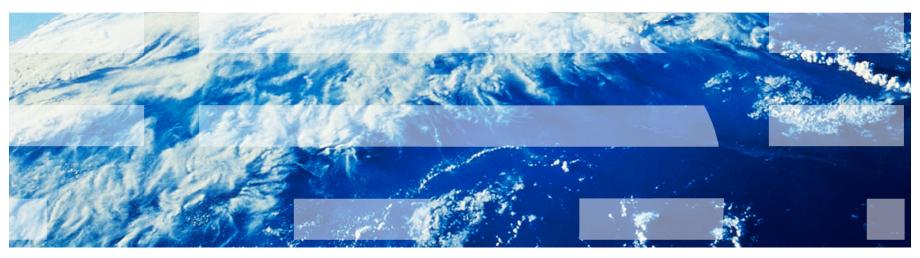# EECS E6893 Big Data Analytics Lecture 1:

## *Overview of Big Data Analytics*

**Ching-Yung Lin**, Ph.D.

Adjunct Professor, Depts. of Electrical Engineering and Computer Science

IEEE Fellow

September 6th, 2024

# Definition and Characteristics of Big Data

"*Big data is high-**volume**, high-**velocity** and high-**variety** information assets that demand **cost-effective**, **innovative** forms of information processing for **enhanced insight and decision making**.*"  -- Gartner

which was derived from:

"*While enterprises struggle to consolidate systems and collapse redundant databases to enable greater operational, analytical, and collaborative consistencies, changing economic conditions have made this job more difficult. E-commerce, in particular, has exploded data management challenges along three dimensions: **volumes, velocity and variety**. In 2001/02, IT organizations much compile a variety of approaches to have at their disposal for dealing each.*" – Doug Laney
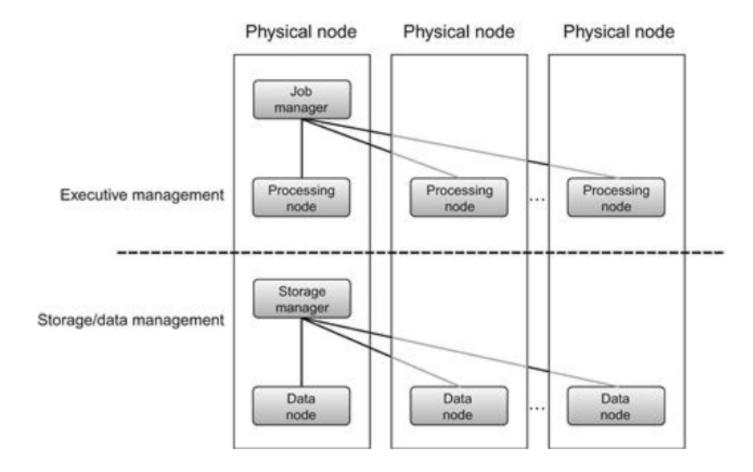
# What made Big Data needed?

"Big Data Analytics", David Loshin

- Processing capability: CPU, processor, or node.
- Memory
- Storage
- Network



Physical node    Physical node    Physical node

Job manager

Executive management    Processing node    Processing node   ...   Processing node

Storage manager

Storage/data management    Data node    Data node   ...   Data node

"Big Data Analytics", David Loshin

# Scalability — Scale Up & Scale Out

- Scale out
  - Use more resources to distribute workload in parallel
  - Higher data access latency is typically incurred
- Scale up
  - Efficiently use the resources
  - Architecture-aware algorithm design



Example: Resource utilization for a large production cluster at Twitter data center



www.stanford.edu/~cdel/2014.asplos.quasar.pdf

- For independent data ==> scale up may not have obvious advantage than scale out
- For linked data ==> utilizing scale up as much as possible before scale out

| Aspect | Typical Scenario | Big Data |
|---|---|---|
| Application development | Applications that take advantage of massive parallelism developed by specialized developers skilled in high-performance computing, performance optimization, and code tuning | A simplified application execution model encompassing a distributed file system, application programming model, distributed database, and program scheduling is packaged within Hadoop, an open source framework for reliable, scalable, distributed, and parallel computing |
| Platform | Uses high-cost massively parallel processing (MPP) computers, utilizing high-bandwidth networks, and massive I/O devices | Innovative methods of creating scalable and yet elastic virtualized platforms take advantage of clusters of commodity hardware components (either cycle harvesting from local resources or through cloud-based utility computing services) coupled with open source tools and technology |
| Data management | Limited to file-based or relational database management systems (RDBMS) using standard row-oriented data layouts | Alternate models for data management (often referred to as NoSQL or "Not Only SQL") provide a variety of methods for managing information to best suit specific business process needs, such as in-memory data management (for rapid access), columnar layouts to speed query response, and graph databases (for social network analytics) |
| Resources | Requires large capital investment in purchasing high-end hardware to be installed and managed in-house | The ability to deploy systems like Hadoop on virtualized platforms allows small and medium businesses to utilize cloud-based environments that, from both a cost accounting and a practical perspective, are much friendlier to the bottom line |

"Big Data Analytics", David Loshin

# Techniques towards Big Data

- Massive Parallelism
- Huge Data Volumes Storage
- Data Distribution
- High-Speed Networks
- High-Performance Computing
- Task and Thread Management
- Data Mining and Analytics
- Data Retrieval
- Machine Learning
- Data Visualization

➔ Techniques exist for years to decades. Why is Big Data hot?

E6893 Big Data Analytics – Lecture 1: Overview

© 2024 CY Lin, Columbia University

- More data are being collected and stored
- Open source code
- Commodity hardware / Cloud

- More data are being collected and stored
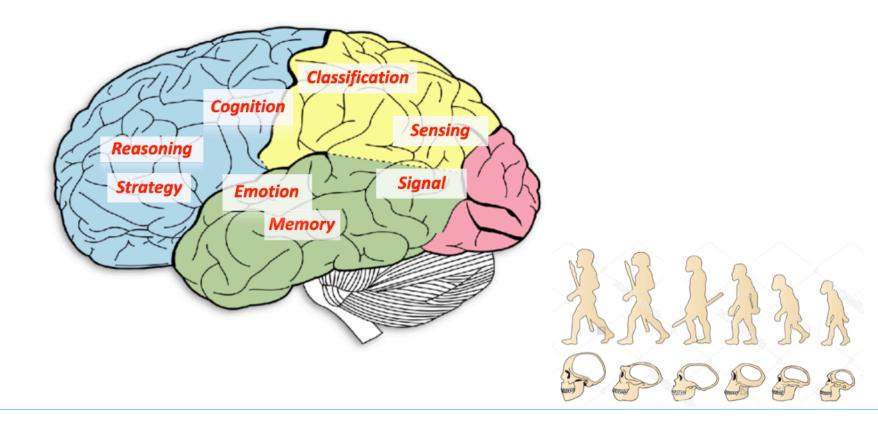- Open source code
- Commodity hardware / Cloud

- <span style="color:green">High-Volume</span>
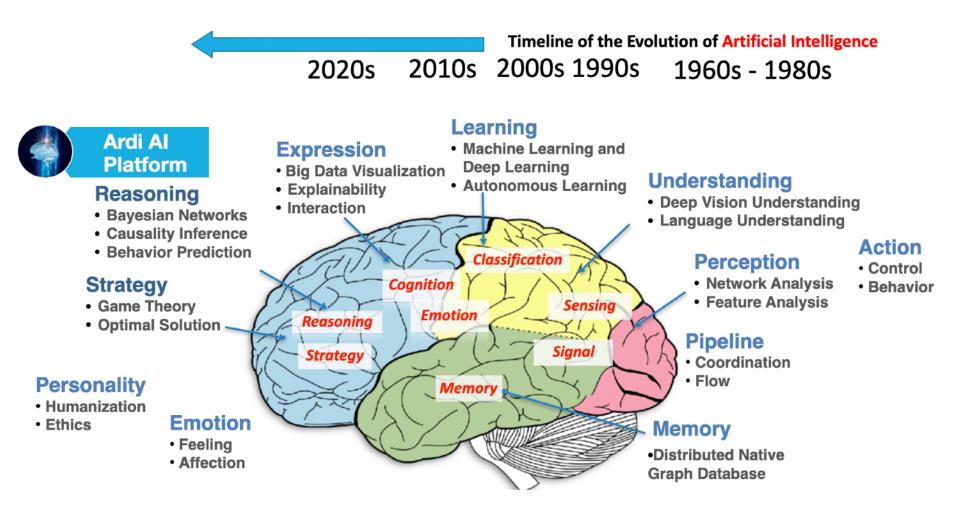➜ - <span style="color:orange">High-Velocity</span>
- <span style="color:red">High-Variety</span>

➜ **Artificial Intelligence**

E6893 Big Data Analytics – Lecture 1: Overview

# Evolution of Intelligence



**Direction of the Evolution of Intelligence**

Classification

Cognition

Sensing

Reasoning

Strategy

Emotion

Signal

Memory

E6893 Big Data Analytics – Lecture 1: Overview

# Evolution of Artificial Intelligence is similar, but much faster



Timeline of the Evolution of **Artificial Intelligence**

2020s    2010s    2000s 1990s    1960s - 1980s

**Ardi AI Platform**

**Reasoning**
- Bayesian Networks
- Causality Inference
- Behavior Prediction

**Strategy**
- Game Theory
- Optimal Solution

**Personality**
- Humanization
- Ethics

**Emotion**
- Feeling
- Affection

**Expression**
- Big Data Visualization
- Explainability
- Interaction

**Learning**
- Machine Learning and Deep Learning
- Autonomous Learning

**Understanding**
- Deep Vision Understanding
- Language Understanding

**Action**
- Control
- Behavior

**Perception**
- Network Analysis
- Feature Analysis

**Pipeline**
- Coordination
- Flow

**Memory**
- Distributed Native Graph Database

*Classification*
*Cognition*
*Emotion*
*Reasoning*
*Strategy*
*Sensing*
*Signal*
*Memory*

https://www.graphen.ai/products/ardi.html

E6893 Big Data Analytics – Lecture 1: Overview

© 2024 CY Lin, Columbia University

# Course Outline

| Class Date | Class Number | Topics Covered |
|------------|--------------|----------------|
| 09/06/24 | 1 | Introduction to Big Data Analytics |
| 09/13/24 | 2 | Big Data Platforms & Algorithms |
| 09/20/24 | 3 | Real-Time Stream Analysis |
| 09/27/24 | 4 | Linked Big Data Analysis |
| 10/04/24 | 5 | End-to-End System Workflow |
| 10/11/24 | 6 | Big Data Visualization |
| 10/18/24 | 7 | Large Language Models |
| 10/25/24 | 8 | GPU-Based Massive Data Analysis |
| 11/01/24 | 9 | GPU-Accelerated Machine Learning |
| 11/08/24 | 10 | *Final Project Proposal Presentation* |
| 11/15/24 | 11 | AI Finance Applications |
| 11/22/24 | 12 | *Final Project Progress Presentation* |
| 11/29/24 |  | *Thanksgiving Holiday* |
| 12/06/24 | 13 | AI Medical Applications |
| 12/13/24 | 14 | *Big Data Analytics Workshop* |

E6893 Big Data Analytics – Lecture 1: Overview

The Apache™ Hadoop® project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

The project includes these modules:
- **Hadoop Common**: The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™)**: A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN**: A framework for job scheduling and cluster resource management.
- **Hadoop MapReduce**: A YARN-based system for parallel processing of large data sets.

http://hadoop.apache.org

# Four distinctive layers of Hadoop

APACHE **Spark** ™

*Lightning-fast unified analytics engine*

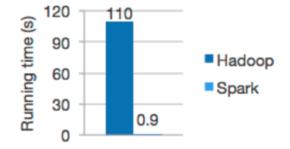| Download | Libraries ▾ | Documentation ▾ | Examples | Community ▾ | Developers ▾ |

**Apache Spark™** is a unified analytics engine for large-scale data processing.
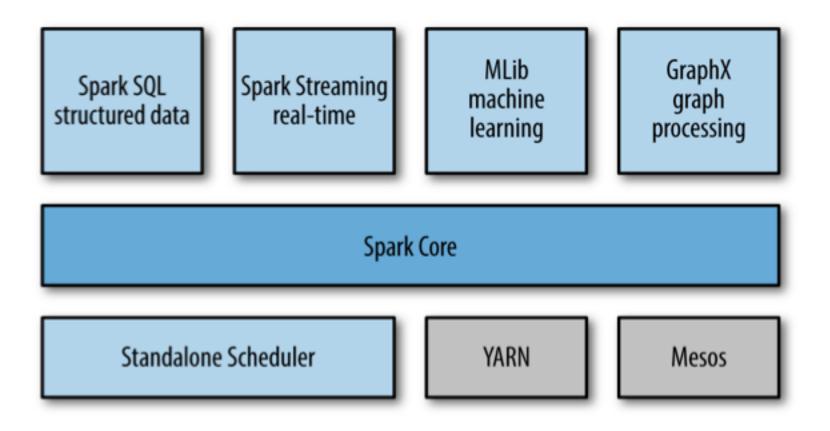
## Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.



Logistic regression in Hadoop and Spark

E6893 Big Data Analytics – Lecture 1: Overview

© 2024 CY Lin, Columbia University

# Main Spark Stack

E6893 Big Data Analytics – Lecture 2: Big Data Platform
**© 2024 CY Lin, Columbia University**

# Course Main Thrust 3: Streaming and Linked Big Data Analytics

# Course Main Thrust 4: Workflow and Analytics Pipeline



- A scheduler, which handles both triggering scheduled workflows, and submitting Tasks to the executor to run.
- An executor, which handles running tasks. In the default Airflow installation, this runs everything *inside* the scheduler, but most production-suitable executors actually push task execution out to *workers*.
- A *webserver*, which presents a handy user interface to inspect, trigger and debug the behaviour of DAGs and tasks.
- A folder of *DAG files*, read by the scheduler and executor (and any workers the executor has)
- A *metadata database*, used by the scheduler, executor and webserver to store state.
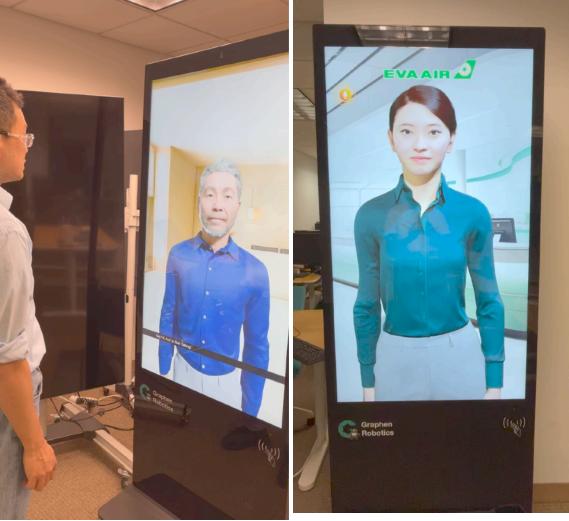
# Course Main Thrust 5: Big Data Visualization

E6893 Big Data Analytics – Lecture 1: Overview

© 2024 CY Lin, Columbia University

# Course Main Thrust 6: Generative AI and Large Language Model



Graphen Aiia — World's First AI Digital Human for Daily Life

...a Analytics – Lecture 1: Overview

# Course Main Thrust 7: GPU-Based Big Data Analysis



Data analytics workflows have traditionally been slow and cumbersome, relying on CPU compute for data preparation, training, and deployment. Accelerated data science can dramatically boost the performance of end-to-end analytics workflows, speeding up value generation while reducing cost.
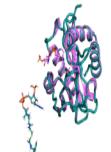
# Course Main Thrust 8: Big Data AI Solutions

- **Big Data and AI for Finance**
- **Big Data and AI for Healthcare**

Graphen Small Mole Drug Dev ➜ 1/27 of the Time; 1/9000 of the Cost, comparing to traditional methods



*"Tools from established companies like **Google**, startups like **Graphen**, and AI chipsets from vendors like **NVIDIA** and **Intel** will help accelerate the speed of drug discovery, development, and testing, allowing pharmaceutical companies and healthcare authorities to combat the pandemic." – ABI research, May 2020*

- **Key Differentiator of this class:** Focusing on building a full-spectrum understanding of the latest Big Data Analytics technologies and using them to build real industry real-world solutions.

- **Sapphire Big Data Analytics Open Source Applications:** Create a Big Data open source toolsets for various industries (and disciplines)



- **Dataset and Use Cases:** Welcome!!

# Course Outline

| Class Date | Class Number | Topics Covered |
|---|---|---|
| 09/06/24 | 1 | Introduction to Big Data Analytics |
| 09/13/24 | 2 | Big Data Platforms & Algorithms |
| 09/20/24 | 3 | Real-Time Stream Analysis |
| 09/27/24 | 4 | Linked Big Data Analysis |
| 10/04/24 | 5 | End-to-End System Workflow |
| 10/11/24 | 6 | Big Data Visualization |
| 10/18/24 | 7 | Large Language Models |
| 10/25/24 | 8 | GPU-Based Massive Data Analysis |
| 11/01/24 | 9 | GPU-Accelerated Machine Learning |
| 11/08/24 | 10 | *Final Project Proposal Presentation* |
| 11/15/24 | 11 | AI Finance Applications |
| 11/22/24 | 12 | *Final Project Progress Presentation* |
| 11/29/24 | | *Thanksgiving Holiday* |
| 12/06/24 | 13 | AI Medical Applications |
| 12/13/24 | 14 | *Big Data Analytics Workshop* |

E6893 Big Data Analytics – Lecture 1: Overview

# Course Information

- Website:

    http://www.ee.columbia.edu/~cylin/course/bigdata/

- Textbook:

    -- None, but reference book(s) and/or articles/papers will be provided each lecture.

## Big Data Analytics

From Strategic Planning to Enterprise Integration with Tools, Techniques, NoSQL, and Graph

MK MORGAN KAUFMANN

David Loshin

Chapter 1: Market and Business Drivers for Big Data Analysis

Chapter 2: Business Problems Suited to Big Data Analytics

Chapter 3: Achieving Organizational Alignment for Big Data Analytics

Chapter 4: Developing a Strategy for Integrating Big Data Analytics into the Enterprise

Chapter 5: Data Governance for Big Data Analytics: Considerations for Data Policies and Processes

Chapter 6: Introduction to High-Performance Appliances for Big Data Management

Chapter 7: Big Data Tools and Techniques

Chapter 8: Developing Big Data Applications

Chapter 9: NoSQL Data Management for Big Data

Chapter 10: Using Graph Analytics for Big Data

Chapter 11: Developing the Big Data Roadmap

# Highly Recommended Reference Book for Lectures 1-4

Isaac Triguero and Mikel Galar

**Large-Scale Data Analytics with Python and Spark**

A Hands-on Guide to Implementing Machine Learning Solutions

**Part I: Understanding and Dealing with Big Data**

Chapter 1: Introduction

Chapter 2: MapReduce

**Part II: Big Data Frameworks**

Chapter 3: Hadoop

Chapter 4: Spark

Chapter 5: Spark SQL and DataFrames

**Part III: Machine Learning for Big Data**

Chapter 6: Machine Learning with Spark

Chapter 7: Machine Learning for Big Data

Chapter 8: Implementing Classic Methods: k-Means and Linear Regression

Chapter 9: Advanced Examples — Semi-Supervised, Ensembles, Deep Learning Model Deployment

# Course Grading

- **5 Homeworks: 50%**

  -- **Individual work**; Language Requirement: Python, JavaScript; Get familiar with Linux

  -- **Report (including description of the work, discussions, experiments, etc) and source code**

  - **HW #0: Big Data Environment Setup and Testing**

  - **HW #1: Analytic Algorithms and System Monitoring**

  - **HW #2: Graph Analysis and Analytics Pipeline**

  - **HW #3: Big Data Visualization and LLM**

  - **HW #4: GPU-Based Big Data Analysis**

- **Final Project: 50%**

    **-- Teamwork: 2 - 3 students per team (on campus); 1 - 3 students per team for CVN**

    - **Proposal** (slides — short presentation in the class)

    - **Progress Presentation** (slides — short presentation in the class)

    - **Progress Report** (report)

    - **Final Report** (paper, up to 10 pages)

    - **Workshop Presentation** (Oral and Demo)

    - **Open Source Codes**

    - **Video Presentation** (on YouTube)

# Assignments and Submissions

| Class Date | Class Number | Assignment | Due |
|---|---|---|---|
| 09/06/24 | 1 | HW #0 Big Data Environment Setup and Testing [assignment][tutorial] | |
| 09/13/24 | 2 | HW #0 Tutorial II | |
| 09/20/24 | 3 | HW #1 Analytics Algorithms and Monitoring [assignment][tutorial] | HW #0 |
| 09/27/24 | 4 | HW #1 Tutorial II | |
| 10/04/24 | 5 | HW #2 Graph Analysis and Analytics Pipeline [assignment][tutorial] | HW #1 |
| 10/11/24 | 6 | HW #2 Tutorial II | |
| 10/18/24 | 7 | HW #3 Big Data Visualization and LLM [assignment][tutorial] | HW #2 |
| 10/25/24 | 8 | HW#3 Tutorial II | |
| 11/01/24 | 9 | HW #4 GPU-based Big Data Analysis [assignment][tutorial] | HW #3 |
| 11/08/24 | 10 | | Proposal Slides |
| 11/15/24 | 11 | HW#4 Tutorial II | |
| 11/22/24 | 12 | | HW #4 |
| 11/29/24 | | | |
| 12/06/24 | 13 | | Progress Report |
| 12/13/24 | 14 | | Final Project Materials |

# Other Issues

- Professor Lin:
  - Office Hours:

    By appointment

  - Contact: c.lin@columbia.edu

- TA:

  - Apurva Patel (amp2365):  Monday 10am-12pm (onsite) and Wednesday 12pm-2pm (online)
  - Linyang He (lh3288): Tuesday 3pm-5pm (onsite) and Thursday 1pm-3pm (online)

  - Location: EE Student Lounge (next to the EE office, Mudd 13th Floor)

# 5 Example Big Data Use Case Categories



**Big Data Exploration**
Find, visualize, understand all big data to improve decision making



**Enhanced 360º View of the Customer**
Extend existing customer views (MDM, CRM, etc) by incorporating additional internal and external information sources



**Security/Intelligence Extension**
Lower risk, detect fraud and monitor cyber security in real-time



**Operations Analysis**
Analyze a variety of machine data for improved business results



**Data Warehouse Augmentation**
Integrate big data and data warehouse capabilities to increase operational efficiency

# Big Data Examples -- Application Use Cases

1. Expertise Location
2. Recommendation
3. Commerce
4. Financial Analysis
5. Social Media Monitoring
6. Telco Customer Analysis
7. Healthcare Analysis
8. Data Exploration and Visualization
9. Personalized Search
10. Anomaly Detection
11. Fraud Detection
12. Cybersecurity
13. Sensor Monitoring (Smarter another Planet)
14. Cellular Network Monitoring
15. Cloud Monitoring
16. Code Life Cycle Management
17. Traffic Navigation
18. Image and Video Semantic Understanding
19. Genomic Medicine
20. Brain Network Analysis
21. Data Curation
22. Near Earth Object Analysis

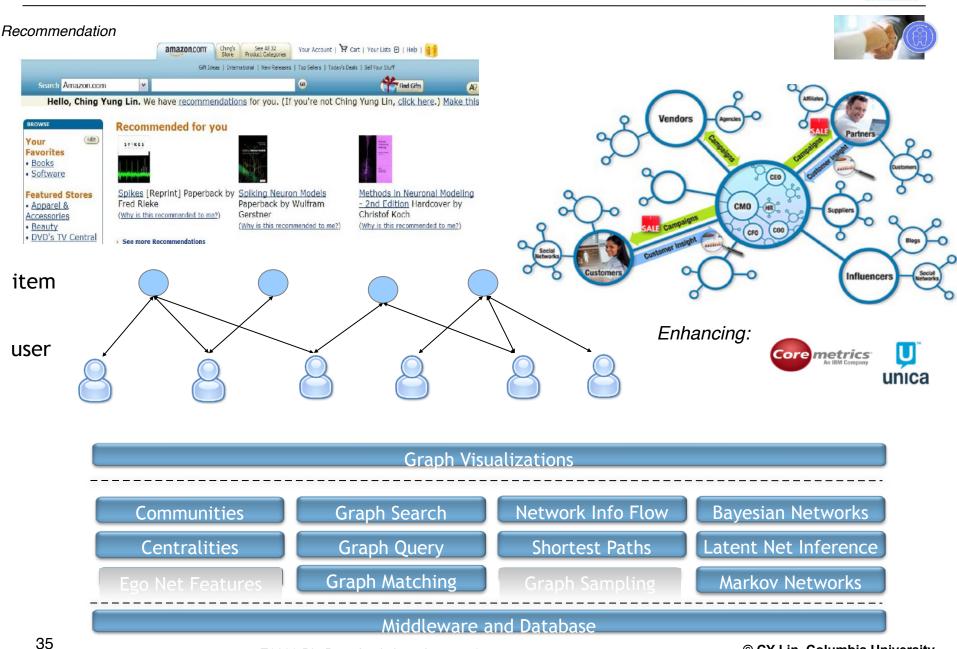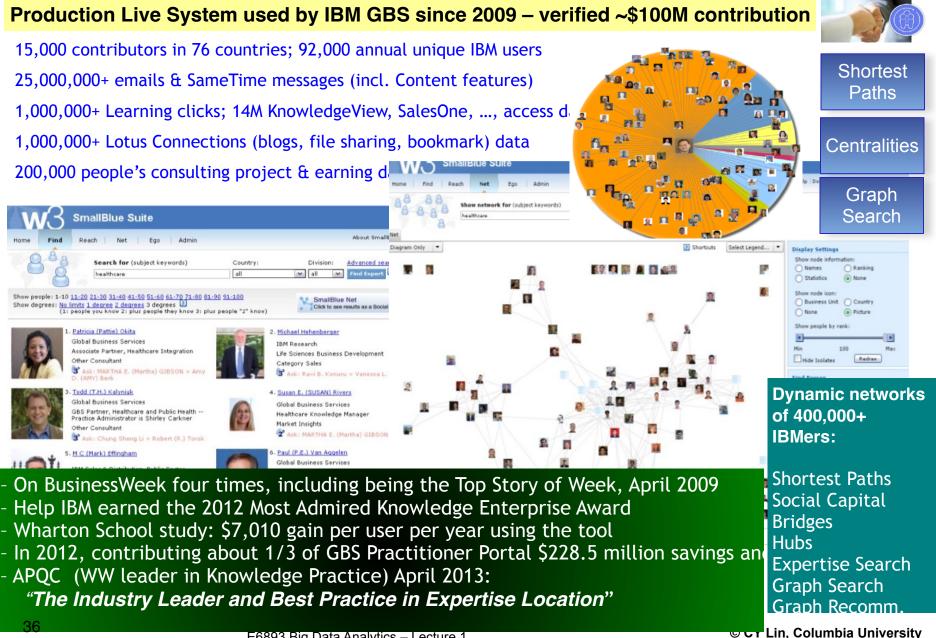# Category 1: 360° View

*Recommendation*

item

user

*Enhancing:*

| Graph Visualizations | | | |
|---|---|---|---|
| Communities | Graph Search | Network Info Flow | Bayesian Networks |
| Centralities | Graph Query | Shortest Paths | Latent Net Inference |
| Ego Net Features | Graph Matching | Graph Sampling | Markov Networks |
| Middleware and Database | | | |

E6893 Big Data Analytics – Lecture 1

# Use Case 1: Social Network Analysis in Enterprise for Productivity

**Production Live System used by IBM GBS since 2009 – verified ~$100M contribution**

15,000 contributors in 76 countries; 92,000 annual unique IBM users

25,000,000+ emails & SameTime messages (incl. Content features)

1,000,000+ Learning clicks; 14M KnowledgeView, SalesOne, …, access d

1,000,000+ Lotus Connections (blogs, file sharing, bookmark) data

200,000 people's consulting project & earning d

Shortest Paths

Centralities

Graph Search



**Dynamic networks of 400,000+ IBMers:**

Shortest Paths
Social Capital
Bridges
Hubs
Expertise Search
Graph Search
Graph Recomm.

– On BusinessWeek four times, including being the Top Story of Week, April 2009
– Help IBM earned the 2012 Most Admired Knowledge Enterprise Award
– Wharton School study: $7,010 gain per user per year using the tool
– In 2012, contributing about 1/3 of GBS Practitioner Portal $228.5 million savings an
– APQC  (WW leader in Knowledge Practice) April 2013:
   *"The Industry Leader and Best Practice in Expertise Location"*

E6893 Big Data Analytics – Lecture 1

| Markov Network | Latent Network | Bayesian Network |
|---|---|---|



How'd those assorted tank tops work out for you?

login → browsing → search → comparing → Checkout

- Behavior Pattern Detection
- Help Needed Detection

**Goal:** *Injecting Network Graph Effects for Financial Analysis. Estimating company performance considering correlated companies, network properties and evolutions, causal parameter analysis, etc.*

- IBM 2003

- IBM 2009



- Data Source:
  - Relationships among 7594 companies, data mining from NYT 1981 ~ 2009

<u>Targets</u>: 20 Fortune companies' normalized Profits

<u>Goal:</u> Learn from previous 5 years, and predict next year

<u>Model:</u> Support Vector Regression (RBF kernel)



profit (R^2 mean)

**Network feature**:
 s (current year network feature),
  t (temporal network feature),
 d (delta value of network feature)
**Financial feature**:
 p (historical profits and

Profit prediction by joint network and financial analysis outperforms network-only by 130% and financial-only by 33%.

39

E6893 Big Data Ana

# Use Case 5: Social Media Monitoring



monitoring categories

Monitoring filter

Live Tweets, Sentiment, Keywords    Dynamic Graph    Zooming / Panning    Real-Time Translation, Loca
Top Retweets

E6893 Big Data Analytics – Lecture 1

© CY Lin, Columbia University

# Use Case 6: Customer Social Analysis for Telco

**Goal:** Extract customer social network behaviors to enable Call Detail Records (CDRs) data monetization for Telco.
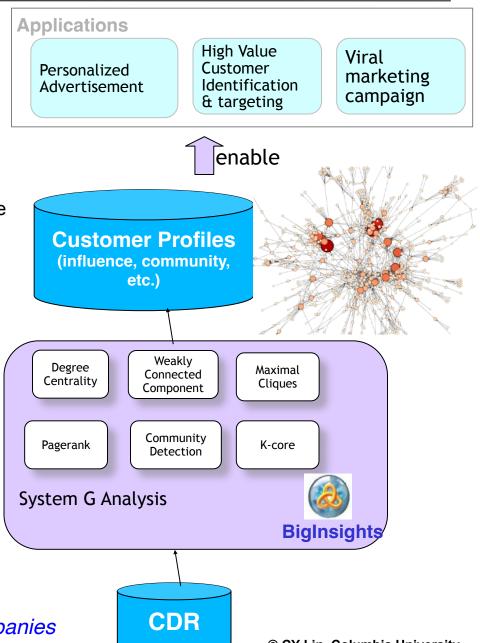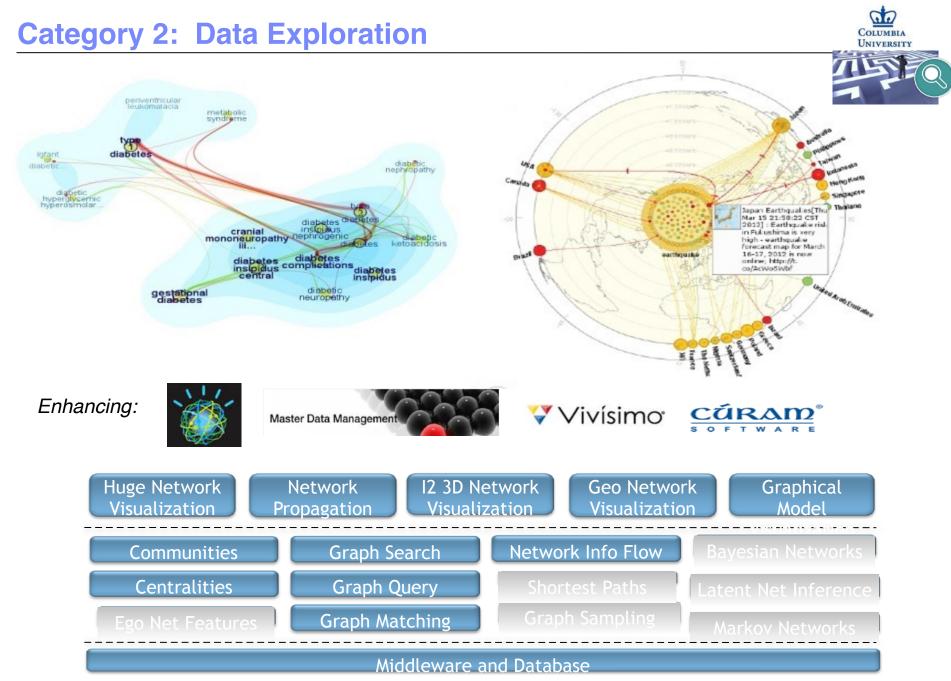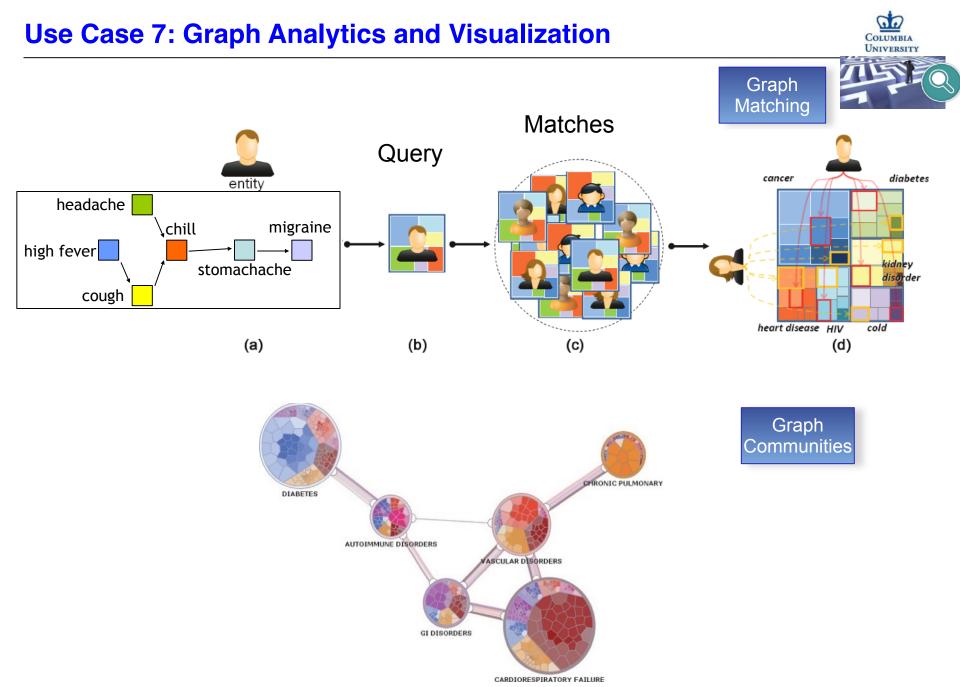
- Applications based on the extracted social profiles
  - Personalized advertisement (beyond the scope of traditional campaign in Telco)
  - High value customer identification and targeting
  - Viral marketing campaign
- Approach
  - Construct social graphs from CDRs based on {caller, callee, call time, call duration}
  - Extract customer social features (e.g. influence, communities, etc.) from the constructed social graph as customer social profiles
  - Build analytics applications (e.g. personalized advertisement) based on the extracted customer social profiles

**Applications**

Personalized Advertisement

High Value Customer Identification & targeting

Viral marketing campaign

enable

**Customer Profiles
(influence, community, etc.)**

Degree Centrality

Weakly Connected Component

Maximal Cliques

Pagerank

Community Detection

K-core

System G Analysis

**BigInsights**

**CDR**

*PoCs with Chinese and Indian Telecomm companies*

*Enhancing:*



| Huge Network Visualization | Network Propagation | I2 3D Network Visualization | Geo Network Visualization | Graphical Model |
|---|---|---|---|---|

| Communities | Graph Search | Network Info Flow | Bayesian Networks |
|---|---|---|---|
| Centralities | Graph Query | Shortest Paths | Latent Net Inference |
| Ego Net Features | Graph Matching | Graph Sampling | Markov Networks |

| Middleware and Database |
|---|

# Use Case 7: Graph Analytics and Visualization



Graph Matching

entity

headache
high fever
chill
cough
stomachache
migraine

(a)

Query

(b)

Matches

(c)

cancer
diabetes
kidney disorder
heart disease    HIV    cold

(d)

Graph Communities

DIABETES
AUTOIMMUNE DISORDERS
CHRONIC PULMONARY
VASCULAR DISORDERS
GI DISORDERS
CARDIORESPIRATORY FAILURE

Whisper : Tracing the information diffusion in Social Media

http://systemg.ibm.com/apps/whisper/index.html



SocialHelix: Visualizaiton of Sentiment Divergence in Social Media

- Belts : communities
- Color : sentiments
- Bars : keywords & hashtags of an event

# Use Case 9: Graph Search



existing search engine

query

Info-Socio networks

index

ranking

re-ranking

Graph analysis

query context

Improved search results

Interest / social network based content recommendations

Graph Search

E6893 Big Data Analytics – Lecture 1

© CY Lin, Columbia University

# Category 3:  Security

**Network Info Flow**

**Ponzi scheme Detection**

**Ego Net Features**

Normal:
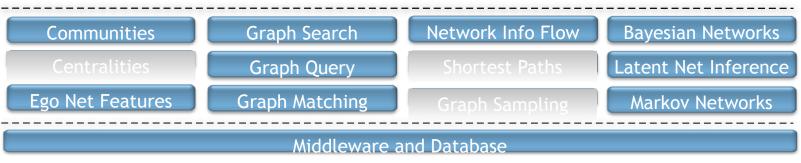(1) Clique-like
(2) Two-way links

Attacker:
Near-Star

## Detecting DoS attack

(a) Single large graph representing TCP SYN and ICMP PING network traffic, with two Denial of Service (DoS) attacks taking place.
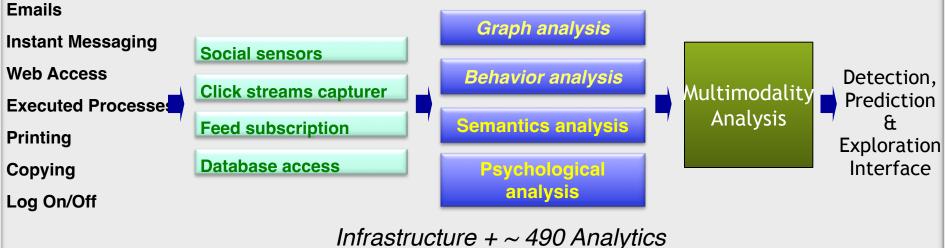
(a) Near-star   (b) Near-clique   (c) Heavy vicinity   (d) Dominant edge

| Graph Visualizations | | | |
|---|---|---|---|
| Communities | Graph Search | Network Info Flow | Bayesian Networks |
| Centralities | Graph Query | Shortest Paths | Latent Net Inference |
| Ego Net Features | Graph Matching | Graph Sampling | Markov Networks |

| Middleware and Database |
|---|

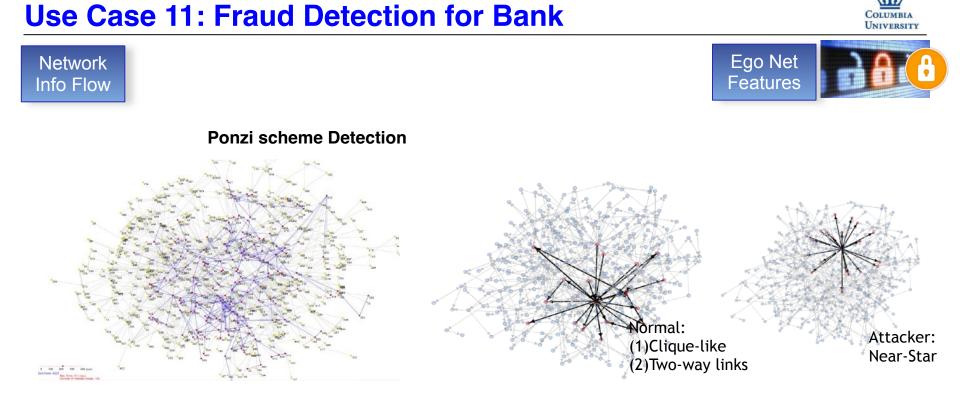E6893 Big Data Analytics – Lecture 1

© CY Lin, Columbia University

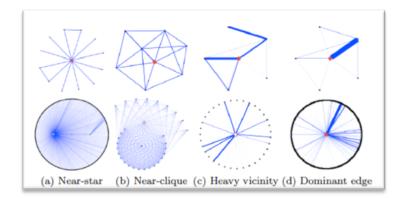**Based on President Executive Order 13587**

**Goal:** System for Detecting and Predicting Abnormal Behaviors in Organization, through large-scale social network & cognitive analytics and data mining, to decrease insider threats such as espionage, sabotage, colleague-shooting, suicide, etc.

THE WALL STREET JOURNAL.

Many Past Espionage Cases Had Links to | U.S. Ups Ante for Spying

**npr**

news › business

To Catch Worker Misconduct, Companies Hire Corporate Detectives

by AILSA CHANG

January 10, 2013   6:25 PM

"Enterprise Information Leakage Impacted economy and jobs" Feb 2013

"What's emerged is a multibillion dollar detective industry"
*npr Jan 10, 2013*

| Emails | | Graph analysis | | |
|---|---|---|---|---|
| Instant Messaging | **Social sensors** | **Behavior analysis** | | Detection, Prediction & Exploration Interface |
| Web Access | **Click streams capturer** | | Multimodality Analysis | |
| Executed Processes | | | | |
| Printing | **Feed subscription** | **Semantics analysis** | | |
| Copying | **Database access** | **Psychological analysis** | | |
| Log On/Off | | | | |

*Infrastructure + ~ 490 Analytics*

Network Info Flow

Ego Net Features

**Ponzi scheme Detection**



Normal:
(1) Clique-like
(2) Two-way links

Attacker: Near-Star



(a) Near-star  (b) Near-clique  (c) Heavy vicinity  (d) Dominant edge

**Detecting DoS attack**
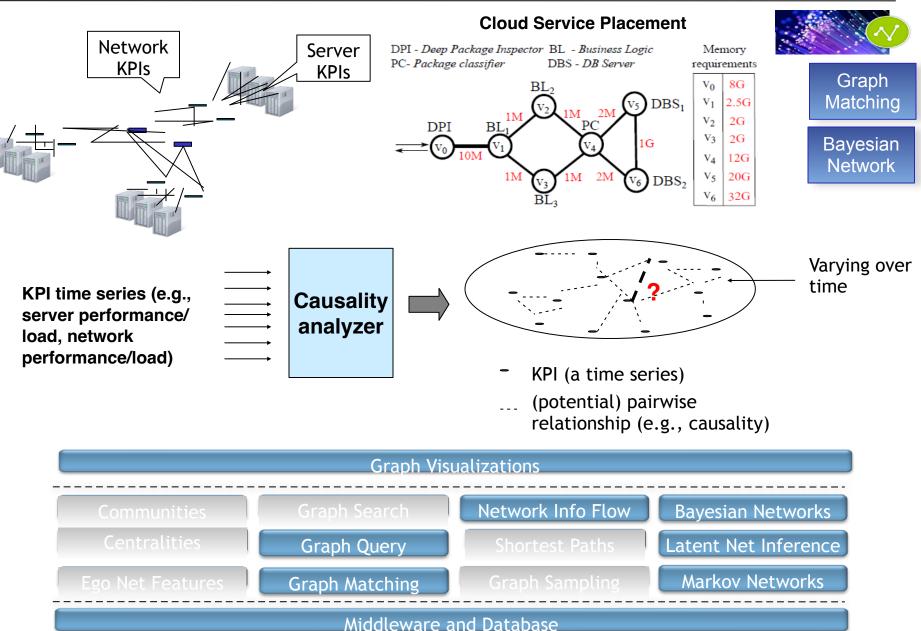


(a) Single large graph representing TCP SYN and ICMP PING network traffic, with two Denial of Service (DoS) attacks taking place.
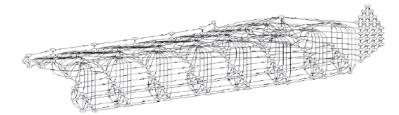
## Cloud Service Placement

Network KPIs

Server KPIs

DPI - *Deep Package Inspector*   BL - *Business Logic*
PC- *Package classifier*   DBS - *DB Server*

| | Memory requirements |
|---|---|
| $V_0$ | 8G |
| $V_1$ | 2.5G |
| $V_2$ | 2G |
| $V_3$ | 2G |
| $V_4$ | 12G |
| $V_5$ | 20G |
| $V_6$ | 32G |

Graph Matching

Bayesian Network

**KPI time series (e.g., server performance/ load, network performance/load)**

**Causality analyzer**

Varying over time

- KPI (a time series)

... (potential) pairwise relationship (e.g., causality)

Graph Visualizations

| Communities | Graph Search | Network Info Flow | Bayesian Networks |
| Centralities | Graph Query | Shortest Paths | Latent Net Inference |
| Ego Net Features | Graph Matching | Graph Sampling | Markov Networks |

Middleware and Database

E6893 Big Data Analytics – Lecture 1

# Use Case 13: Smarter *another* Planet

**Goal:** Atmospheric Radiation Measurement (ARM) climate research facility provides *24x7 continuous field observations* of cloud, aerosol and radiative processes. **Graphical models** can automate the validation with improvement efficiency and performance.

**Approach:** BN is built to represent the dependence among sensors and replicated across timesteps. BN parameters are learned from over *15 years* of ARM climate data to support distributed climate sensor validation. Inference validates sensors in the connected instruments.

Bayesian Network



**Dynamic Bayesian Network - Execution Time**
**(Dist. Sensor Network - 21 Measurements x 3 Timesteps)**



- State Graph
- Clique Scheduling

Exec Time (s) / Processor Count

## Bayesian Network
* 3 timesteps        * 63 variables
* 3.9 avg states   * 4.0 avg indegree
* 16,858 CPT entries
## Junction Tree
* 67 cliques
* 873,064 PT entries in cliques

E6893 Big Data Analytics – Lecture 1

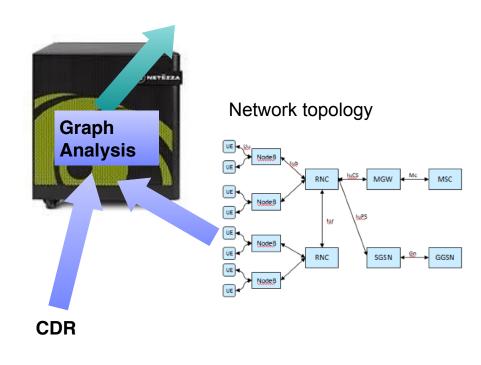# Use Case 14: Cellular Network Analytics in Telco Operation

**Goal:** Efficiently and uniquely identify *internal* state of Cellular/Telco networks (e.g., performance and load of network elements/links) using probes between monitors placed at selected network elements & endhosts

- Applied Graph Analytics to telco network analytics based on CDRs (call detail records): estimate traffic load on CSP network with low monitoring overhead

    (1) CDRs, already collected for billing purposes, contain information about voice/data calls

    (2) Traditional NMS* and EMS** typically lack of end-to-end visibility and topology across vendors

    (3) Employ graph algorithms to analyze network elements which are not reported by the usage data from CDR information

- Approach

    – Cellular network comprises a hierarchy of network elements

    – Map CDR onto network topology and infer load on each network element using graph analysis

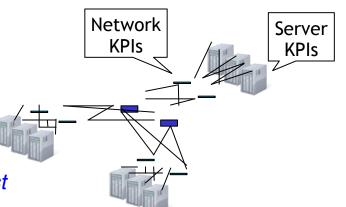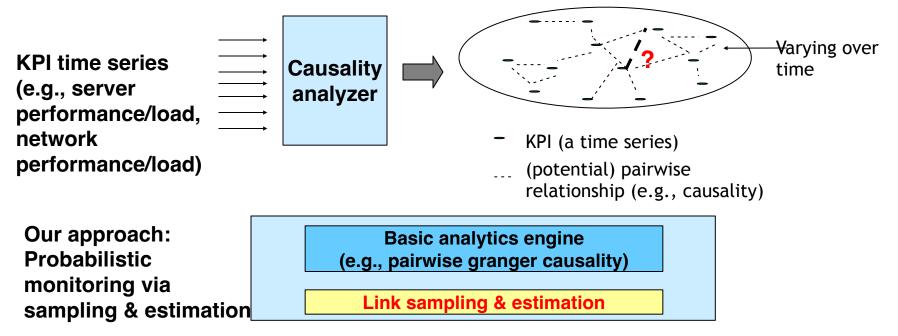    – Estimate network load and localize potential problems

Network load level report

Graph Analysis

Network topology

CDR

# Use Case 15: Monitoring Large Cloud

**Goal:** *M*onitoring technology that can track the time-varying state (e.g., causality relationships between KPIs) of a large Cloud when the processing power of monitoring system cannot keep up with the scale of the system & the rate of change

Network KPIs

Server KPIs

- *Causality relationships (e.g., Granger causality) are crucial performance monitoring & root cause analysis*
- *Challenge: easy to test pairwise relationship, but hard to test multi-variate relationship (e.g., a large number of KPIs)*



**KPI time series (e.g., server performance/load, network performance/load)**

**Causality analyzer**

Varying over time

– KPI (a time series)

··· (potential) pairwise relationship (e.g., causality)

**Our approach: Probabilistic monitoring via sampling & estimation**

**Basic analytics engine (e.g., pairwise granger causality)**

**Link sampling & estimation**

*Select KPI pairs (sampling)*→ *Test link existence* → *Estimate unsampled links based on history*

53 → *Overall graph*

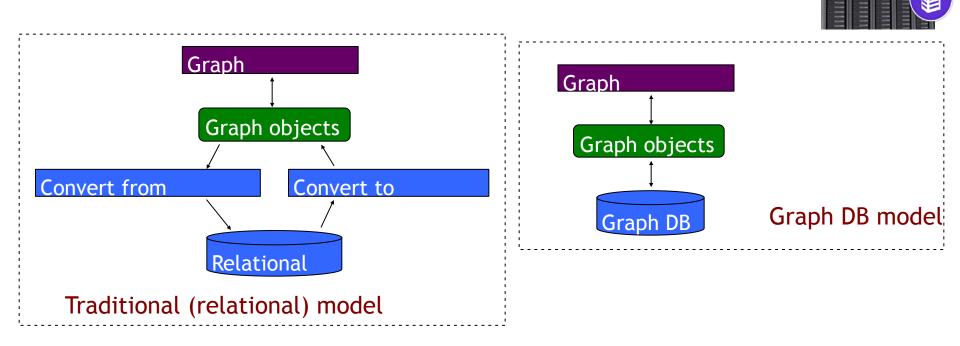E6893 Big Data Analytics – Lecture 1

Traditional (relational) model

Graph DB model

- Advantages of working directly with graph DB for graph applications
  - (1) Smaller and simpler code
  - (2) Flexible schema → easy schema evolution
  - (3) Code is easier and faster to write, debug and manage
  - (4) Code and Data is easier to transfer and maintain
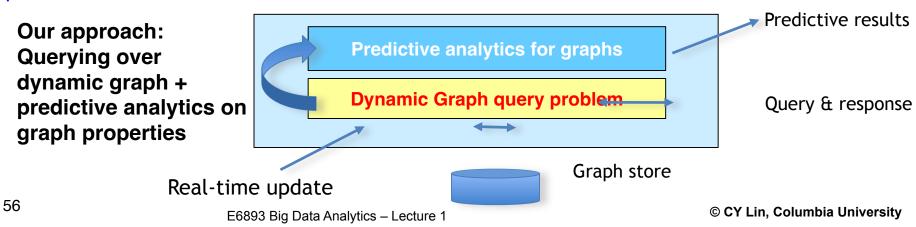
# Use Case 17: Smart Navigation Utilizing Real-time Road Information

**Goal:** *Enable unprecedented level of accuracy in **traffic scheduling** (for a fleet of transportation vehicles) and navigation of individual cars utilizing the **dynamic real-time information** of changing road condition and predictive analysis on the data*

• *Dynamic graph algorithms implemented in System G provide **highly efficient graph query computation** (e.g. shorted path computation) on time-varying graphs (order of magnitudes improvement over existing solutions)*
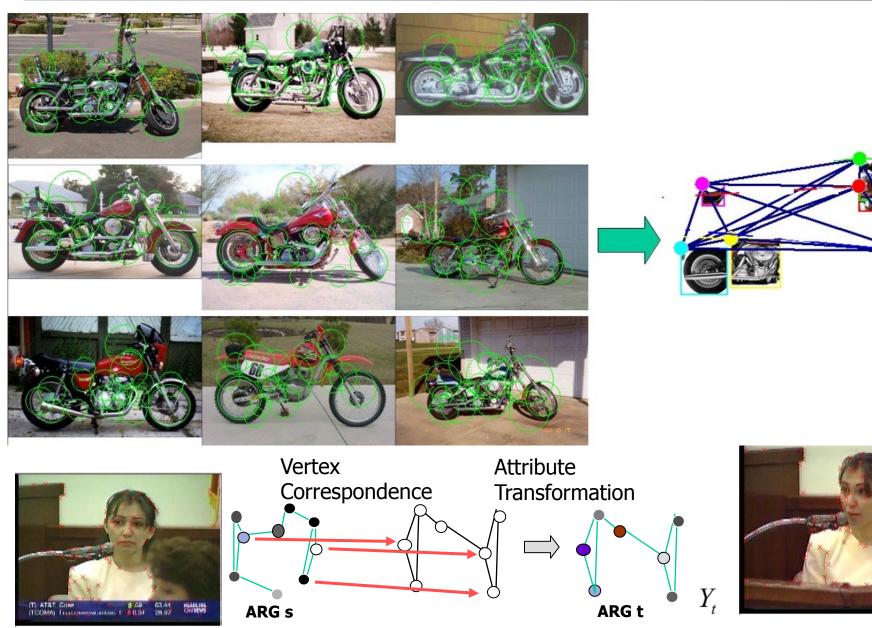
• *High-throughput **real-time predictive analytics** on graph makes it possible to estimate the future traffic condition on the route to make sure that the decision taken now is optimal overall*



**Our approach: Querying over dynamic graph + predictive analytics on graph properties**

**Predictive analytics for graphs**

**Dynamic Graph query problem**

Predictive results

Query & response

Real-time update

Graph store

E6893 Big Data Analytics – Lecture 1

© CY Lin, Columbia University

Vertex Correspondence

Attribute Transformation

**ARG s**
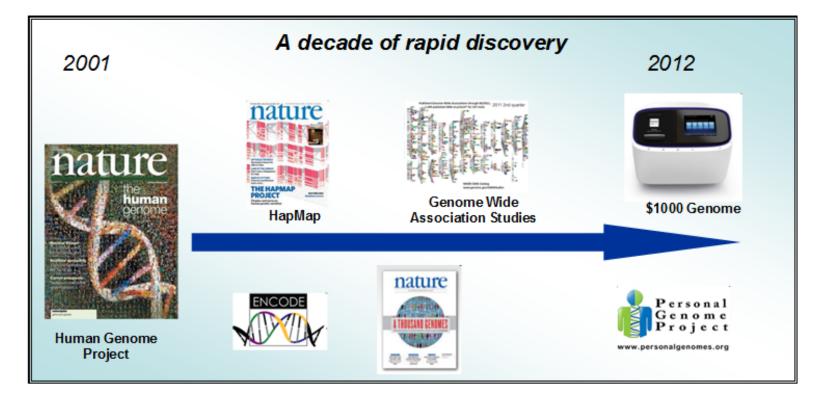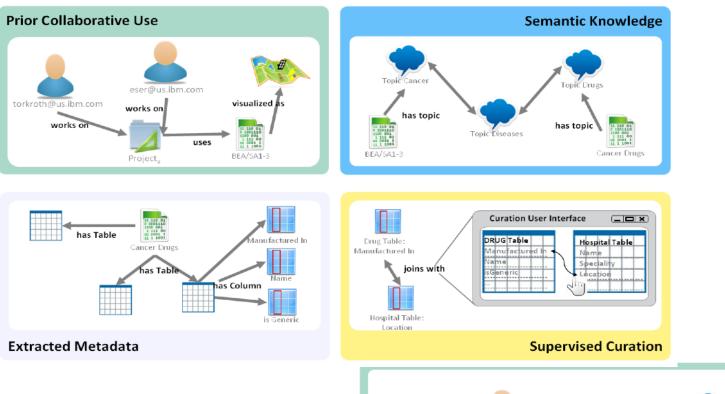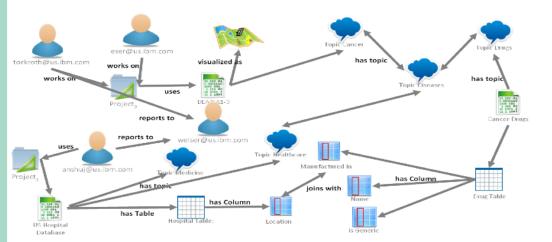
**ARG t**

$Y_t$

E6893 Big Data Analytics – Lecture 1

**Figure 1: Since the Human Genome Project, various projects have started to reveal the mysteries of genomes and the $1000 Genome is almost reality.**
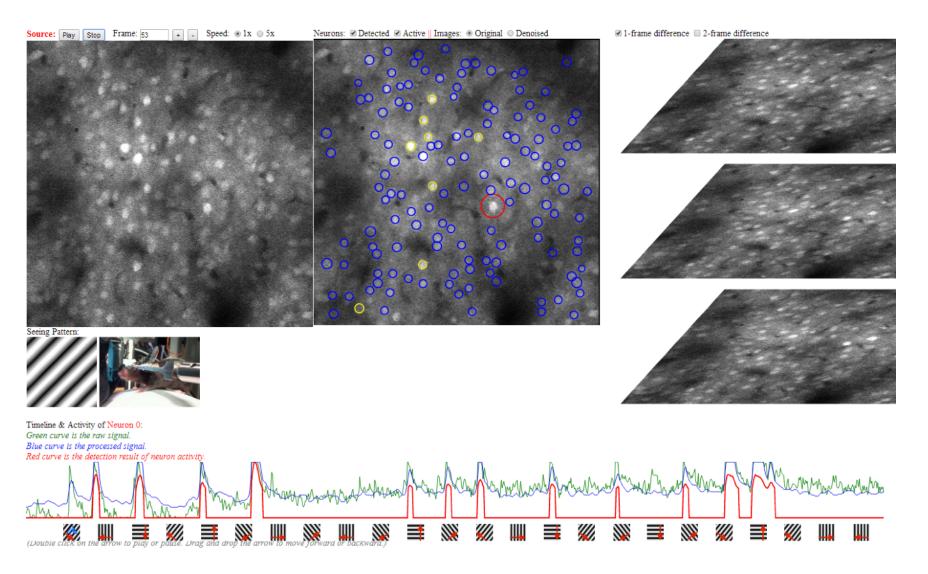
E6893 Big Data Analytics – Lecture 1

- Big Data on Large-Scale Sky Monitoring



Photograph by Rob Ratkowski for the PS1SC

**Dangers from space**

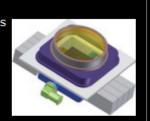Learn about the threat to Earth from asteroids & comets and how the Pan-STARRS project is designed to help detect these NEOs. Learn more...

**1,400,000,000 pixels**

Pan-STARRS has the world's largest digital cameras.

Read about them here...

**The PS1 Prototype**

PS1 goes operational and begins science mission

PS1 Science Consortium formed...

PS1SC Blog

PS1 image gallery

E6893 Big Data Analytics – Lecture 1

# Homework #0: Big Data Environment Setup and Test (due September 20, 5pm)

1. Warm-Up Exercises:
   - Setup Google Cloud account and environment
   - Install Google Cloud SDK
   - Create a Spark cluster
   - Word Count using Google Cloud Storage and Spark
   - Hive and BigQuery

2. Data Analysis — NYC Bike Expert:
   - Load data to a Cloud Storage
   - Simple Analyses through BigQuery

3. Data Analysis — Understanding Shakespeare:
   - Load data to a Cloud Storage
   - Simple Analyses through Word Counts
   - Analyses after running Natural Language Toolkit

E6893 Big Data Analytics – Lecture 1: Big Data Introduction © 2024 CY Lin, Columbia University

# Homework Late Submission Policy

5pm: submission deadline
Next Day midnight: 10% penalty
Two Days late midnight: 20% penalty
Three Days late midnight: 30% penalty
Any late submission more than 3 days will not be accepted.

Please do your each homework as early as possible!! They are all quite 'heavy'.