

Big Data Infrastructures

Philippe Cudré-Mauroux

Fall 2018

Lecture 1 - Class Introduction

Outline

- Introduction to Big Data
- Class overview

Instant Quiz

- **Foreign-Key?**
- **Normal Forms?**
- **Two-Phase Commit?**
- **ACID? BASE?**
- **CAP?**
- **Hadoop? Yarn? Spark?**

Big Data & Me

- My lab @ unifr: eXascale Infolab (<http://exascale.info/>)
 - Previously: M.I.T. (Stonebraker's lab), EPFL (Best PhD Award), U.C. Berkeley
 - Industry also (IBM Watson Research, HP, Microsoft Research Asia, Microsoft CISL, Scigility, Dashcom)
 - Teach Big Data at **Swiss Joint MSc in CS**, U. Lucern, Royal Institute of Tech. (Sweden), IIMT, EPFL, HEC Lausanne

→ How to store and manage Big Data

- 2M € ERC, SNF, Haslerstiftung, H2020
- Verisign, SAP, Microsoft, Amazon, Google, ArmaSuisse



eXascale Infolab



Course Assistants

- Several labs and exercises will be given by Ph.D. students and Senior Researchers



Goals of the Class

- Study new principles of data management
 - Data properties, data independence, scale-out, etc.
- Study key DBMS design issues
 - Transactions, parallel data processing, etc.
 - Consistency, Availability, network Partition
- Hands on some current, *hot* Big Data topics
 - Get exposed to recent trends / technologies / start-ups
- Ensure that
 - You are comfortable using a DBMS
 - You have an idea about (distributed) data management
 - You know a bit about current research topics in data management
 - You know a bit about the current market and industry trends

Exascale Data Deluge

- Web companies
 - Google
 - Ebay
 - Yahoo
- Science
 - Biology
 - Astronomy
 - Remote Sensing
- Financial services
 - retail companies
 - governments, etc.



New data formats
New machines
Peta & exa-scale datasets
Obsolescence of traditional information infrastructures



The privacy debate
A scramble for information threatens individual freedoms | Page 5

Big data
BusinessTechnology
The power behind decisions



How big data put Obama back in the White House | Pages 8 and 9

Big Data
BusinessTechnology
The power behind decisions



D I V E R S E D I S C I P L I N E S , O N E
Biomedical Computing
Published by Simbiots, an NIH National Center for Biomedical Computing

Big Data
Analytics
in Biomedical
Research

PLUS:
Privacy and Biomedical Research:
Building a Trust Infrastructure

Big data can generate significant financial value across sectors



US health care

- \$300 billion value per year
- ~0.7 percent annual productivity growth



Europe public sector administration

- €250 billion value per year
- ~0.5 percent annual productivity growth



Global personal location data

- \$100 billion+ revenue for service providers
- Up to \$700 billion value to end users



US retail

- 60+% increase in net margin possible
- 0.5–1.0 percent annual productivity growth



Manufacturing

- Up to 50 percent decrease in product development, assembly costs
- Up to 7 percent reduction in working capital

SOURCE: McKinsey Global Institute analysis

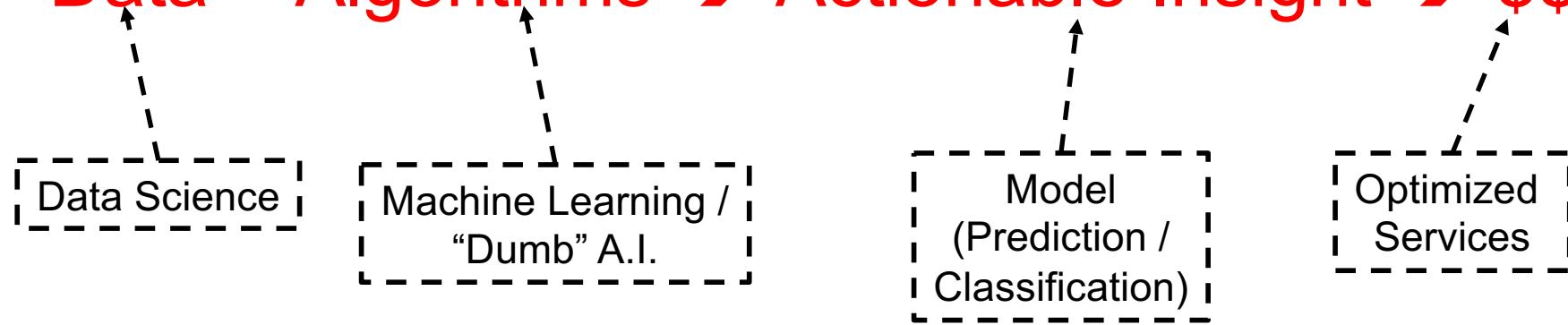
Big Data as a New Class of Asset

- The Age of Big Data (NYTimes Feb. 11, 2012)
<http://www.nytimes.com/2012/02/12/sunday-review/big-datas-impact-in-the-world.html>

“Welcome to the Age of Big Data. The new megarich of Silicon Valley, first at Google and now Facebook, are masters at harnessing the data of the Web — online searches, posts and messages — with Internet advertising. At the World Economic Forum last month in Davos, Switzerland, Big Data was a marquee topic. A report by the forum, “Big Data, Big Impact,” declared **data a new class of economic asset, like currency or gold.**”

Data is the new Oil

- Data + Algorithms → Actionable Insight → \$\$



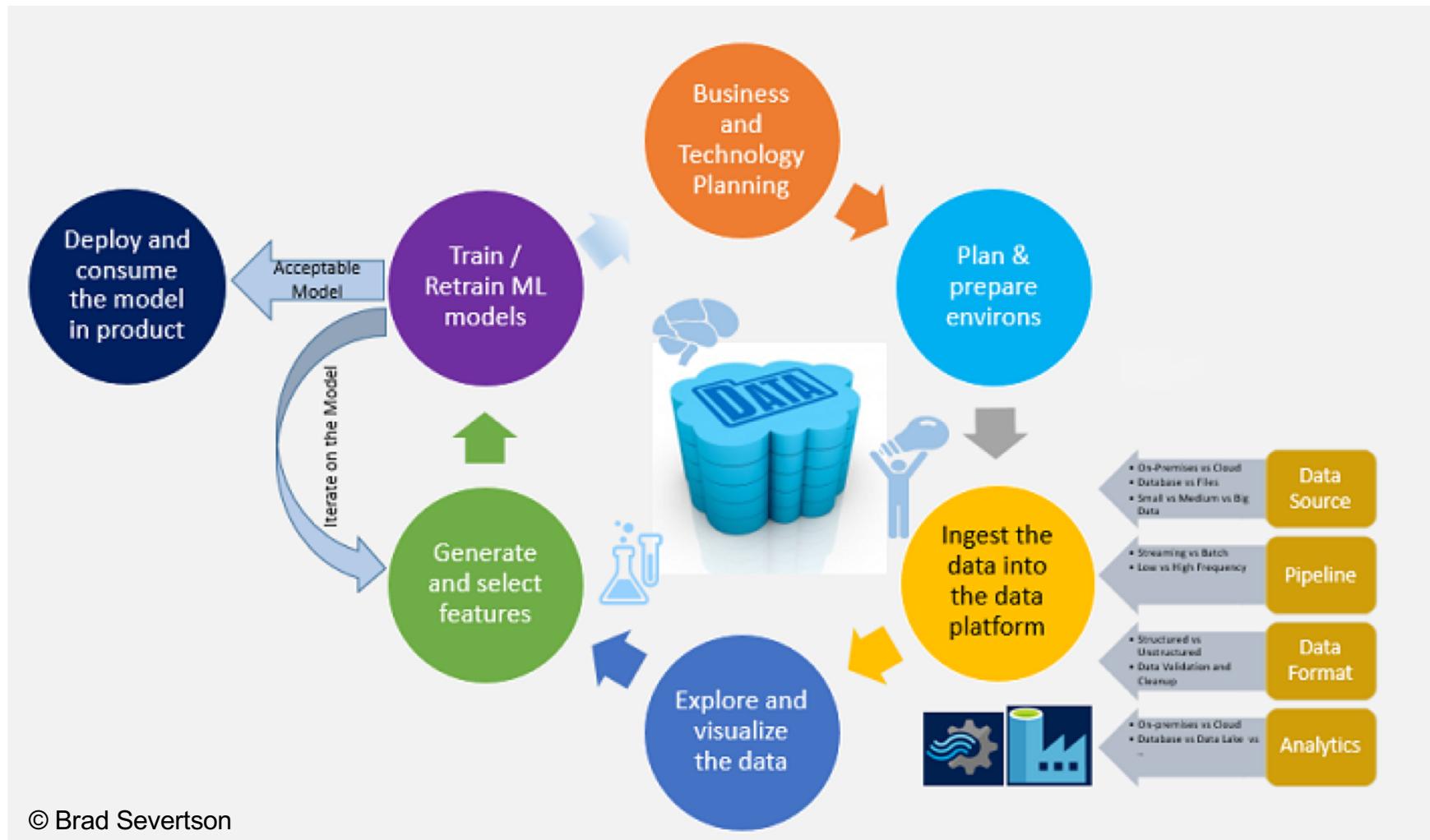
The 3-Vs of Big Data

- **V**olume
 - Amount of data
- **V**elocity
 - speed of data in and out
- **V**ariety
 - range of data types and sources
- [Gartner 2012] "*Big Data are high-volume, high-velocity, and/or high-variety information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization*"

What can you do with the data

- Reporting
 - Post Hoc
 - Real time
- Monitoring (fine-grained)
- Exploration
- Finding Patterns
- Root Cause Analysis
- Closed-loop Control
- Model construction
- Prediction
- ...

Data Science Lifecycle

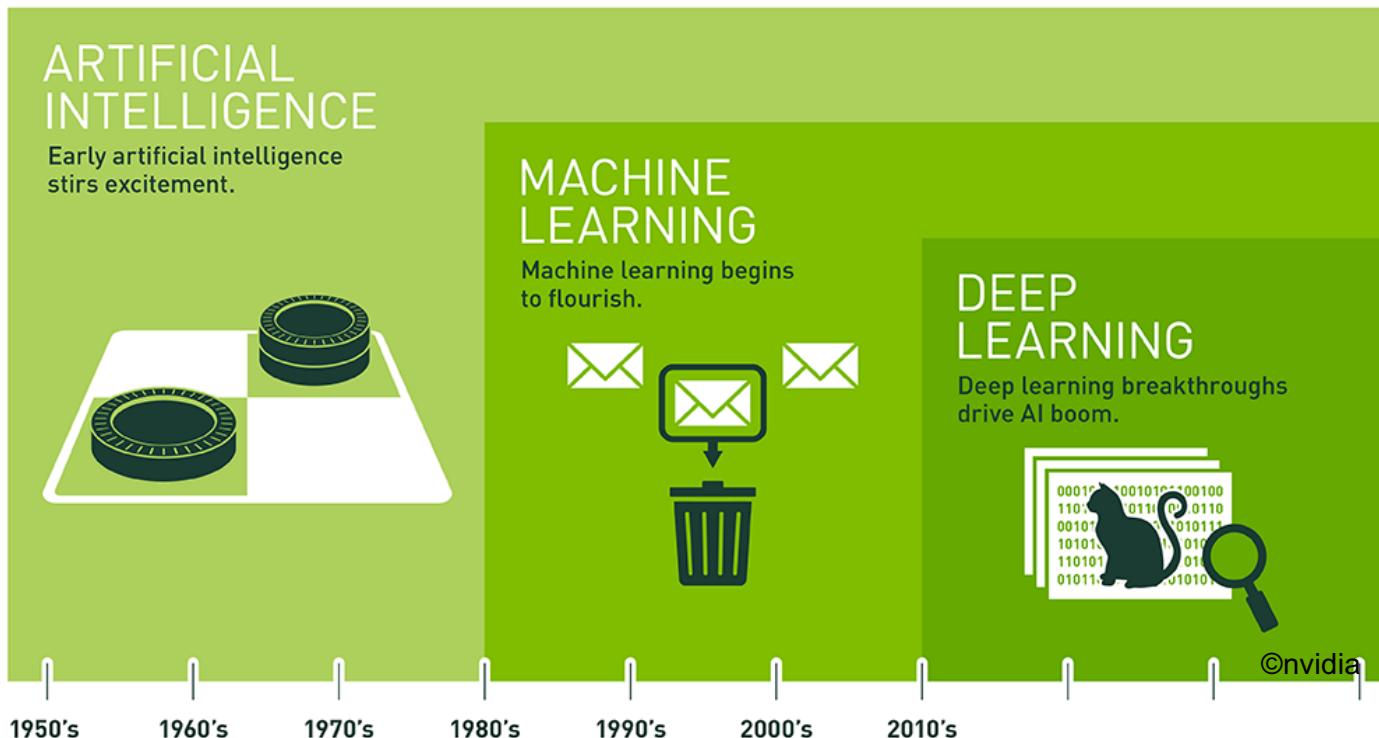


Three Generations of Models

- Descriptive models (yesterday)
- Predictive models (today)
- Prescriptive models (tomorrow)

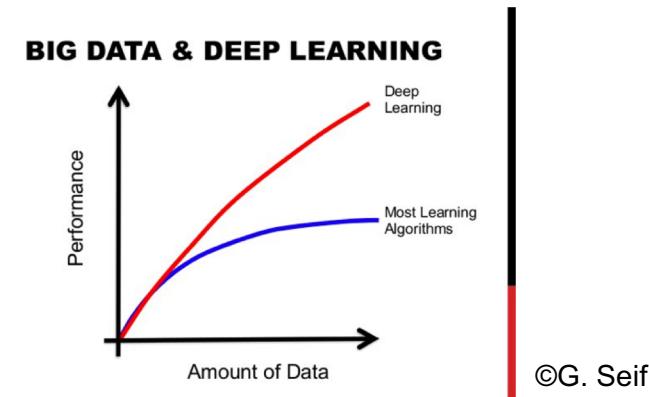
Current Model Development

- Today, models are dominated by **Deep Learning** techniques that are pushed by leading American and Chinese IT companies

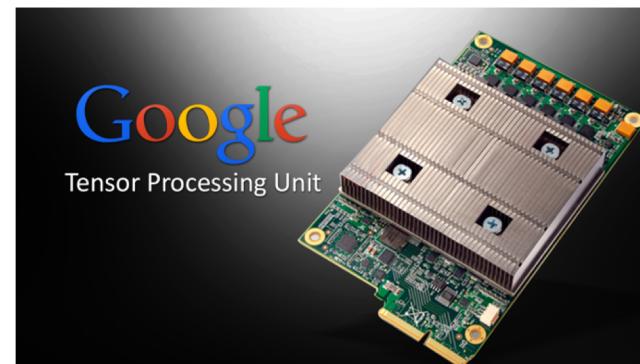


Deep Learning Feeds on Big Data

- Requires *gigantic* amounts of **annotated data** (examples)

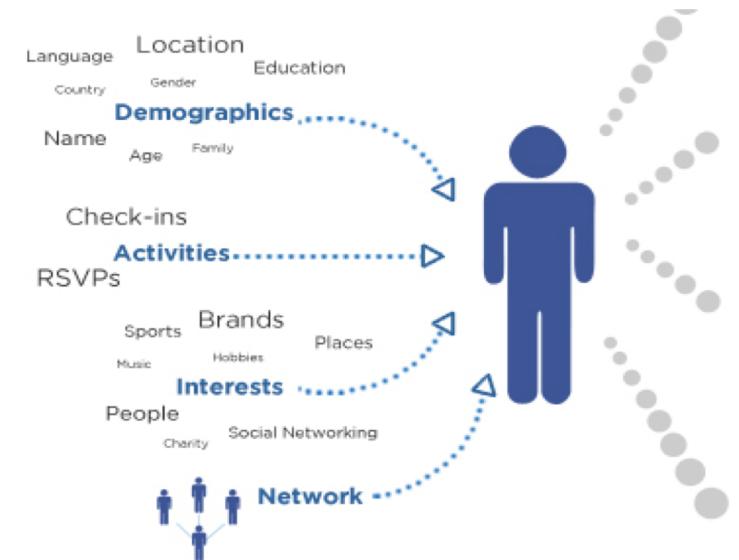


- Requires **enormous computation** (GPUs, TPUs)



Typical Big Data Success Story

- **Modeling users** through Big Data
 - Online ads sale / placement [e.g., Facebook]
 - Personalized Coupons [e.g., Target]
 - Product Placement [Walmart]
 - Content Generation [e.g., NetFlix]
 - Personalized learning [e.g., Duolingo]
 - HR Recruiting [e.g., Gild]



More Data => Better Answers?

- Not that easy...
- More Rows: Algorithmic complexity kicks in
- More Columns: Exponentially more hypotheses
- Another formulation of the problem:
 - Given an inferential goal and a fixed computational budget, provide a guarantee that the quality of inference will increase monotonically as data accrue (without bound)
- In other words:
=> **Data should be a resource, not a load**

What's wrong with my old DBMS?

- Managing Big Data is hard...
 - ... extremely hard
 - Traditional DBMSs are 30 years old, were not meant for Big Data
 - One user, one CPU, one type of queries
 - Obsolete physical model (n-ary storage, B-trees, etc.)
 - Impractical logical guarantees (transactions, ACID)



BIG DATA LANDSCAPE 2017



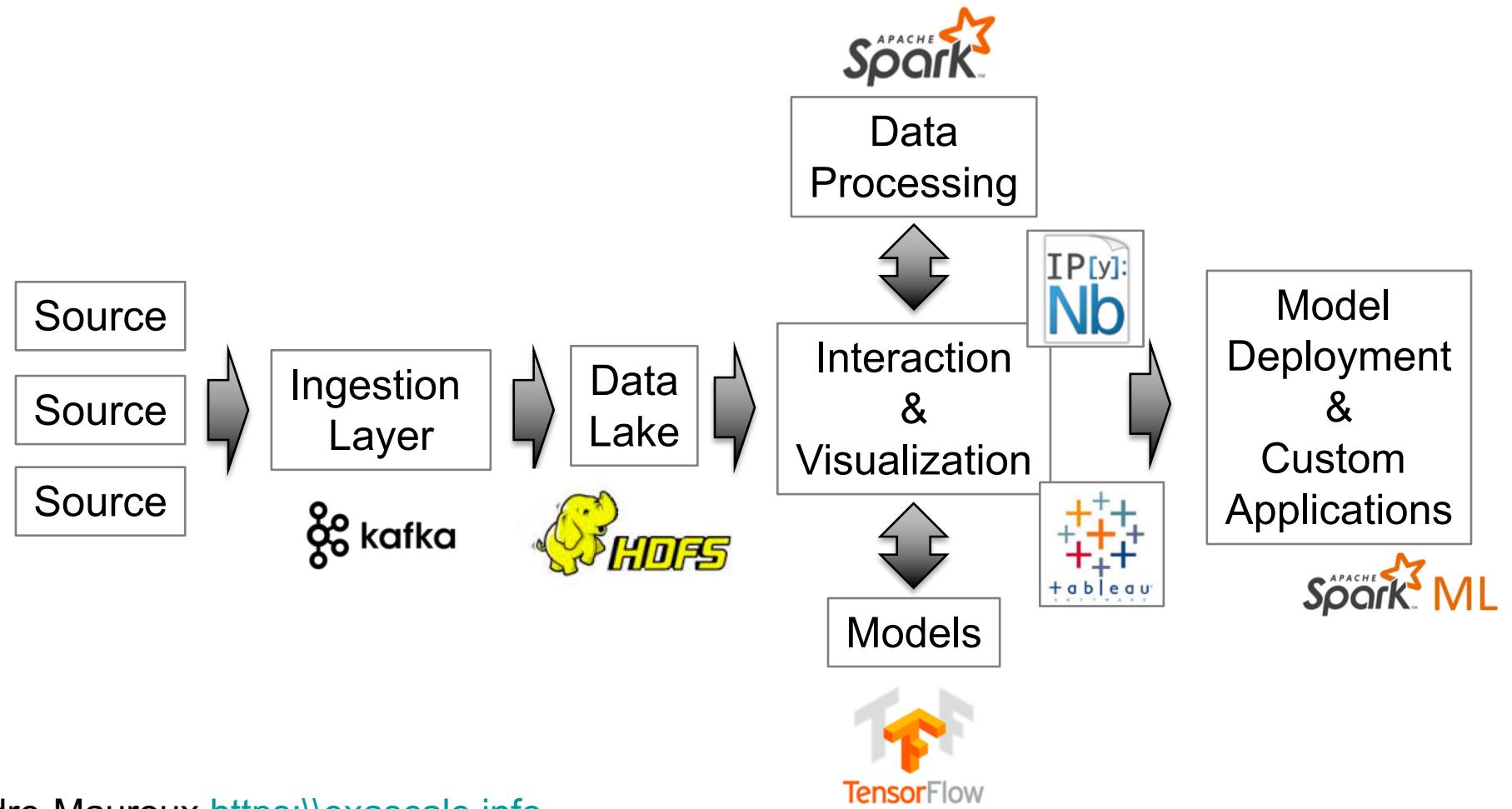
V2 – Last updated 5/3/2017

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

mattturck.com/bigdata2017

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

Data Science Infrastructure (circa 2018)



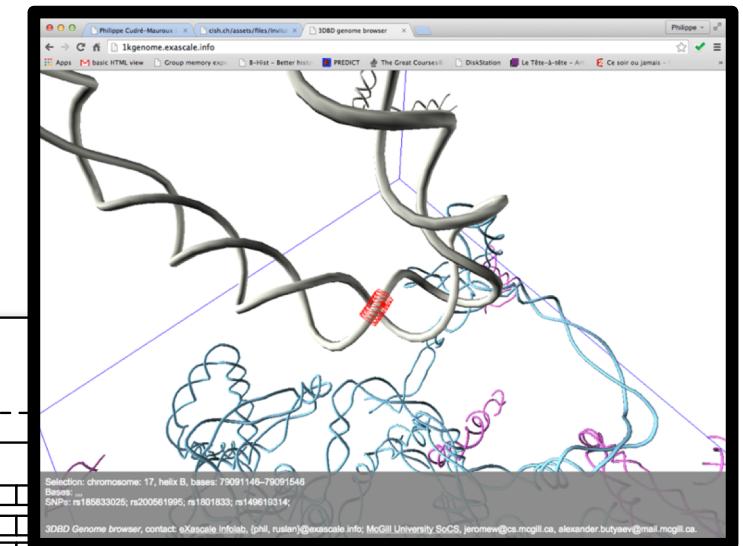
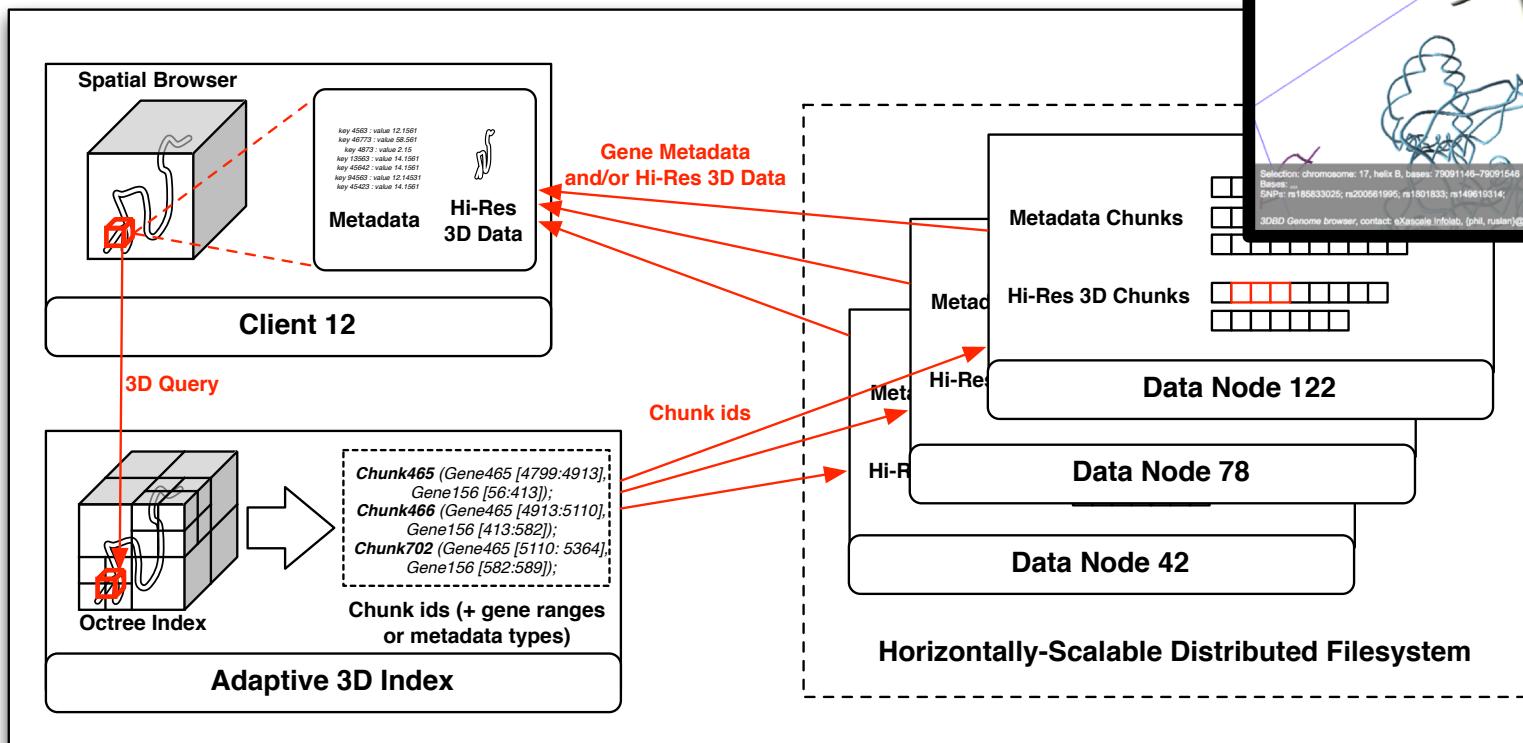
Three Examples from XI

- Three examples from the eXascale Infolab:
 - Volume: 3D Genome Browser
 - Velocity: Detecting anomalies in real-time in smart-cities
 - Variety: Integrating Log Data on Azure

Volume

The 3D Genome Browser

- Graceful scale-out
- Very low latency



Velocity

- Smart(er) Cities!



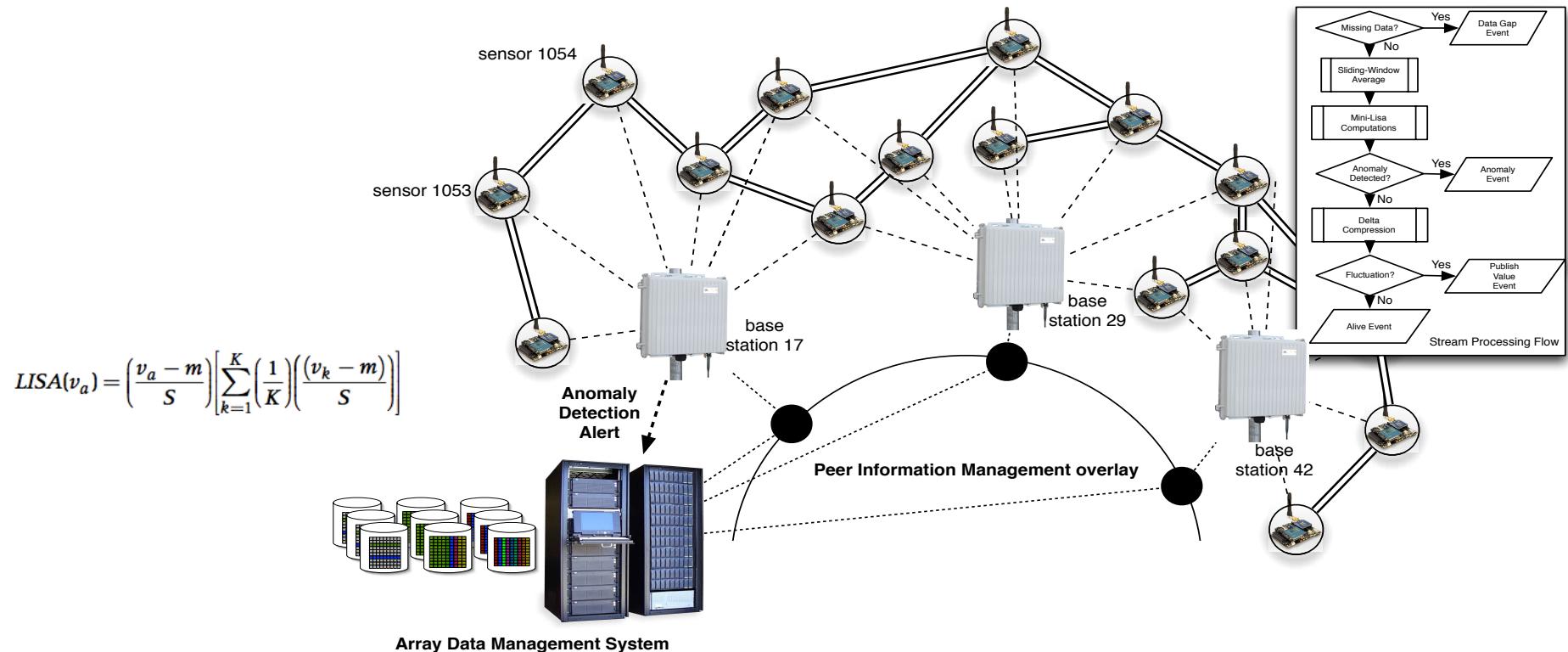
Velocity

- Detecting leaks / pipe bursts / contamination in real-time for water distribution networks

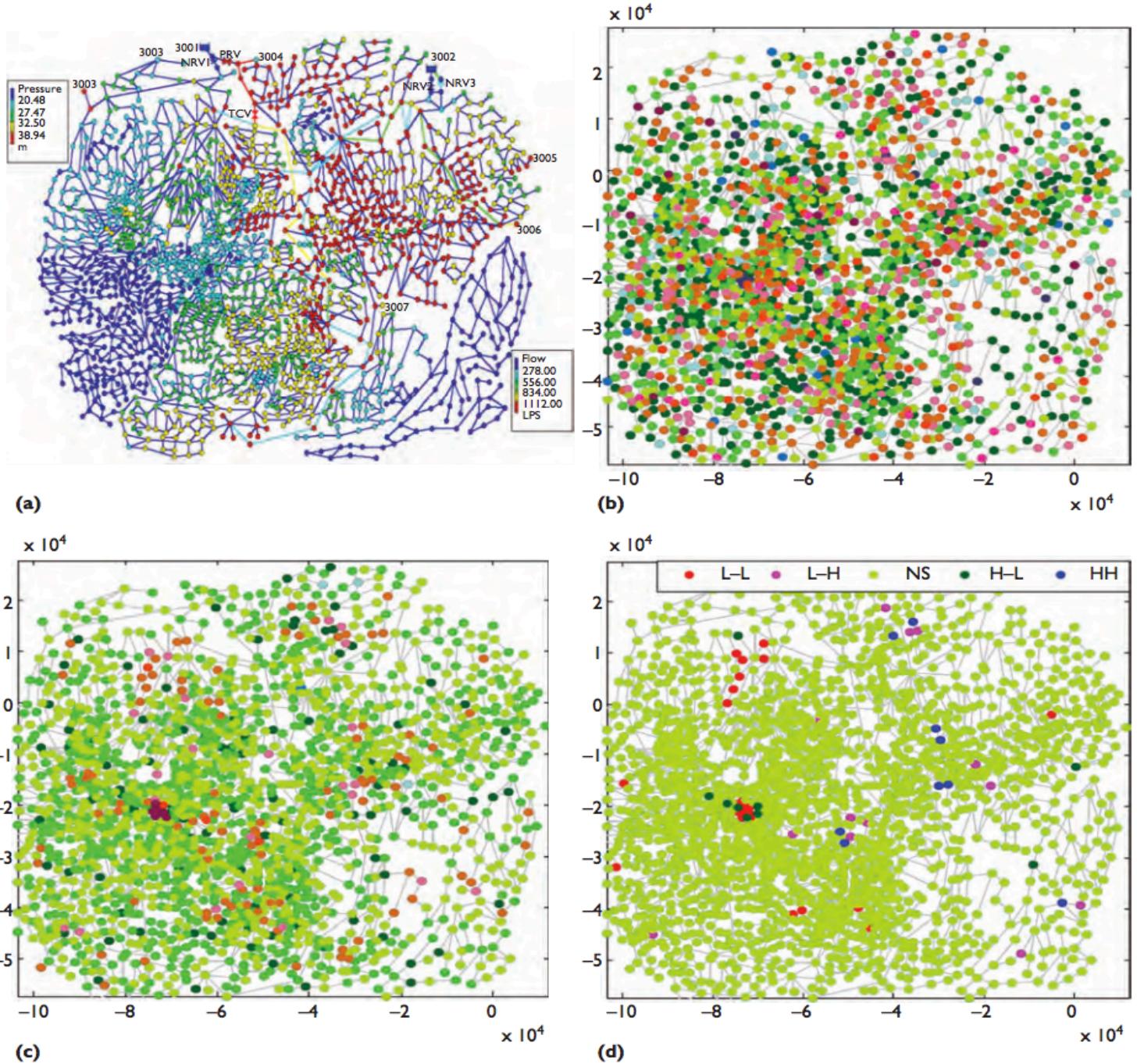


Data at each Vertex!

- Spatial + temporal statistical processing (mini-Lisas)
- Stream processing (Storm) + Array processing (SciDB)

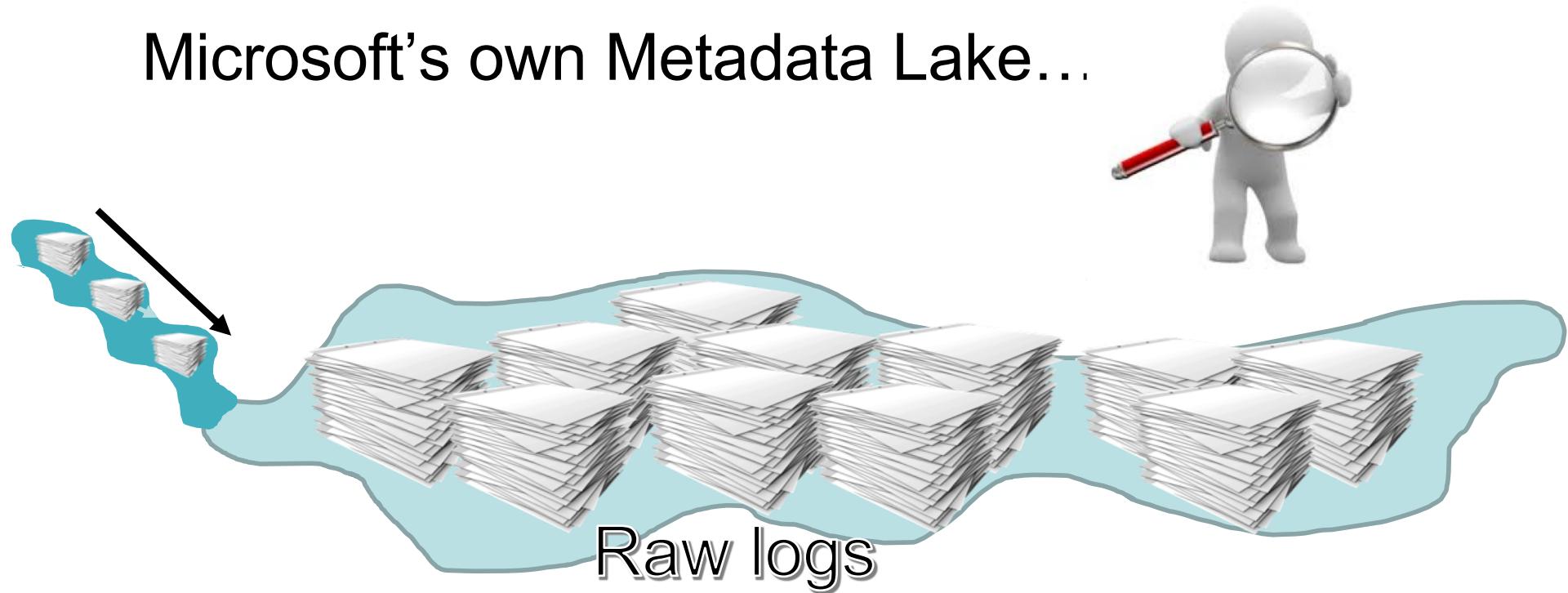


Results (anomalies Detected)



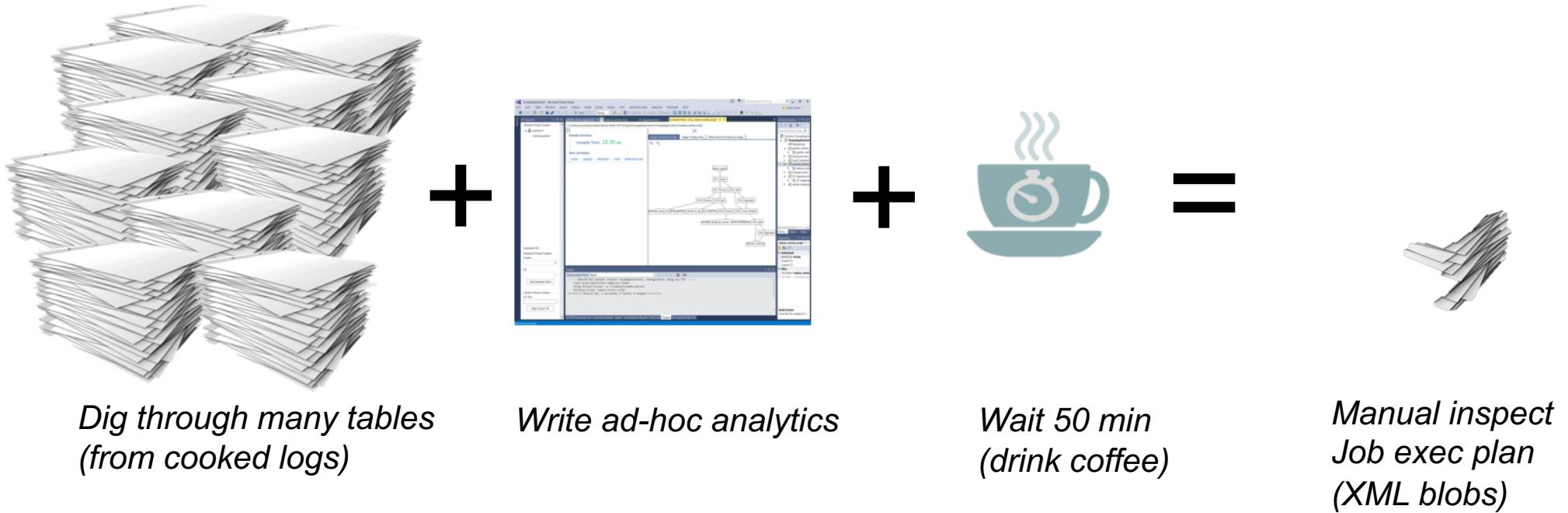
Variety

Microsoft's own Metadata Lake...

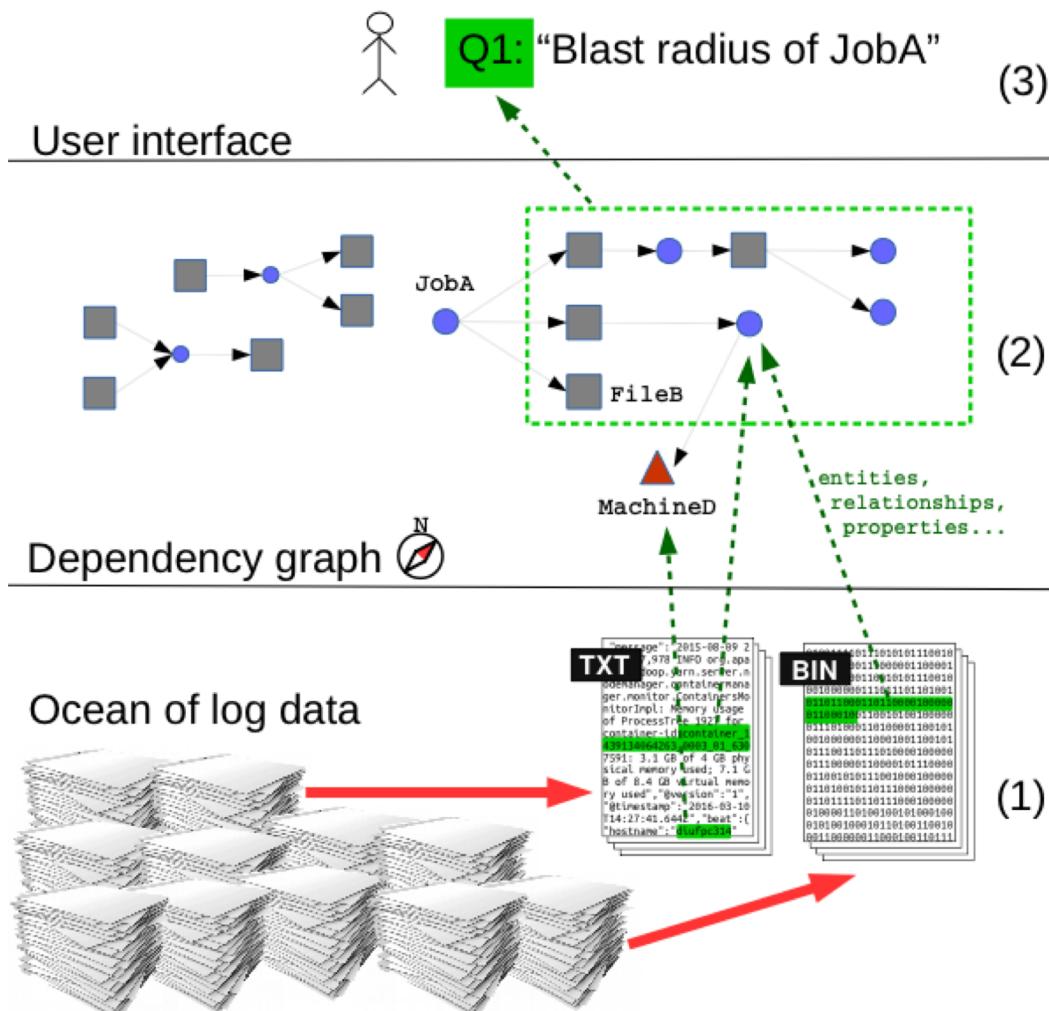


Example 1: Job pipeline analysis

- *User: “I need help with my ML experiment processing Clicklogs”*
- *Ops / Engineer:*



Our Solution: Guider

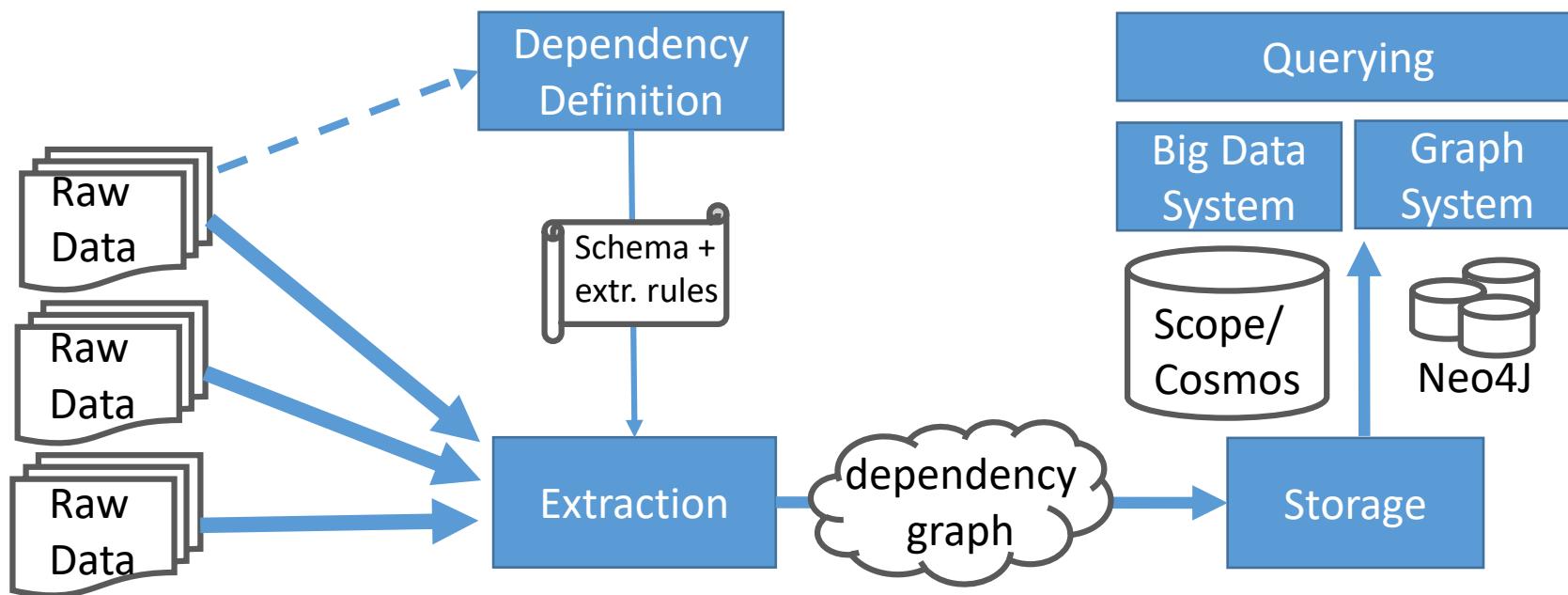


(3) User-level queries return bytes of aggregated data.

(2) Entity graph that represents a lightweight “skeleton” of the logs used for navigation

(1) Petabytes of daily logs

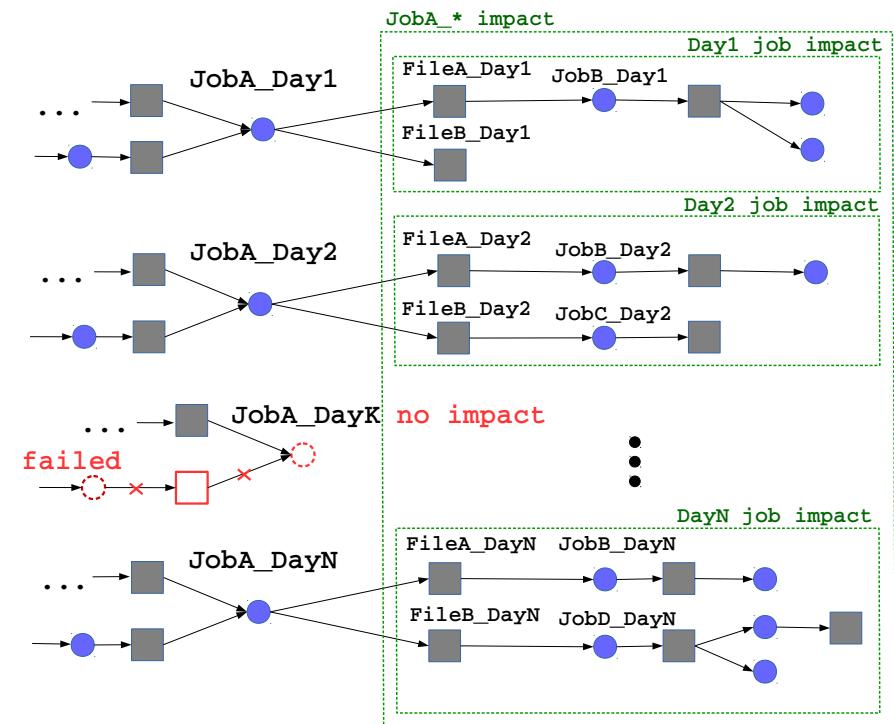
Guider Architecture



Dependency-Driven Analytics: a Compass for Uncharted Data Oceans.
Ruslan Mavlyutov, Carlo Curino, Boris Asipov, and Phil Cudre-Mauroux. CIDR 2017

Guider Use-Cases

1. Auditing and Compliance [in production]
2. Job Scheduling [Morpheus]
3. Global Job Ranking
4. Datacenter migration



Big Data Today in



- Big Data is not a new technology: it's a fact;
 - Deal with it → POCs and productized in most banks, insurance companies, retailers
- Largely behind US (and Asia)
- Leader in EU landscape
 - Research
 - Large Companies
 - SMEs
 - Administrations

Outline

- Introduction to Big Data
- Class overview

Class Format

- One lecture per week: Thursday, 9:15am-12am, **E230**
- A number of **labs** (on machine/laptop)
 - Not systematically every week
 - Groups of 1-2 students
 - Each group uploads their solution on Ilias by Wednesday evening (following week)
- Lectures + discussion
 - Don't be afraid to express yourselves and ask questions

One Note About the Slides

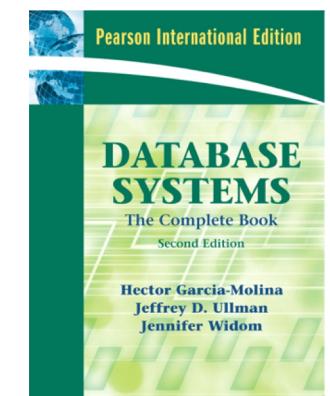
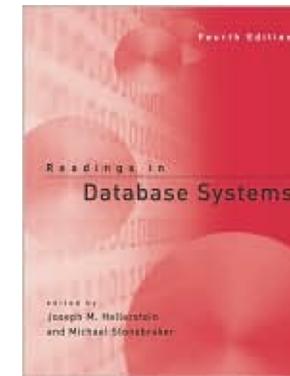
Some slides adapted from
Prof. Magda Balazinska's
fantastic course @ UWash
with her kind permission

Thanks Madga!



Readings and Notes

- Some readings are based on papers
 - Mix of old seminal papers and new papers
 - Papers will be available online on website
 - Many come from the “red book” [no need to get it but available at the Perolles 90 library]
 - Three types of readings
 - Mandatory, additional resources, and optional
- Background readings from the following book
 - **Database Systems: The Complete Book**, Garcia-Molina (available at Perolles 90)
 - Intro to databases (Stanford online course)
(or simply some of their notes: <https://cs.stanford.edu/people/widom/cs145/>)
- Lecture notes (the pptx slides)
 - Posted on class website before each lecture



Tentative Menu (1/2)

- 20.09 BDI 1: intro
- 27.09 BDI 2: SQL lab (**homework 1**)
- 04.10 BDI 3: Storage+Indexing
- 11.10 BDI 4: Query Execution
- 18.10 BDI 5: Column Stores + Query Execution lab
(**homework 2**)
- 25.10 BDI 6: Transactions + CAP
- 01.11 **Toussaint (holiday)**

Tentative Menu (2/2)

- 08.11 BDI 7: Mid-Term
- 15.11 Dies academicus (no course)
- 22.11 BDI 8: NoSQL + HDFS+HBase lab ([homework 3](#))
- 29.11 BDI 9: GraphDB + lab ([homework 4](#))
- 06.12 BDI 10: Hadoop+Yarn+Map/Reduce + lab
([homework 5](#))
- 13.12 BDI 11: Logging & Recovery
- 20.12 BDI 12: Spark + lab ([homework 6](#))

Evaluation

- Final Grade = 25% labs + 25% mid-term + 50% final exam
- Labs grade:
 - 6 “pass”: 6
 - 5 “pass”: 5.5
 - 4 “pass”: 5.0
 - 3 “pass”: 4.0
 - 2 “pass”: 3.0
 - 1 “pass”: 2.0
 - 0 “pass”: 1.0