

# Lecture 1: Introduction to Data Analytics

Benjamin Simeon Harvey (Ben)  
Adjunct Faculty, EMSE Department  
The George Washington University  
E-mail: [bsharve@nsa.gov](mailto:bsharve@nsa.gov)  
Web: TBD

# Agenda

- I. Introductions
- II. Overview of Syllabus and Course Expectations
- III. Course Core Topics and Lectures
  - I. Big Data Foundations
  - II. Big Data and Technology Concepts
  - III. Big Data Analysis and Science
  - IV. Advanced Big Data Analysis and Science
- IV. Motivational Examples
- V. In-class discussion:
  - I. What is Data Science (EDA and Data Science Process)?
  - II. What is Big Data?
- VI. Big Data Foundations
- VII. Summary & Case Study

# I. Introductions

# Introduction: "Cool Things"

- I work for the National Security Agency as a Cryptologic Computer Scientist
  - EOD data was 10/09 and started in the Computer Science Development Program (CDP)
- Intelligence Community National Counterintelligence and Security Professional Award
  - Data Scientist for the Corporate Disclosures Action Group
- Office of the Director National Intelligence (ODNI) Award for Human Capital
  - NSA Day of Cyber
- I played Division I football and basketball at Mississippi Valley State University
  - The alma mater of “Jerry Rice” the greatest wide receiver of all time

# Introduction: My background



- Undergraduate B.S Degree in Computer Science
  - Mississippi Valley State University, Itta-Bena, MS
- Post-baccalaureate in Bioinformatics and Integrative Genomics (BIG) and a fellow at NIH Department of Clinical Research Informatics (DCRI)
  - MIT-Harvard HST, Cambridge, MA
  - NIH Clinical Center, Bethesda, MD
- M.Sc. In Computer Science, and Ph.D. in Computer Science with a focus in Bioinformatics
  - Bowie State University, Bowie, MD
- Research Assistant at Bowie State University's Biomedical Computing Lab
  - Bowie State University, Bowie, MD
- Mentors include computer scientists, enterprise architects, cyber security specialists, system engineers, & biostatisticians (Vince Carey, Ph.D. Creator of R/Bioconductor)

# Introduction (your turn!)



- Name
- Academic background
- Work experience
- What are your goals for this class?
- How do your experiences shape these goals?
- What does "Big Data" mean to you? (why are you here?)
- Something “Cool” about yourself

## II. Syllabus and Course Expectations

# Course Outcomes



- Core Topic 1: Fundamentals of Big Data
  - **Explain Big Data from a business and technology perspective**, along with an overview of common benefits, challenges, and adoption issues.
- Core Topic 2: Big Data Analysis and Technology Concepts
  - **Apply contemporary analysis practices, technologies and tools within Big Data environments** through programming and at a conceptual level
  - **Know the common analysis functions and features offered by Big Data software solutions**, as well as a high-level understanding of the **back-end components** that enable these functions.

# Course Outcomes

- Core Topic 3: Big Data Analysis and Science
  - **Apply the learned topic areas and analysis techniques to Big Data with an emphasis on how analysis and analytics need to be carried out individually and collectively in support of the distinct characteristics, assess requirements and challenges associated with Big Data datasets.**
- Core Topic 3: Advance Big Data Analysis and Science
  - **Knowledge of probability & statistics, modeling, and analysis techniques for data patterns, clusters, classification, and text analytics, as well as the identification of outliers and errors that affect the significance and accuracy of predictions made on Big Data datasets.**

# Assignments - Portfolio Project (50%)



- In this class you will create a (1) Data Science portfolio and (2) “ToolKit”
  - The goal is to consolidate learning materials, assignments, and Data Science tools within a portfolio for displaying your skillset to colleagues and potential employers and future tool use.
- Assignment 1 – Creating a Portfolio: Intro to GitHub and R/Python and EDA
- Assignment 2 - Statistical Inference
- Assignment 3 - Machine Learning
- Assignment 4 – Use the Data Science Process (DSP) to analyze your selected Big Data set and then visualize, interpret, and communicate your results in your Portfolio. Essentially, document your research project in your portfolio
  - Students will find a dataset based upon their research interests, things that inspire them, or something we discussed in class.
  - Detail your steps in developing your solution, including how you collected the data, alternative solutions you tried, describing statistical methods you used, and the insights you obtained in your portfolio.

# Student Final Project (50%)

- Final projects are to be done in either teams or as an individual
- Each team will present a talk (e.g., PowerPoint or other technique)
- Each team will produce a written document
  - 2-5 pages long, formatted according to IEEE style
    - Points will be withheld if the document is not formatted appropriately
- Peer evaluations will be factored in the determination of the term project grade.
- Final projects should use at least one of the techniques covered in class and from the assignments
- Final projects will be placed in your final portfolio (Assignment 4)
- This is a significant part of your grade. Choose your problem early. Can also be done as a group.

# Student Final Project Format (50%)

## Research Proposal

- Overview, and Motivation: Provide an overview of the project goals and the motivation for it. Consider that this will be read by people who did not see your project proposal.
- Related Work: Anything that inspired you, such as a paper, a web site, or something we discussed in class.
- Initial Questions: What questions are you trying to answer? How did these questions evolve over the course of the project? What new questions did you consider in the course of your analysis?
- Data: Source, scraping method, cleanup, storage, etc.

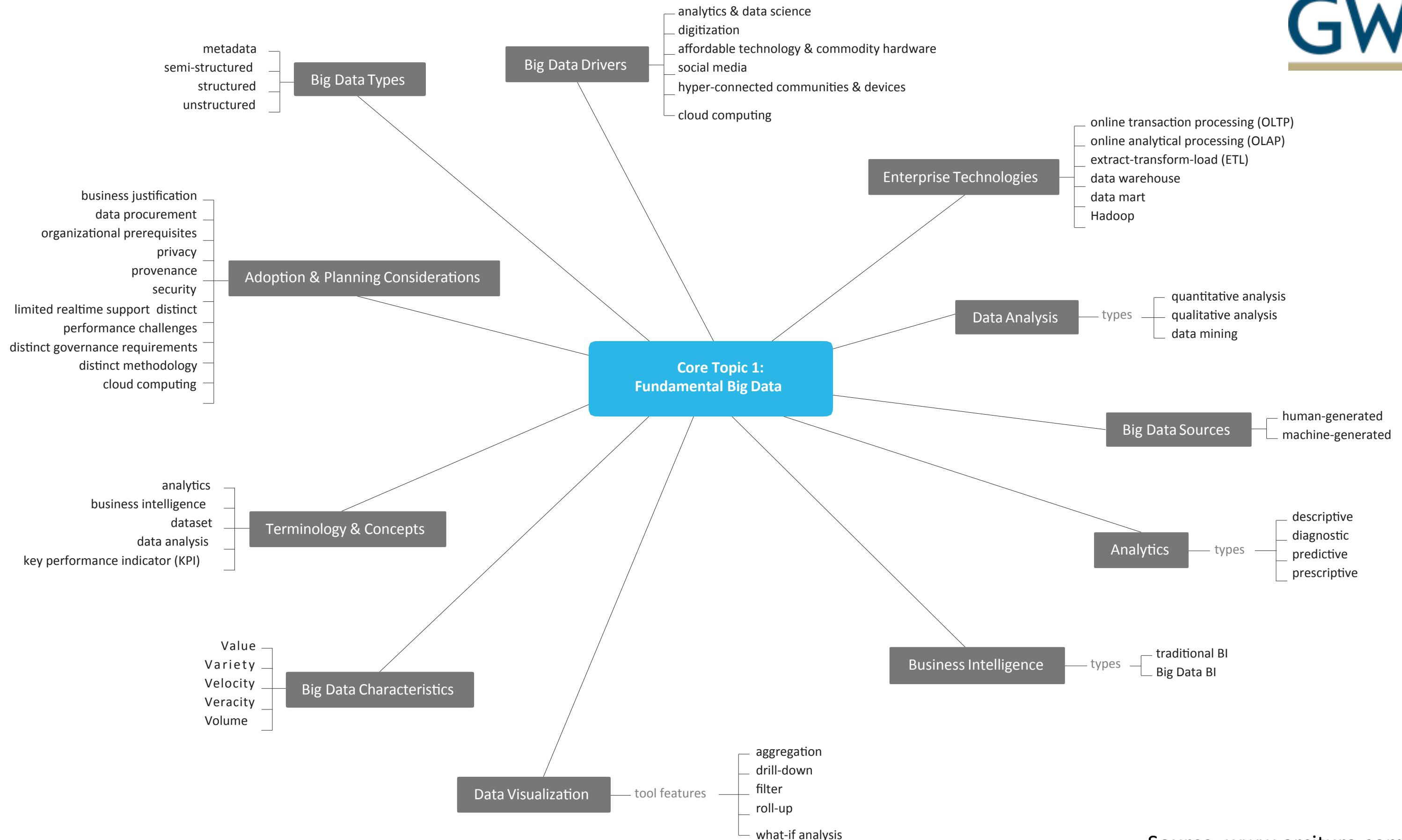
Final Paper (also includes Abstract, Overview, Motivation, Related Work, Initial Questions and Data from the Research Proposal):

- Exploratory Data Analysis: What visualizations did you use to look at your data in different ways? What are the different statistical methods you considered? Justify the decisions you made, and show any major changes to your ideas. How did you reach these conclusions?
- Final Analysis: What did you learn about the data? How did you answer the questions? How can you justify your answers?
- Presentation: Present your final results in a compelling and engaging way using text, visualizations, images, and videos on your project web site.

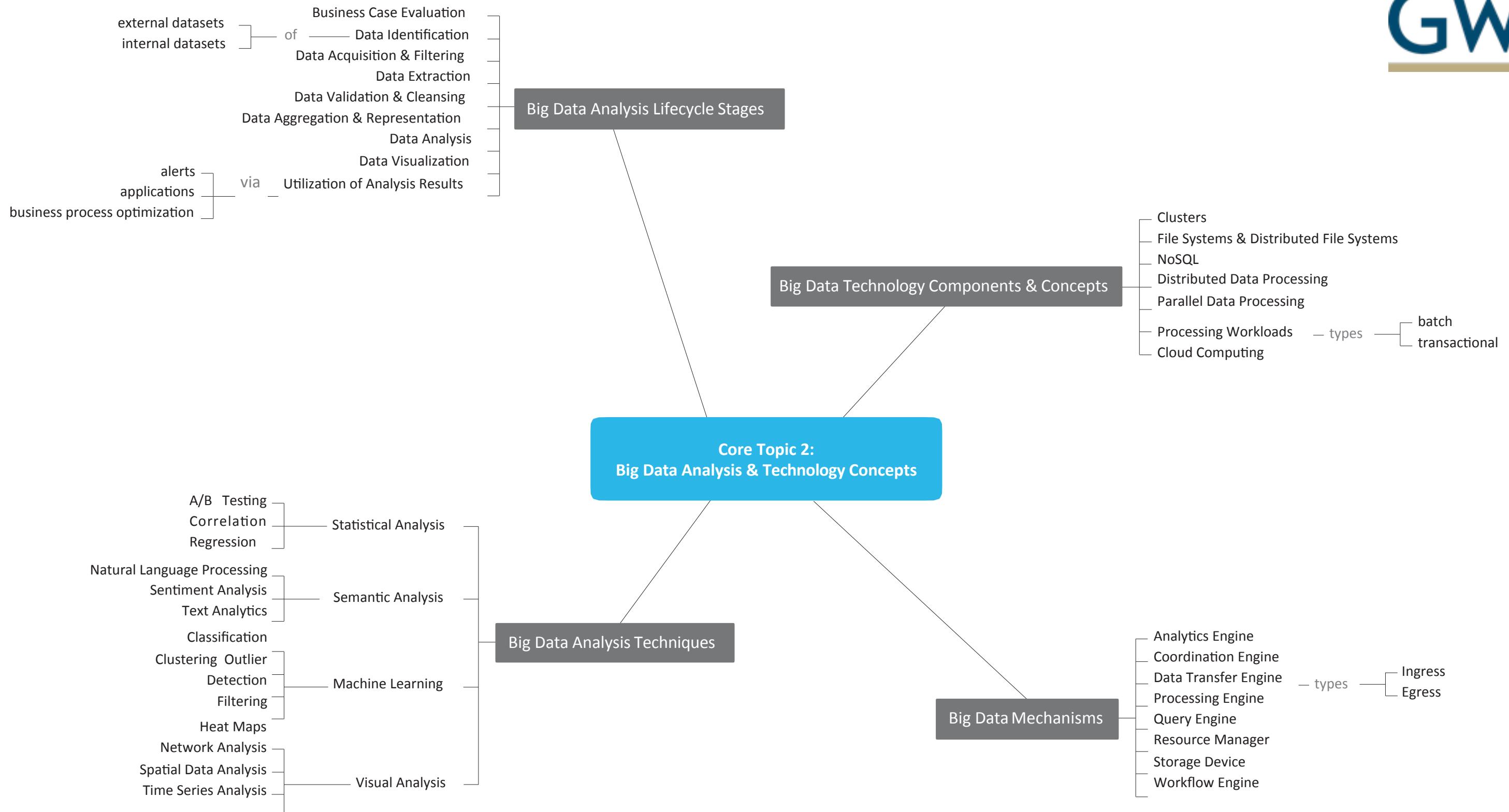
# Academic Integrity

- Plagiarism is a direct violation of the GW Code of Academic Integrity which you and I have both agreed to.
- Course policy:
  - 1<sup>st</sup> violation of academic integrity will result in a warning consisting of a loss of points equal to the value of the plagiarized assignment.
    - A plagiarized homework assignment will be given a score of -5% of your total grade (equivalent to getting zero on two assignments)
    - A plagiarized exam is equivalent to a zero on two exams.
    - A plagiarized final project is equivalent to -40% of your grade and is a practical failure of the course.
  - 2<sup>nd</sup> violation of academic integrity will result in a complaint filed in the Office of Academic Integrity with a recommended MINIMUM sanction of failure of the course.
  - No exceptions
  - If you are unsure if your work is plagiarized, ask me before you turn it in.
- Plagiarism will not be tolerated.
  - If you feel that you have been incorrectly accused of plagiarism, you may appeal to the GW Committee on Academic Integrity and a hearing will be held.

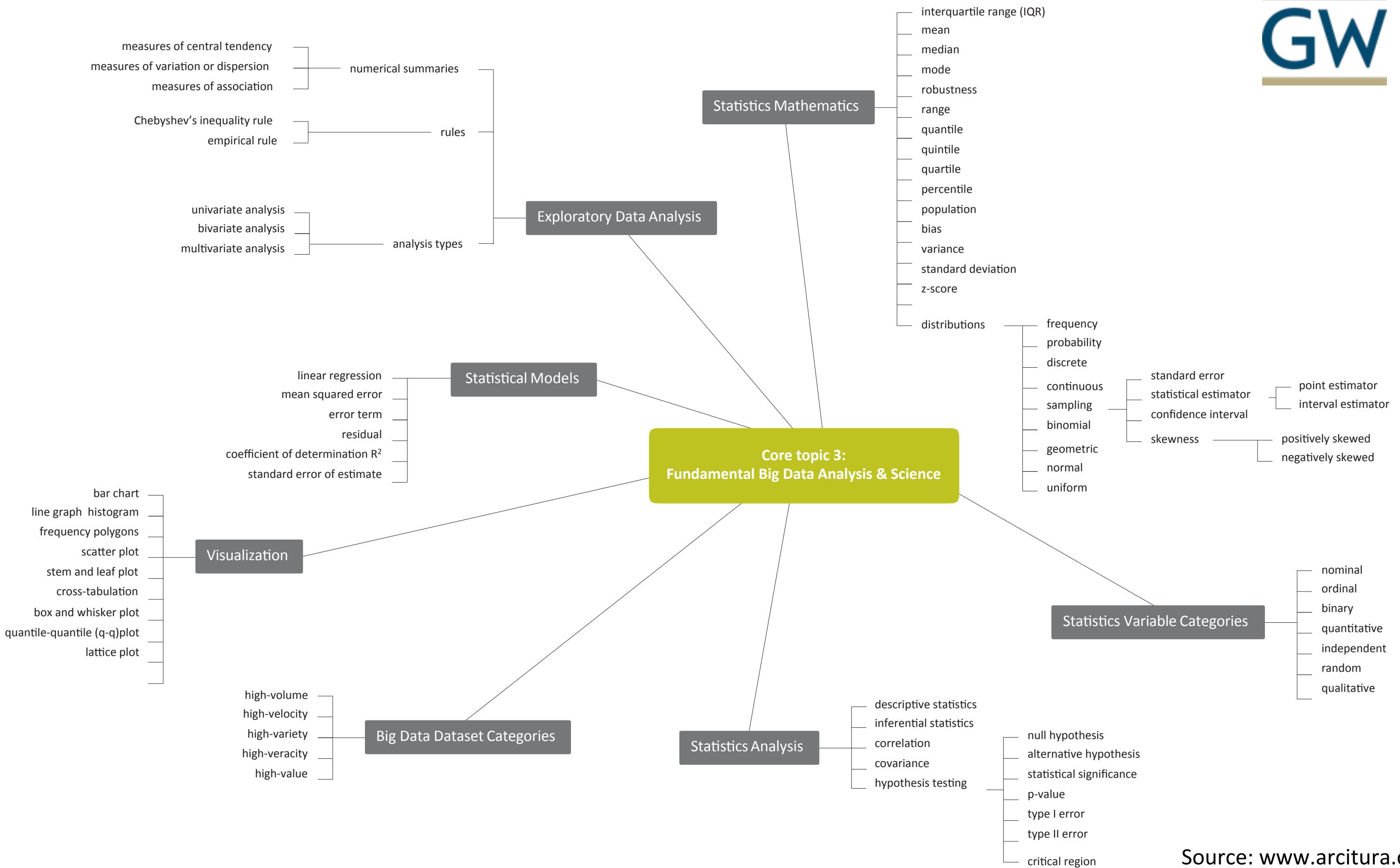
## II. Course Material and Lectures



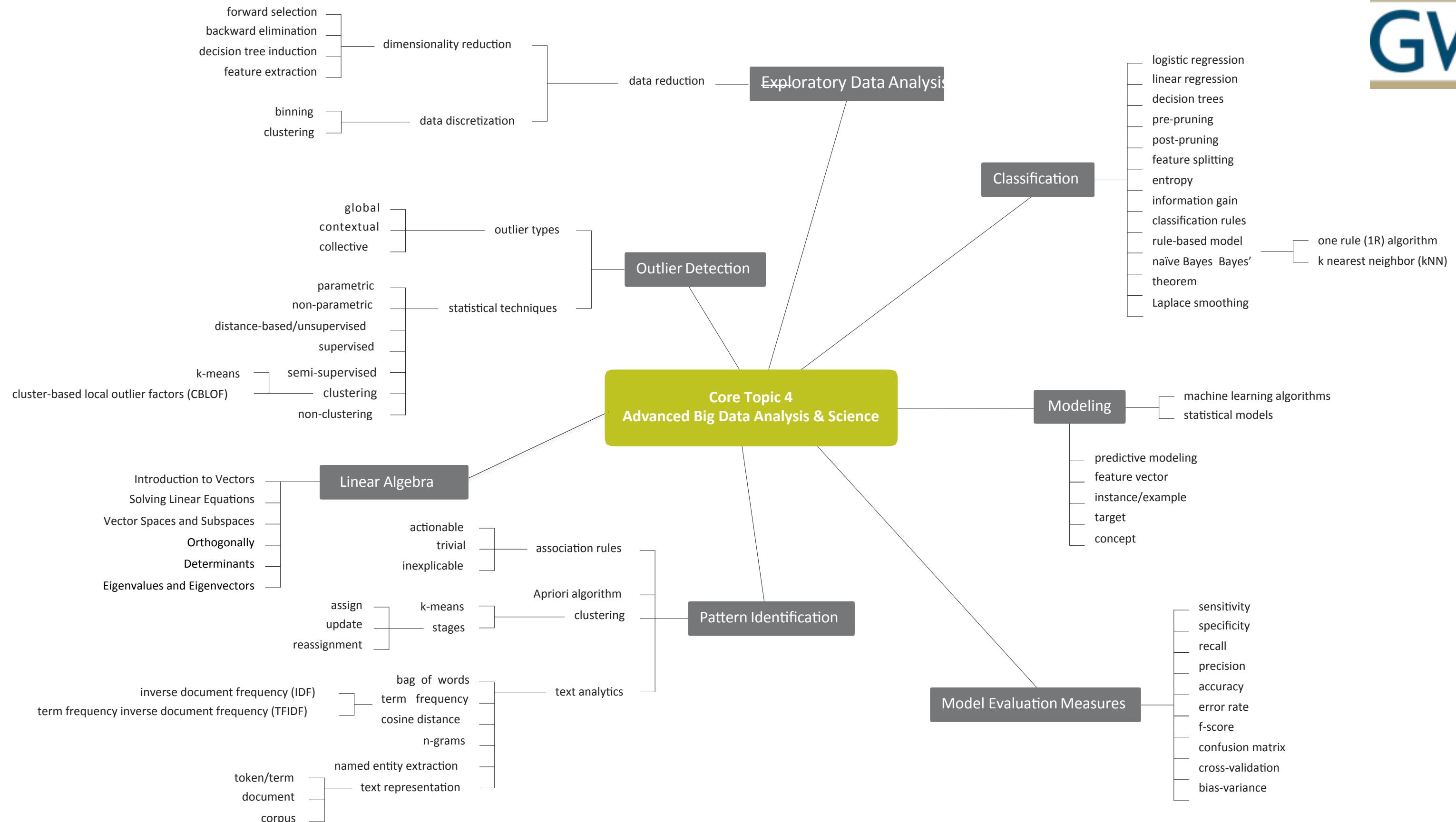
Source: [www.arcitura.com](http://www.arcitura.com)



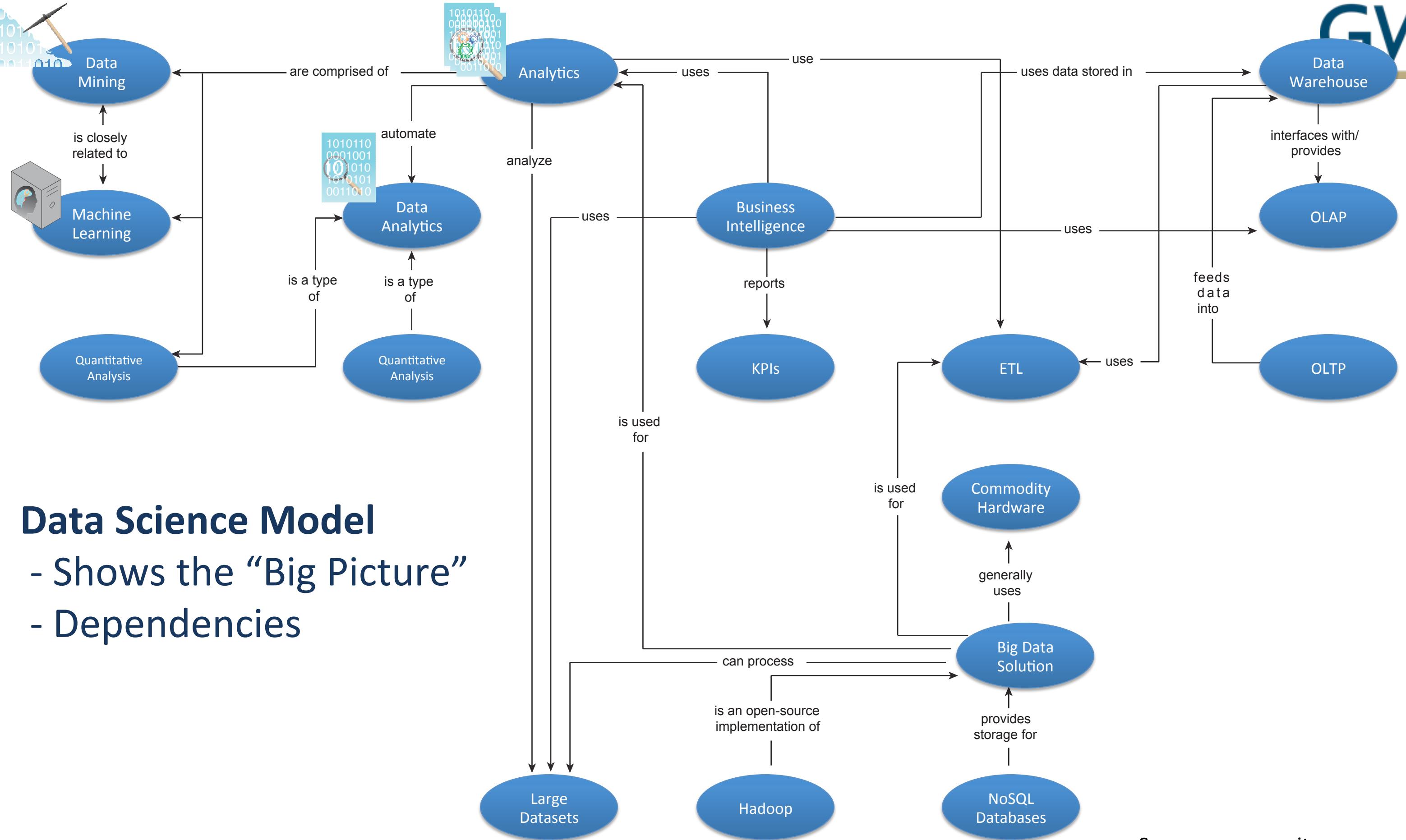
Source: [www.arcitura.com](http://www.arcitura.com)



Source: [www.arcitura.com](http://www.arcitura.com)



Source: www.arcitura.com



## Data Science Model

- Shows the “Big Picture”
- Dependencies

Source: [www.arcitura.com](http://www.arcitura.com)

# Data Analysis Has Been Around for a While

1935: "The Design of Experiments"

R.A. Fisher



1939: "Quality Control"

W.E.  
Demming

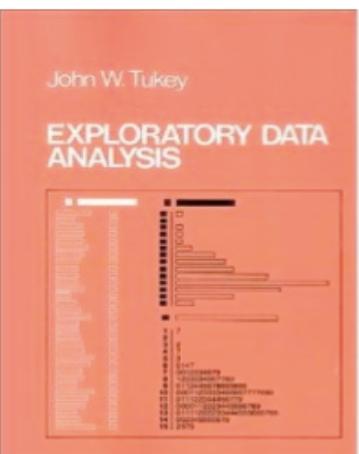


1958: "A Business Intelligence System"

Peter Luhn



1977: "Exploratory Data Analysis"



Howard  
Dresner

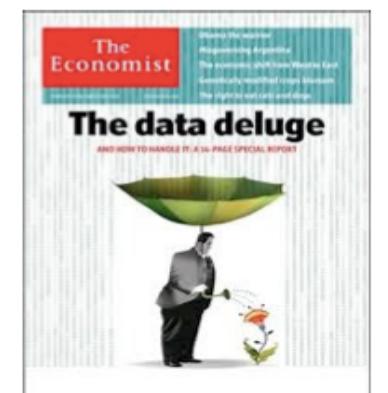


1989: "Business Intelligence"

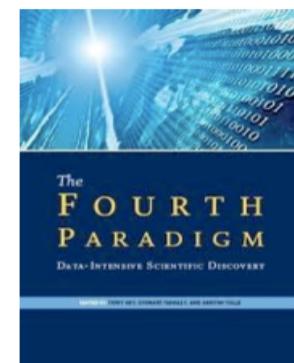
1997: "Machine Learning"



2010: "The Data Deluge"

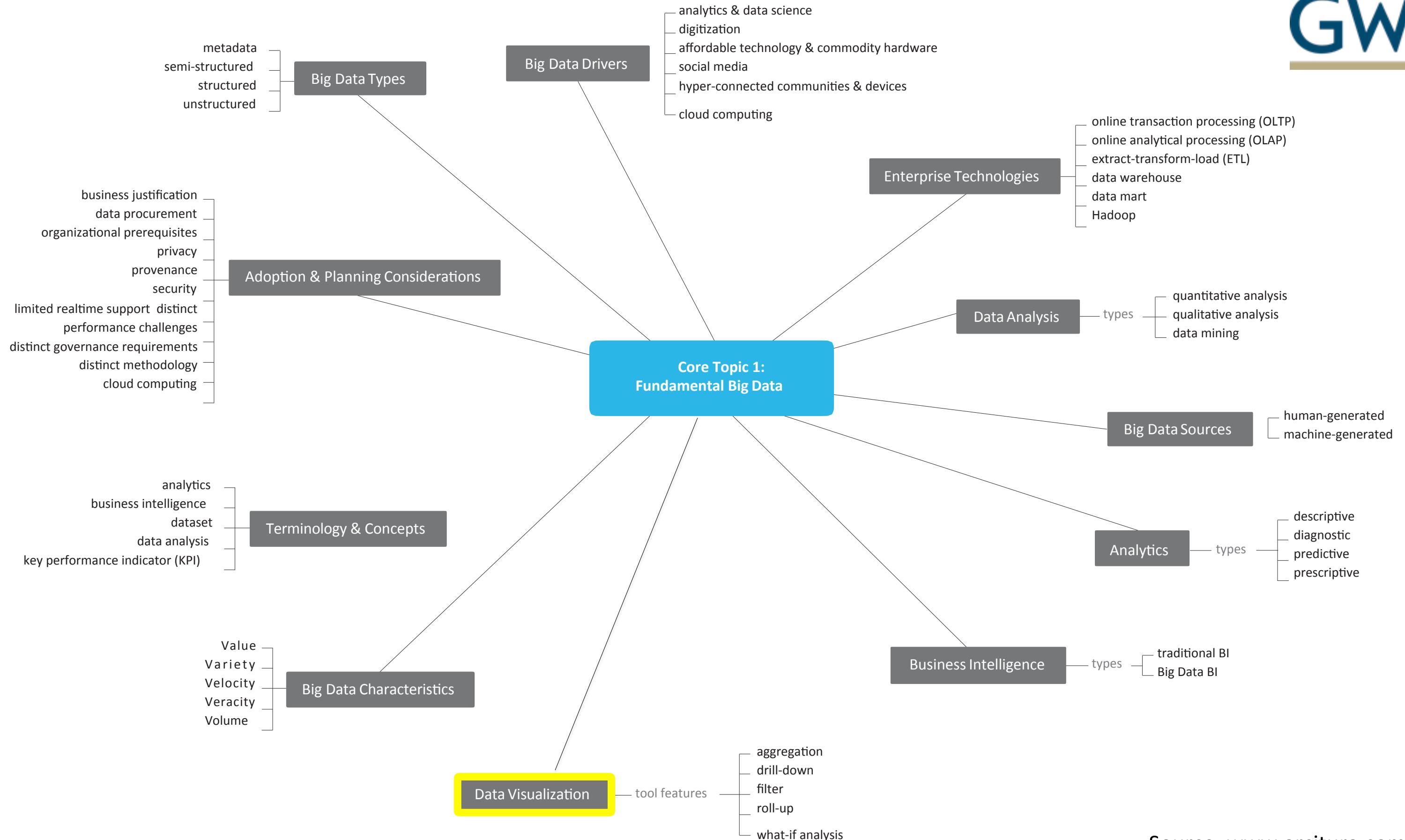


1996: Google



Abridged Version of Jeff Hammerbacher's timeline for CS 194, 2012

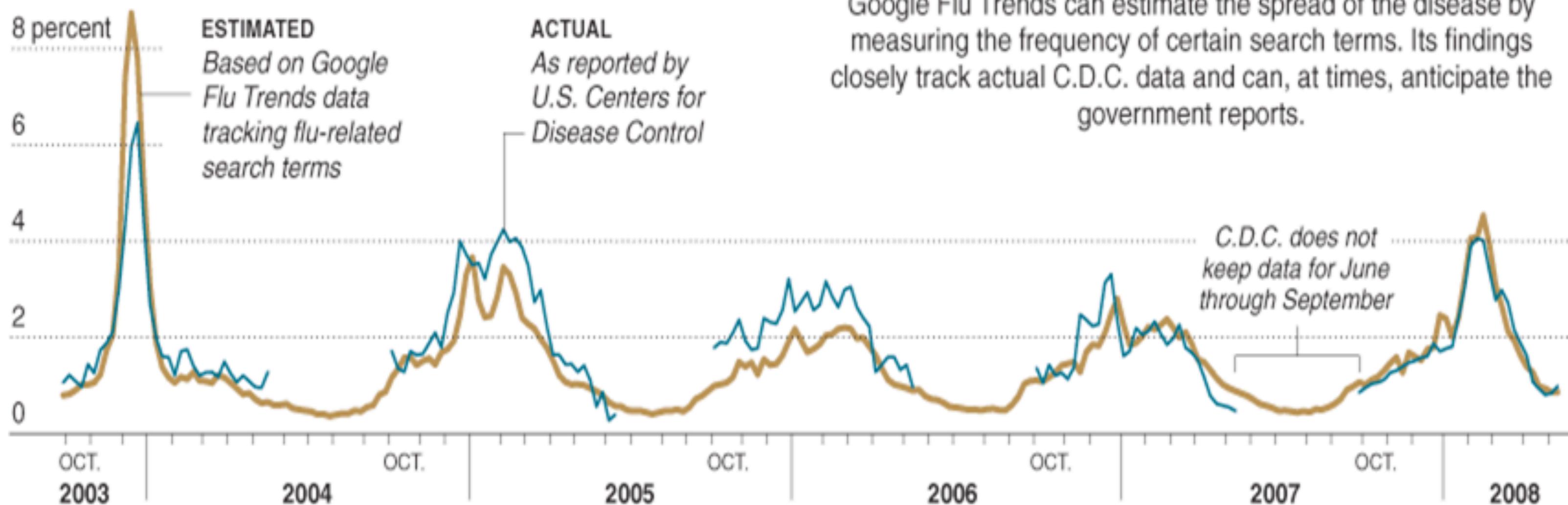
## III. Motivational Examples



Source: [www.arcitura.com](http://www.arcitura.com)

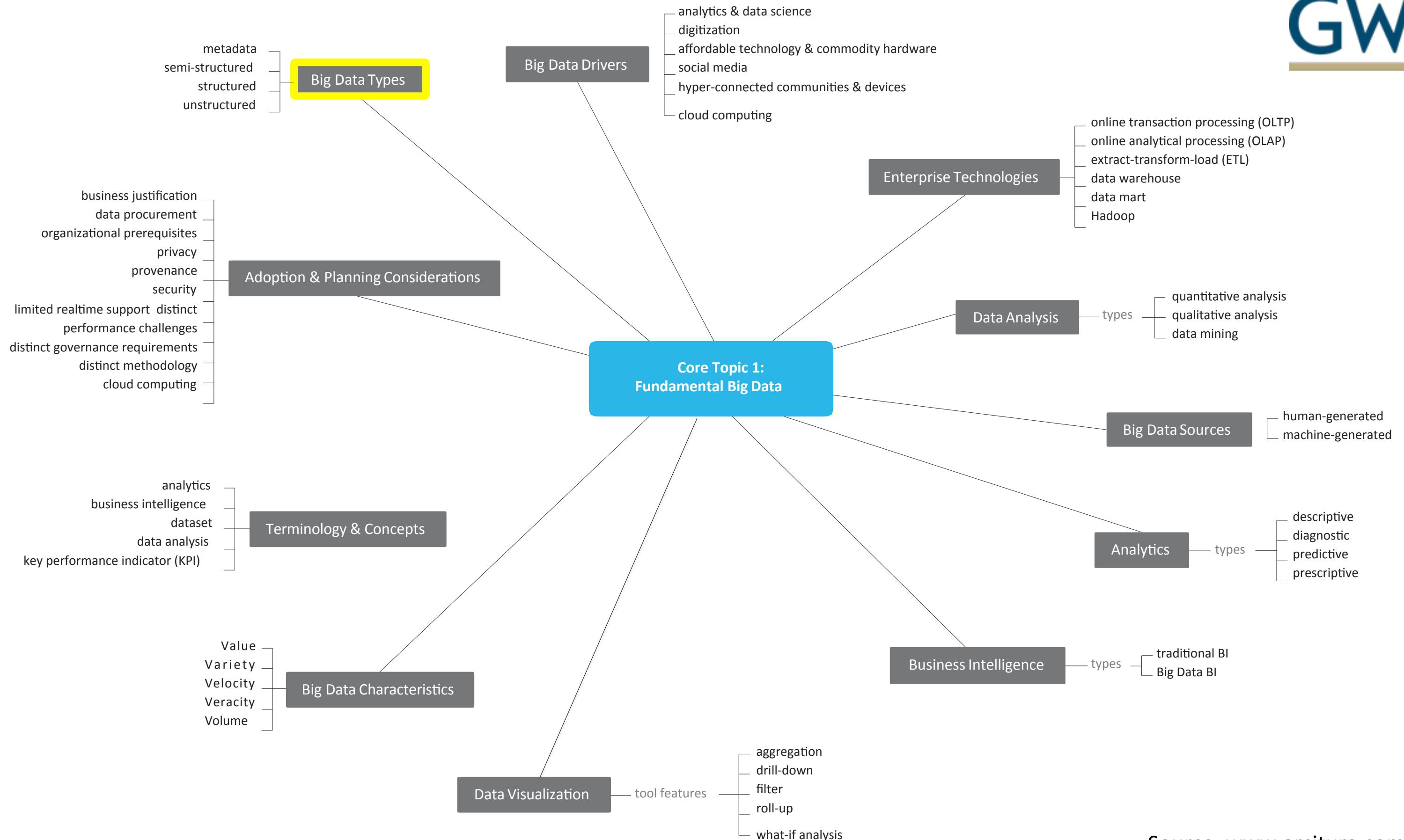
# Using Google to Monitor the Flu

PERCENT OF HEALTH VISITS FOR FLU-LIKE SYMPTOMS *Mid-Atlantic region*



Sources: Google; Centers for Disease Control

THE NEW YORK TIMES



Source: [www.arcitura.com](http://www.arcitura.com)

# Large-Scale Cancer Genomics

(10MB per Sample – 10's of GBs per dataset)

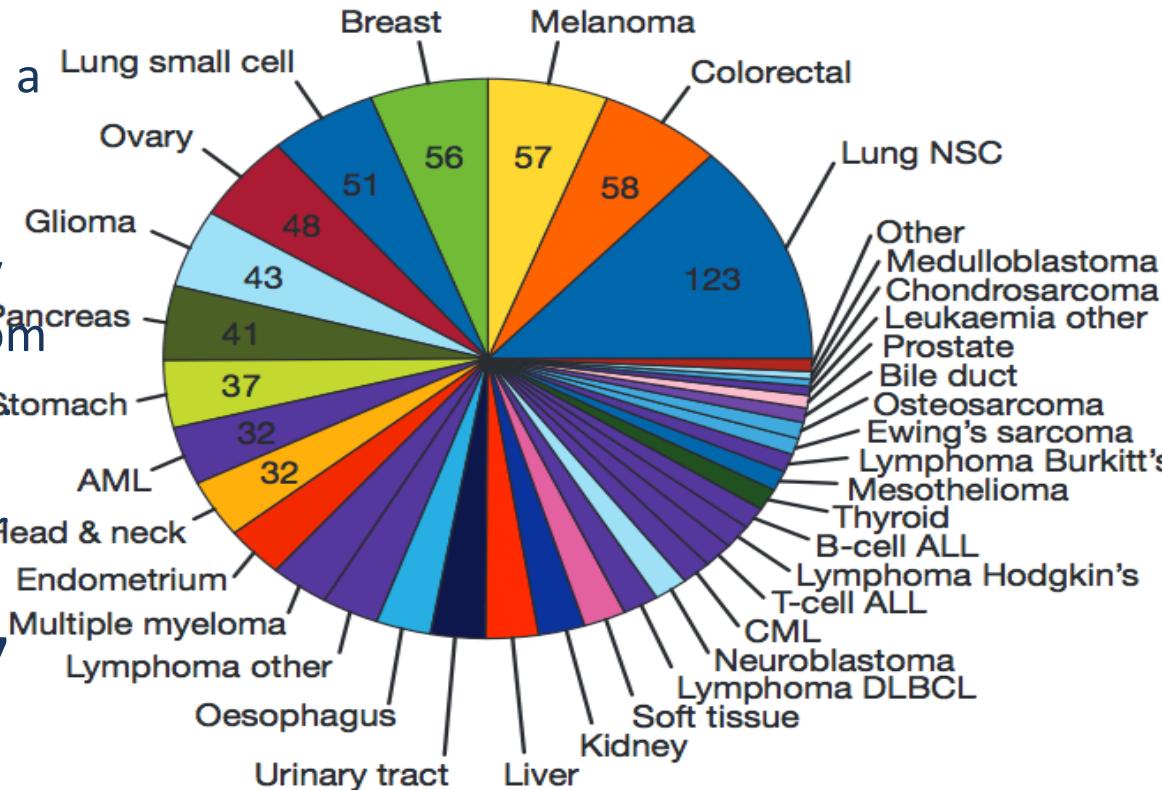
**GW**

**Global Cancer Map (GCM) Data:** 218 tumor samples, spanning 14 common tumor types, and 90 normal tissue samples to oligonucleotide microarray gene expression analysis:

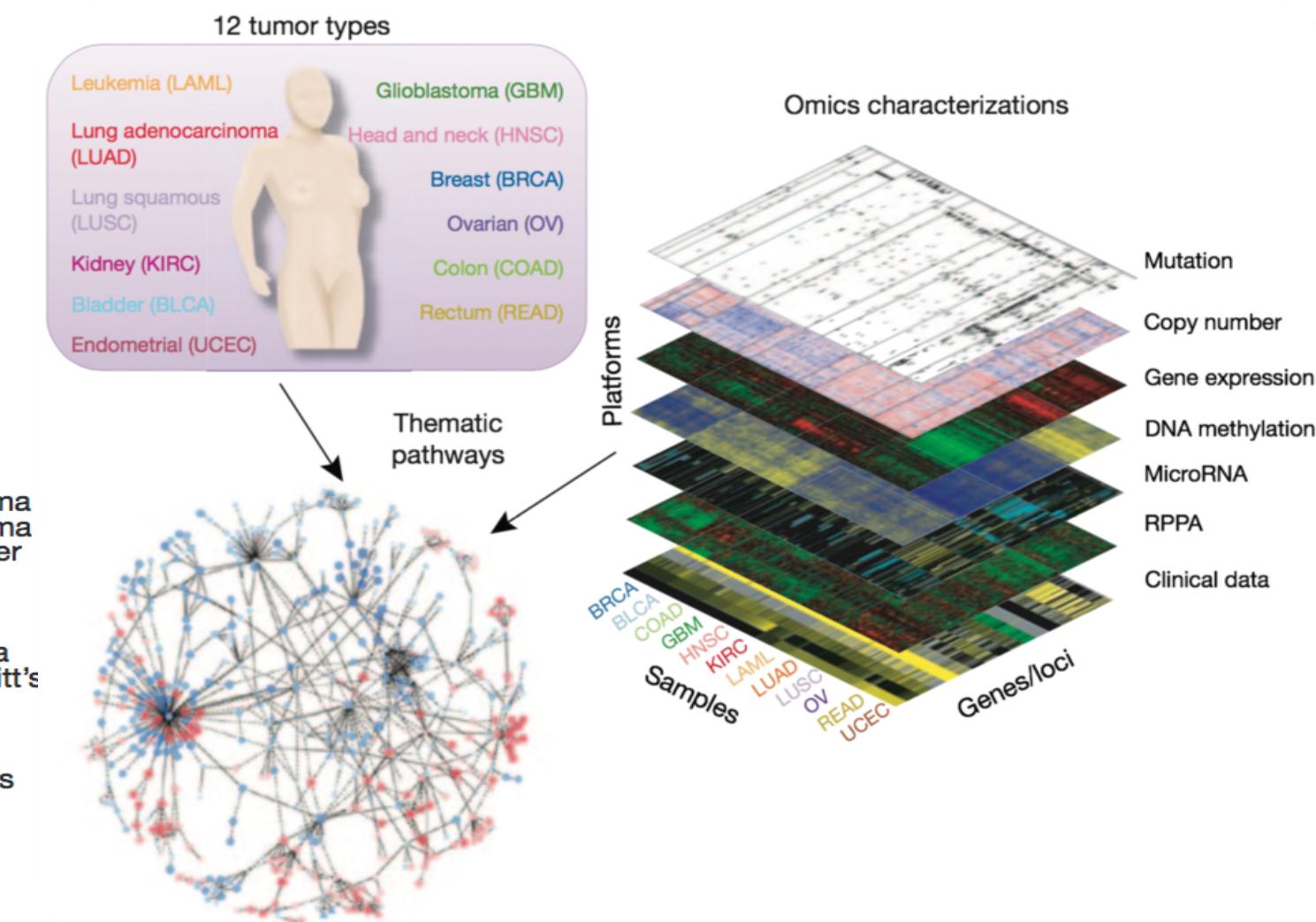
**16,063 genes**

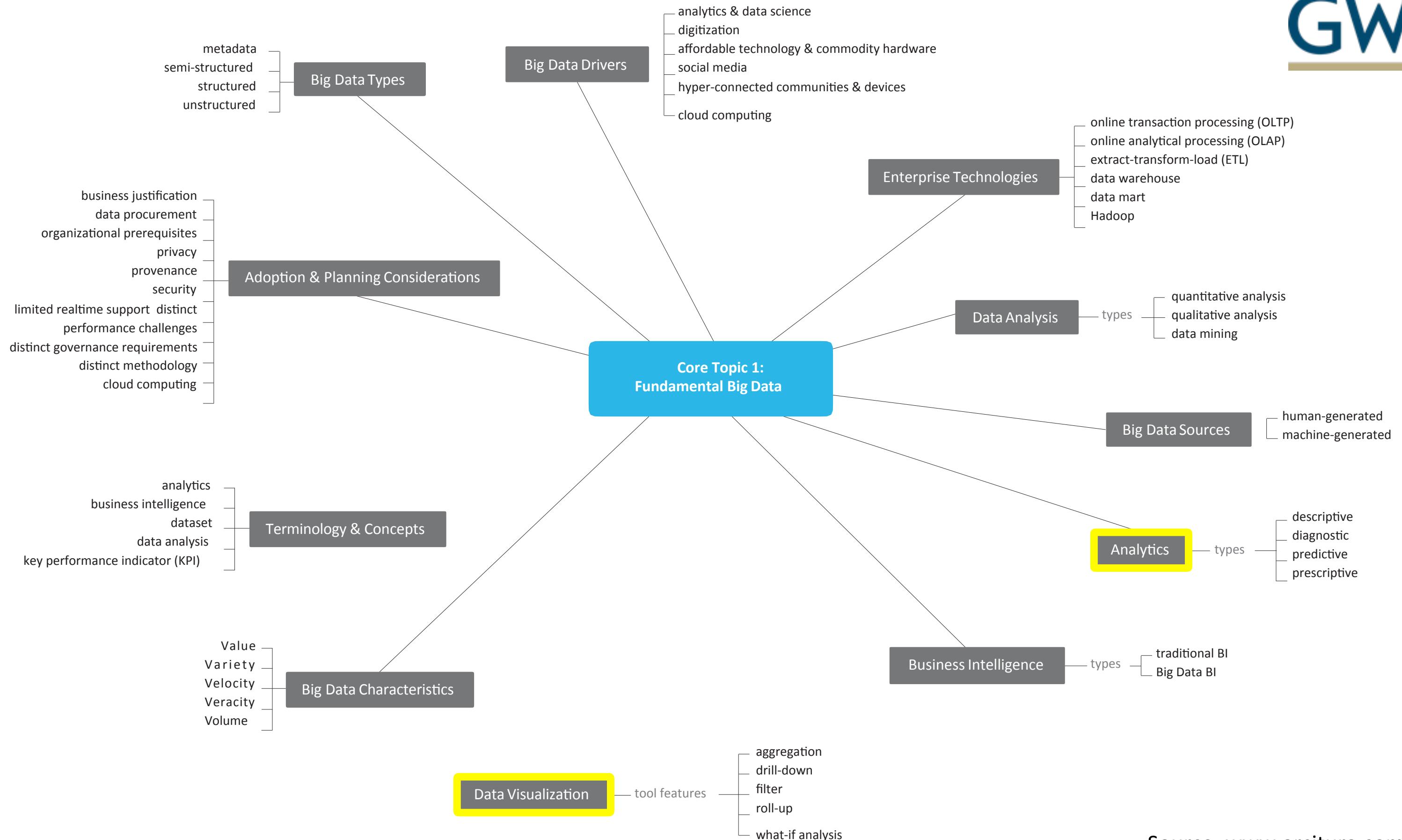


**Cancer Cell Line Encyclopedia (CCLE)**: a compilation of gene expression, chromosomal copy number and massively parallel sequencing data from 947 human cancer cell lines coupled with pharmacological profiles for 24 anticancer drugs across 479 of the cell lines: **18,897 genes**



**The Cancer Genome Atlas:** DNA methylation is an epigenetic mark which can be associated with transcriptional inactivity when located in promoter regions. Ovarian cancer study for gene expression and methylation correlated values for **22,000 genes across 598 samples**





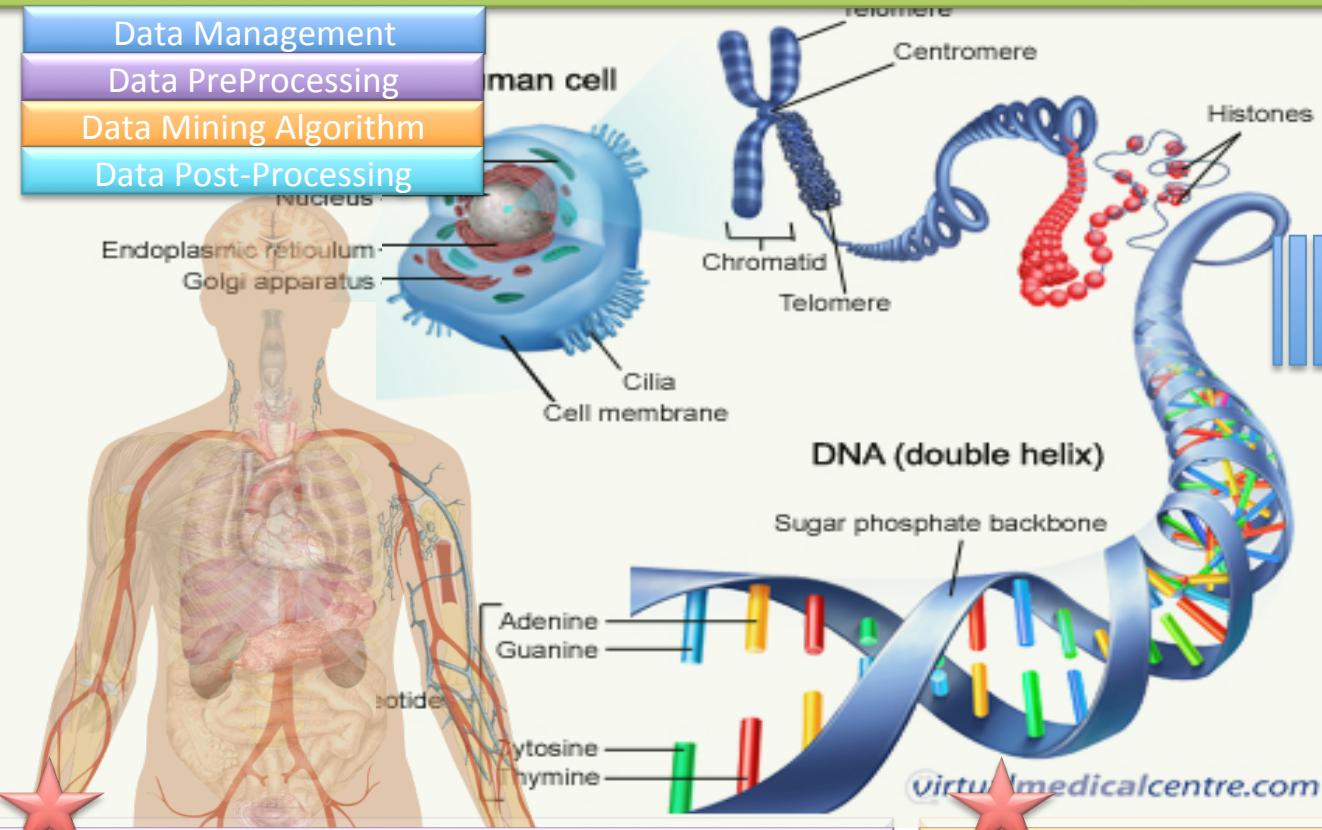
Source: [www.arcitura.com](http://www.arcitura.com)

# GENOMIC SIGNAL PROCESSING ANALYSIS

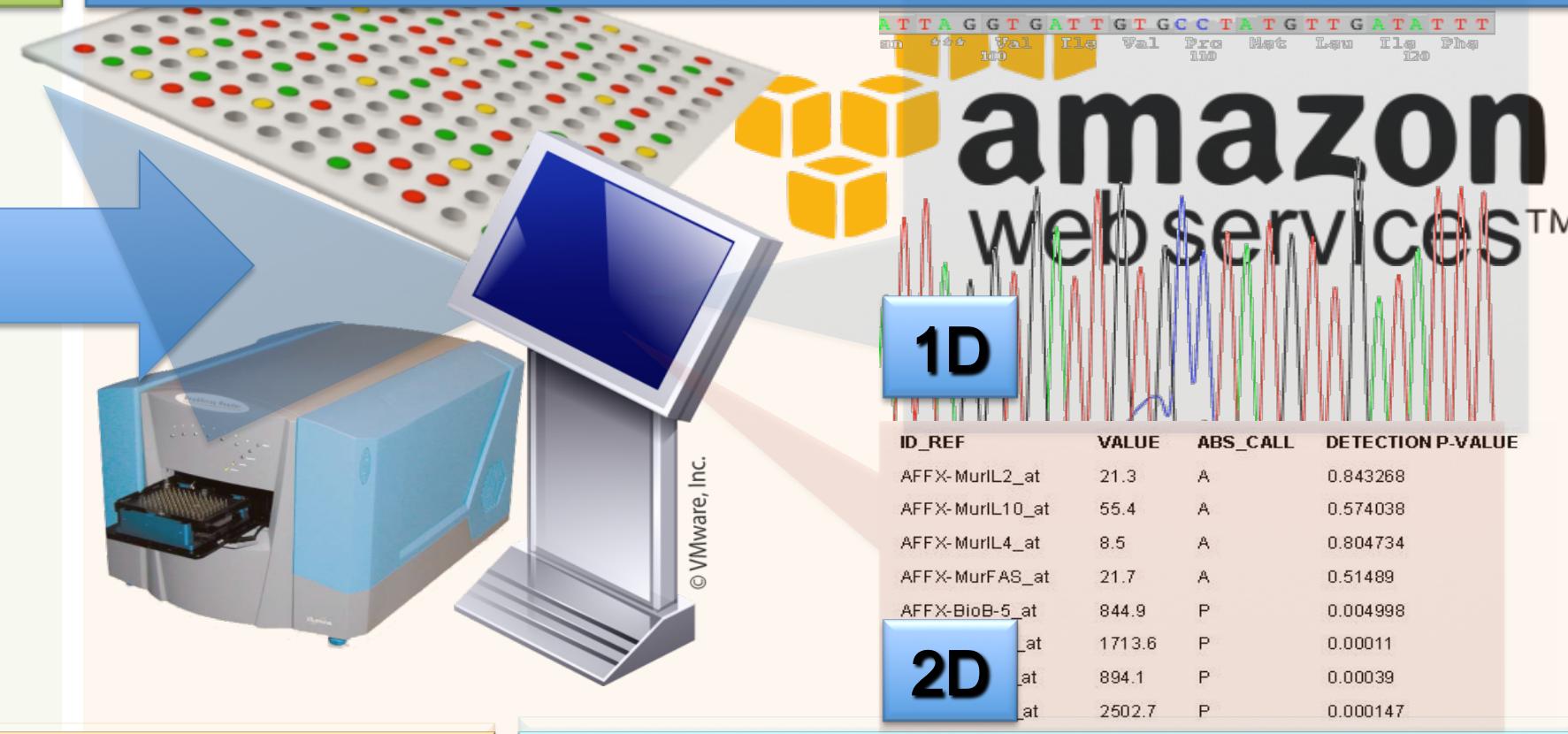
**GW**

**Genomic Signal Processing:** is a genomic data mining process which involves uncovering patterns, associations, anomalies, and statistically significant structures and events in genetics and genomics data

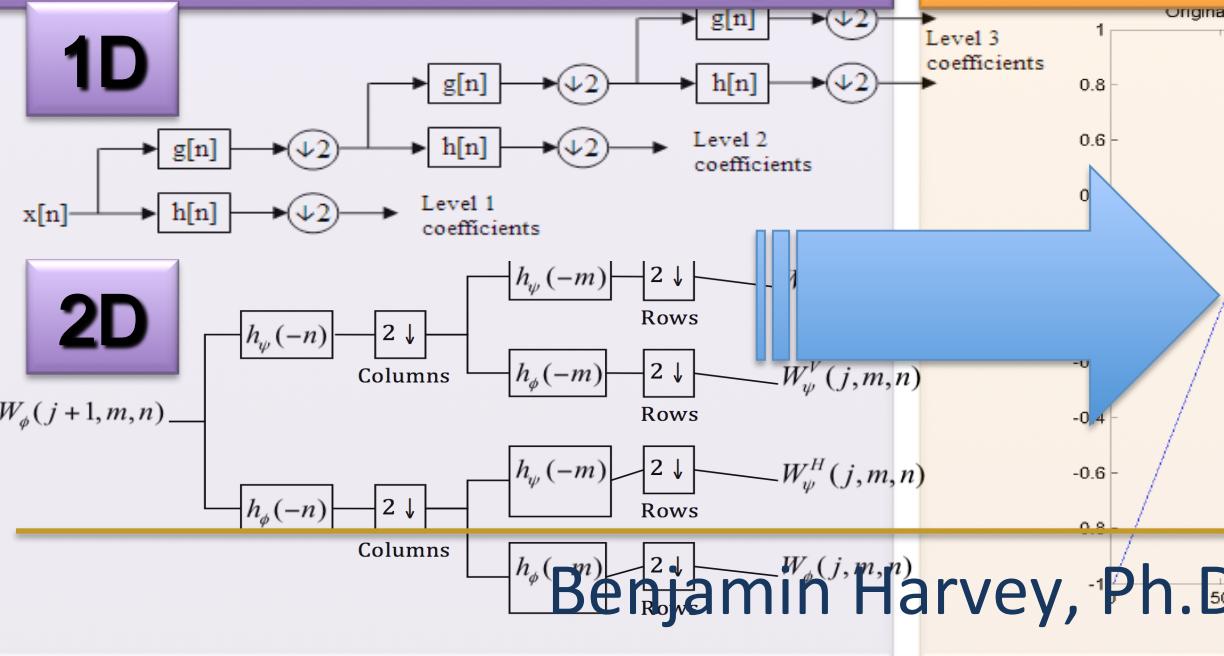
## Data Fusion and Sampling



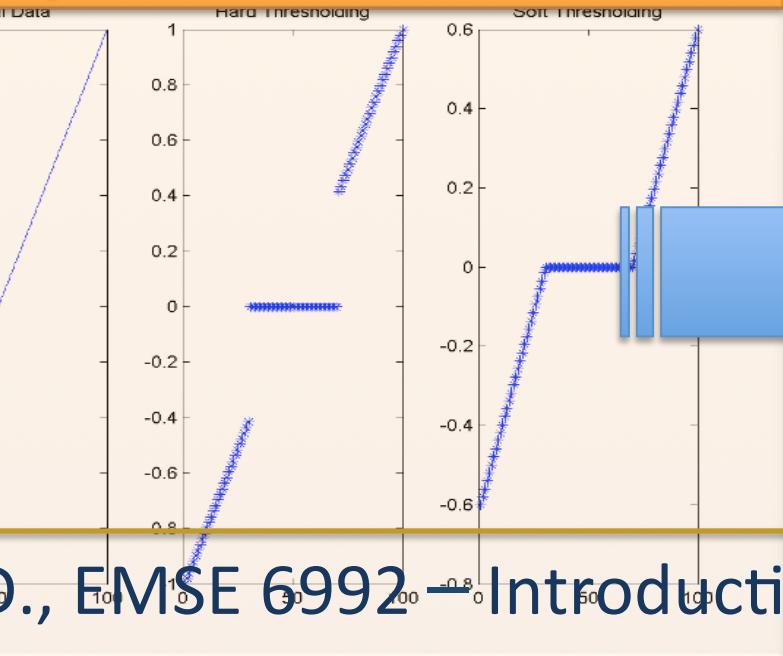
## Data Storage, Integration, and Representation



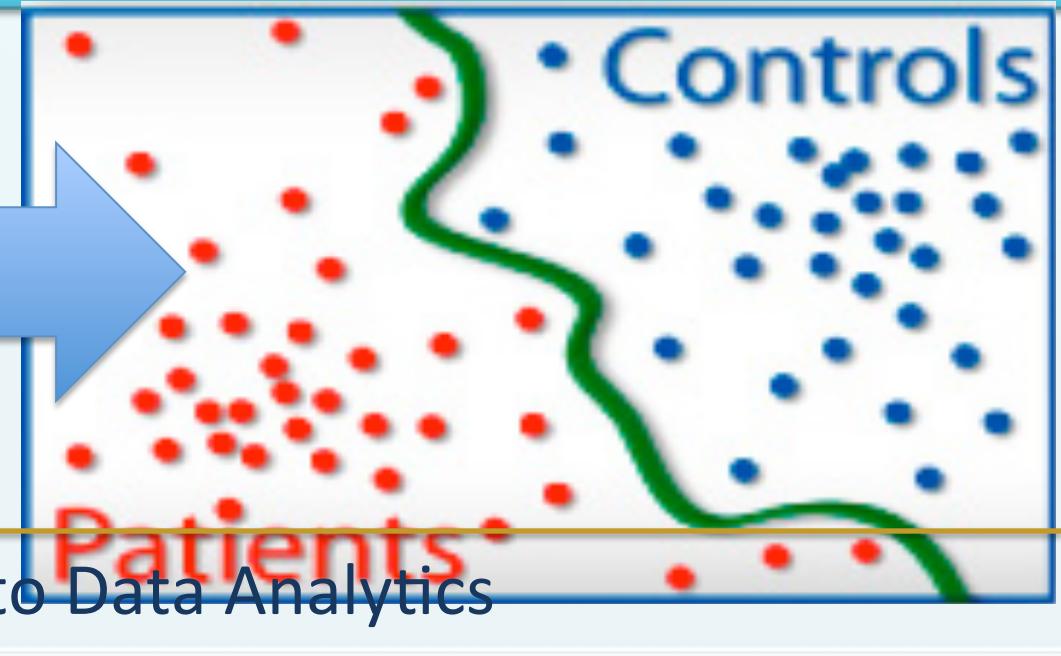
## Signal Pre-Processing

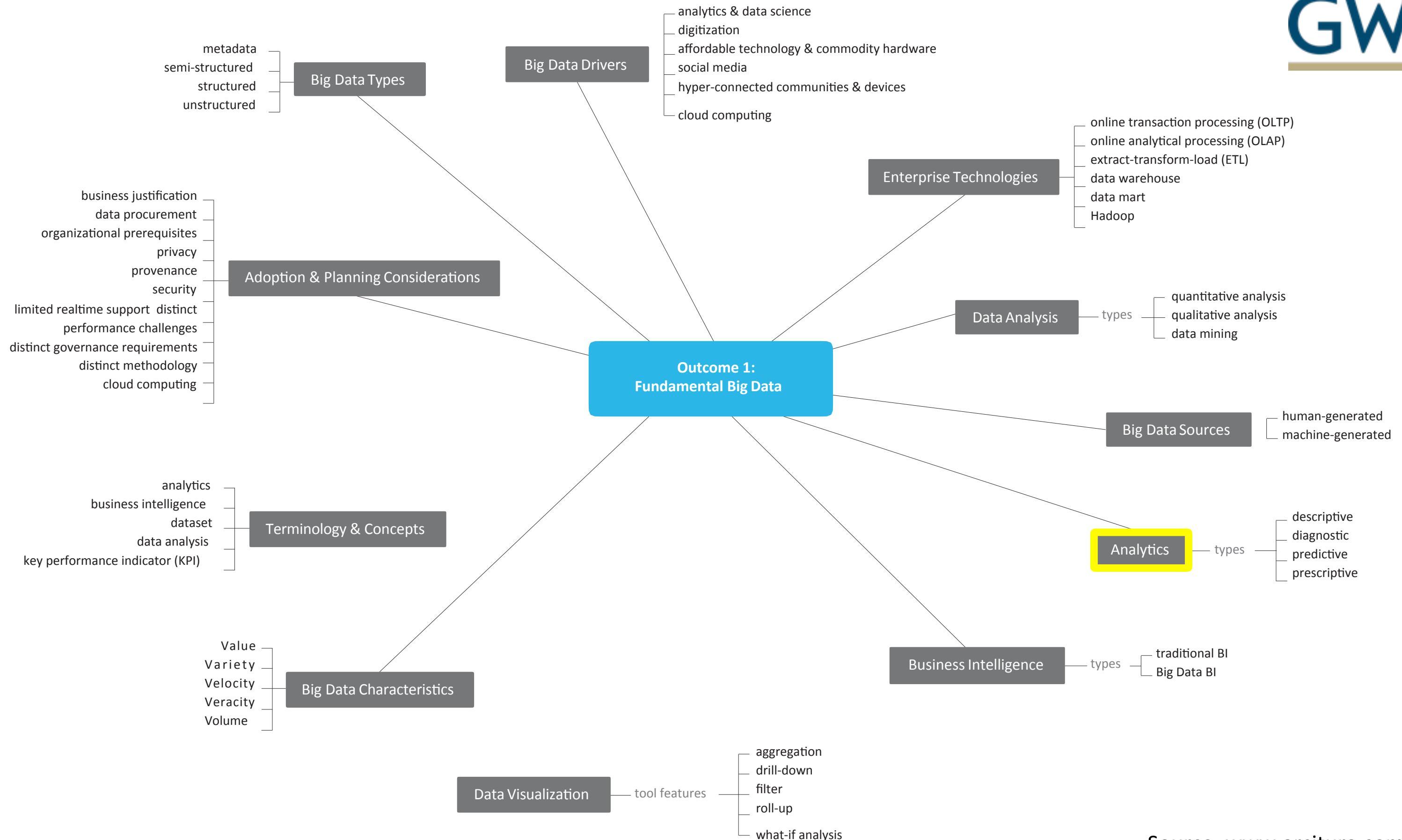


## Denoising & Feature Extraction



## Classification and Validation

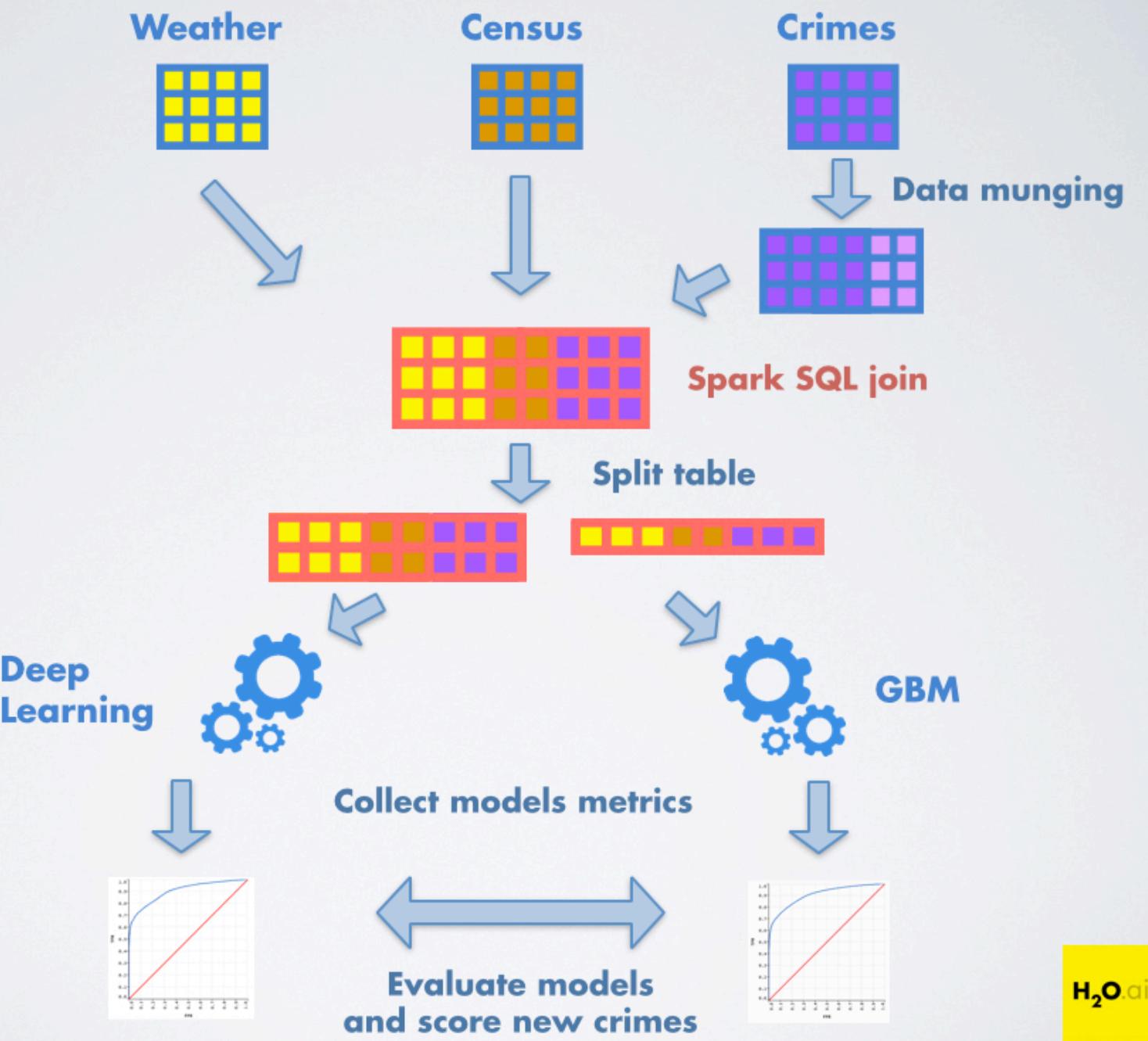




Source: [www.arcitura.com](http://www.arcitura.com)

# H<sub>2</sub>O ai and Crime Prediction

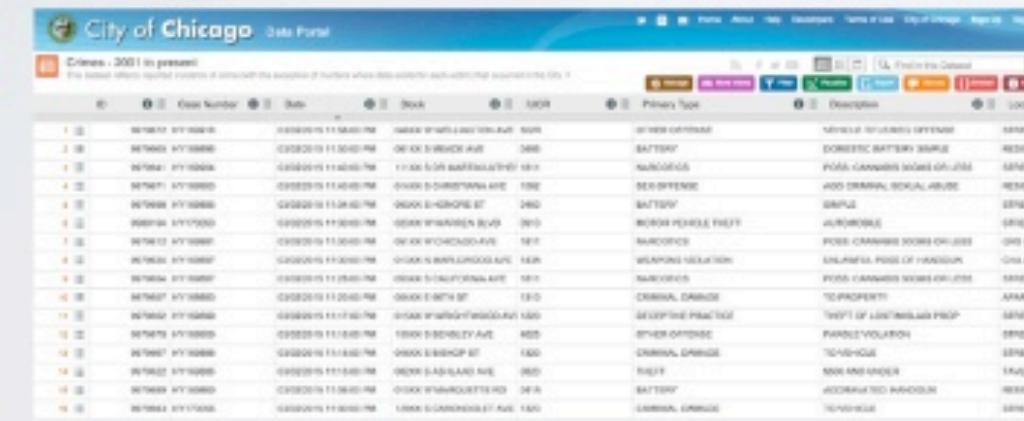
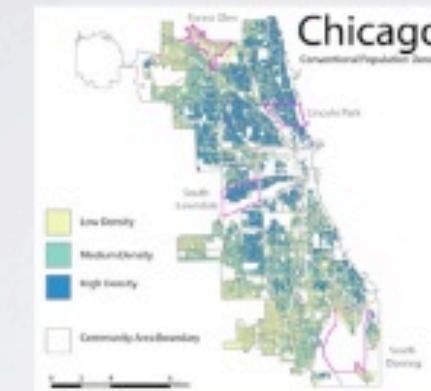
GW



## OPEN CITY, OPEN DATA

“...my kind of town” - F. Sinatra

U.S. Census  
Data ?



Weather Data ?



Crime Data

~4.6 Million rows of crimes from 2001, updated weekly\*  
External data source considerations???

H2O

deep learning

machine learning

H2O.ai

## H2O.ai Raises \$20M For Its Open Source Machine Learning Platform

Posted Nov 9, 2015 by Frederic Lardinois (@fredericl)

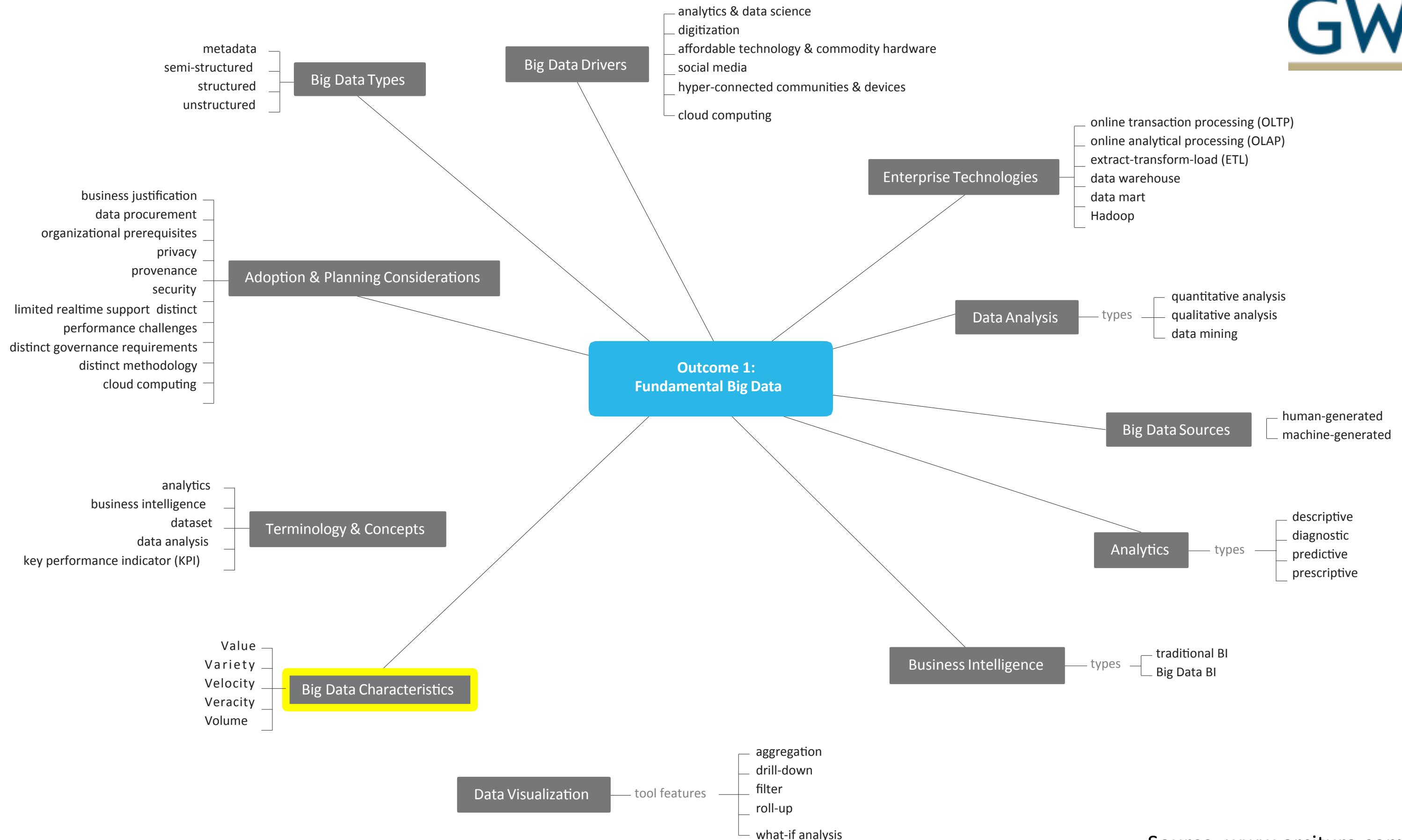
Popular Posts

Benjamin Harvey, Ph.D., EMSE 6992 – Introduction to Data Analytics

Next Story

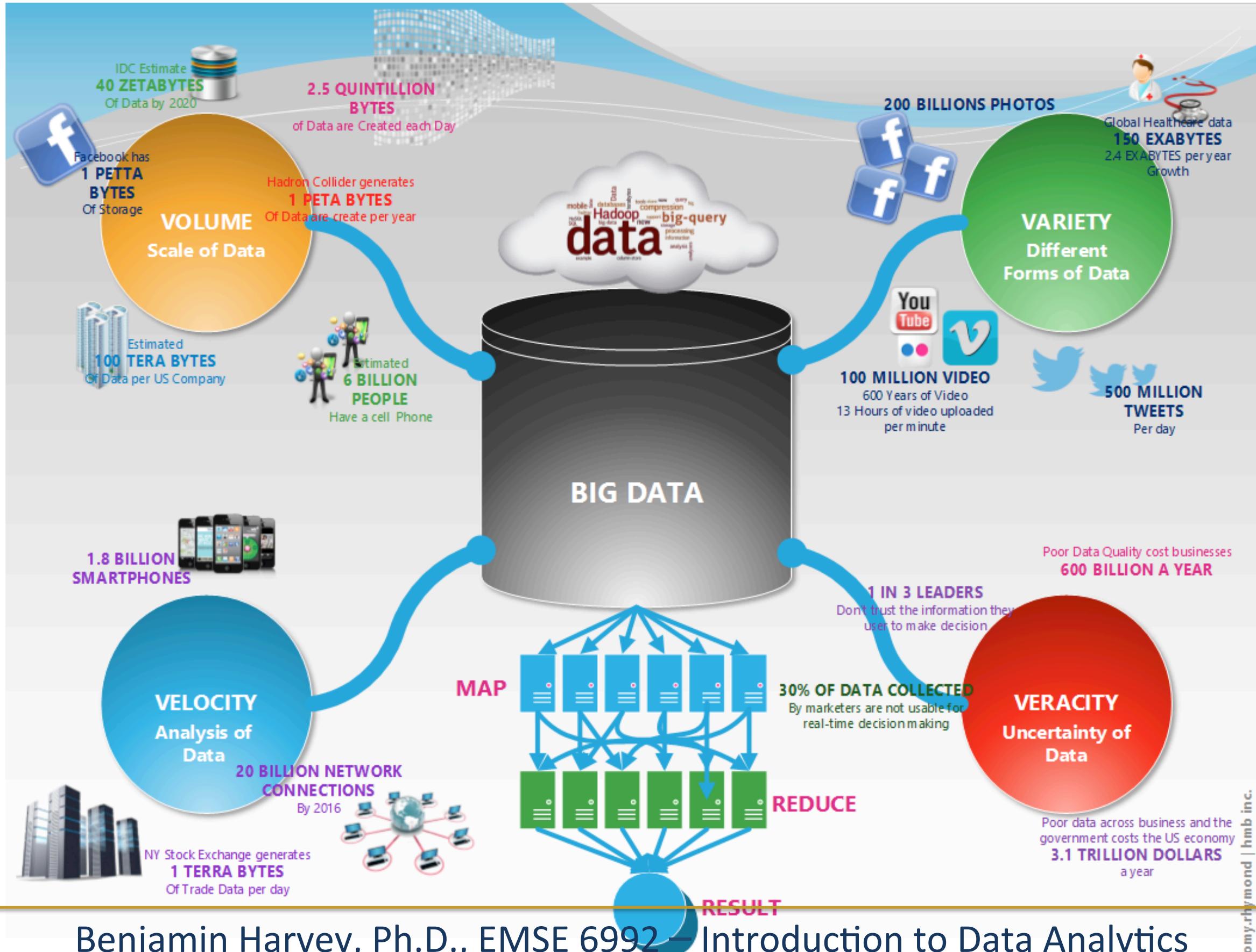
Mark Cuban

# IV. What is Big Data?

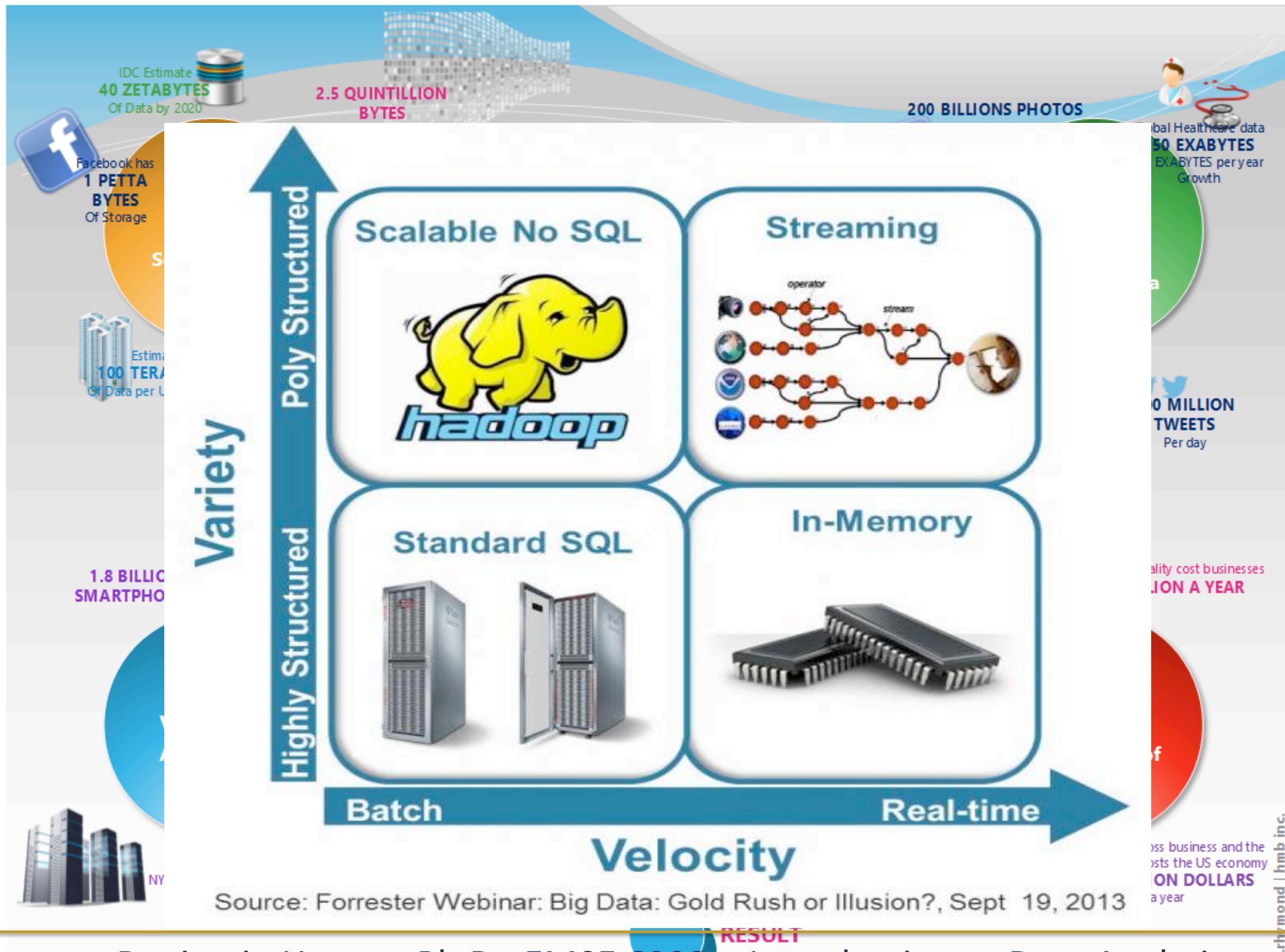


Source: [www.arcitura.com](http://www.arcitura.com)

# What is Big Data?



# What's NOT Big Data?



*There's certainly a lot of it!*



# Data, data everywhere...

1 Zettabyte

800 EB

1.8 ZB

8.0 ZB

Data produced each year

1 Exabyte

5 EB

161 EB

IBM builds 120 petabyte cluster out of 200,000 hard drives



Share This Article

By Sebastian Anthony on August 26, 2011 at 6:18 am | 16 Comments  
Smashing all known records by a multiple of 10, IBM Research Almaden, California, has developed hardware and software technologies that will allow it to strap together 200,000 hard drives to create a single storage cluster of 120 petabytes — or 120 million gigabytes. The drive collective, when it is complete, is expected to store one trillion files — or to put it in Apple terms, two billion hours of MP3 music.

120 PB

1 Petabyte

Human brain's capacity

2002

2006 2009 2011 2015

1 Petabyte == 1000 TB

1 TB = 1000 GB

60 PB

14 PB

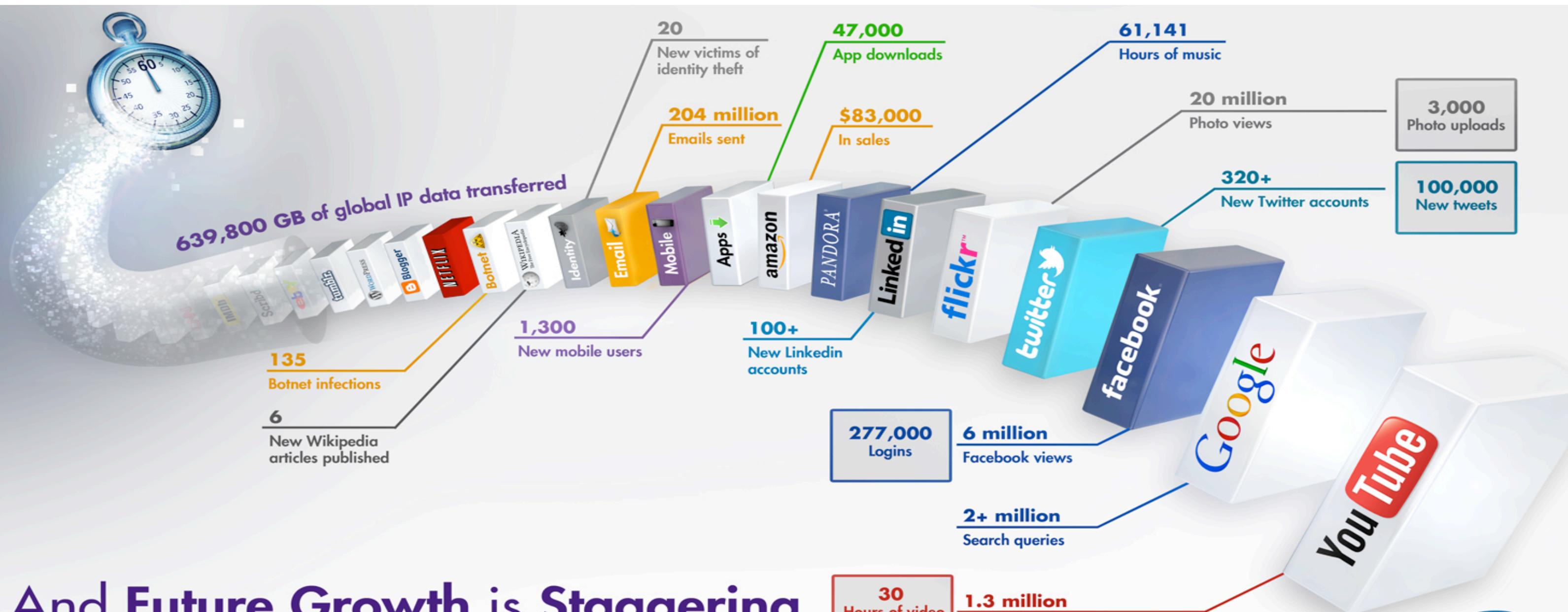
100-years of HD video + audio

## References

- (2015) 8 ZB: <http://www.emc.com/collateral/analyst-reports/idc-extracting-value-from-chaos-ar.pdf>
- (2011) 1.8 ZB: <http://www.emc.com/leadership/programs/digital-universe.htm>
- (2009) 800 EB: <http://www.emc.com/collateral/analyst-reports/idc-digital-universe-are-you-ready.pdf>
- (2006) 161 EB: <http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>

- (2002) 5 EB: <http://www2.sims.berkeley.edu/research/projects/how-much-info-2003/execsum.htm>
- (life in video) 60 PB: in 4320p resolution, extrapolated from 16MB for 1:21 of 640x480 video (w/ sound) – almost certainly a gross overestimate, as sleep can be compressed significantly!
- (brain) 14 PB: <http://www.quora.com/Neuroscience-1/How-much-data-can-the-human-brain-store>

# 60 Second Data Analysis

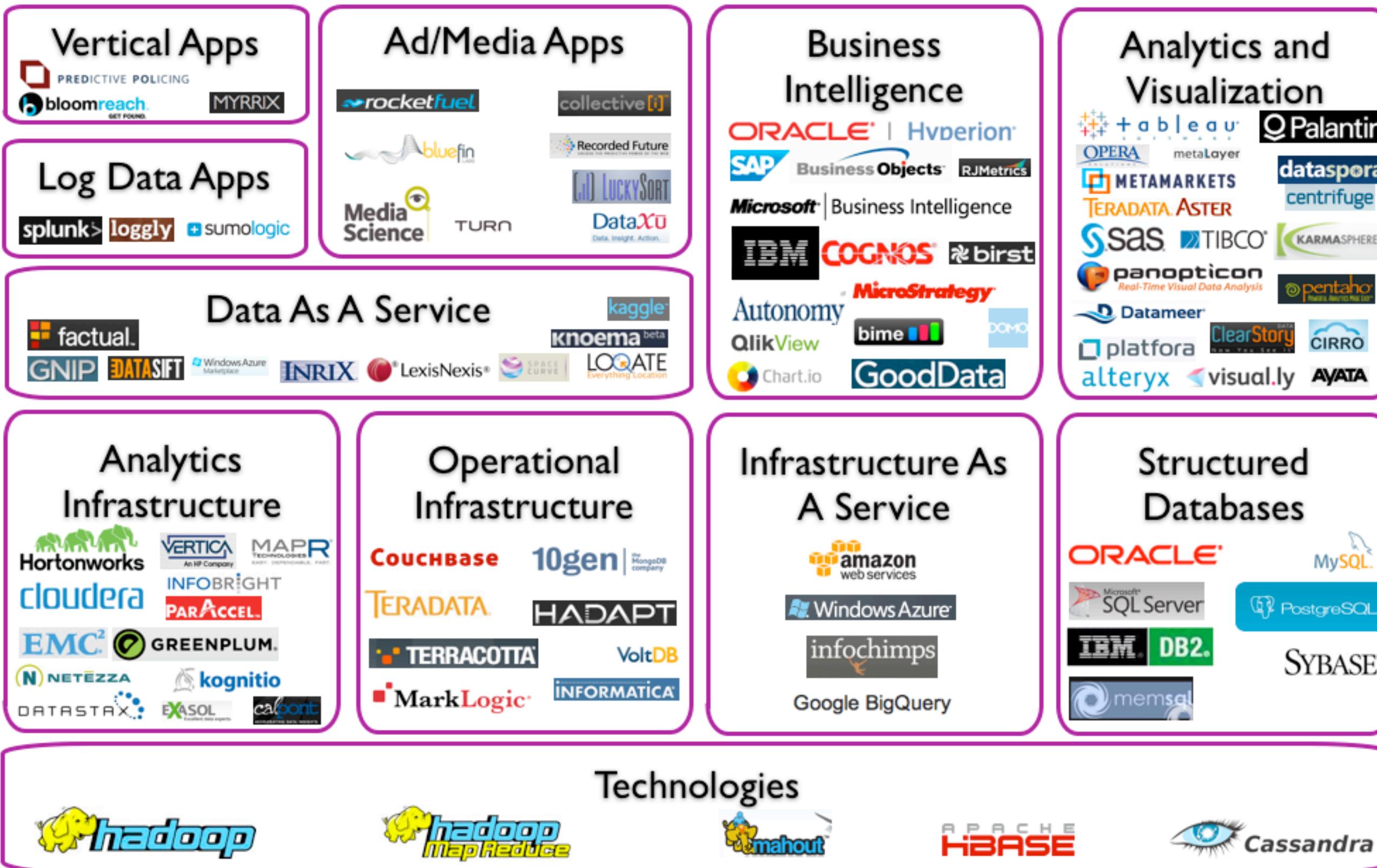


And Future Growth is Staggering



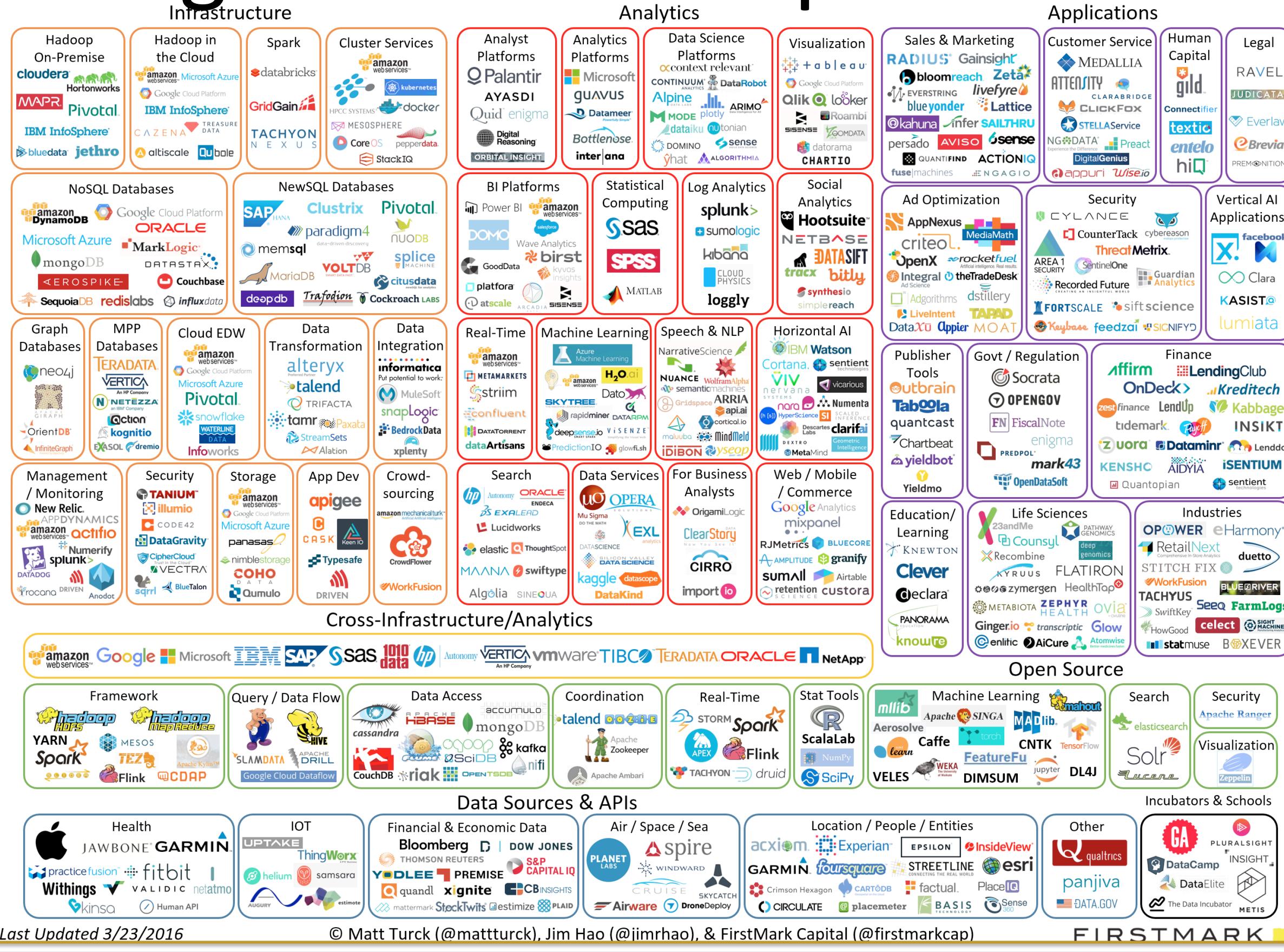
# Big Data Landscape - 2012

**GW**



# Big Data Landscape - 2016

GW



Last Updated 3/23/2016

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)

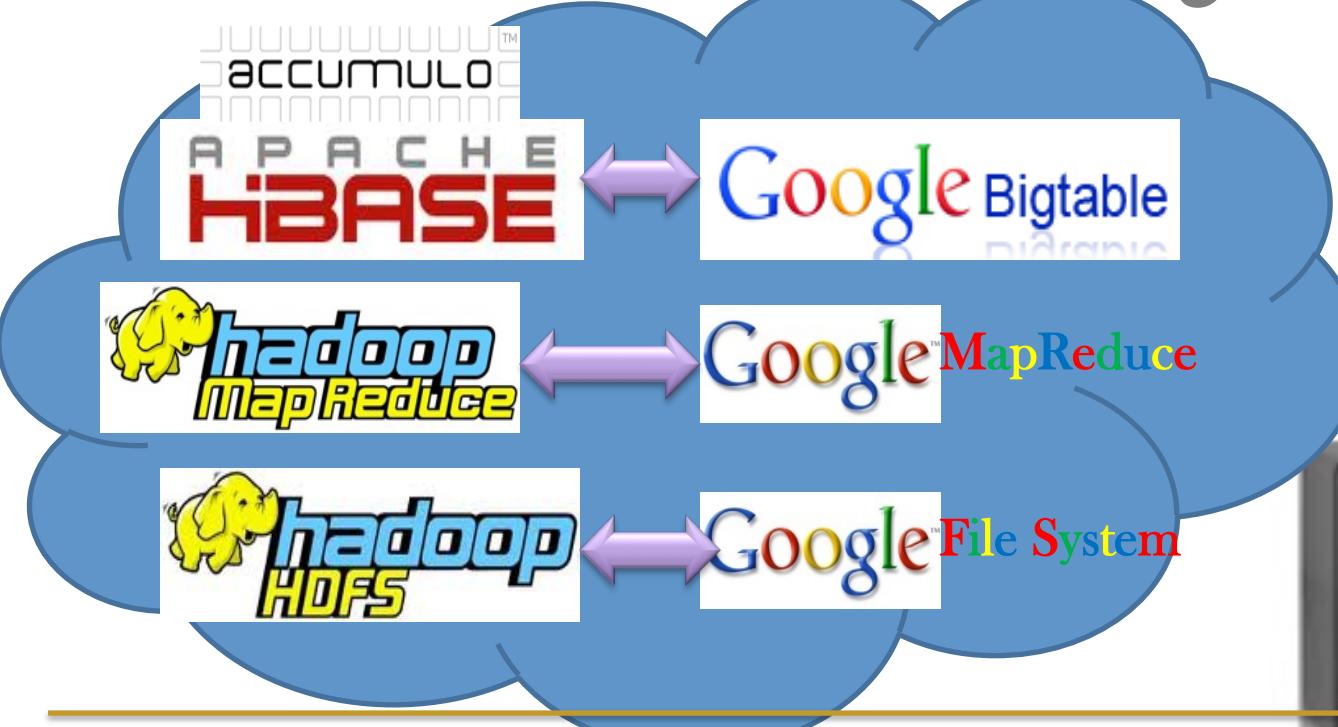
FIRSTMARK

Benjamin Harvey, Ph.D., EMSE 6992 – Introduction to Data Analytics



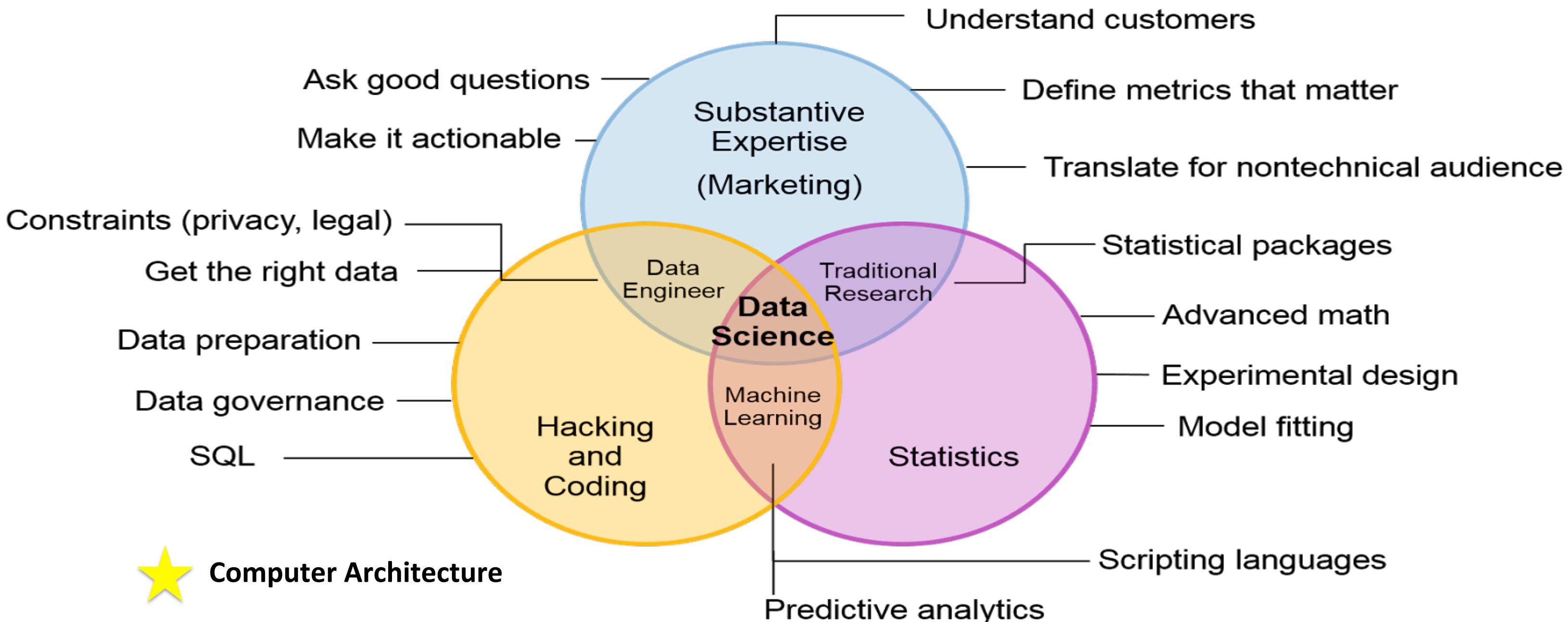
**APACHE  
HBASE**

## Data Cloud History

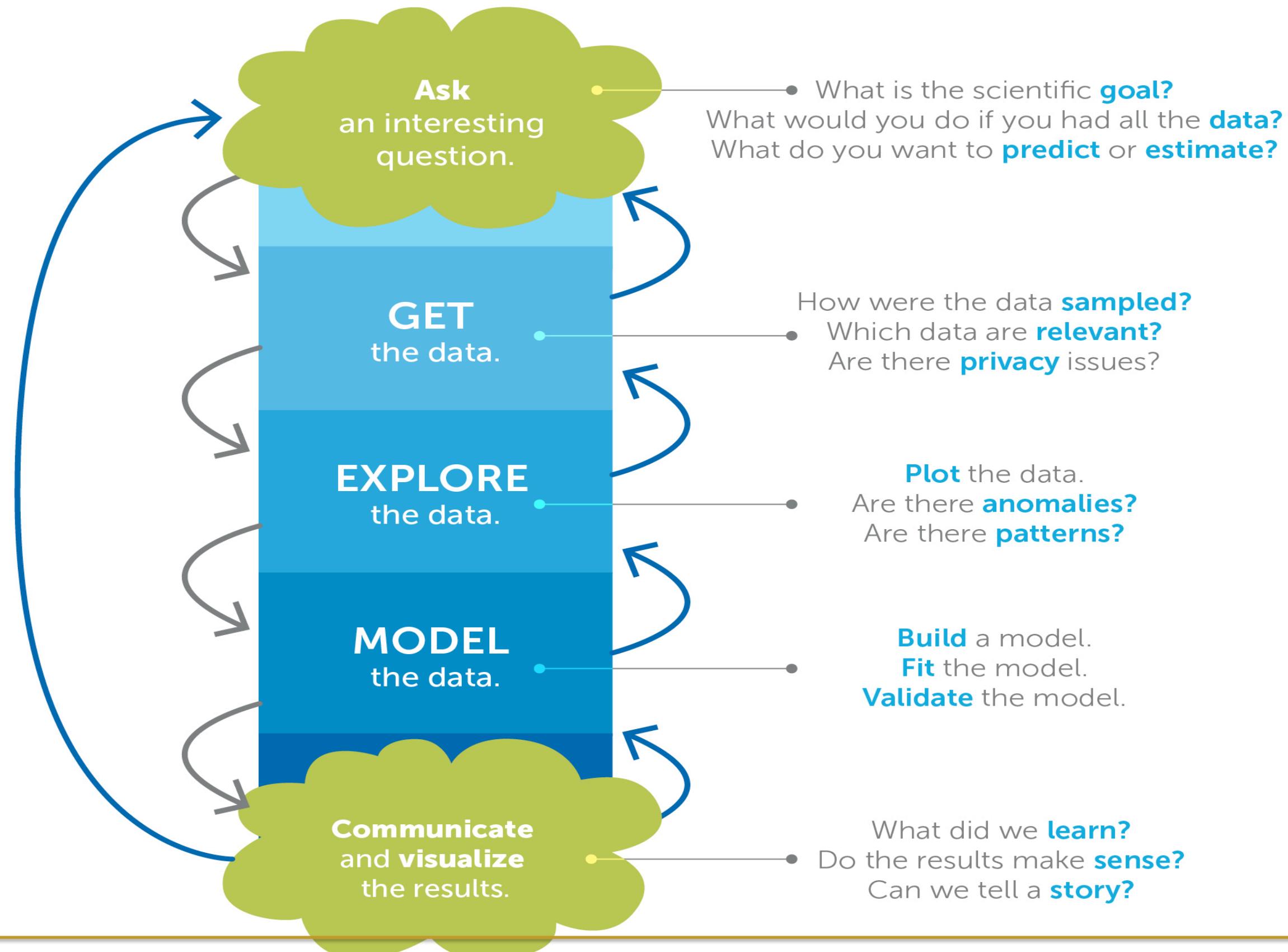


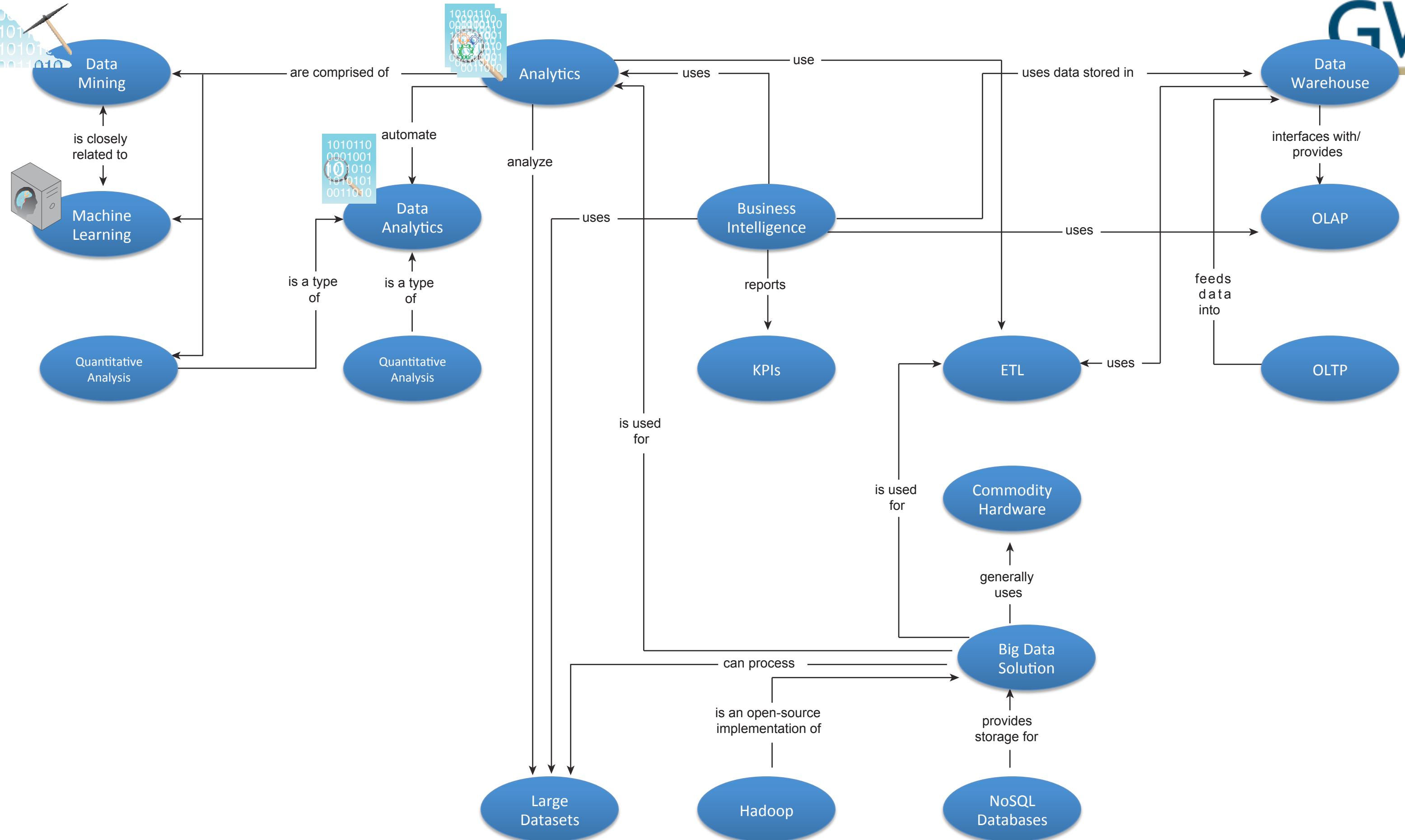
## IV. What is Data Science

# Data Science Venn Diagram



# Data Science Process





# Discussion

- What is Data Science to you?
- Provide an example.
- What is “Big Data” to you?
- Provide an example.

# Summary: What is Data Science? **GW**

## What is Big Data?

Data Science: To gain insights through computation, statistics, and visualization



- In this course, you will learn four core topics that have been associated with Data Science
  - Big Data Foundations, Big Data and Technology Concepts, Big Data Analysis and Science, Advanced Big Data Analysis and Science
- The definition of Big Data is subjective and really depends of the following:
  - Volume, Velocity, Veracity, Variety of data
  - Current CPU capacity

## IV. Foundations of Big Data

# Agenda

## Concepts and Terminology

- Datasets
- Data Analysis
- Types of Data Analytics
  - Descriptive Analytics, Diagnostic Analytics, Predictive Analytics, Prescriptive Analytics
- Business Intelligence (BI)
- Key Performance Indicators (KPI)

## Big Data Characteristics

- Volume, Velocity, Variety, Veracity, Value

## Different Types of Data

- Structured Data
- Unstructured Data
- Semi-structured Data
- Metadata

## Case Study Background

# Concepts and Terminology



## Data Analysis

- Data analysis is the process of examining data to find facts, relationships, patterns, insights and/or trends. The overall goal of data analysis is to support better decision making.

## Data Analytics

- Data analytics is a discipline that includes the management of the complete data lifecycle, which encompasses collecting, cleansing, organizing, storing, analyzing and governing data.

## Big Data Analytics

- The lifecycle generally involves identifying, procuring, preparing and analyzing large amounts of raw, unstructured data to extract meaningful information that can serve as an input for identifying patterns, enriching existing enterprise data and performing large-scale searches.

Data analytics enable data-driven decision-making with scientific backing so that decisions can be based on factual data and not simply on past experience or intuition alone.

- descriptive analytics
- diagnostic analytics
- predictive analytics
- prescriptive analytics

# Concepts and Terminology

## Descriptive Analytics

- Descriptive analytics are carried out to answer questions about events that have already occurred.



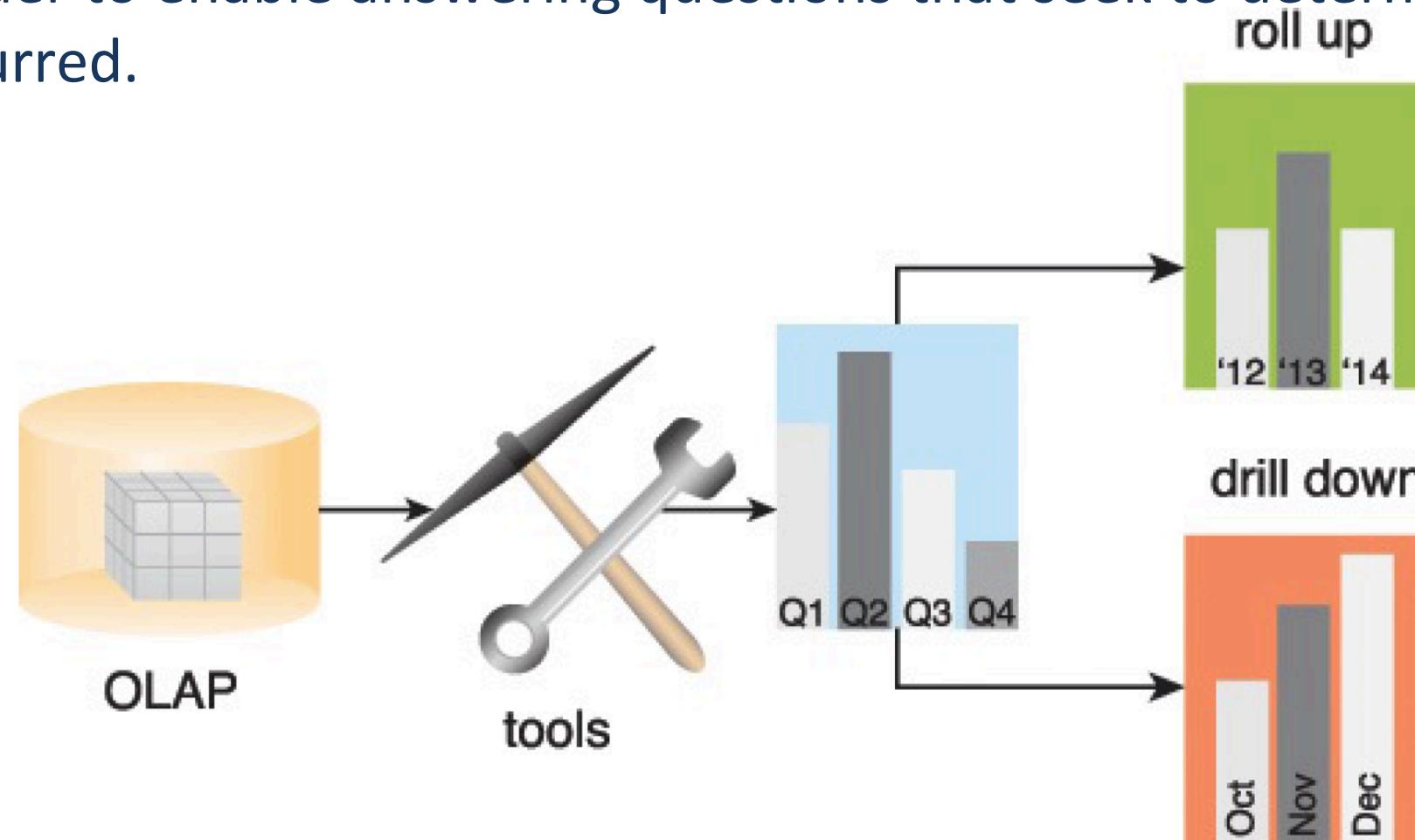
**Figure 1.5** The operational systems, pictured left, are queried via descriptive analytics tools to generate reports or dashboards, pictured right.

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

# Concepts and Terminology

## Diagnostic Analytics

- Diagnostic analytics aim to determine the cause of a phenomenon that occurred in the past using questions that focus on the reason behind the event.
- The goal of this type of analytics is to determine what information is related to the phenomenon in order to enable answering questions that seek to determine why something has occurred.



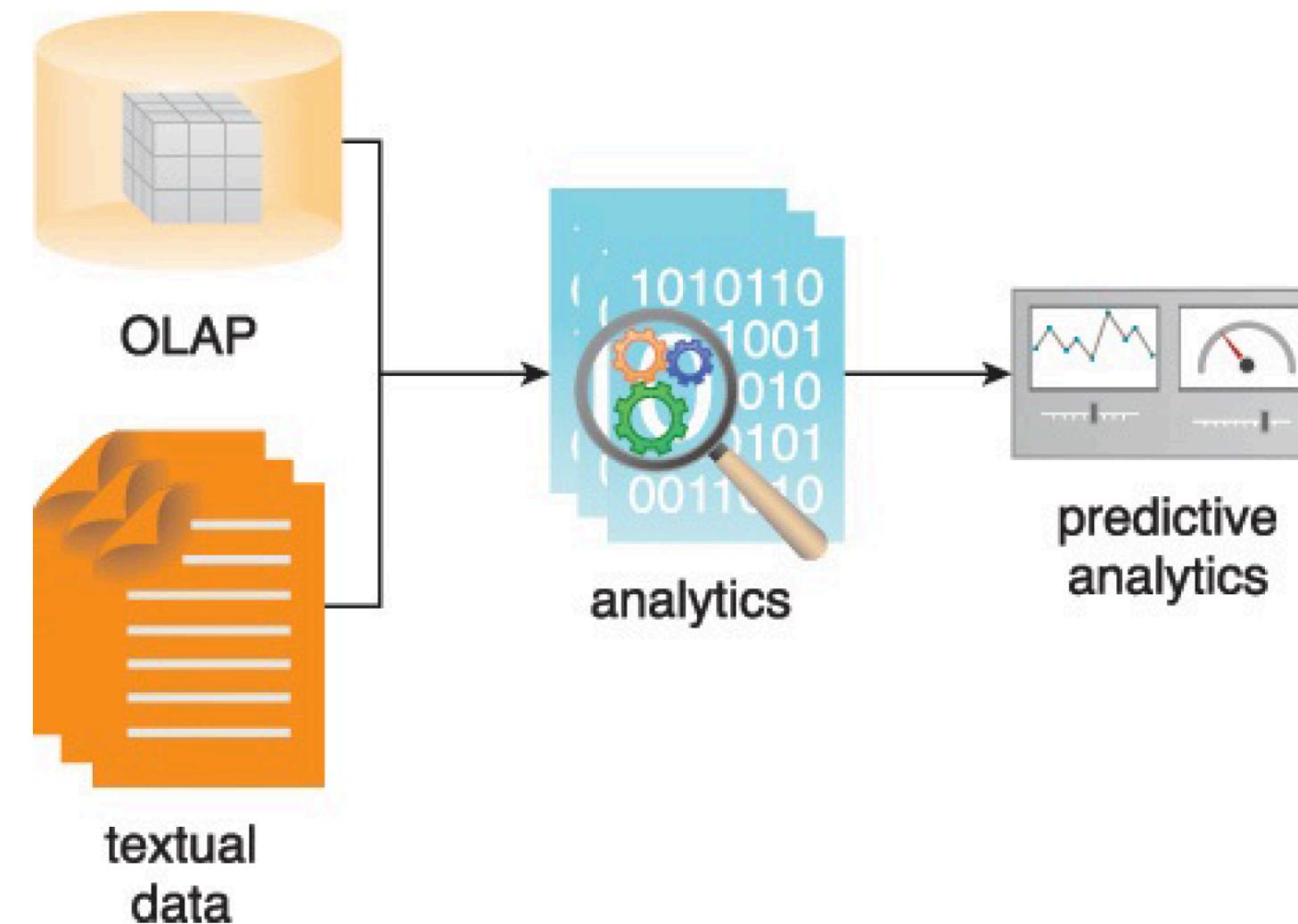
**Figure 1.6** Diagnostic analytics can result in data that is suitable for performing drill-down and roll-up analysis.

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

# Concepts and Terminology

## Predictive Analytics

- Predictive analytics are carried out in an attempt to determine the outcome of an event that might occur in the future. Information is enhanced with meaning to generate knowledge that conveys how that information is related.



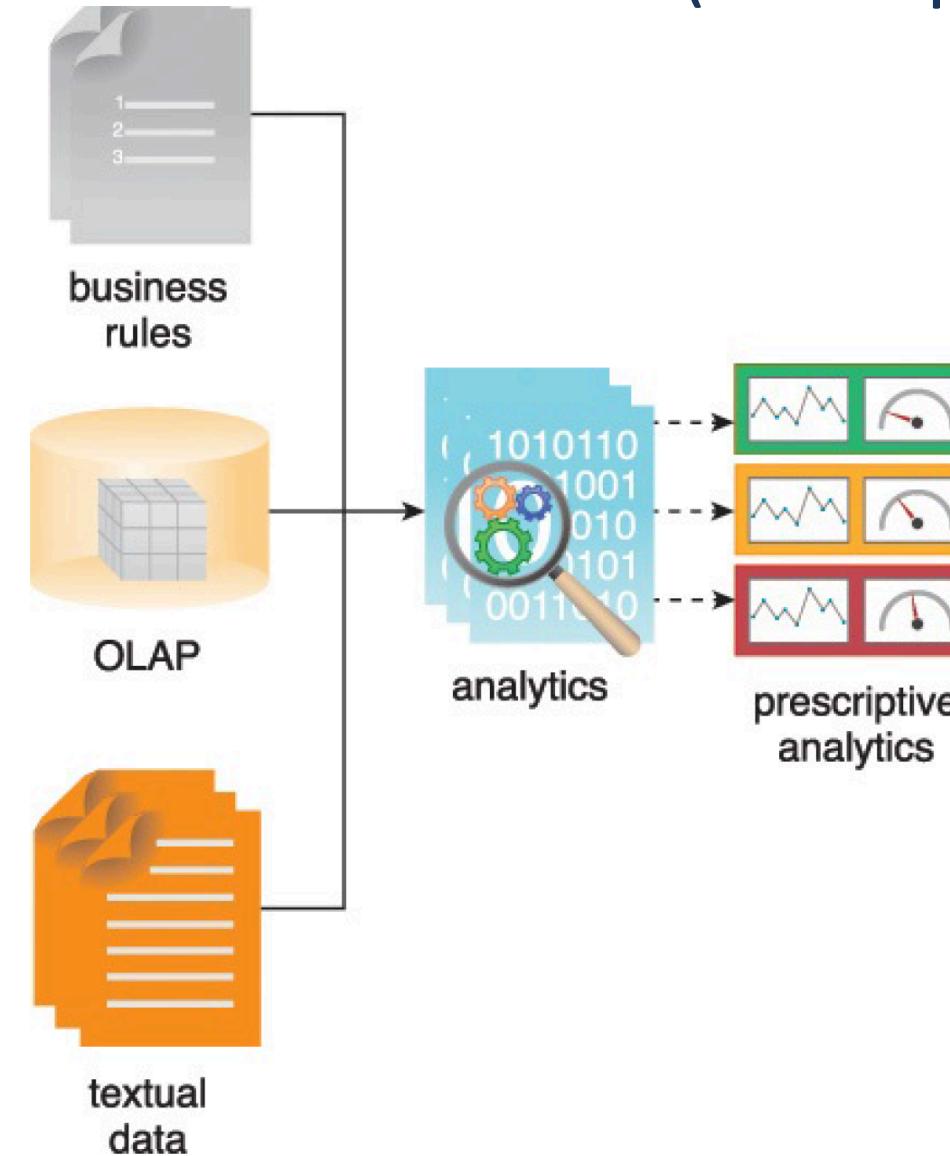
**Figure 1.7** Predictive analytics tools can provide user-friendly front-end interfaces.

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

# Concepts and Terminology

## Prescriptive Analytics

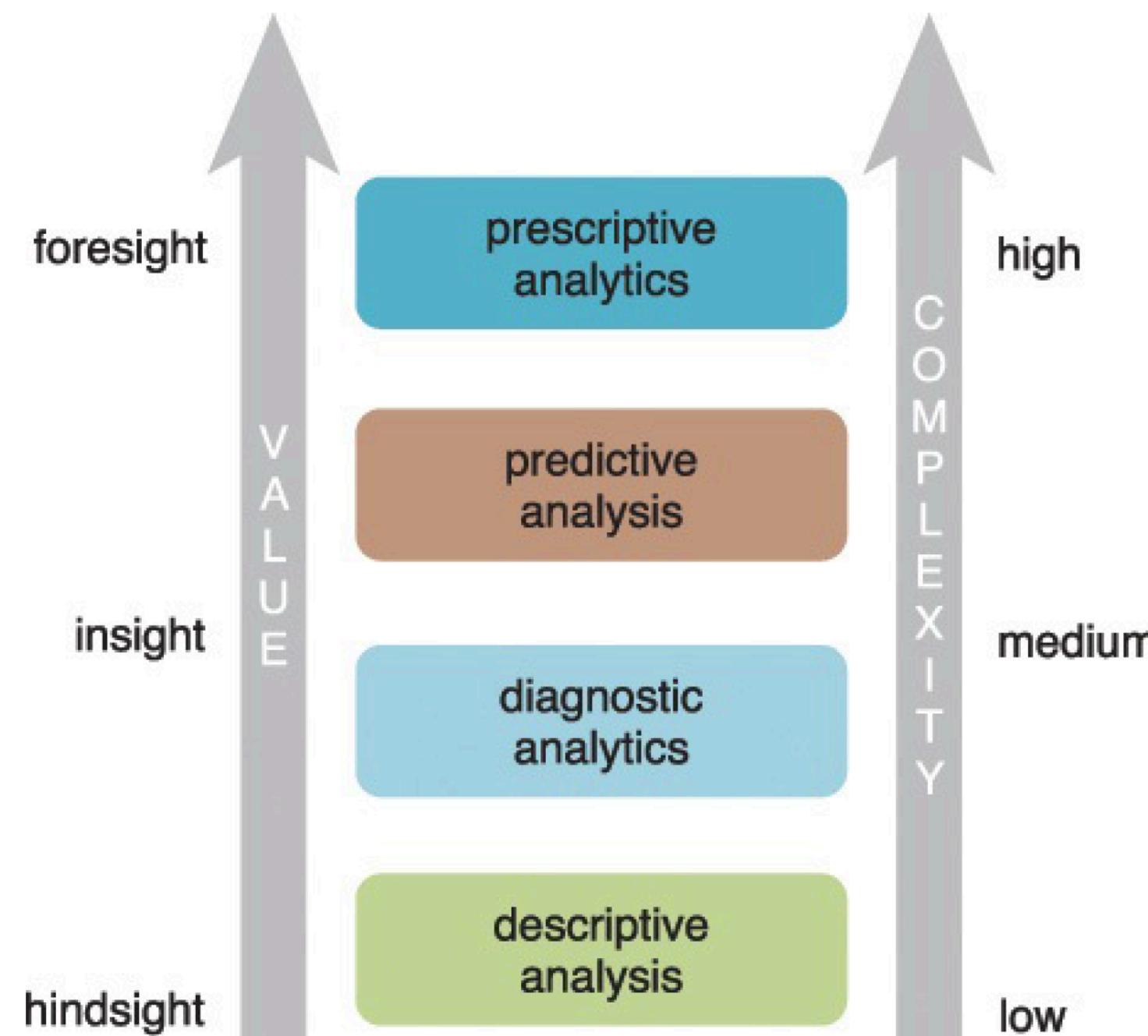
- Prescriptive analytics build upon the results of predictive analytics by prescribing actions that should be taken (best option and why).



**Figure 1.8** Prescriptive analytics involves the use of business rules and internal and/or external data to perform an in-depth analysis.

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

# Concepts and Terminology



**Figure 1.4** Value and complexity increase from descriptive to prescriptive analytics.

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

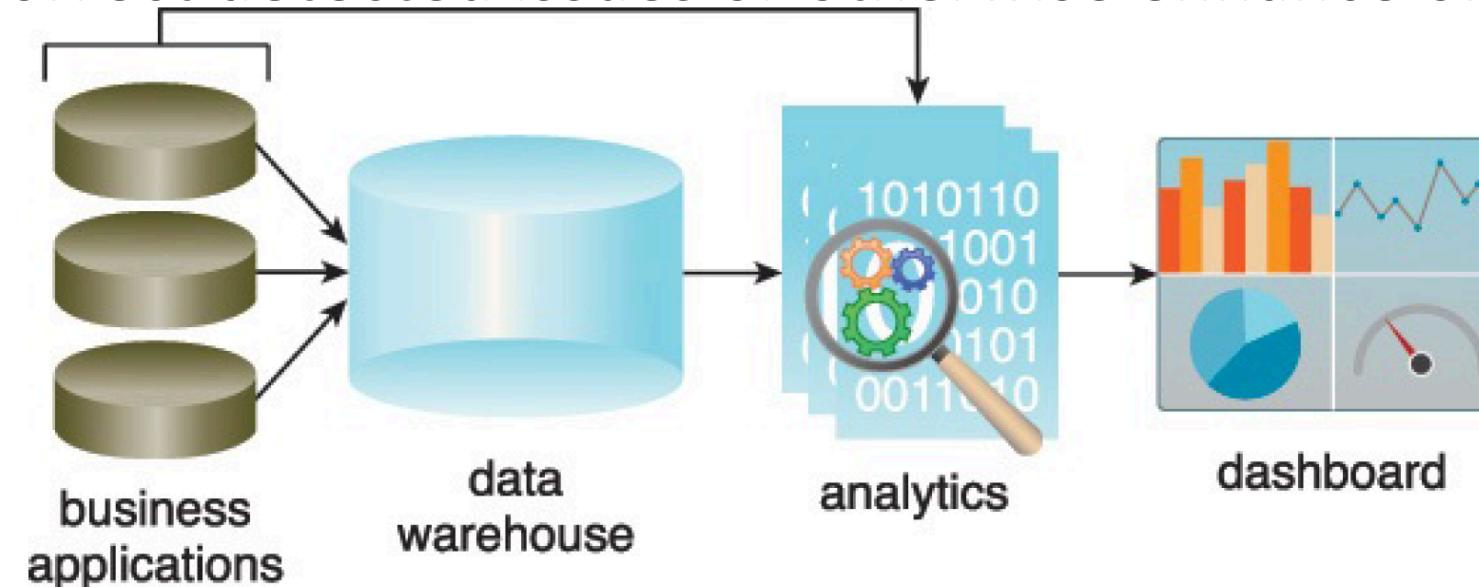
# Concepts and Terminology

## Key Performance Indicators (KPI)

- A KPI is a metric that can be used to gauge success within a particular business context.
- KPIs are linked with an enterprise's overall strategic goals and objectives.
- They are often used to identify business performance problems and demonstrate regulatory compliance.

## Business Intelligence (BI)

- BI enables an organization to gain insight into the performance of an enterprise by analyzing data generated by its business processes and information systems.
- The results of the analysis can be used by management to steer the business in an effort to correct detected issues or otherwise enhance organizational performance



**Figure 1.9** BI can be used to improve business applications, consolidate data in data warehouses and analyze queries via a dashboard.

# Concepts and Terminology

## Structured Data

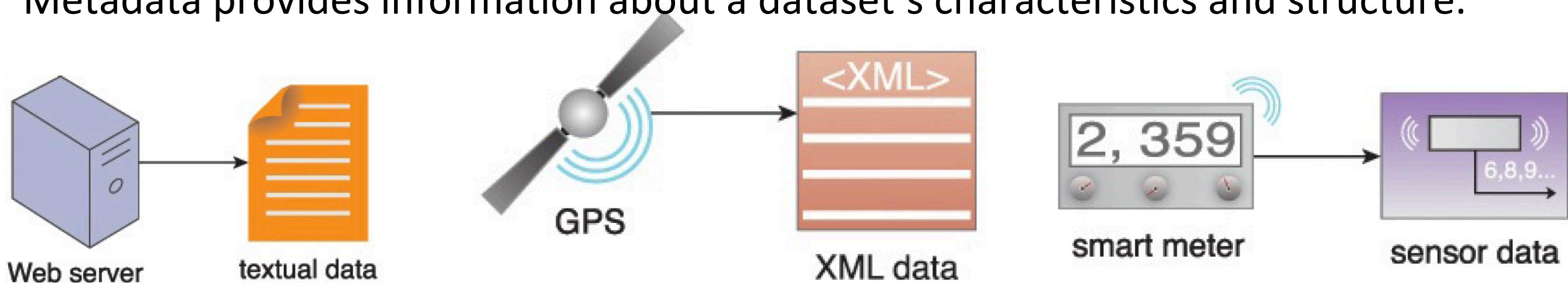
- Structured data conforms to a data model or schema and is often stored in tabular form. It is used to capture relationships between different entities and is therefore most often stored in a relational database.

## Semi-structured Data

- Semi-structured data has a defined level of structure and consistency, but is not relational in nature. Instead, semi-structured data is hierarchical or graph-based.

## Unstructured Data

- Data that does not conform to a data model or data schema is known as unstructured data.
- Metadata provides information about a dataset's characteristics and structure.



**Figure 1.17** Examples of machine-generated data include web logs, sensor data, telemetry data, smart meter data and appliance usage data

Erl, Thomas, Wajid Khattak, and Paul Buhler. Big data fundamentals: concepts, drivers & techniques. Prentice Hall Press, 2016.

# Case Study: Ensure to Insure

## Identifying Types of Data

- The IT team members go through a categorization exercise of the various datasets that could be identified:
  - Categorize them as either: Structured, Unstructured and Semi-Structured

## Identifying Data Characteristics

- The IT team members want to gauge different datasets that are generated inside ETI's boundary as well as any other data generated outside ETI's boundary that may be of interest to the company in the context of volume, velocity, variety, veracity and value characteristics.
- The team members take each characteristic in turn and discuss how different datasets manifest that characteristic.

Do you think the new Big Data strategy will work?