displayed for sales have the matched clothes from the previous seasons which are no longer available in the shop.



Expand

A Big Data Problem...

- Their own data did not and could not solve their sales problems.
- Because they did not have the necessary data to solve their sales problems.



An airline case study...

- Top level management of an airline company would like to know the service provided at the operational levels – to find out if their customers are happy or not
- Top level management is not able to see the operational levels. They only see the high level management.
- They can analyze their own data, and will not find anything.
- A Big Data Problem their own data will not provide the solution to their problems.



An airline case study...

- So, they conducted a "sentiment analysis" from social media.
- To find out how the feeling from people toward the airline
- No questionnaires so it is more objective
- Based on the results they get from sentiment analysis, they are able to rank the airlines based on the satisfactory from passengers through social media.
- Input (twitter feeds), output (ranking of airlines)



An airline case study...

- Sentiment analysis uses some quantitative values
- Long queue -1, nice food +1, delay -1, etc.
- Challenges:
 - It is supposed a negative comment, but the comment does not reflect the negativity explicitly
 - Use of punctuation (e.g. "good" ③)



Big Data problem has been here long ago...

since 1800s !!!


since 1800s !!!

- There was **Cholera epidemic** in London (1831-1832).
- A lot of patients data has been collected. Experts said that cholera, like other diseases, was transmitted through inhalation of contaminated vapors.
- Cholera struck England again in 1854, even with more deaths; that proved the experts wrong.



since 1800s !!!

- There was **Cholera epidemic** in London (1831-1832).
- A lot of patients data has been collected. Experts said that cholera, like other diseases, was transmitted through inhalation of contaminated vapors.
- Cholera struck England again in 1854, even with more deaths; that proved the experts wrong.

Was it because they didn't have the technology? NO

Was it because they didn't have The data? YES

Motivating Example – 3

Big Data problem has been here long ago.

since 1800s !!!

- There was CLA vra Condemic in London (1981) 832).
- A lot of patients data by Oblighteeted. Experts said that cholera, like other disea points onsmitted through inhalation of contaminated valors.
- Cholera struck England again in 1854, even with more deaths.

Was it because they didn't have the technology? NO

Was it because they didn't have The data? YES



since 1800s !!!

- John Snow was an apprenticed to a surgeon living at Newcastle-on-Tyne during the first cholera outbreak (1831-32). He then became a medical student, and later became a surgeon (MD, University of London, 1844).
- Dr Snow published a paper on "mode of communication of cholera", suggesting that the "Cholera Poison" reproduced in the human body, and was spread through the contamination of food or water.
- Experts disagreed. Dr Snow could not prove them wrong, because the patient data they had did not indicate anything. They simply did not have the data. Patient data was not enough!!!



since 1800s !!!

- Then during the second Cholera outbreak in 1854, Dr Snow persisted, by collecting non-health related data of patients.
- He surveyed the location of the deaths and their water pumps.





http://www.youtube.com/watch?v=TUTmg1iVX8E

+

Big data is all about data that we don't have...



- 1. Misconceptions
- 2. Motivating Examples
- 3. Challenges

Big Challenges 1. Understanding Business Needs

One very important question:

Where is the data???

Also...

What data I need?

Who owns the data?

How to get the data?



Ella @MrsEllaBella Mar 7 Last night I discovered that @Cinnabon can be delivered throughout the UK. I dropped so many hints to hubby. Bet he didn't get it at all. Expand



Lady Hannah @mardycow Mar 7 @MrsEllaBella @cinnabon @garethemery has promised to get me one of these when I visit him in LA! Expand



Ella @MrsEllaBella Mar 7 @mardycow @cinnabon @garethemery they are amazing. I always get one at Chicago airport when transferring to/from Austin - love them!!

12:34 p.m. - Mar 7, 2013 · Details



View photo

Cinnabon @Cinnabon @MrsEllaBella @mardycow We look forward to seeing you both! Until

then, something to remember us by :) pic.twitter.com/DXysPJjzF9

Mar 7

Big Challenges 1. Understanding Business Needs

- Data is coming at a rate faster than what we can absorb.
- Data is coming from everywhere:
 - Social networks
 - Mobile phones
 - Financial data
 - Sensors
 - · · · ·
 - Ambient data: data is everywhere, and yet we didn't see it.
 - BUT... those who can react quickly wins.
 - The importance is in the speed of the feedback loop, taking data from input to decision.





Data Integration, or Data Linkage

- M Brodie (AINA2010 keynote): "Data integration accounts for ~40% software project costs".
- Data Linkage for:
 - Health records
 - Spatial data linkage
 - **...**

Big Challenges 2. Data Integration



Data Cleaning and Data Preparation

- Raw data, not ready for any kind of processing and analytics.
- Not many work in this area.
- Perhaps not scientific enough, or considered dirty work,...
- ...although data cleaning and data preparation account for 80% of the entire data processing and analytics. It is a serious business.
- There are actually many available techniques: Parsing, data transformation, duplicate elimination, statistical methods, etc.

Big Challenges 3. Data Storage and Management

Data from social media, data streams... 1.



Ella @MrsEllaBella Last night I discovered that @Cinnabon can be delivered throughout the UK. I dropped so many hints to hubby. Bet he didn't get it at all. Expand



Lady Hannah @mardycow Mar 7 @MrsEllaBella @cinnabon @garethemery has promised to get me one of these when I visit him in LA! Expand



Ella @MrsEllaBella Mar 7 @mardycow @cinnabon @garethemery they are amazing. I always get one at Chicago airport when transferring to/from Austin - love them!!

12:34 p.m. - Mar 7, 2013 · Details



Cinnabon @Cinnabon

@MrsEllaBella @mardycow We look forward to seeing you both! Until then, something to remember us by :) pic.twitter.com/DXysPJjzF9 View photo



Mar 7

Big Challenges 3. Data Storage and Management

- 1. Data from social media, data streams...
- 2. DBMS, cloud storage, data storage technologies

Memory is the new disk

- Cost of main memory is plunging, so main memory solution becomes very practical.
- If the data not quite fits into main memory, you might look at flash as another technology.
- SAP's Main Memory Database push application logic close to the database to leverage the sophisticated optimization and parallelization techniques.

Big Challenges 4. Analytics

Google Cloud Platform

Home	Pro	ducts	Solutions	Pricing	Support	Cust	omers	Partners	Resources	
App Engi	ne	Comput	te Engine	Cloud Storag	ge BigQue	ery (Cloud SQL	. More Pr	oducts	

Google BigQuery

Analyze Big Data in the cloud using SQL and get real-time business insights in seconds using Google BigQuery. Use a fully-managed data analysis service with no servers to install or maintain.



New! Try the BigQuery tour to see it for yourself. Start the BigQuery tour >>



Reliable & Secure

Complete peace of mind as your data is automatically replicated across multiple sites and secured using access control lists.



Scale infinitely

You can store up to hundreds of terabytes, paying only for what you use.



Blazing fast

Run ad hoc SQL queries on multi-terabyte datasets in seconds.

AP	The Best-Run E	Businesses Ru	n SAP	Inited States	Newsletter +1-800-872-1
Big Data			lig Data Solutions	>	
Bigdata		c	ustomer Testimo	>	
		c	community Expert	>	
AP.con	n Solutions	Big Data			

Make Big Data real: Real time, real world, real results

Transform the way you do business with Big Data offerings from SAP. Transact, analyze, and act on massive volumes of data – instantly – with our real-time platform. Create new business models and revenue streams with Big Data apps and analytics for the real world. And get the support you need to drive real results – with our Big Data services.

Turn Big Data into big progress >





Need enterprise level support? Contact sales

Get an inside look at BigQuery. Download the BigQuery technical white paper.

Big Challenges 4. Analytics

The Best-Run Business	es Run SAP United St	ites Newsletter +1-800- <mark>872-</mark> 1
Big Data	Big Data Solutions	>
	Customer Testimonials	>
	Community Experts	x
SAP.com > Solutions > Big Date	ta	

Google Cloud Platform

Turn Big Data into big progress >

Share 💟 🖪 in 💱 🖂

Is data integration a separate process?

Will the integrated data remain decentralized, or by Data of the second real-time platform. Create new business models and revenue streams with become centralized? Big Data apps and analytics for the real world. And get the support you need to drive real results - with our Big Data services.

Will it be real-time analysis?

Analyze Big Data if the cloud using SQL and get real-time business insights in seconds using Google BigQuery. Use a fully-managed data analysis service wit

servels the imported data preserve privacy?



App Engine Compute

Big Data, Big Challenges

- 1. Understanding Business Needs
- 2. Data Integration
- 3. Data Storage and Management
- 4. Analytics

Big Data is all about data that we don't have.

Other than this, we just do our research normally – and it is not Big Data.





Big Data is all about data that we don't have.

Other than this, we just do our research normally – and it is not Big Data.