# Big Data Integration

Philippe Cudré-Mauroux

eXascale Infolab, University of Fribourg
Switzerland

**eXascale Infolab**

# Instant Quiz

- n-Vs of Big Data?
- Yarn?
- Spark
- Knowledge Graph?

# eXascale Infolab (XI)

- New lab @ U. of Fribourg–Switzerland
- <span style="color:red">Big Data/AI infrastructures</span> for social / semantic / scientific data
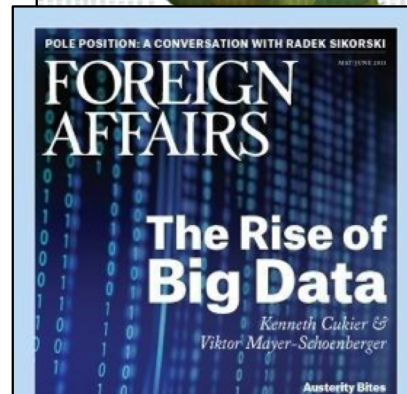


https://exascale.info/

# 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

# Data is the new Oil

- Data + Algorithms ➜ Actionable Insight ➜ $$

Big Data / Data Science

Machine Learning / "Dumb" A.I.

Model (Prediction / Classification)

Optimized Services

# Big data can generate significant financial value across sectors

**US health care**
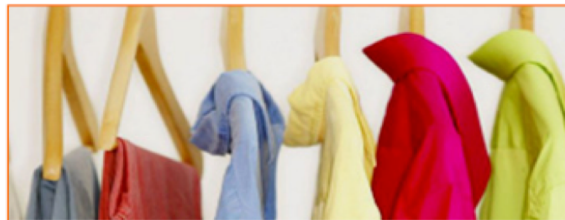- $300 billion value per year
- ~0.7 percent annual productivity growth

**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
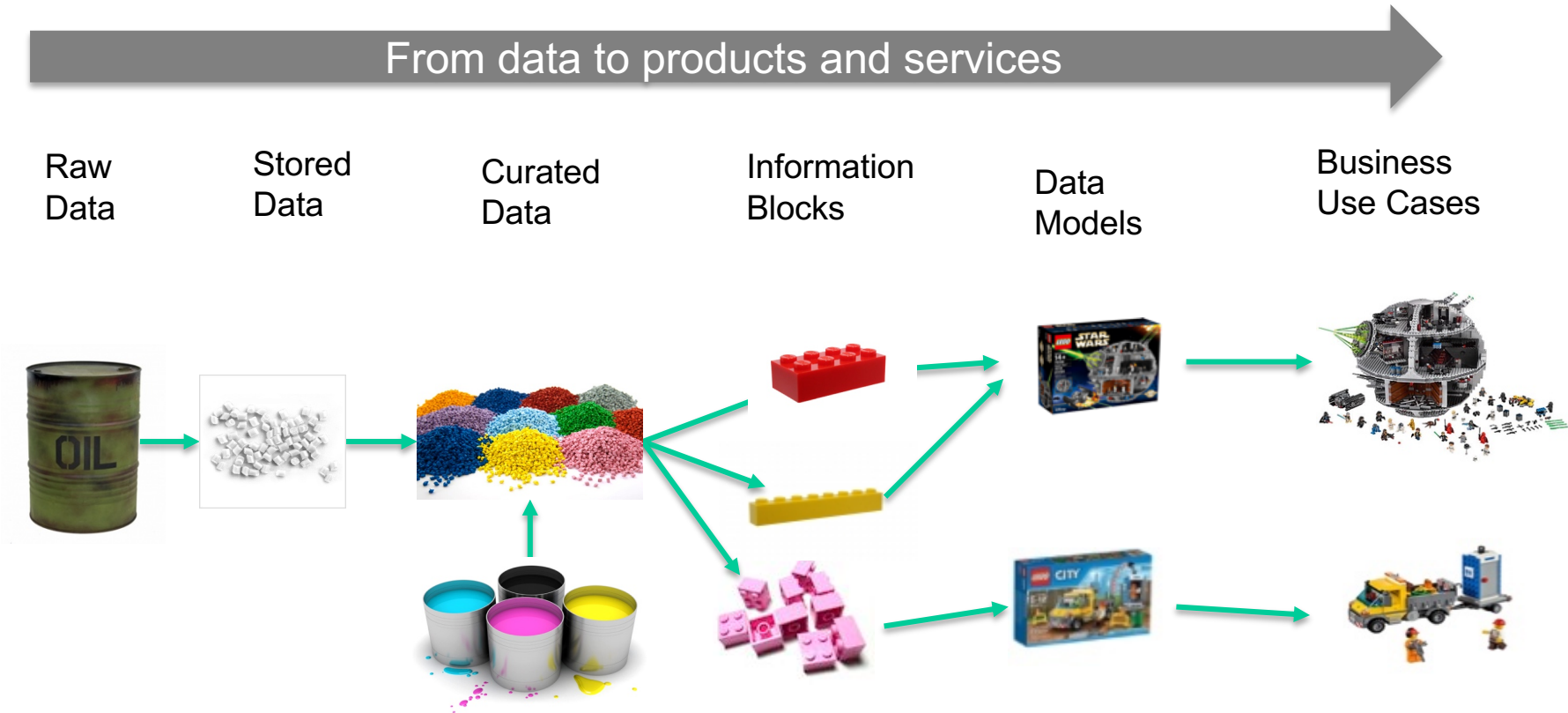
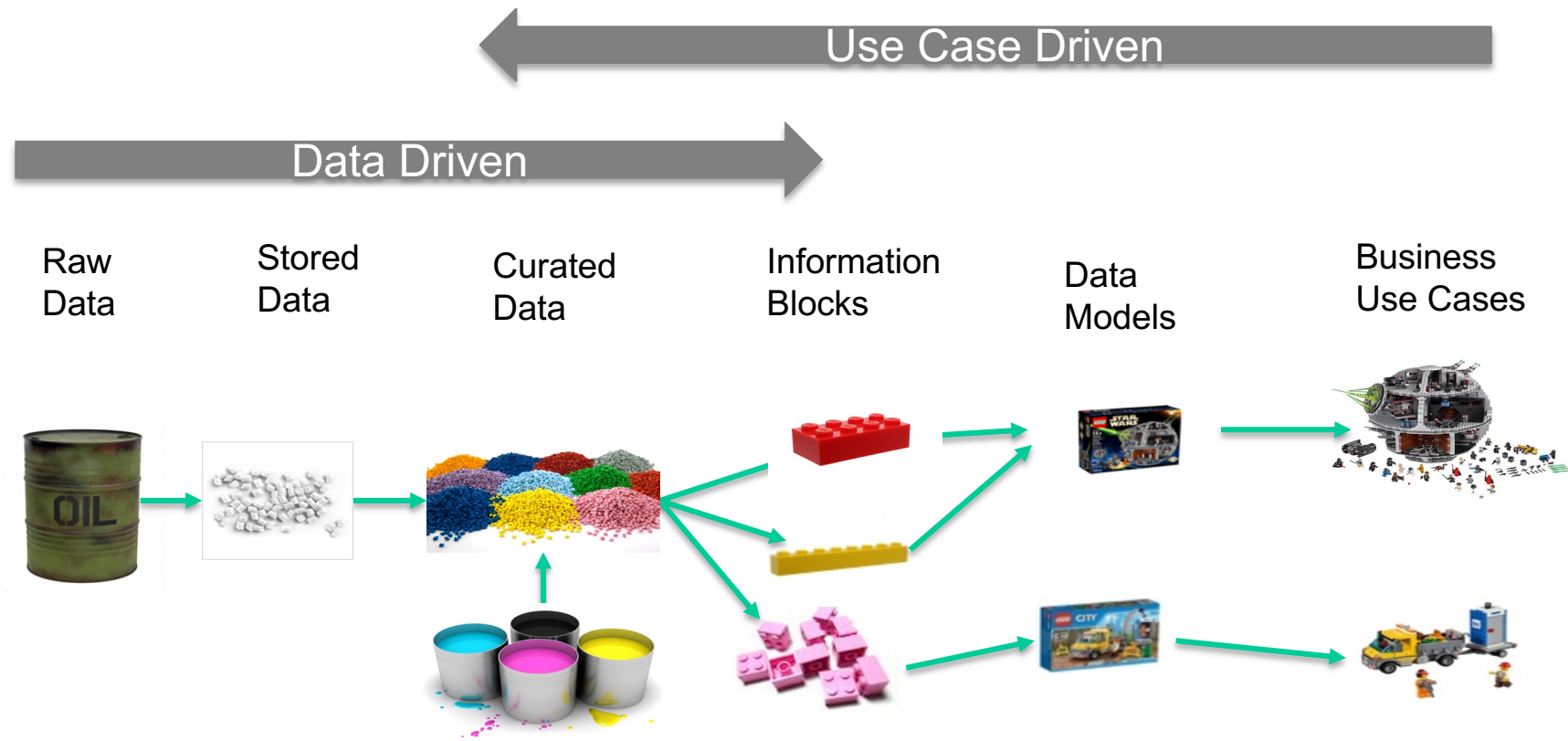SOURCE: McKinsey Global Institute analysis

# The *n*-Vs of Big Data

- Volume
  - amount of data (scale *out* not *up*)

- Velocity
  - speed of data in and out

- Variety
  - 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"*
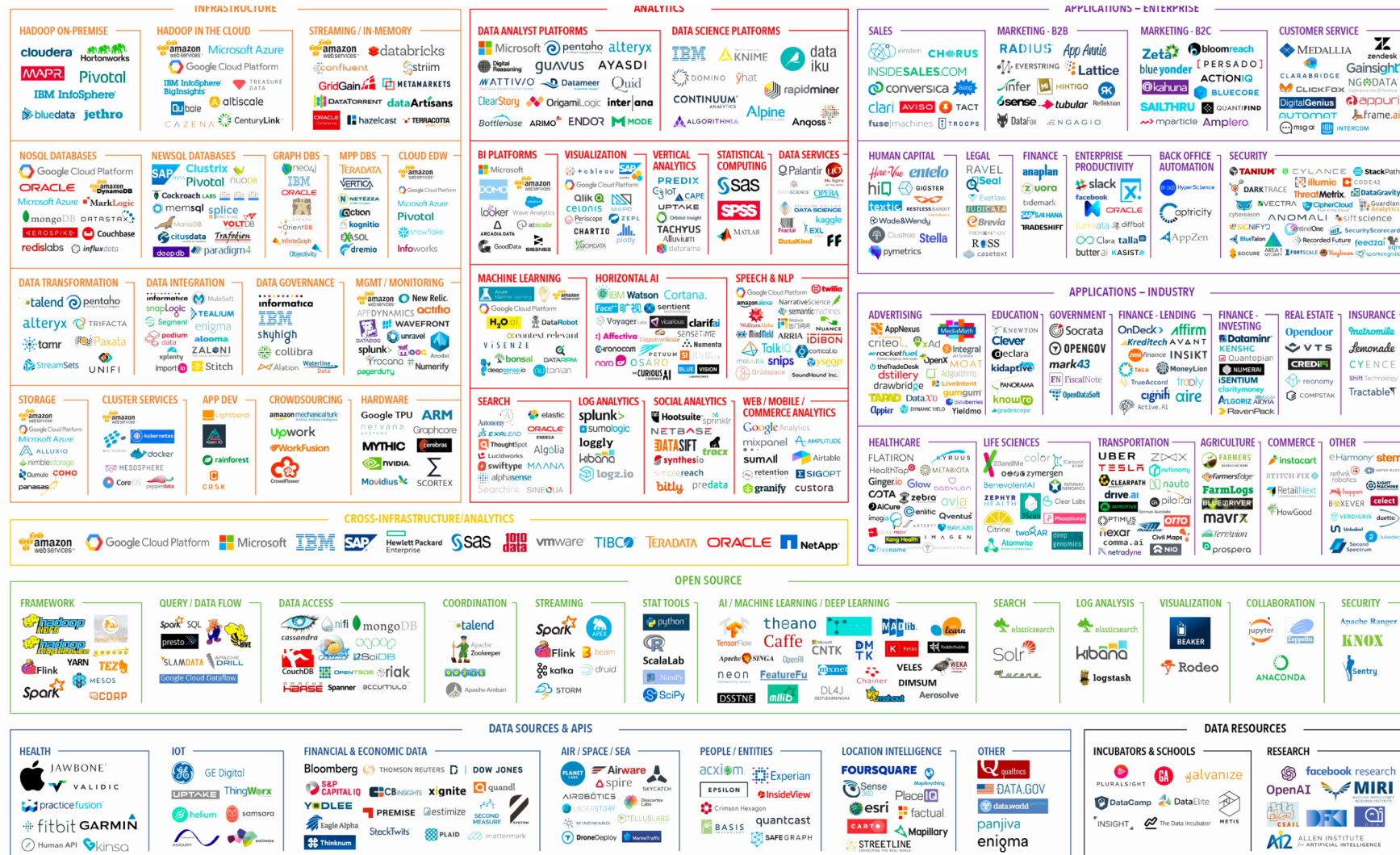
# Data vs. Traditional Assets

From data to products and services
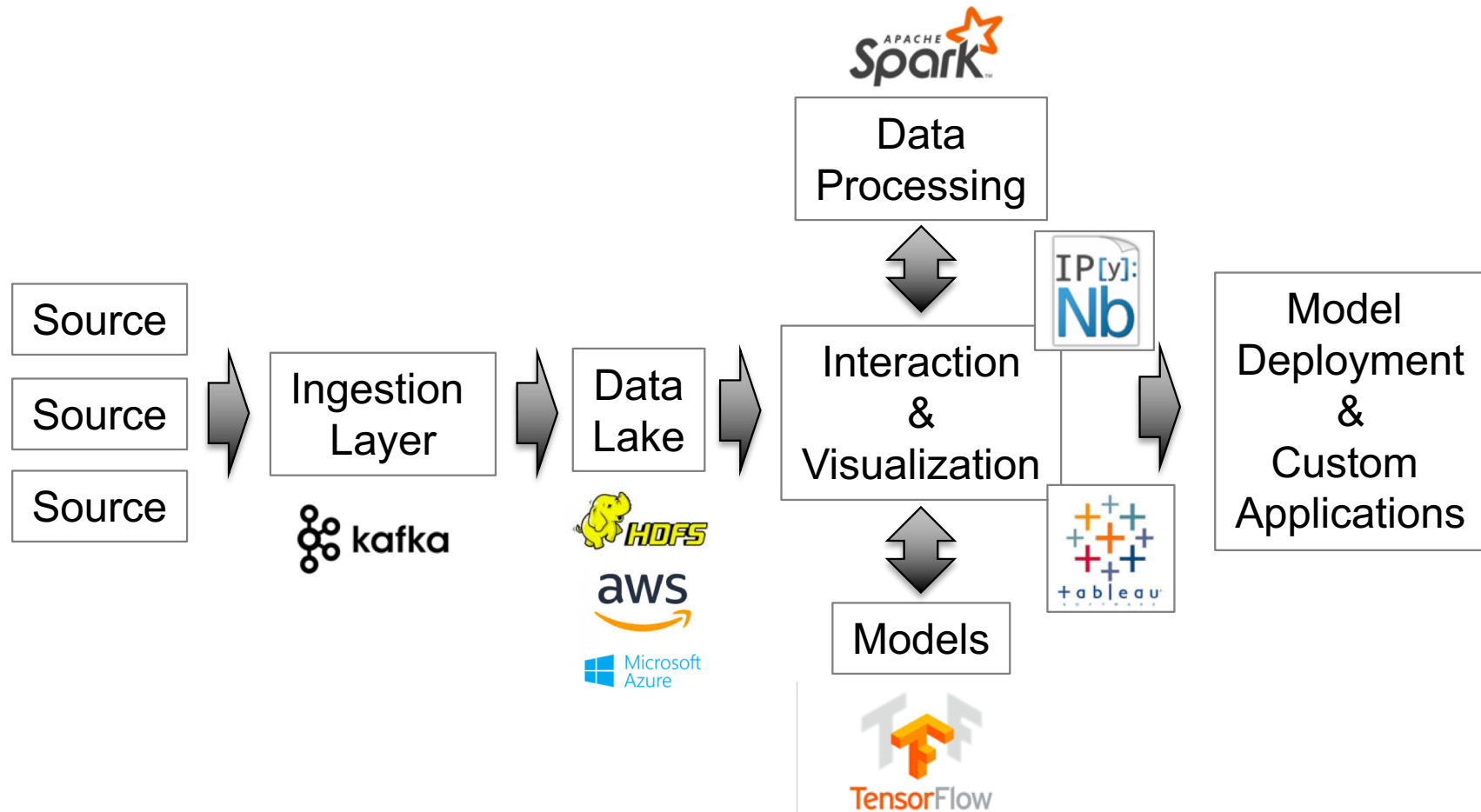
Raw Data

Stored Data

Curated Data

Information Blocks

Data Models

Business Use Cases

# Use Case or Data Driven



Use Case Driven

Data Driven

| Raw Data | Stored Data | Curated Data | Information Blocks | Data Models | Business Use Cases |
|---|---|---|---|---|---|

# Big Data Landscape

# Typical Big Data Architecture (circa 2018)

# The *n*-Vs of Big Data

- Volume
  - amount of data


- Velocity
  - speed of data in and out


- Variety (*fusing* n data sources as an input to a model)
  - range of data types and sources

# Entity-Centric Data Fusion

Higher-level apps

Captures both direct and indirect relationships

Knowledge Graph

# Three Big Data Fusion Applications

1. Anomaly Detection for Smart Cities
2. Crime Prediction using Data Fusion
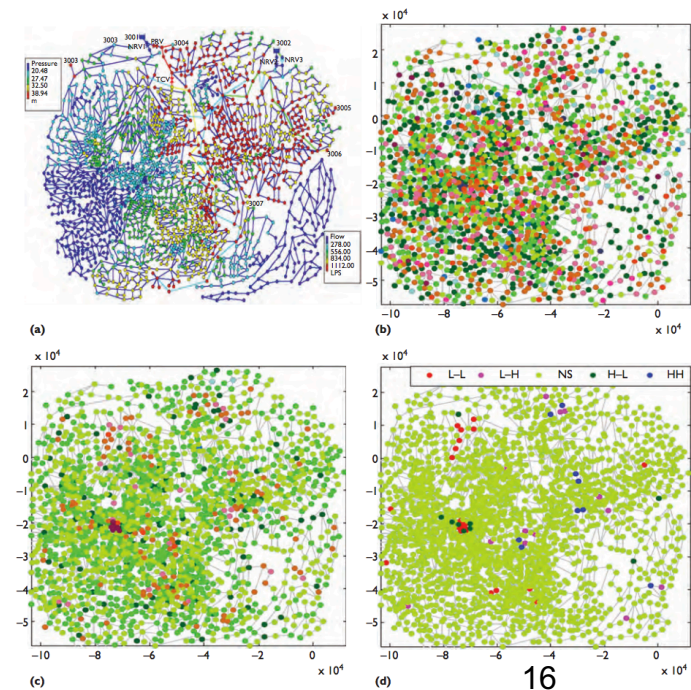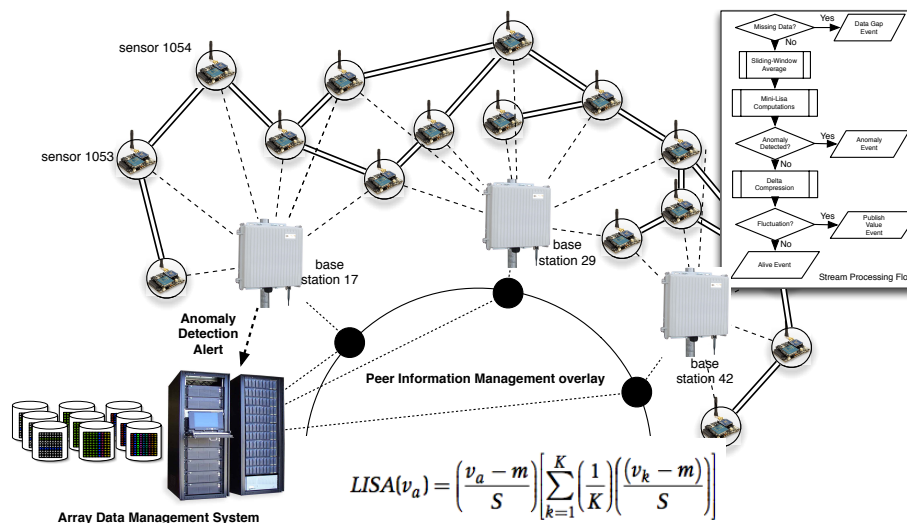3. Infrastructure Monitoring

# 1. Anomaly Detection for Smart Cities

- Detecting leaks / pipe bursts / contamination in real-time for water distribution networks
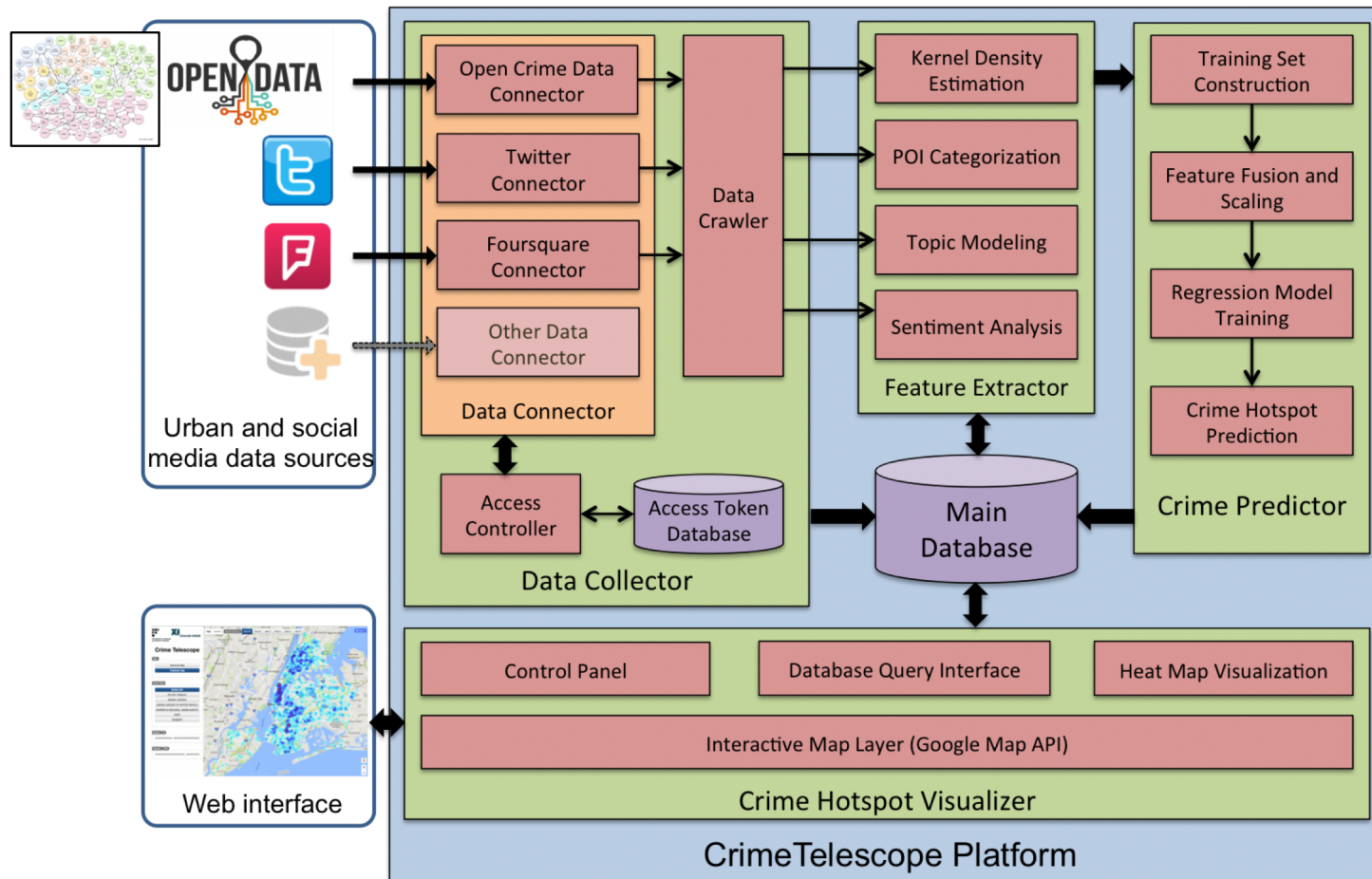
# Sensors installed in the water pipes!

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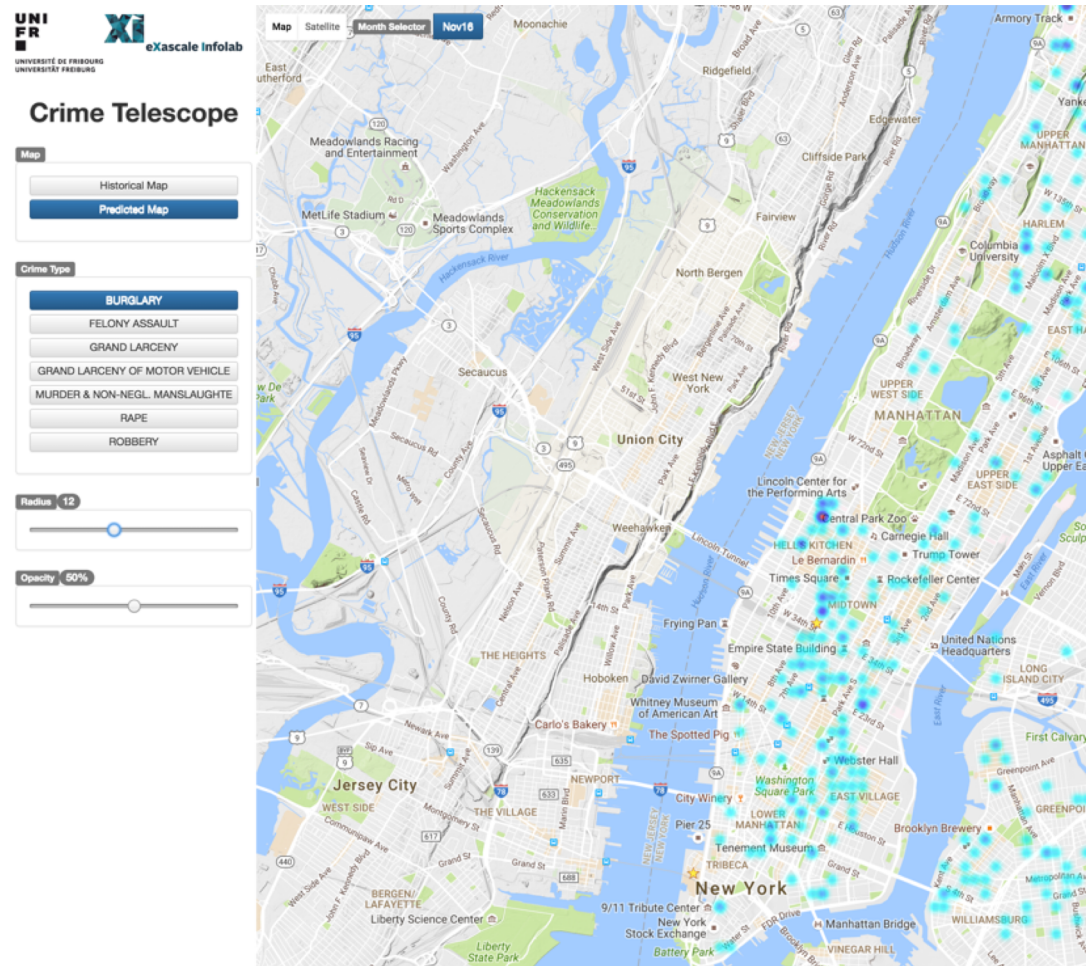- Stream processing (Storm) + Array processing (SciDB)

# 2. Crime Prediction using Big Data Fusion



- How to predict crime hotspots more accurately?

- Fusion of historical, urban & social data

CrimeTelescope: Crime Hotspot Prediction based on Urban and Social Media Data Fusion. D.Yang, T. Heaney, A. Tonon, L. Wang, P. Cudre-Mauroux. WWWJ 2017.
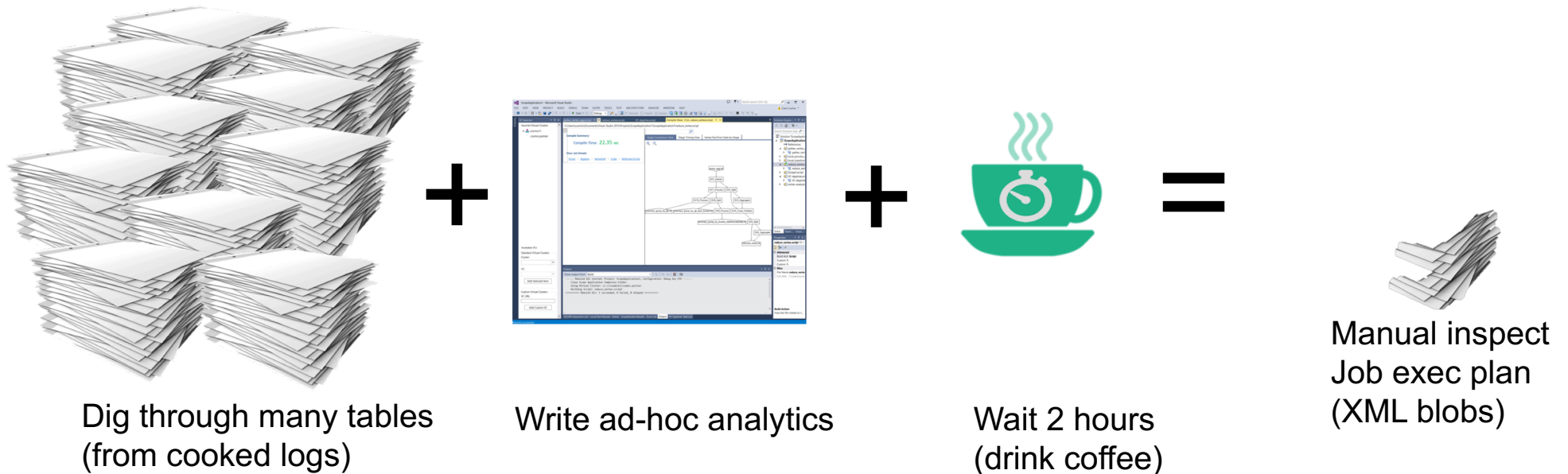
# Crime Prediction using Data Fusion

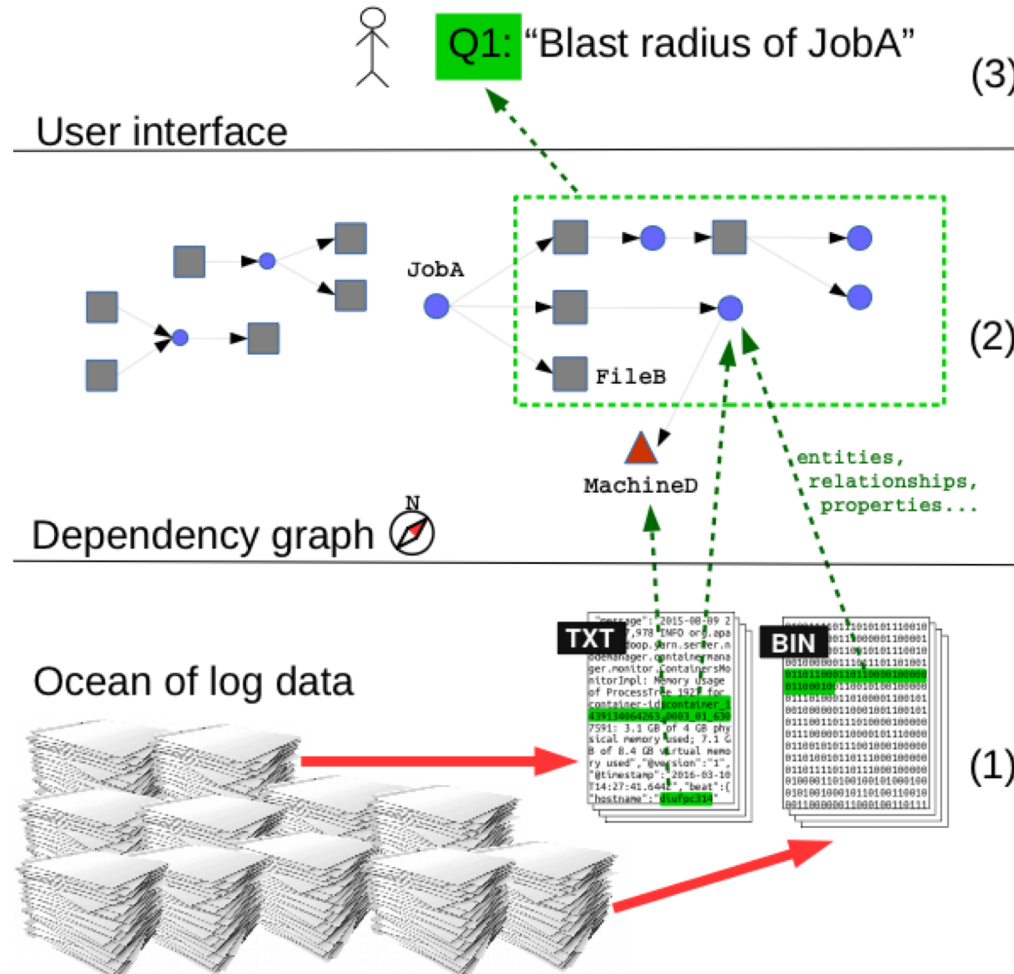# 3. Infrastructure Monitoring

Microsoft's own Metadata Lake…

Raw logs

# Example 1: Job pipeline analysis (state-of-the-art)

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



Dig through many tables
(from cooked logs)

**+**

Write ad-hoc analytics

**+**

Wait 2 hours
(drink coffee)

**=**

Manual inspect
Job exec plan
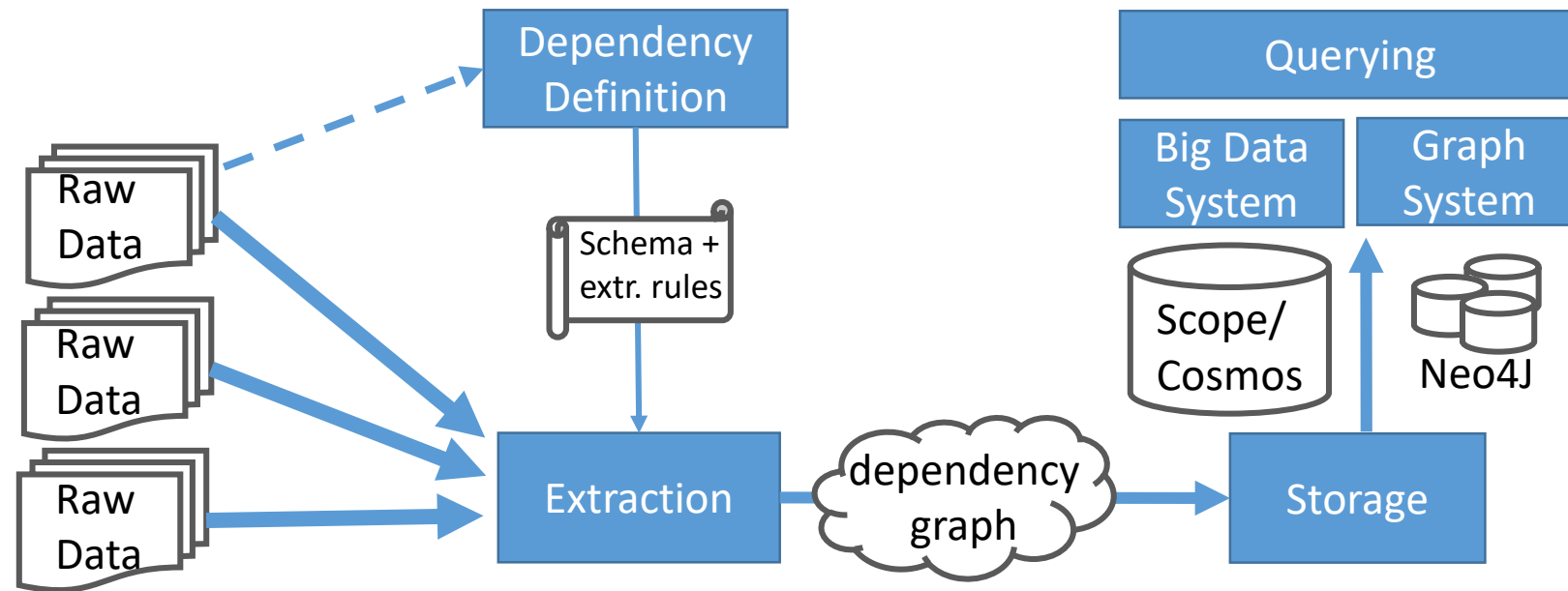(XML blobs)

# 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**
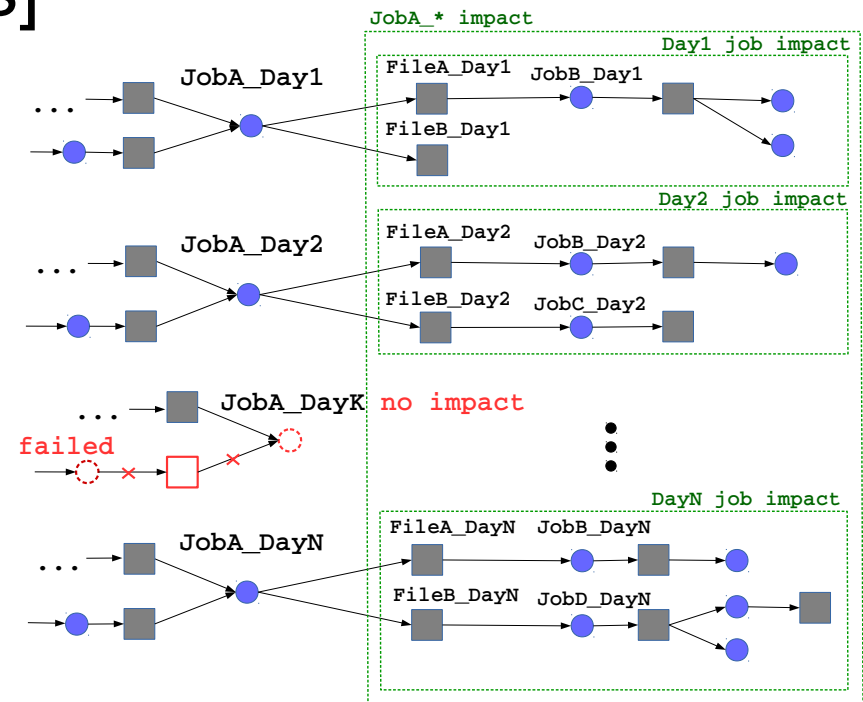
# Guider Architecture



Dependency-Driven Analytics: a Compass for Uncharted Data Oceans.
Ruslan Mavlyutovm, Carlo Curinom, Boris Asipovm, 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

# Thanks for your Attention!



https://exascale.info