

Big Data Systems and Analytics: Opportunities and Challenges

Ling Liu
Professor

Distributed Data Intensive Systems Lab
School of Computer Science
Georgia Institute of Technology



Theme of this Course



Big Data Systems* *Big Data Analytics



Scope of the Course



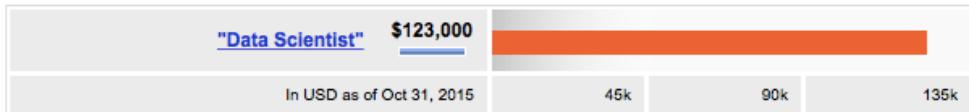
- How various parts of a big data computer system (hardware, software and applications) are put together?
- What are the appropriate approaches to realize high-performance, scalability, reliability, and security in practical big data systems?
- How to effectively manage very large amounts of data and extract value and knowledge from them

Administravia

- Course Website:
<http://www.cc.gatech.edu/~lingliu/courses/8803/2016Fall/cs8803-BDS.html>
- Office Hours
 - Thursdays 11am -12noon, or by appointment
- Grading
 - ◆ Class Participation 15%
 - ◆ Homework Assignments (4) 20%
 - ◆ Project 50%
 - ◆ Final (Technology Review) 15 %
- Announcements
 - ◆ In class and on T-Square

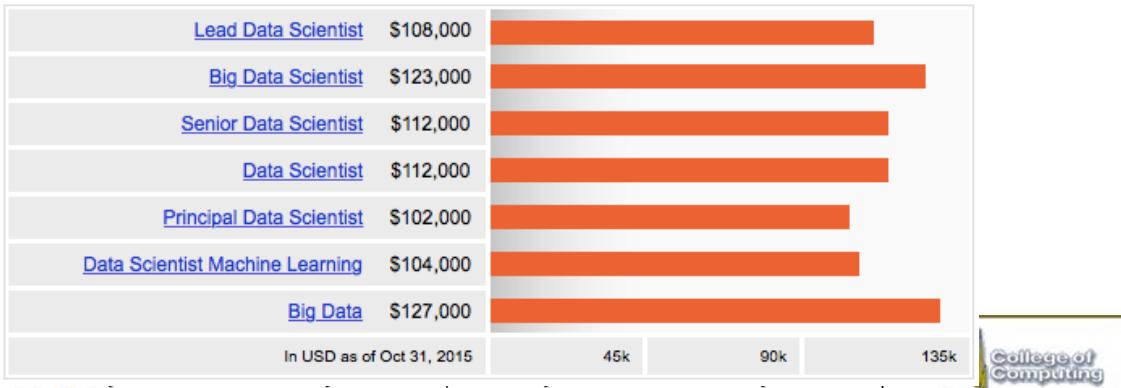
From Computer Programming to Data Scientists

Average Salary of Jobs Matching Your Search



Average "Data Scientist" salaries for job postings nationwide are 113% higher than average salaries for all job postings nationwide.

Average Salary of Jobs with Related Titles



Outline

- Big Data and Big Data Computing
 - Big Data Challenges
 - Innovation Opportunities in Big Data Computing
 - Programmable algorithm abstractions
 - Parallel processing abstractions
 - Leveraging Multiple Learning Models
 - Challenges of Big Data
 - Data Privacy
 - Data Quality



Characteristics of Big Data Sets

- Huge
- Distributed
 - Dispersed over many servers
- Dynamic
 - Items add/deleted/modified continuously
- Heterogeneous
 - Many agents access/update data
- Noisy
 - Inherent
 - Unintentional
 - Malicious
- Unstructured / semi-structured
 - No database schema

Big Data Puzzles

What is Big Data?

What makes data, “Big” Data?

Why Big Data now?



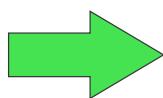
What is Big Data: A CS Perspective

- Big data refers to datasets
 - that are *beyond the ability of legacy approaches* to manage at an acceptable level of quality and / or
 - that *exceed the capacity of conventional systems* (hardware and/or software) to process within an acceptable elapsed time.
- Definition of big data: Subjective & evolving
 - As technology advances over time, the size of datasets that qualify as big data will also increase.
 - The definition is varying by sector, depending on
 - ◆ what kinds of software tools are commonly available and
 - ◆ what sizes of datasets are common in a particular industry or science domain.

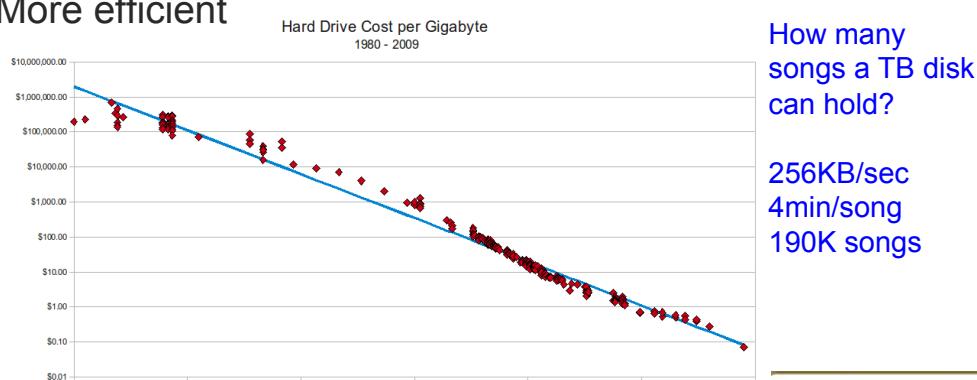


Why Big Data now: Better storage technology

- Storage & disks
 - Cheaper
 - More volume
 - Physically smaller
 - More efficient



Large data sets are
affordable



Why Big Data now: Cloud Computing: Pay per use



- Pay-as-you-go
- Elasticity
- Multi-tenancy
- Economics of Scale



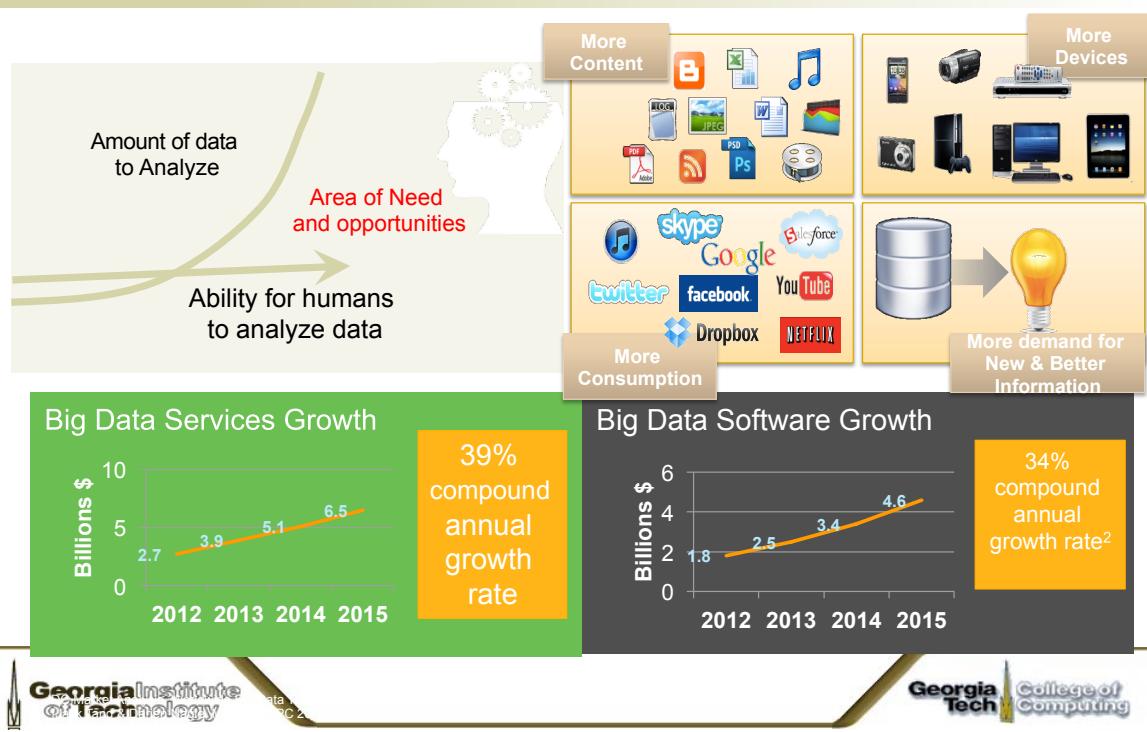
More **affordable** to perform
big data analytics



11

Georgia Tech College of Computing

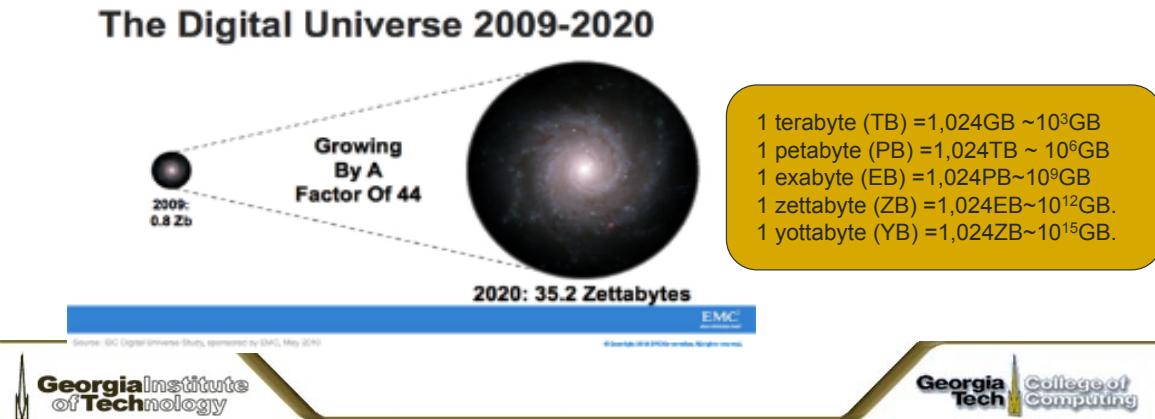
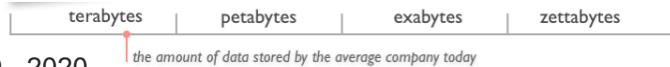
Why now → Data grows faster than intelligence



Characteristics of Big Data (1): Scale (Volume)

■ Data Volume

- 44x increase from 2009 - 2020
 - From 0.8 zettabytes to 35zb
- Data volume is increasing exponentially



Big Data: The Volume Challenge

- Challenges for **sensing, collection, generation** of more data, more meaningful data, more useful data
- Challenges for **computation** on data bigger in volume
- Challenges for **communication** on data bigger in volume
- Challenges for **analysis (algorithms)** on data bigger in volume
- How would we design **volume adaptive** big data processing **framework, systems and services?**

Characteristics of Big Data (2): Speed (Velocity)



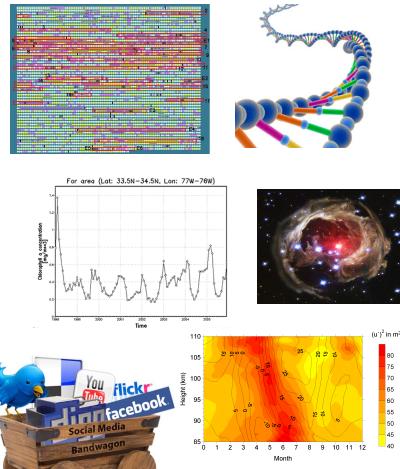
- Data are generated fast and need to be processed fast
- Online Data Analytics + Real time Analytics
- Delayed decisions → missing opportunities
- **Examples**
 - **E-Promotions:** Based on your current location, your purchase history, what you like → send promotions right now for store next to you
 - **Healthcare monitoring:** sensors monitoring your activities and body → any abnormal measurements require immediate reaction

Big Data: The Velocity Challenge

- Challenges for collecting sensor data in real time (fast rates of data generations and consumption, regarding data input sources, data output rate, varying rates to sense and collect data, the varying rates of data output)
- Challenges for analyzing sensor data in real time: how fast to analyze, and its implication/impact.
- **What would be velocity adaptive big data processing framework, systems and services?**

Characteristics of Big Data (3): Complexity (Variety)

- Input data related complexity
 - Various formats, types, and structures
 - Text, numerical, images, audio, video, sequences, time series, social media data, multi-dim arrays, etc...
 - Static data vs. streaming data
 - A single application may be generating/collecting many types of data



Characteristics of Big Data (3): Complexity (Variety)

- Algorithms related complexity
 - Different algorithms may derive different insights/outputs from the same collection of datasets
 - What analysis does make sense?
- Metrics (Quality and Performance measurement)
 - Different measures can be used to measure performance, effectiveness, and utility of big data
 - **What metrics do make sense?**

To leverage different algorithms and metrics
→ derive high quality value and data utility



Many Sources Generate Big Data



Scientific instruments
(collecting all sorts of data)



Mobile devices
tracking all objects all the time)



Sensor technology and networks
(measuring all kinds of data)

- The progress and innovation are no longer hindered by the ability to collect data; but by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion

The Evolution of Data Generation Models

- The Model of Generating/Consuming Data has Changed**

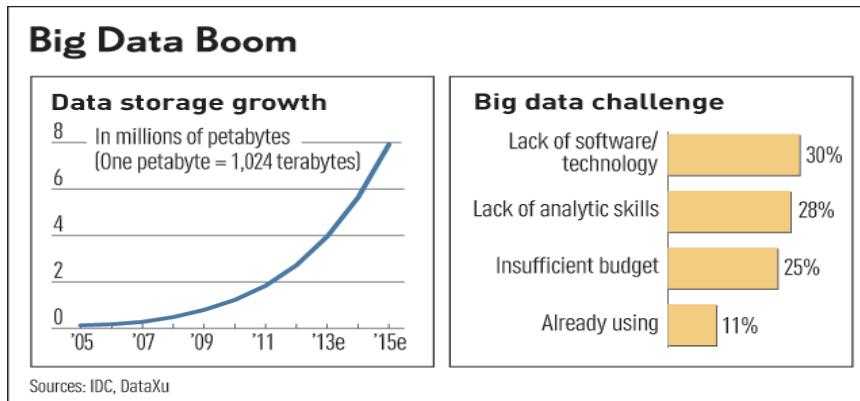
Old Model: Few companies are generating data, all others are consuming data



New Model: everyone is generating data and consuming data



Challenges in Handling Big Data



- **The Bottleneck is in technology**
 - New architecture, algorithms, techniques are needed in data collection, storage, processing, data center management, data security and privacy ...
- **Also in technical skills**
 - Experts in using the new technology and dealing with big data



21

The ‘Big Data’ Phenomenon

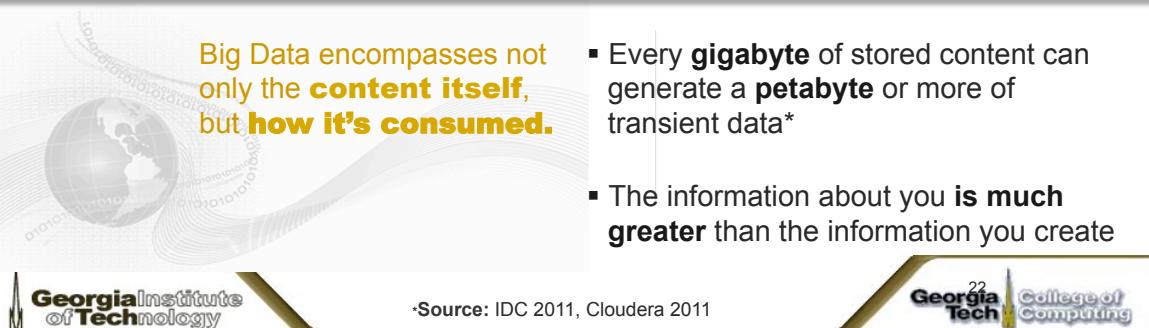
Big Data Drivers:

- The proliferation of data capture and creation technologies
- Increased “interconnectedness” drives consumption (creating more data)
- Inexpensive storage makes it possible to keep more, longer
- Innovative software and analysis tools turn data into information



Big Data encompasses not only the **content itself**, but **how it's consumed**.

- Every **gigabyte** of stored content can generate a **petabyte** or more of transient data*
- The information about you is **much greater** than the information you create



*Source: IDC 2011, Cloudera 2011

22 Georgia Tech College of Computing

Big Data Challenges

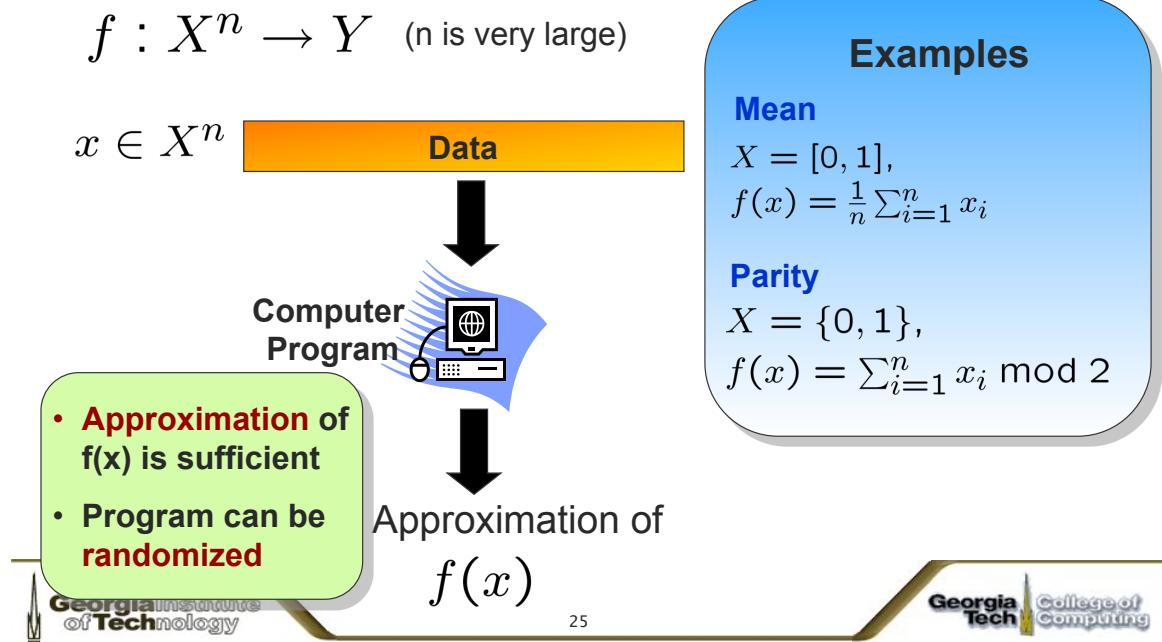
It's not just about "big"

- Cost-effectively managing the **volume, velocity and variety** of data
- Deriving value across **structured and unstructured** data
 - Flexible and effective inference/analytics/correlation analysis techniques
 - Time series, histograms, graphs
- Adapting to **context changes** and integrating **new data sources and types**
 - Seamlessly adaptive, highly available
 - Failure tolerant, highly reliable

New challenges The Scalability Dilemma

- State-of-the Art Machine Learning techniques do not scale to large data sets.
 - Memory constraints
 - Network bandwidth constraints
- Data Analytics frameworks can't handle lots of incomplete, heterogeneous, dirty data.
- Processing architectures struggle with increasing diversity of programming models and job types.
- Goal driven feature extractions demand a lot more innovations on learning and modeling techniques
- ...

Abstract model of computing



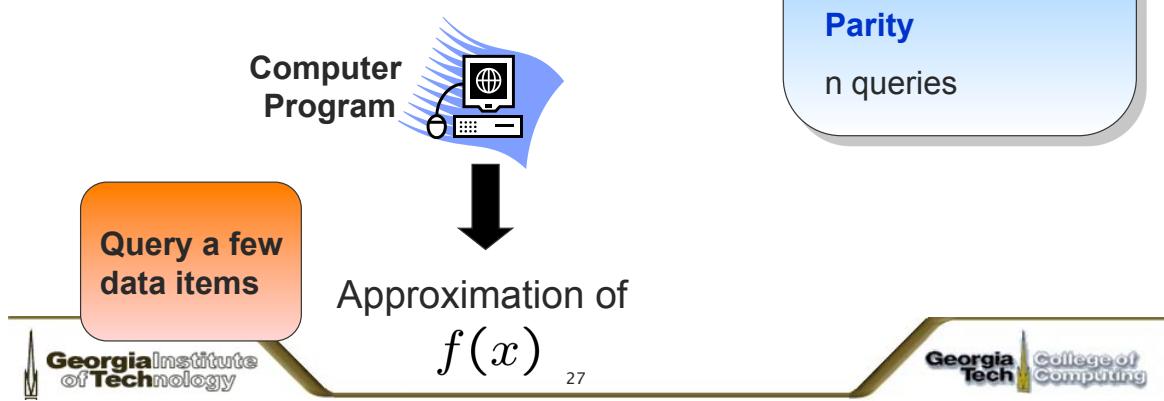
Models for Computing over Big Data

- Optimization with
 - Random sampling
 - Data Streams
 - Sketching
 - Discrete Optimization with submodular function
 - Correlation Analysis (if time permits)
- Parallel Data/Computation Partition
 - Horizontal partition (range/row partition)
 - Vertical partition (column partition)
 - Hash partition
 - Graph partition

Random Sampling

$$f : X^n \rightarrow Y \quad (\text{n is very large})$$

$x \in X^n$ Data



Examples

Mean

O(1) queries

Parity

n queries

Random Sampling

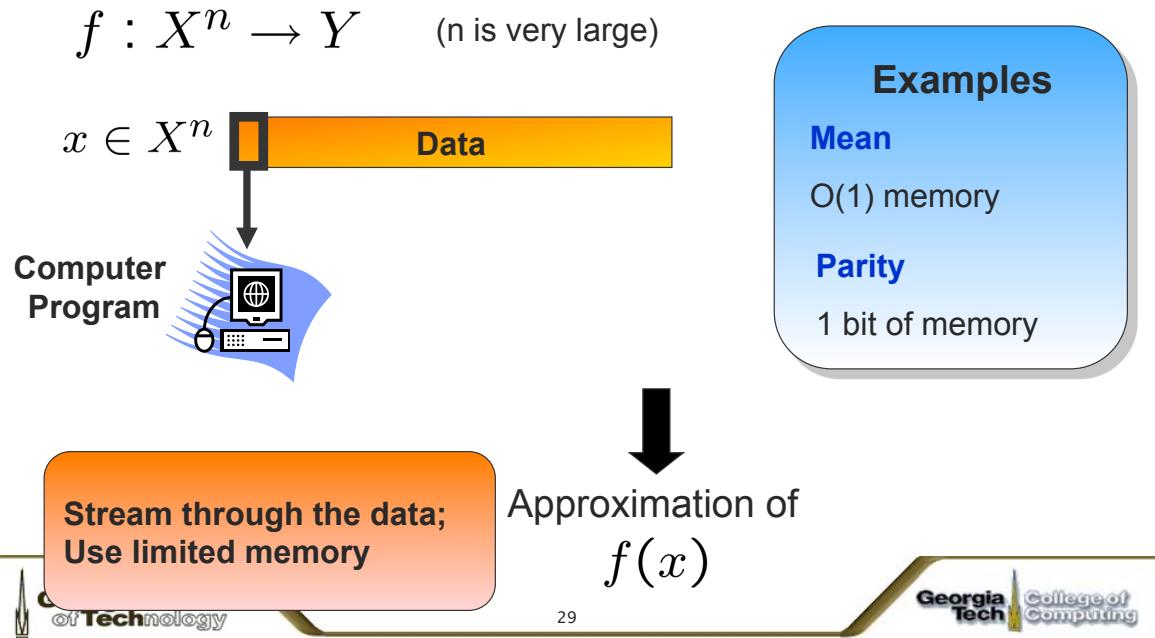
■ Advantages

- Ultra-efficient
 - ◆ Sub-linear running time & space (could even be independent of data set size)

■ Disadvantages

- May require random access
 - Doesn't fit many problems
 - Hard to sample from disorganized data sets

Data Streams



Data Streams

- Advantages
 - Sequential access
 - Limited memory
- Disadvantages
 - Running time is at least linear
 - Too restricted for some problems

Sketching

$$f : X^n \times X^n \rightarrow Y \quad (\text{n is very large})$$

Sketch1

$$x \in X^n$$

Sketch2

$$y \in X^n$$

Examples

Equality

$O(1)$ size sketch

Hamming distance

$O(1)$ size sketch

L_p distance ($p > 2$)

$\Omega(n^{1-2/p})$ size sketch

Compress each data segment into a small “sketch”

Compute over the sketches



Approximation of
 $f(x, y)$

31



Sketching

Advantages

- Appropriate for distributed data sets
- Useful for “dimension reduction”

Disadvantages

- Too restricted for some problems
- Usually, at least linear running time

Discrete Optimization

Many discrete optimization problems can be formulated as follows:

Given **finite set V** , we want to **select a subset A** (subject to some **constraints**) maximizing **utility $F(A)$** .

$$\max_{A \subseteq V} F(A)$$

Source: Guestrin, Krause@IJCAI 2009



Opportunities for Big Data Analytics



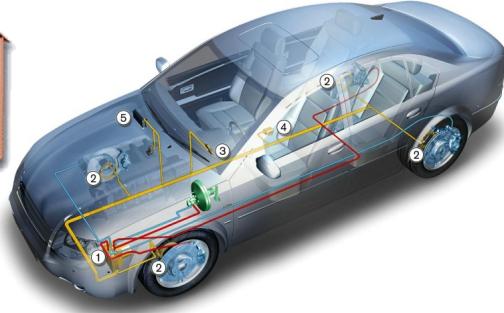
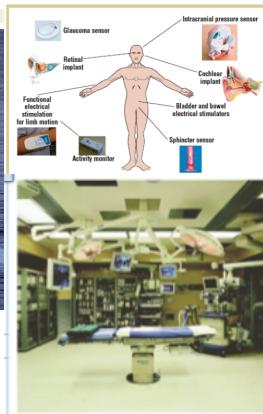
reliance on non-renewable resources



pollution



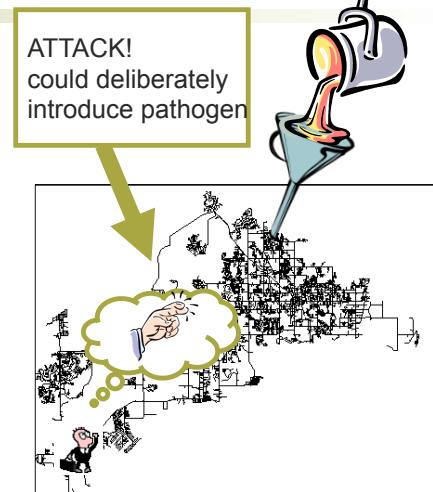
global warming



Georgia Institute of Technology

Water distribution networks

- Water distribution in a city → very complex system
- Pathogens in water can affect thousands (or millions) of people
- Currently: Add chlorine to the source and hope for the best



Source: Guestrin, Krause@IJCAI 2009



Sensor placement

Given: finite set V of locations, sensing quality F

Want: $\mathcal{A}^* \subseteq V$ such that

$$\mathcal{A}^* = \operatorname{argmax}_{|\mathcal{A}| \leq k} F(\mathcal{A})$$

NP-hard!

Greedy algorithm:

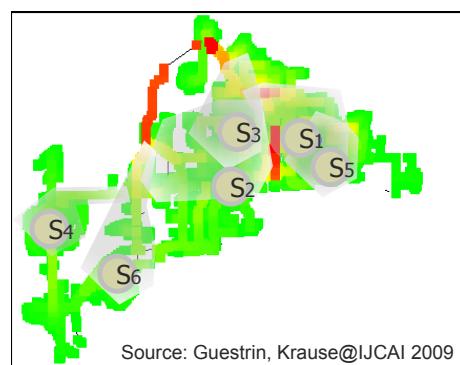
Start with $A = \{\}$

For $i = 1$ to k

$$s^* := \operatorname{argmax}_s F(A \cup \{s\})$$

$A := A \cup \{s^*\}$

M

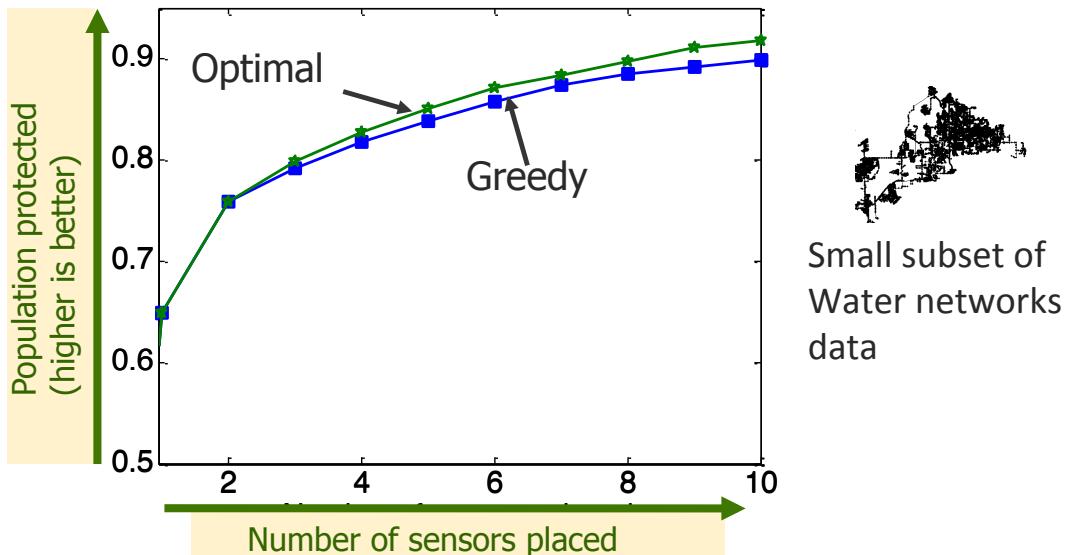


Source: Guestrin, Krause@IJCAI 2009

How well can this simple heuristic do?



Performance of greedy algorithm



Greedy score empirically close to optimal. Why?



Source: Guestrin, Krause@IJCAI 2009

37



One reason submodularity is useful

Theorem [Nemhauser et al '78]

Suppose F is *monotonic* and *submodular*. Then greedy algorithm gives constant factor approximation:

$$F(A_{\text{greedy}}) \geq \underbrace{(1 - 1/e)}_{\sim 63\%} \max_{|A| \leq k} F(A)$$

- Greedy algorithm gives near-optimal solution!
- In general, guarantees best possible unless $P = NP$!

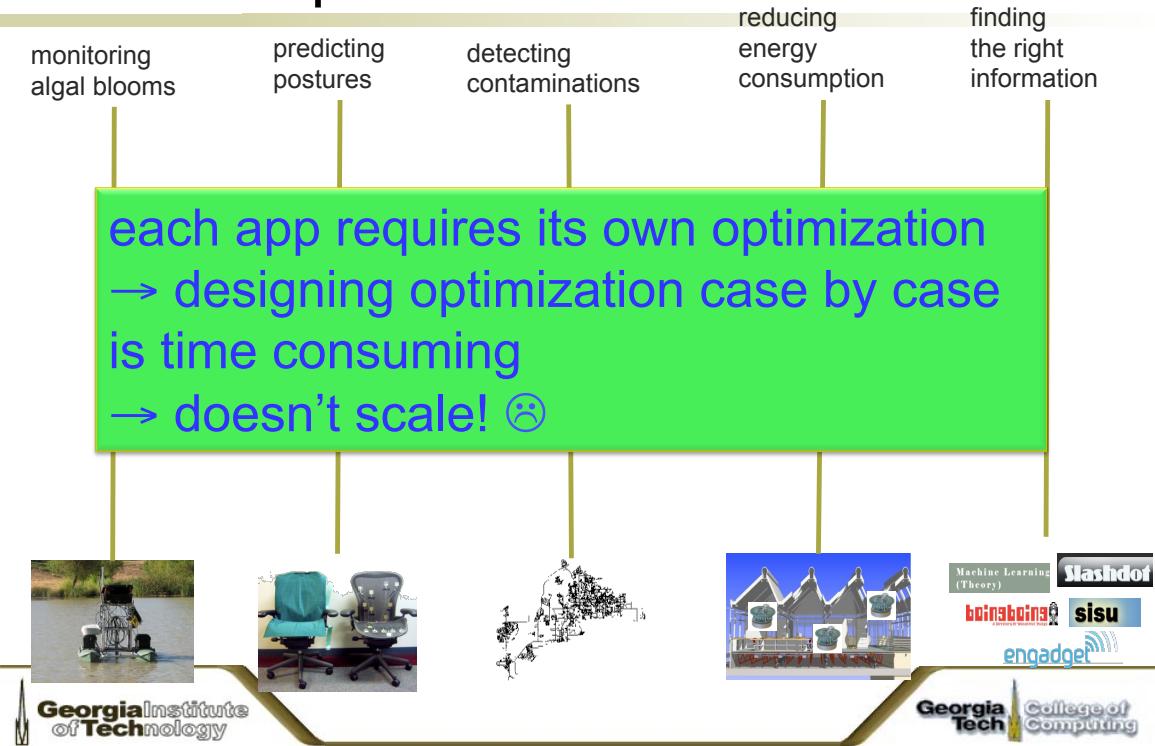


Source: Guestrin, Krause@IJCAI 2009

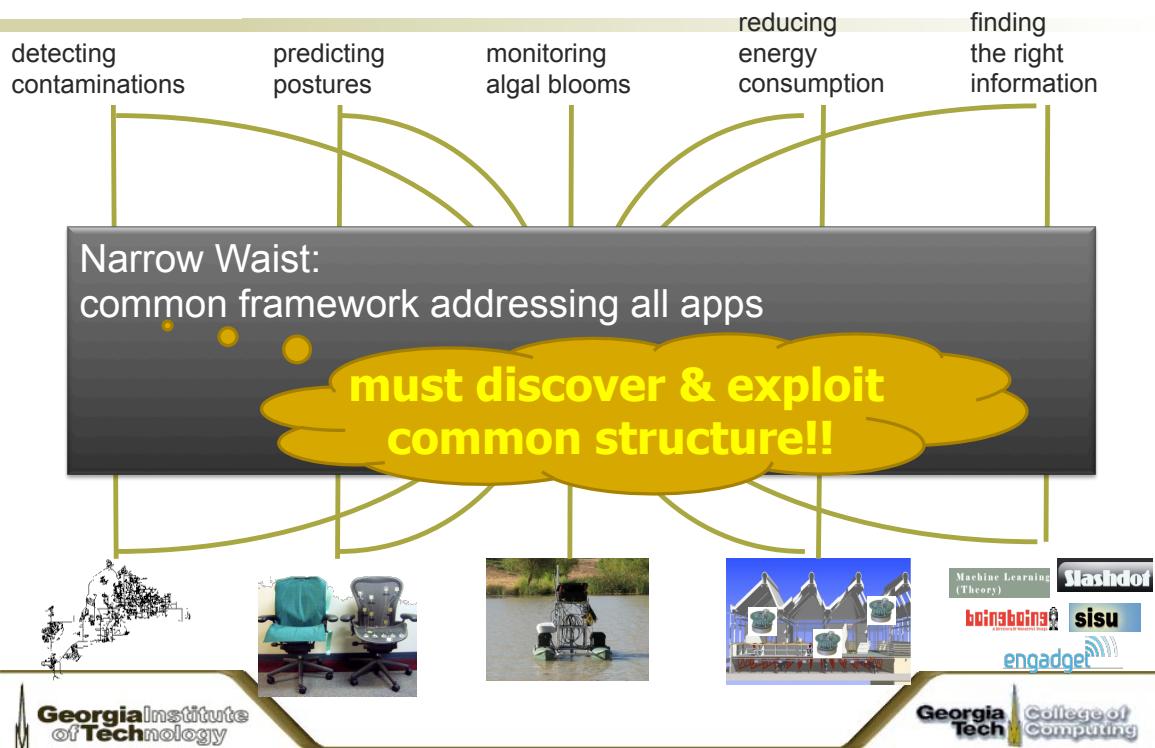
38



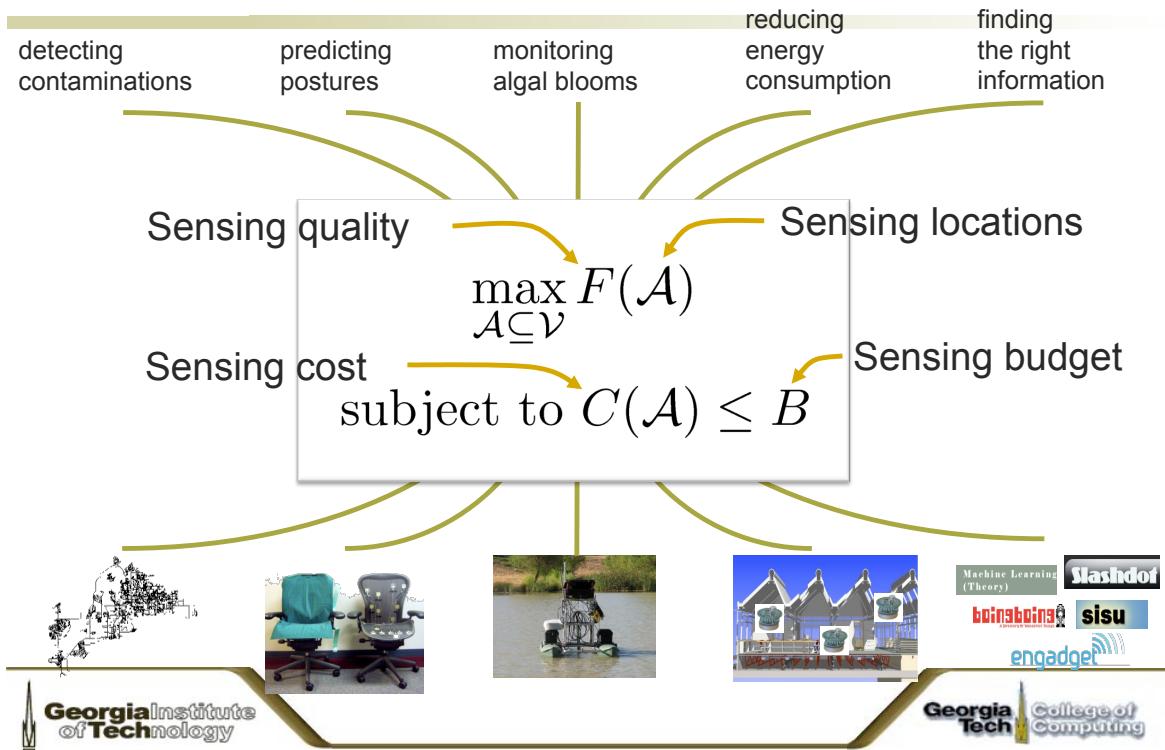
Many big data applications use discrete optimization...



The quest for the optimization narrow waist



The quest for the optimization narrow waist



More Examples

- **Finding top k most influential people in a large social network**
- **Finding top k most interesting blogs from sea of blogs**
- **Finding top k best matches of keyword search**

Social Influence: Problem Setting

- Given
 - a limited budget B for initial advertising (e.g. give away free samples of product)
 - estimates for influence between individuals
- Goal
 - trigger a large cascade of influence (e.g. further adoptions of a product)
- Question
 - Which set of individuals should B target at?
- Application besides product marketing
 - spread an innovation
 - detect stories in blogs



Keyword search is not enough

still too much information

Google scholar "online learning" Search Adv Sch

Scholar How to find my personalized top responses necessarily need

book Facilitated C G Collison, E ... ED448684 Facilitating Cited by 337 - Related articles - Cached - All 4 versions

[pdf] ► Building learning communities in cyberspace RM Palloff, K Pratt - 1999 - online2.org ... By Rena Palloff and Keith Pratt. Online learning: ... Some questions involving online learning ??? How do we know when a student is engaged with the subject matter? ... Cited by 1760 - Related articles - View as HTML - BL Direct - All 7 versions

Adaptive online learning algorithms for blind separation: maximum entropy and minimum ... ► kfupm.edu.sa [PDF] HH Yang, S Amari - Neural computation, 1997 - MIT Press ... Adaptive Online Learning Algorithms for Blind Separation: Maximum Entropy and Minimum Mutual Information ... Page 3. Adaptive Online Learning Algorithms 1459 ... Cited by 285 - Related articles - BL Direct - All 9 versions

book Teaching & learning online: Pedagogies for new technologies J Stephenson - 2001 - books.google.com ... when online 3 Shirley Alexander and David Bond 2 Learning technology and learning relationships 16 Terry Mayes 3 Problems with online learning are systemic ... Cited by 159 - Related articles - All 3 versions

[pdf] ► Examining social presence in online courses in relation to students' perceived learning and ...

doesn't highlight any structure in the space

Georgia Tech College of Computing

Approximation Optimization: Analytics of the Web Data

77% read blogs
184M blogs
[Universal McCann, '08]

TIME
IN PARTNERSHIP WITH CNN
Thursday, Nov. 20, 2008
How Many Blogs Does the World Need?
By Michael Kinsley

Slashdot **engadg** **boingboing**

CELEBRITY DOG WATCHER.COM
Dug up celebrity dog news

you have 10 minutes a day to read blogs / news

which of the million blogs should you read?

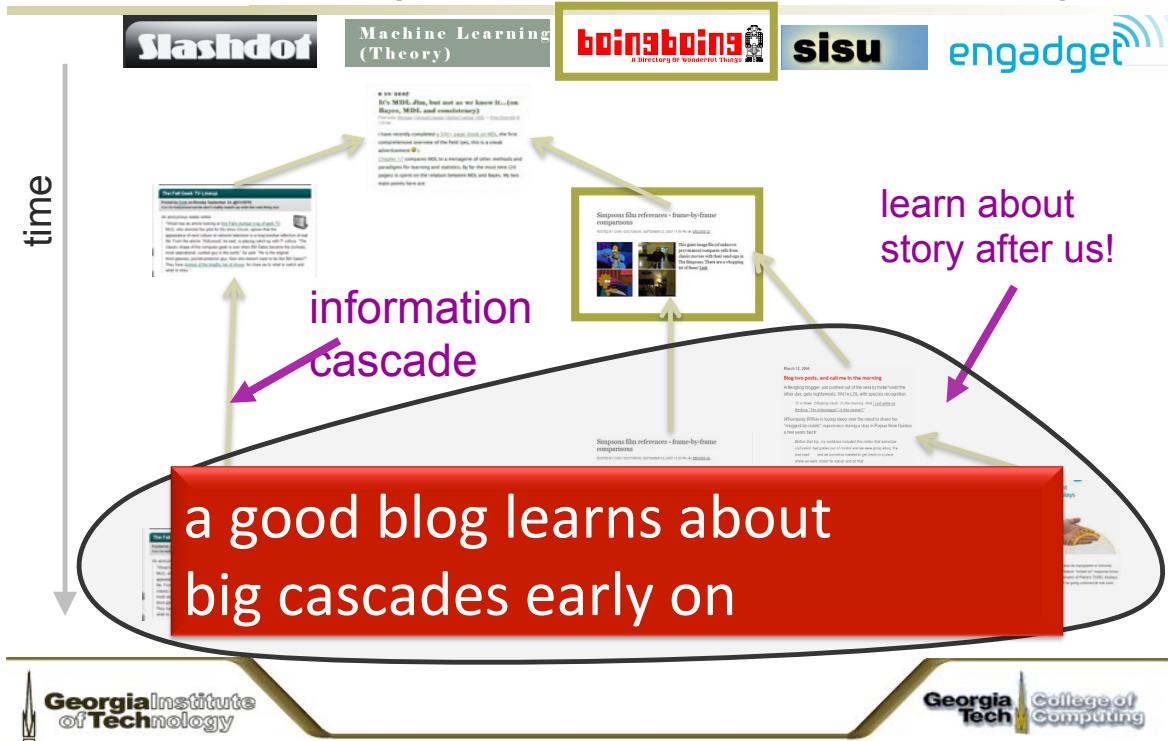
information overload!!! 😞

Homer Simpson

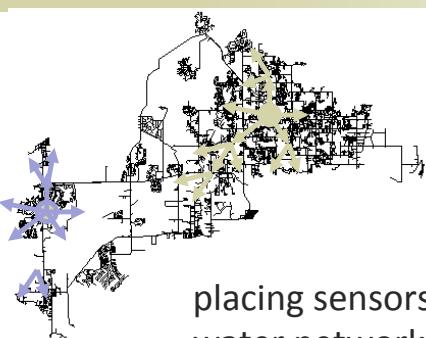
Georgia Institute of Technology

Information Cascades

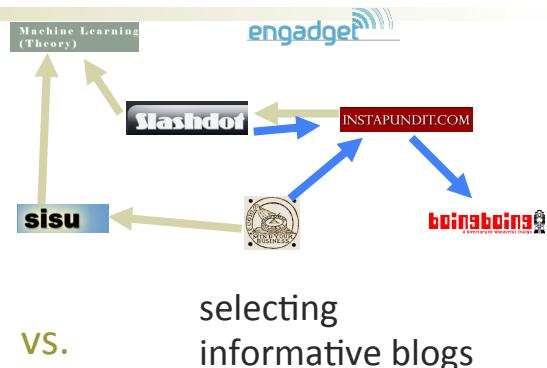
[Leskovec, Krause, G., Faloutsos, VanBriesen, Glance '07]



Water vs. Web

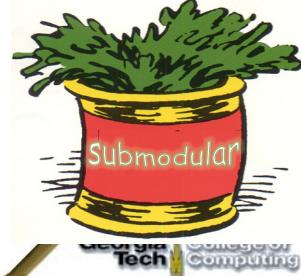


placing sensors in
water networks



want to pick nodes to **detect big cascades early**

in both apps, utility
functions submodular ☺



Georgia Institute
of Technology

Georgia Tech College of Computing

The power of the efficient narrow waist

constrained maximization of submodular functions

placing k sensors

robust sensing

complex
constraints

sequential sensing

finding needles
in haystack

better
generalization

better use of
“attention”

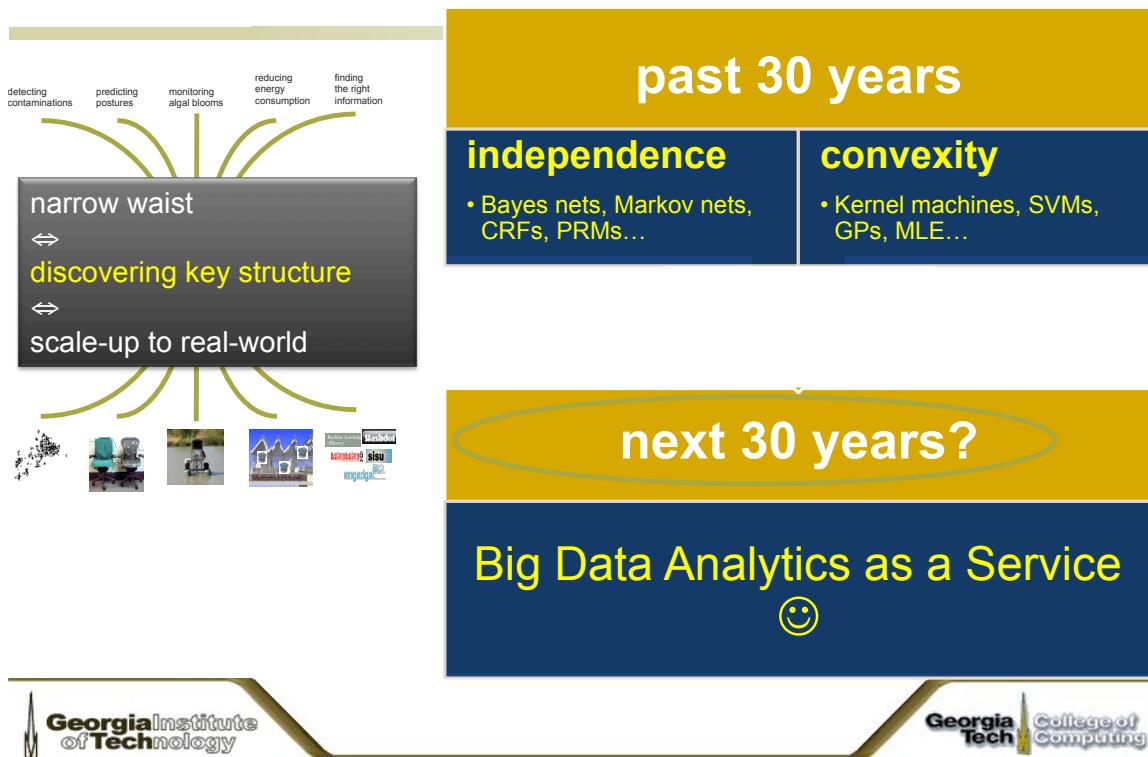
personalization



same algorithms
new insights



Structural insights: challenges of next decade



The basic foundations of data analytics
are changing

**new computing
architectures**

**new notions of
“sensing”**

unique new challenges → new structure
new intelligence

- distributed algorithms for learning over massive datasets on huge computer clusters
- parallel probabilistic inference
- Correlation analysis and inference of multiple time series
- ...

Data Analysis as a Service: Open Issues



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."

- **Privacy Protection** at both Individual and Corporate/Enterprise Level (handling the risk of unwanted privacy intrusion)
- **Noise / Quality Control** on the input, during the computation and on the output phases of data analytics (handling the risk of misuse of BIG DATA)

Questions

