

# Intelligent Knowledge Lakes

*The Age of Artificial Intelligence and Big Data*



**University of Malaya**

**8 January 2020**

**Dr. Amin Beheshti**

**Director, AI-enabled Processes (AIP) Research Centre**

**Director, Data Analytics Research Lab**

**Web:** <https://aip-research-center.github.io/>



**MACQUARIE**  
University  
SYDNEY · AUSTRALIA

# Introduction to **Big Data**

# Introduction to Big Data

3

**Application**



**Set of Related Programs**



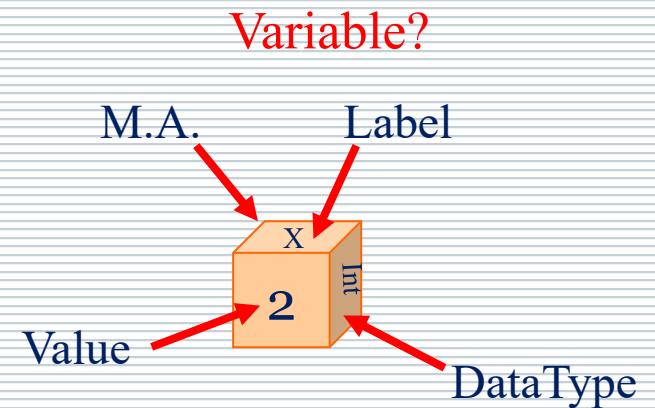
**Set of Related Functions**



**Set of Related Statements**



- Assignment
- Selection
- Iteration
- Jump
- ...



X ← 2

# Introduction to Big Data

4

**Application**



**Set of Related Programs**



**Set of Related Functions**



**Set of Related Statements**

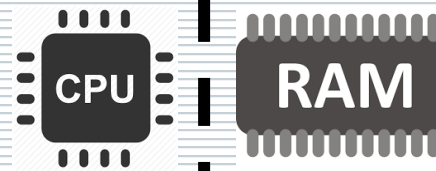
- Assignment
- Selection
- Iteration
- Jump
- ...

**Data Structures?**

Variable  
Array  
Record  
Class  
.  
.  
.  
ADT

Abstraction

Reusability



# Introduction to Big Data

5

**Application**

Execute on

**Computing Platform**

Hardware

Software



**Platform Independent Application**

# Introduction to Big Data

6

Application



Platform Independent Application  
(e.g. Web Applications)



Tim Berners-Lee

# Introduction to Big Data

7

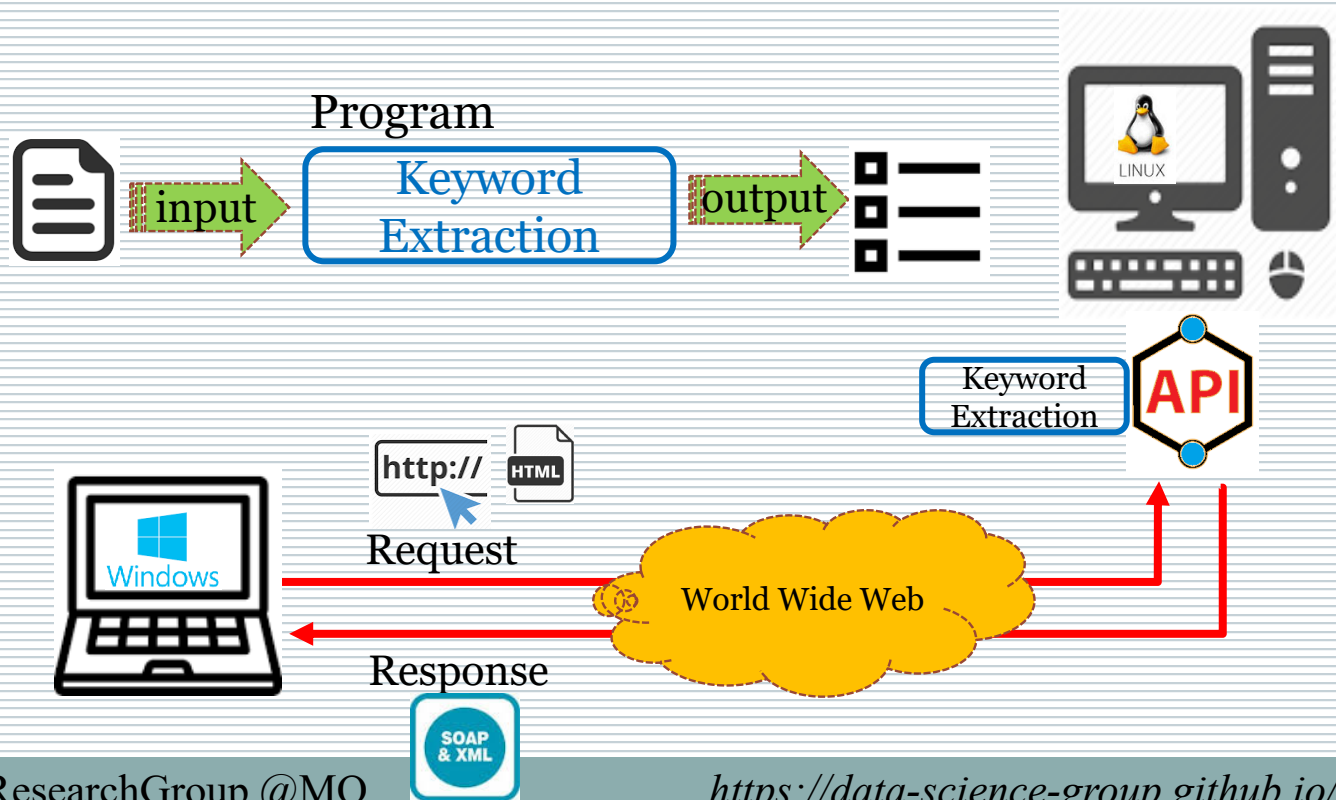
Application



Platform Independent Application  
(e.g. Web Applications)



Web Services



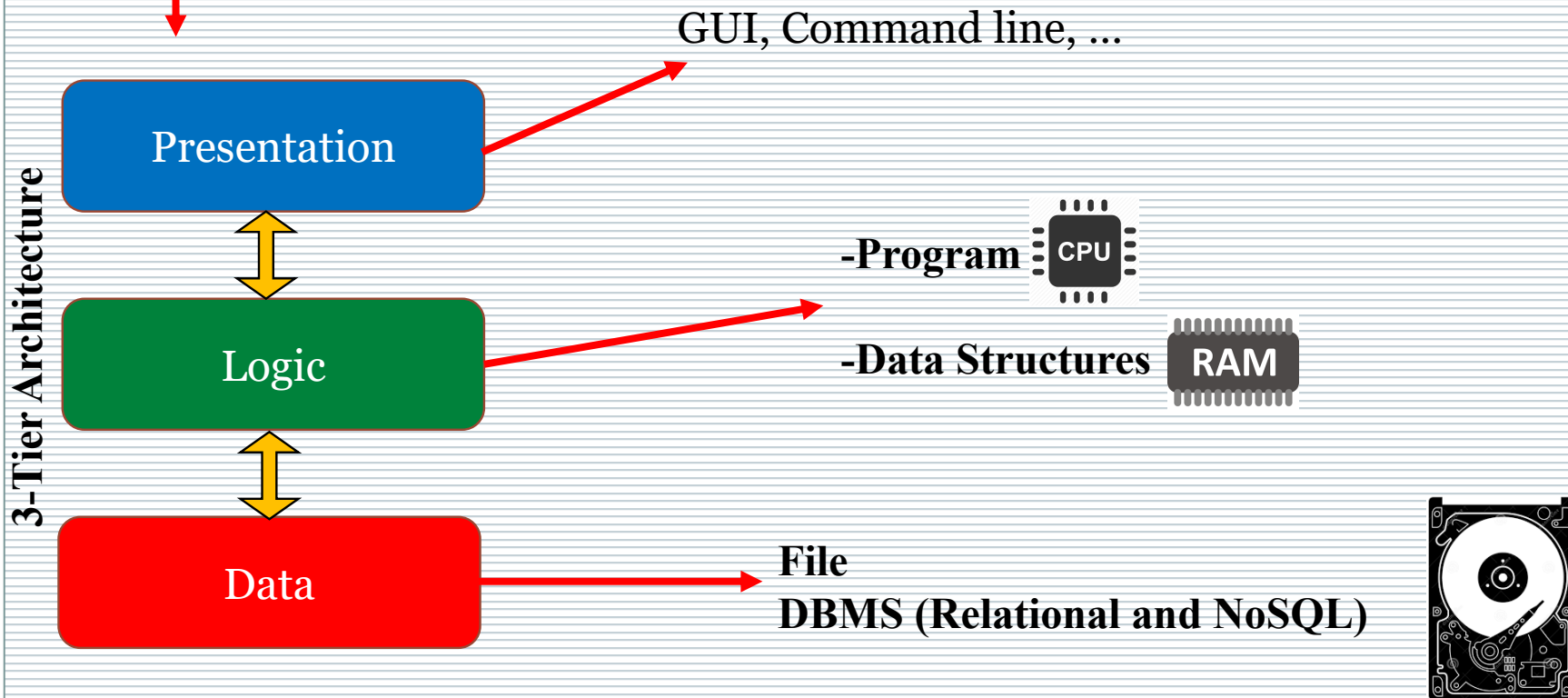
-API Engineering  
-Microservices

# Introduction to Big Data

8

## Application

Architecture





# What is **Data** ?

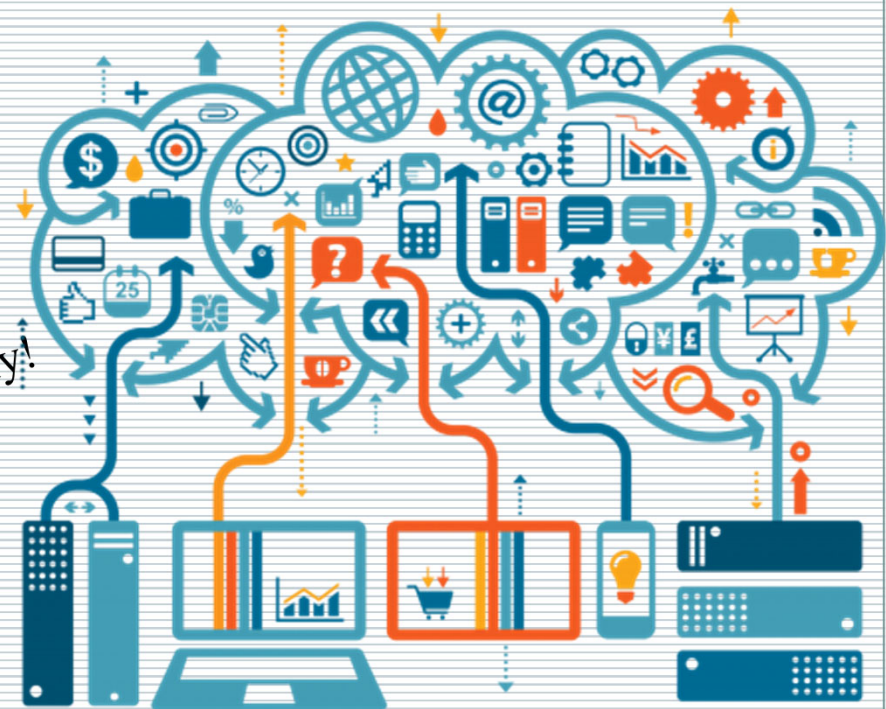
# What is Data ?

10

Every day, we create **2.5 quintillion** bytes of data.

- posts to social media sites
- sensors used to gather climate information
- digital pictures and videos
- purchase transaction records
- cell phone GPS signals
- ...

- 500 Million Tweets sent each day!
- 5.75 BILLION Facebook likes every day.
- 3.6 Billion Instagram Likes each day.
- 4.3 BILLION Facebook messages posted daily!
- 6 BILLION daily Google Searches!
- ...

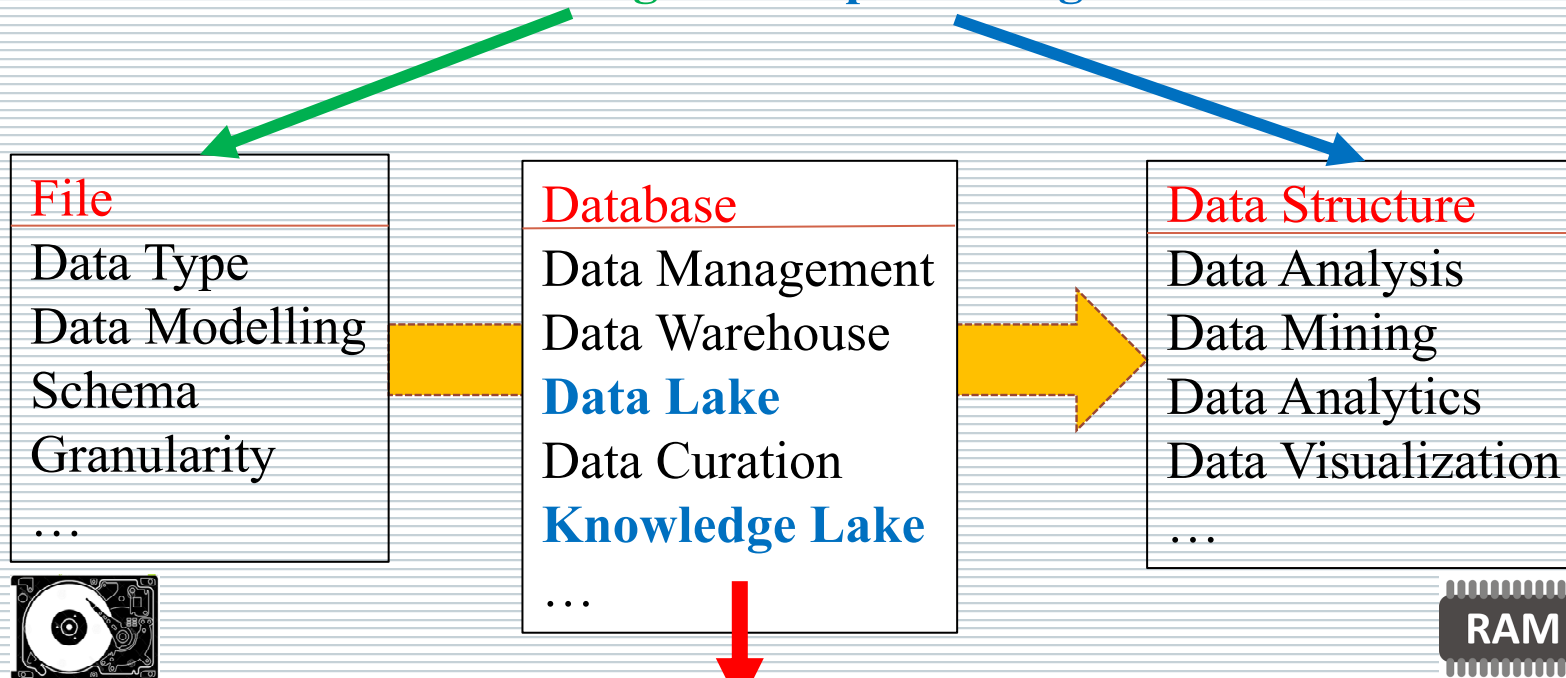


<https://www-01.ibm.com/>; <http://www.internetlivestats.com/>; <http://semeon.com/>

# What is Data ?

11

In computing, **data** is information that has been translated into a form that is efficient for **storage** and/or **processing**.



Preparing the Data for Processing & Analytics  
(**Organizing** and **Curating**)

# What is Metadata ?

12

We are **Tracing** everything:

- What is happening?
- Who is doing that?
- Where it is happening?
- When?
- Why?
- How?
- ...

**Cross-Cutting  
Aspects**

- Provenance
- Versioning
- Privacy
- Security
- ...

• **Smart Phones**, tracks:

- Our location,
- Our speed,
- What apps we are using,
- What music we listen to,
- ...

• **Smart TVs**, tracks:

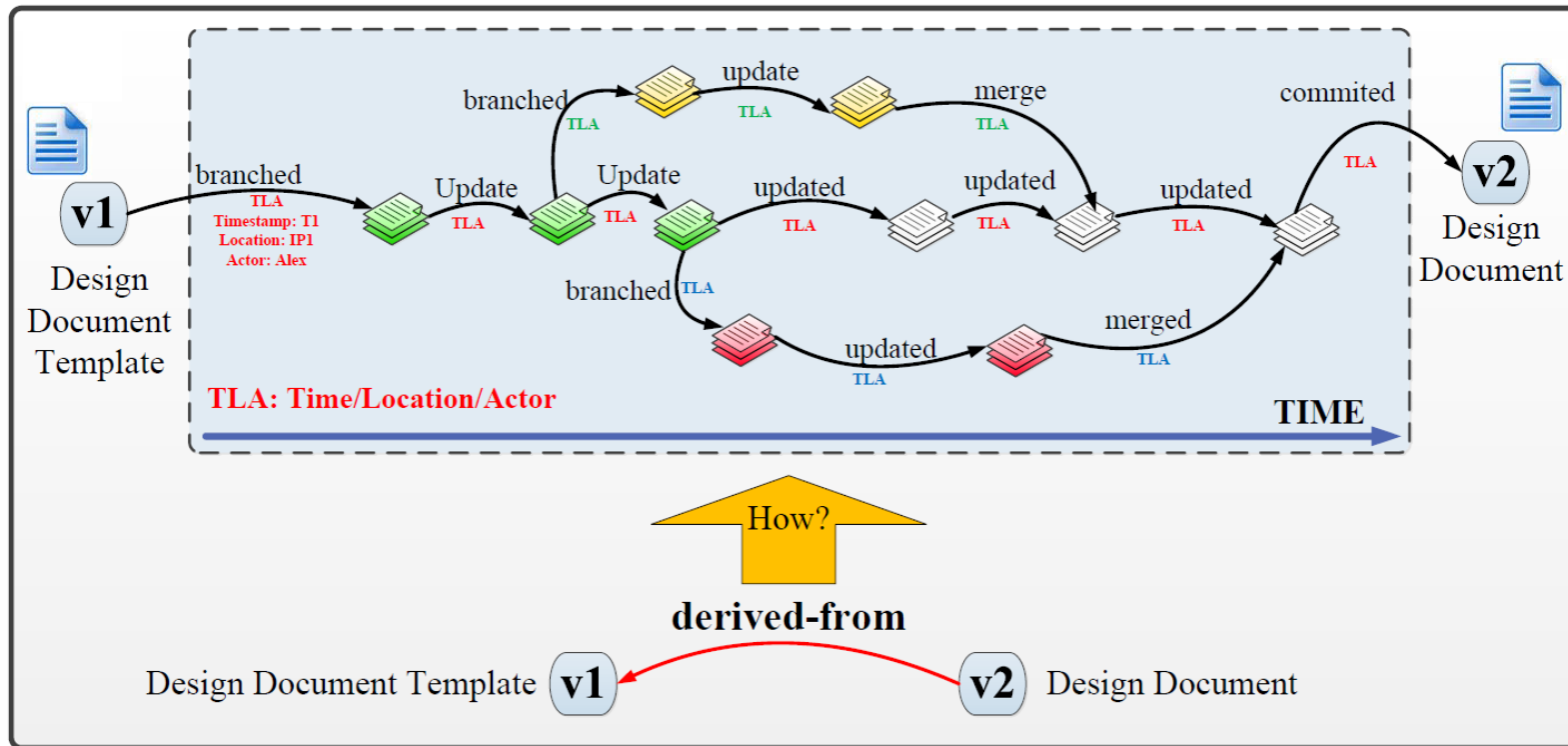
- Channels we watch,
- Time and duration,
- Apps we use,
- ...

• **Smart Watches**, tracks:

- Our health signs,
- Our activity,
- Location,
- ...

# What is Metadata ?

**Provenance**, a kind of metadata, refers to the documentation of an object's lifecycle. This documentation (often represented as a graph) should include all the information necessary to reproduce a certain piece of data or the process that led to it.



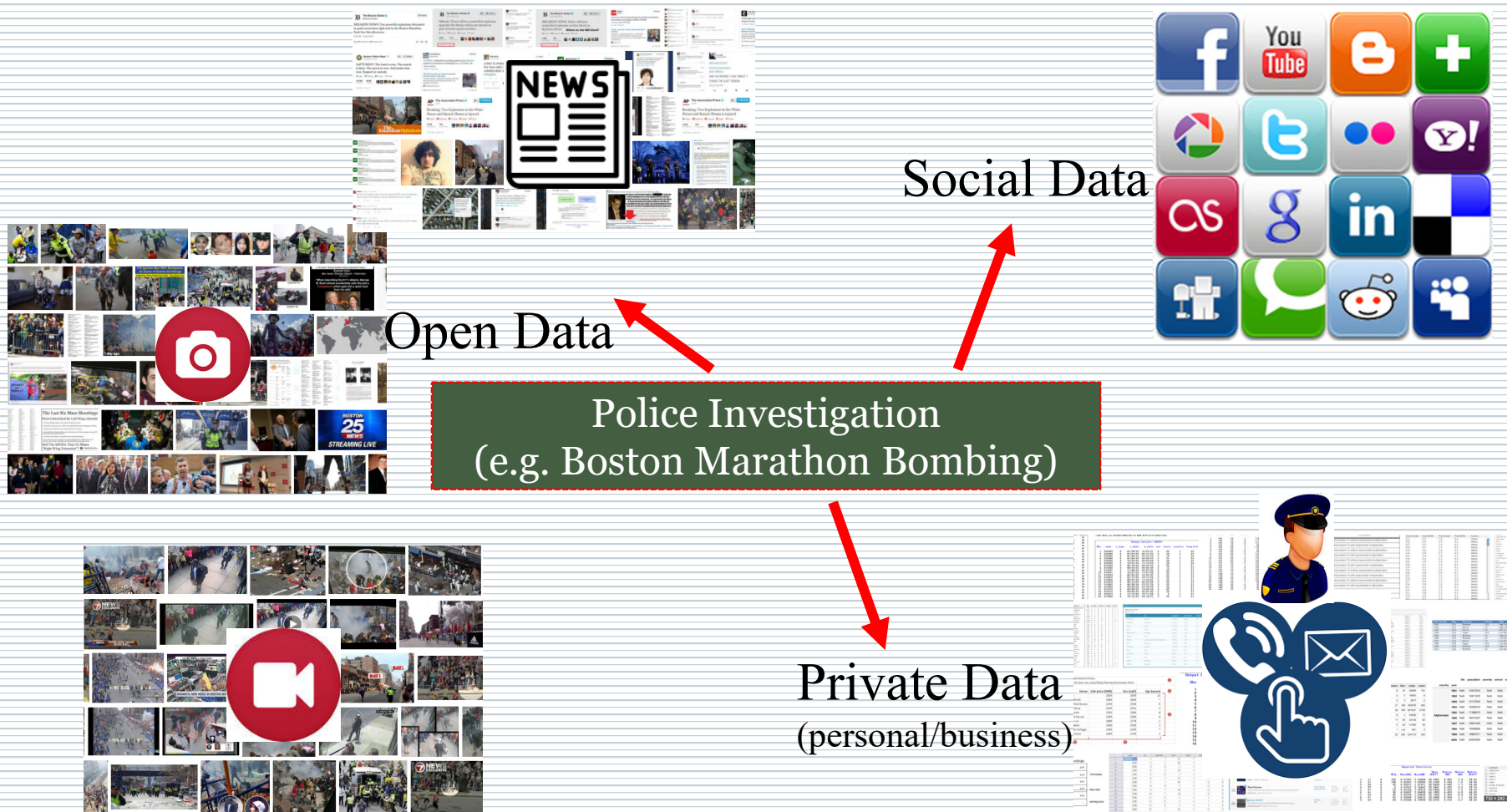
Beheshti et al. "Temporal provenance model (TPM): model and query language, 2012  
[https://www.w3.org/2005/Incubator/prov/wiki/What\\_Is\\_Provenance](https://www.w3.org/2005/Incubator/prov/wiki/What_Is_Provenance)

# What is **Big Data** ?

Example

# What is Big Data ?

15



Beheshti et al. "ProcessAtlas: A scalable and extensible platform for business process analytics", Software: Practice and Experience, 2018

Copyright © DataAnalyticsResearchGroup @MQ

<https://data-science-group.github.io/>

Example

# What is Big Data ?

16

- Typical properties of the big data:
- wide physical distribution
  - diversity of formats
  - non-standard data models
  - independently-managed
  - heterogeneous semantics

Private Data  
(personal/business)



# What is **Big Data** ?

17

- **Big data** refers to our ability to collect and analyse the ever expanding amounts of **data** and **meta-data** that we are generating every second!
- **Big data** can be seen as a massive number of small **data islands** from Private (Personal/Business), Open and Social Data.

**Organizing**, **Curating**, **Analysing** and **Presenting**  
this data is *challenging* and of high interest.

# Organizing Big Data

# Organizing Big data

19

- How to store vast amount of noisy data (varying from structured entities to unstructured documents) being generated on a continuous basis ?

## The **Four V's** of **Big Data**

### Volume

the vast amounts of data generated every second.

### Variety

the increasingly different types of data.

### Velocity

the speed at which new data is generated and moves around.

### Veracity

the reliability and predictability of imprecise data types.

# Big data - Volume

# Big data - Volume

21

**Volume**, the quantity of data to be stored, is a key characteristic of Big Data.

STORAGE CAPACITY UNITS		
TERM	CAPACITY	ABBREVIATION
Bit	0 or 1 value	b
Byte	8 bits	B
Kilobyte	1024* bytes	KB
Megabyte	1024 KB	MB
Gigabyte	1024 MB	GB
Terabyte	1024 GB	TB
Petabyte	1024 TB	PB
Exabyte	1024 PB	EB
Zettabyte	1024 EB	ZB
Yottabyte	1024 ZB	YB

\* Note that because bits are binary in nature and are the basis on which all other storage values are based, all values for data storage units are defined in terms of powers of 2. For example, the prefix *kilo* typically means 1000; however, in data storage, a kilobyte =  $2^{10}$  = 1024 bytes.

Database Systems, Design, Implementation, & Management, 13<sup>th</sup> Edition, Carlos Coronel – Steven Morris

# Big data - Volume

22

**Volume**, the quantity of data to be stored, is a key characteristic of Big Data.

How to deal with storing large volume of data ?

## Scale Up:



Keep the same number of Systems, but migrating each system to a larger System.

e.g. Changing from a server with 16 CPU cores and 1 TB storage system to a server with 64 CPU cores and a 100 TB storage system.

## Scale Out:



When the workload exceeds the capacity of a server, the work load is spread out across a number of servers.

This is also referred to as **Clustering**.

## Notice:

It is cheaper to buy ten 100 TB storage systems than it is to buy a single 1 PB storage system

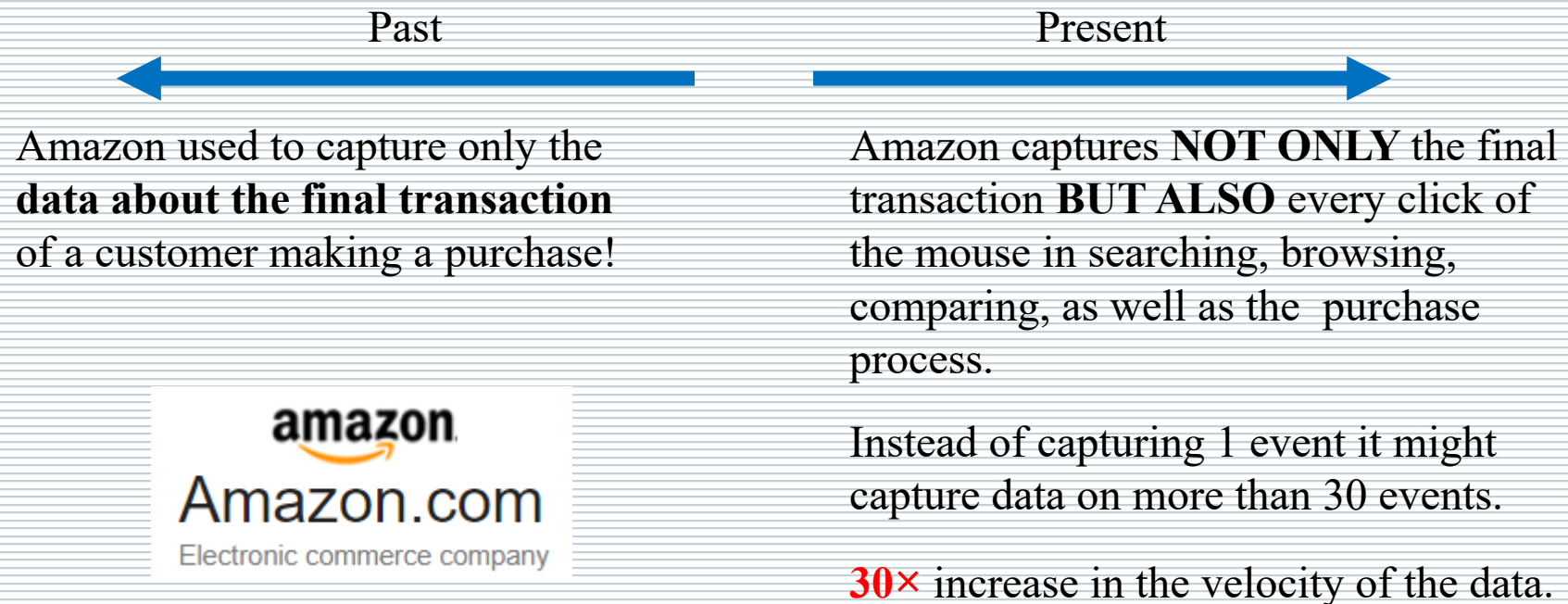
# Big data - Velocity

# Big data - Velocity

24

**Velocity**, refers to the **rate at which new data enters the system** as well as the **rate at which the data must be processed**.

Example:





# Big data - Velocity

25

**Velocity**, refers to the **rate at which new data enters the system** as well as the **rate at which the data must be processed**.



The velocity of processing can be broken down into: **Stream** and **Feedback Loop Processing**

**Stream Processing**, requires analysis of the data stream as it enters the system.

(Focus on the INPUT)

*Example:*

CERN Large Hadron Collider (the largest and most powerful particle accelerator in the world) experiments produce about 600 TB per second of raw data.



All this data can not be processes, accordingly scientists created algorithms to decide ahead of time which data will be kept; and to **filter the data down** to only about 1 GB per second.

# Big data - Velocity

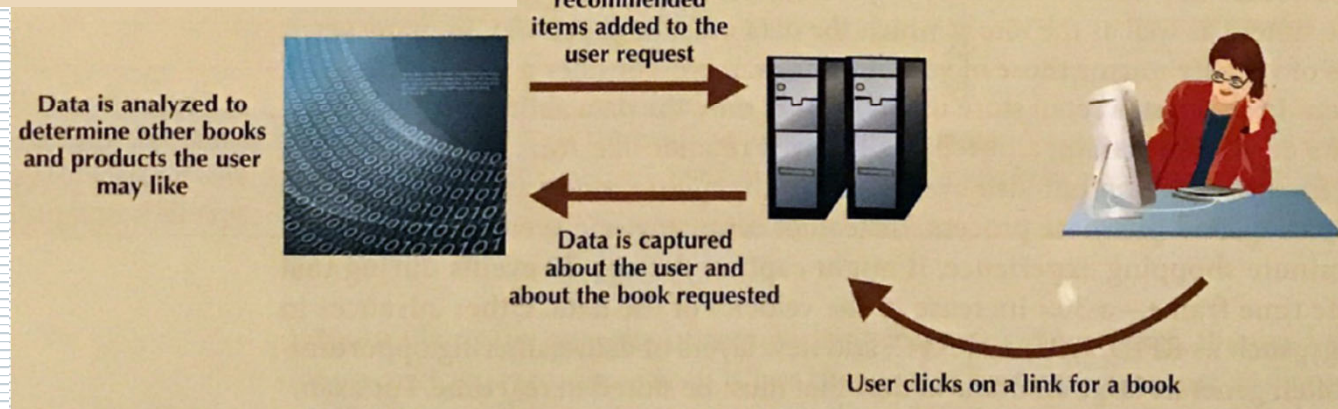
26

**Velocity**, refers to the rate at which new data enters the system as well as the rate at which the data must be **processed**.



The velocity of processing can be broken down into: **Stream** and **Feedback Loop** Processing

**Feedback Loop Processing**, refers to the analysis of the data to produce actionable results. (Focus on the OUTPUT)



# Big data - Variety

# Big data - Variety

28

**Variety**, refers to the vast array of **formats and structures in which the data may be captured**: structured, unstructured and semi-structured.

**Structured Data**, is data that has been organized to fit a predefined data model.

**Unstructured Data**, is data that is not organized to fit into a predefined data model.

**Semi-structured Data**, combines elements of both Structured and Unstructured.

# Data Persistence

# Data Persistence

30

## Various related Data Islands:

From open to private and social data.

## Various Technologies to persist the big data:

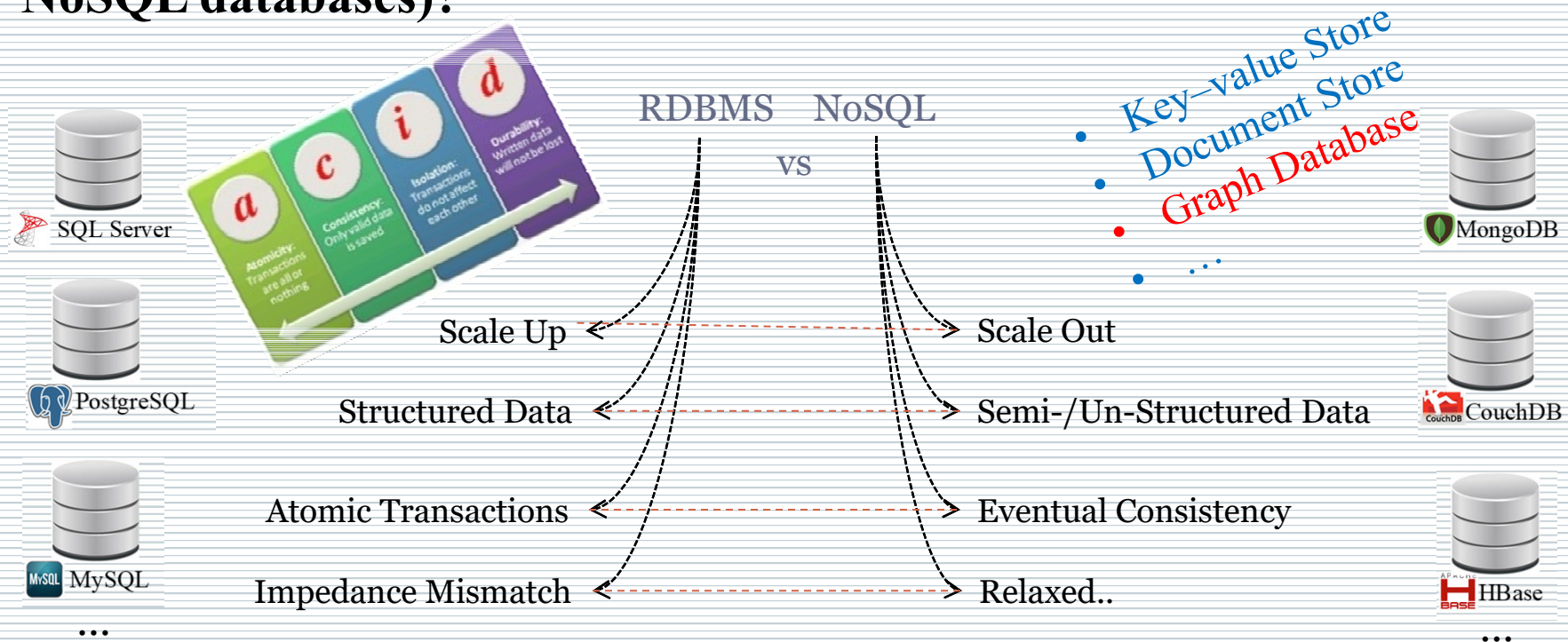
From Relational to NoSQL



# Data Persistence

31

- How to store vast amount of noisy data (varying from structured entities to unstructured documents) being generated on a continuous basis ?
- **What technology to use for persisting the data (from Relational to NoSQL databases)?**



# Introduction to Data Lakes

32

## DBMS and Relational DBMS



# What Is a Database Management System (DBMS) ?

33

A collection of files that store the data

+

A big C program written by someone else that accesses and updates those files for you!



# What is a Relational DBMS?

34

A **Relational Database Management System (RDBMS)** is a database management system (DBMS) based on the **Relational Model** invented by *Edgar F. Codd* at IBM's San Jose Research Laboratory.

In the **Relational Model**, all data must be stored in relations (tables), and each relation consists of rows and columns.

## Where are RDBMS used ?

- Backend for traditional “database” applications
- Backend for large Websites
- Backend for Web services

NoSQL  
Not only SQL!

# NoSQL

36

**NoSQL**, is a new generation of database management systems that is not based on the traditional Relational Database Model.

## NoSQL DATABASES

NoSQL CATEGORY	EXAMPLE DATABASES	DEVELOPER
Key-value database	Dynamo Riak Redis Voldemort	Amazon Basho Redis Labs LinkedIn
Document databases	MongoDB CouchDB OrientDB RavenDB	MongoDB, Inc. Apache OrientDB Ltd. Hibernate Rhinos
Column-oriented databases	HBase Cassandra Hypertable	Apache Apache (originally Facebook) Hypertable, Inc.
Graph databases	Neo4J ArangoDB GraphBase	Neo4j ArangoDB, LLC FactNexus

Database Systems, Design, Implementation, & Management, 13<sup>th</sup> Edition, Carlos Coronel – Steven Morris

# NoSQL – Key-Value Databases

37

**Key-Value Databases**, are conceptually the simplest of the NoSQL data models.

Data will be stored as a collection of Key-Value pairs.



An **identifier** for a value

The value can be anything such as:

- Text
- Document (XML/JSON)
- Image
- etc.

The Database does not attempt to understand the content of the value!

(It is the role of the application to analyse and understand the content)

Database Systems, Design, Implementation, & Management, 13<sup>th</sup> Edition, Carlos Coronel – Steven Morris

# NoSQL – Document Databases

38

**Document Databases**, are conceptually similar to Key-Value Databases, and can be considered as a sub-type of Key-Value Databases.

The Value component can only contain Documents !

The Document can be in any encoded format, such as:

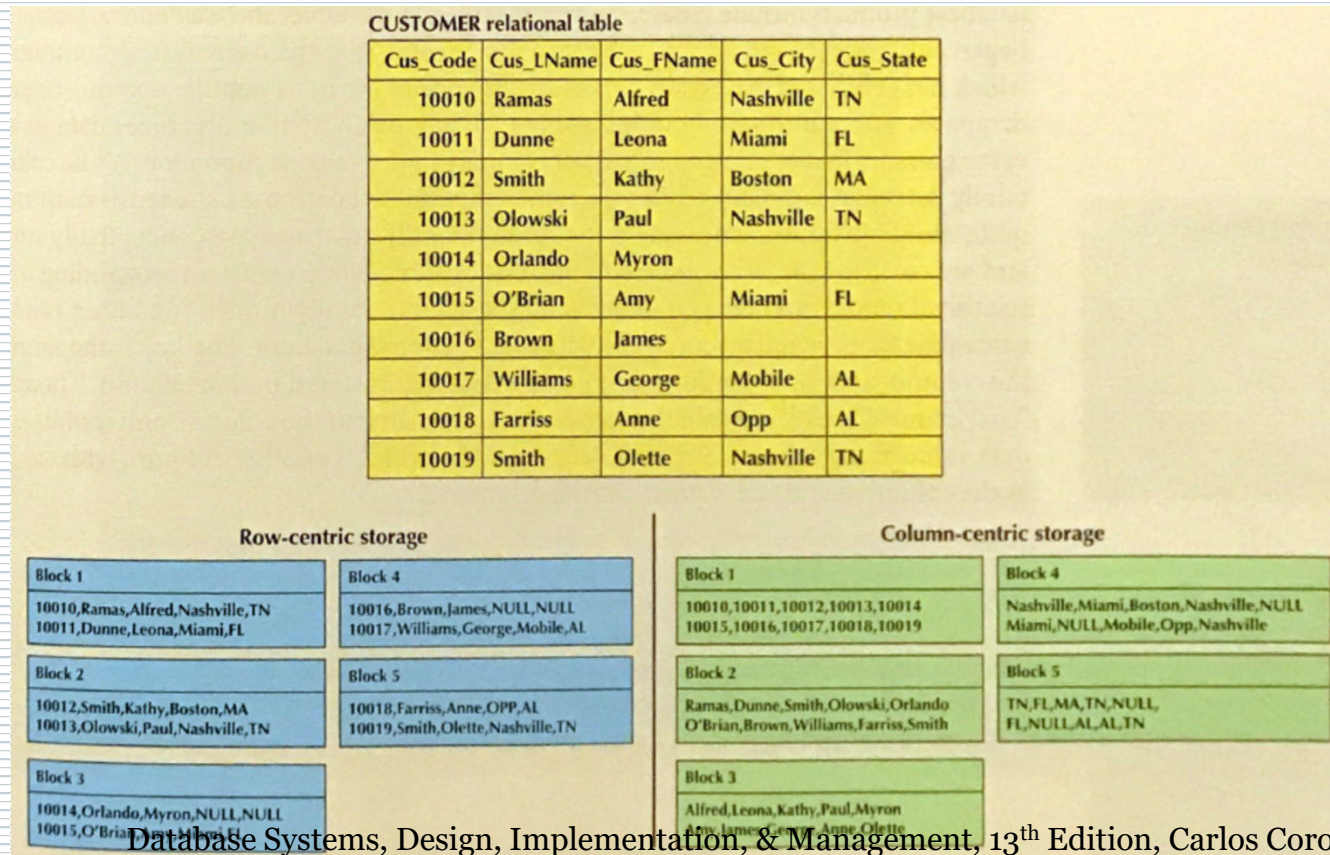
- XML
- JSON
- BSON (Binary JSON)
- etc.

Unlike Key-Value Databases, the Document Database do attempt to understand the content of the Value!

# NoSQL – Column-Oriented Databases

39

**Column-Oriented Database**, stores the data in **blocks by column** instead of by rows.



Database Systems, Design, Implementation, & Management, 13<sup>th</sup> Edition, Carlos Coronel – Steven Morris

# NoSQL – Column-Oriented Databases

40

**Column-Oriented Database**, stores the data in **blocks by column** instead of by rows.

This type of database:

- works very well for databases that are primarily used to run queries over few columns but many rows, as is done in many reporting systems and data warehouses.
- Would be inefficient for processing transactions since Insert, Update and Delete activities would be very disk intensive.

Example: HBase, HyperTable, Cassandra.



Developed by Facebook; one of the most popular Column-Oriented DBs.



# NoSQL – Graph Databases

41

**Graph Database**, is a NoSQL DB based on Graph theory to store data about relationship-rich environments.

A Mathematical and Computer Science field that models relationships (edges) among objects called nodes.

Modelling and storing data about relationships is the focus of Graph Databases.

Interest in Graph Databases originated in the area of **social networks**.



Beheshti et al., "Galaxy: A Platform for Explorative Analysis of Open Data Sources", **EDBT**, 2016.

<http://www.cse.unsw.edu.au/~sbeheshti/EDBT16/>

# Indexing Big Data

42

## Search NoSQL Documents:

- **Elasticsearch** can be used to search all kinds of documents.
- Elasticsearch uses **Lucene** (an indexing and search library) and tries to make all its features available through the JSON and APIs.

<https://www.elastic.co/products/elasticsearch>  
<https://lucene.apache.org/>

## Database Service:

- Dozens of new DBs! how do we choose which DB to use?
- Solution:
  - Manage multiple database technologies and weave them together at the app layer..
  - Make this service accessible through a single API

# Database as a Service

43

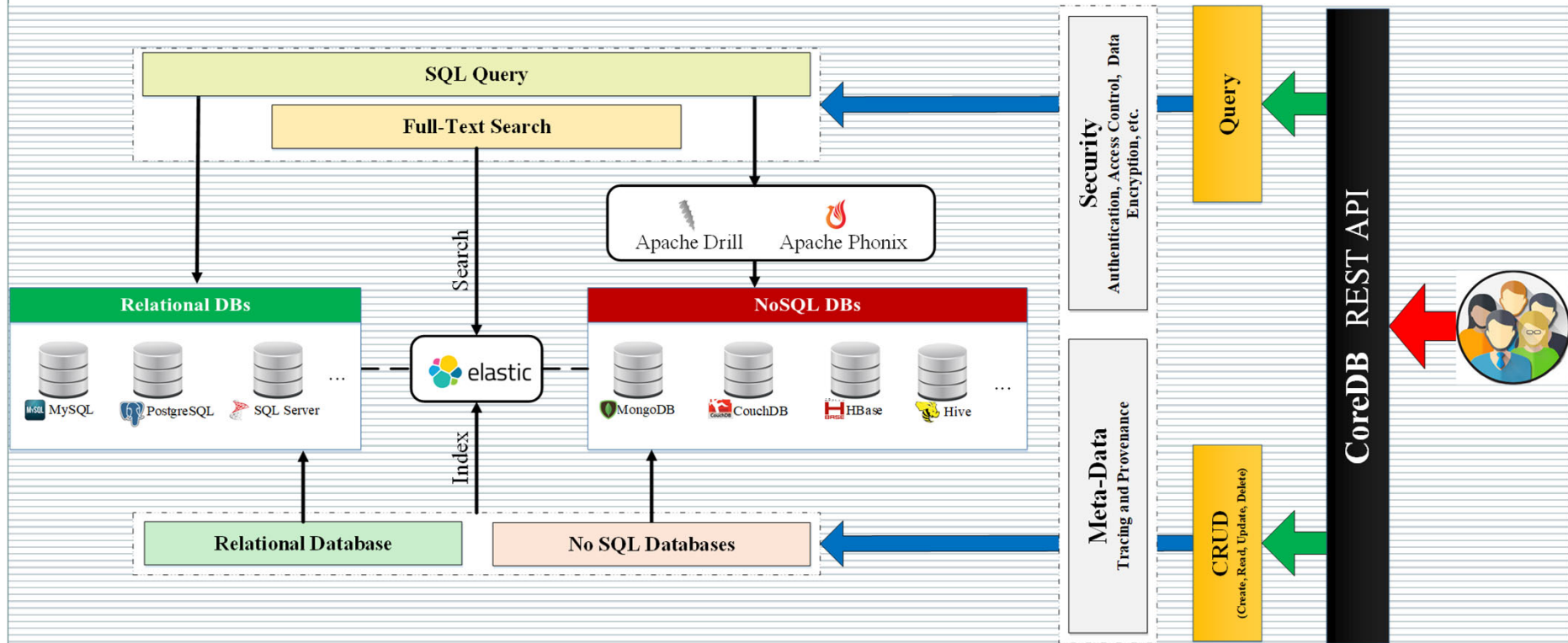
- Database as a service (DBaaS) is a cloud computing service model that provides users with some form of access to a database without the need for setting up physical hardware, installing software or configuring for performance.  
<https://www.contentful.com/>  
<https://orchestrate.io>
- Dozens of new DBs! how do we choose which DB to use?
- Solution: **Data Lakes**
  - Manage multiple database technologies and weave them together at the app layer.
  - Make this service accessible through a single API to support CRUD and Query data.

# Data Lake

# Data Lake

45

A **Data Lake** is a storage repository that holds a vast amount of raw **data** in its native format, including structured, semi-structured, and unstructured **data**.



Beheshti et al., **CoreDB: a Data Lake Service**, CIKM 2017; <https://github.com/unsw-cse-soc/CoreDB>

# Data Lake vs. Data Warehouse

46

## Data Lake

Raw Data

Structured, Semi-Structured, Unstructured

Schema-on-Read

Designed for Low Cost Storage

Highly Agile, Configured and re-Configure as needed

Maturing

Data Scientists and Analysts

Data

Variety

Processing

Volume

Agility

Security

Users

## Data Warehouse

Processed Data

Structured

Schema-on-Write

Expensive for Large Data Volumes

Less Agile, Fixed Configuration

Mature

Business Analysts

A **data warehouse** is a system used for reporting and data analysis, and is a central repository of integrated data from one or more disparate sources.

# Curating **Big Data**

# Curating Big Data

48

Data curation has been defined as the active and on-going management of data through its lifecycle of interest and usefulness.

Data Curation is the **process** of transforming raw data into **Contextualized Data**.

Data curation includes all the **tasks** needed for principled and controlled data **creation, maintenance, and management**, together with the capacity to **add value** to data.

Freitas et al., "Big data curation". In New Horizons for a Data-Driven Economy, 2016.

Arocena et al., "Benchmarking data curation systems". IEEE Data Eng. Bull., 2016.



# Curating Big Data

49

Big Data Curation involves:

- **Identifying** relevant data sources,
- **Ingesting** data and knowledge,
- **Cleaning**,
- **Integration**,
- **Transformation** (Normalization and aggregation),
- **Adding Value** (Preparing Raw Data for Analytics):
  - **Extraction**
  - **Enrichment**,
  - **Linking**,
  - **Summarization.**

# Identifying relevant data sources

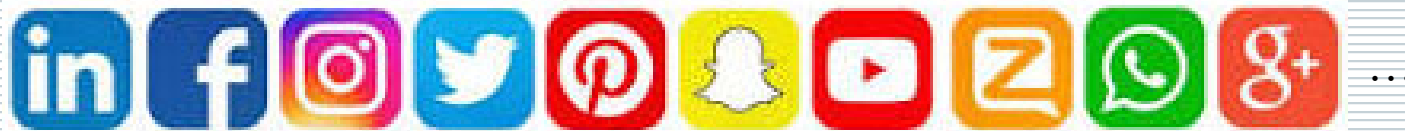
50

With more data repositories constantly being published every day, choosing appropriate data sources for a specific analyst GOAL becomes very important.

## Private Personal/Business Data:



## Social Data:



## Open Data:

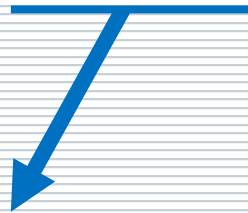


# Ingesting data and knowledge

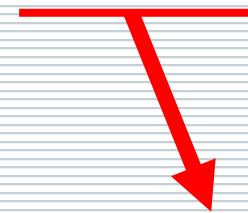
51

**Data ingestion** is the process of obtaining and importing data for immediate use or storage in a database.

Data can be **streamed** in real time or ingested in **batches**.



When data is ingested in real time, each data item is imported **as it is** emitted by the source.



When data is ingested in batches, data items are imported in **discrete chunks** at periodic intervals of time.

**Hortonworks Dataflow (HDF)**

**Makes Big Data Ingest Easy!**



<https://hortonworks.com/solutions/data-ingestion/>



# Big Data **Cleaning**

53

## **Data Cleaning:**

- also known as data cleansing and data scrubbing.
- is the process of amending or removing data in a database that is incorrect, incomplete, improperly formatted, or duplicated.
- is the number one problem in data warehousing

## Other data problems which requires data cleaning

- duplicate records,
- incomplete data,
- inconsistent data

# Big Data **Integration**

54

**Data integration**, combines data from multiple sources.

## **Issues during data integration:**

### ➤ **Schema integration**

- integrate metadata (about the data) from different sources
- Entity identification problem: identify real world entities from multiple data sources. E.g. Change of Name issue.

### ➤ **Detecting and resolving data value conflicts**

- for the same real world entity, attribute values from different sources are different, e.g., different scales (metric vs. British units)

### ➤ **Removing duplicates and redundant data**

- An attribute can be derived from another table (annual revenue)
- Inconsistencies in attribute naming. e.g., A.lastName vs. B.familyName (same attribute?)

[http://www.cs.ccsu.edu/~markov/ccsu\\_courses/datamining-3.html](http://www.cs.ccsu.edu/~markov/ccsu_courses/datamining-3.html)

# Big Data Transformation

55

**Data Transformation** is usually used to smooth the noisy data, summarize, generalize, or normalize the data scale falls within a small, specified range.

- Smoothing: remove noise from data (binning, clustering, regression)
- Normalization: scaled to fall within a small, specified range such as  $-1.0$  to  $1.0$  or  $0.0$  to  $1.0$
- Attribute/feature construction
  - New attributes constructed / added from the given ones
- Aggregation: summarization or aggregation operations apply to data
- Generalization: concept hierarchy climbing
  - Low level/ primitive/raw data are replace by higher level concepts  
(Granularity !)

[http://www.cs.ccsu.edu/~markov/ccsu\\_courses/datamining-3.html](http://www.cs.ccsu.edu/~markov/ccsu_courses/datamining-3.html)

# Curation: Tasks for **Adding Value**

56

- **Extraction,**
- **Enrichment,**
- **Linking,**
- **Summarization.**

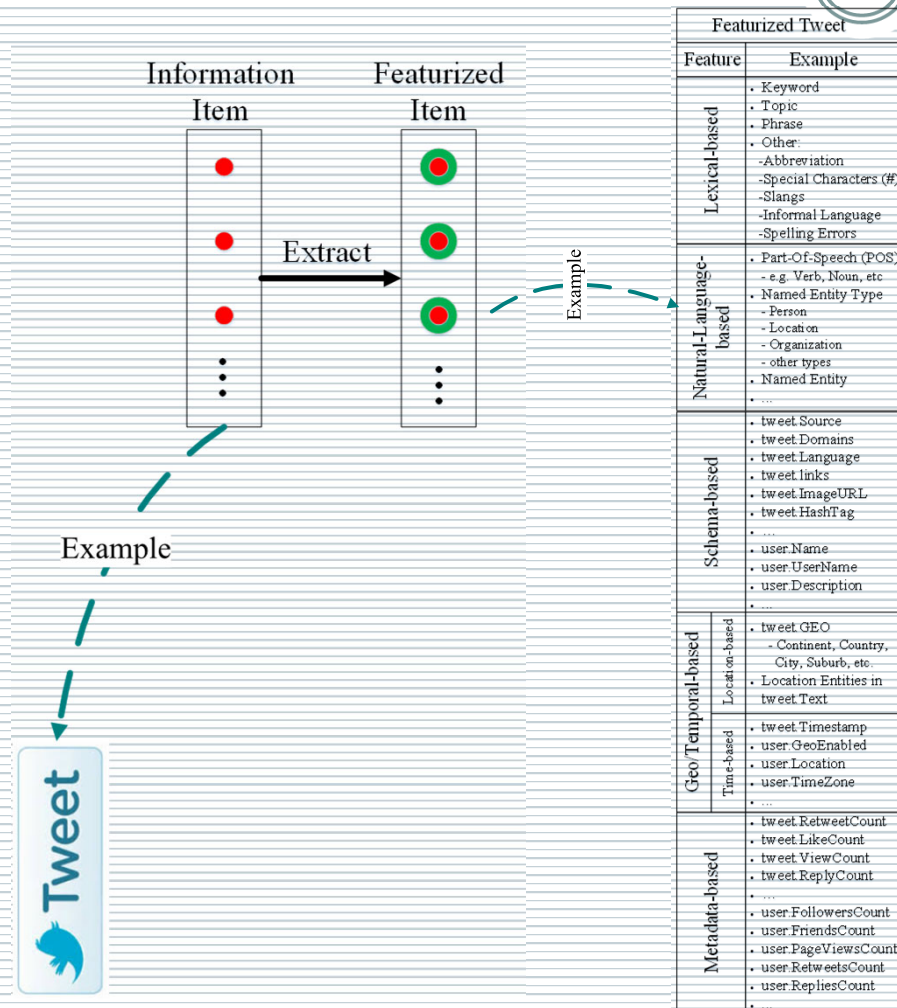
(Preparing Raw Data for Analytics)

Beheshti et al., "**DataSynapse: A Social Data Curation Foundry**". Distributed and Parallel Databases (DAPD) Journal, 2018



# Extraction – Featurized Item

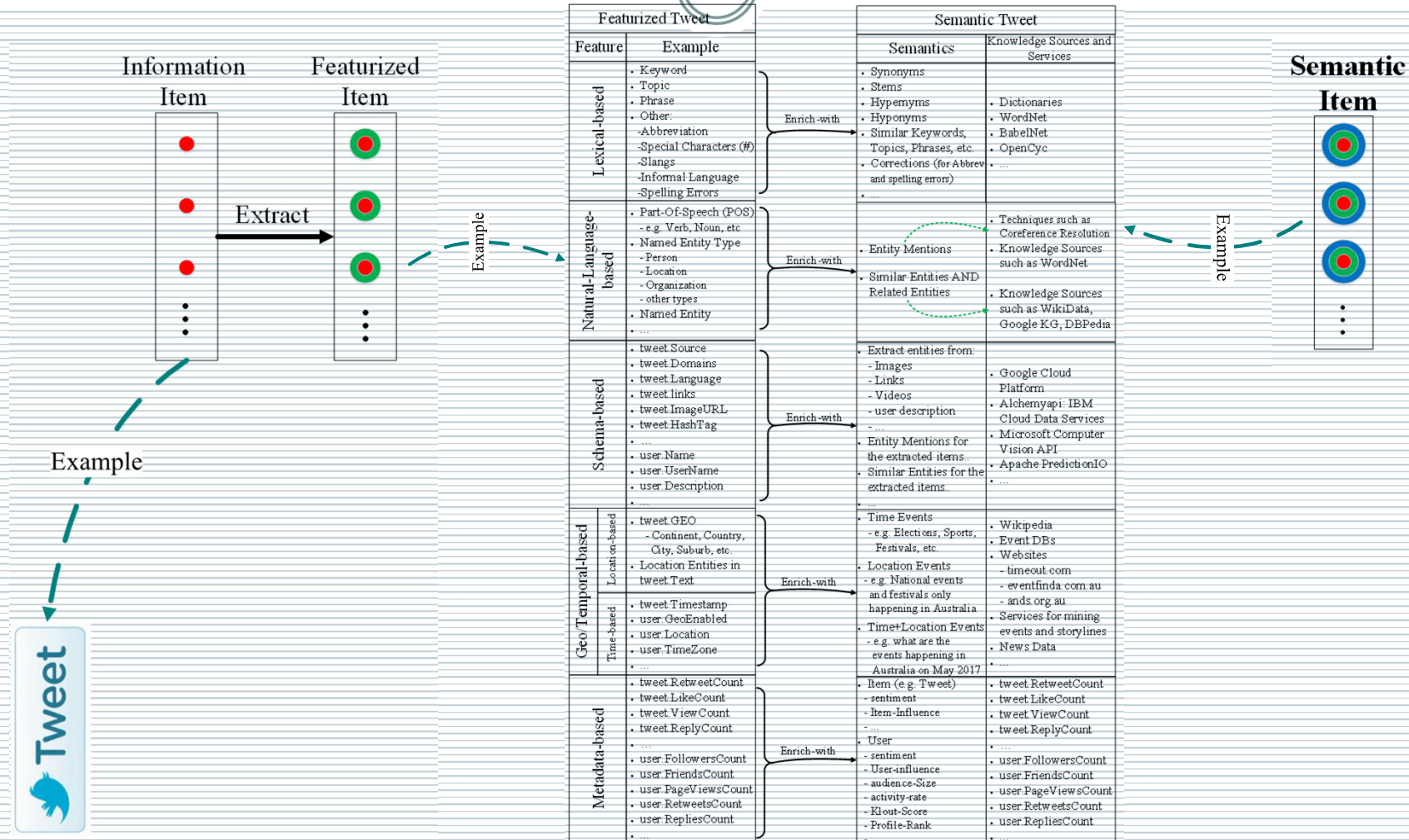
57



Beheshti et al., "DataSynapse: A Social Data Curation Foundry". Distributed and Parallel Databases (DAPD) Journal, 2018

# Enrichment – Semantic Item

58

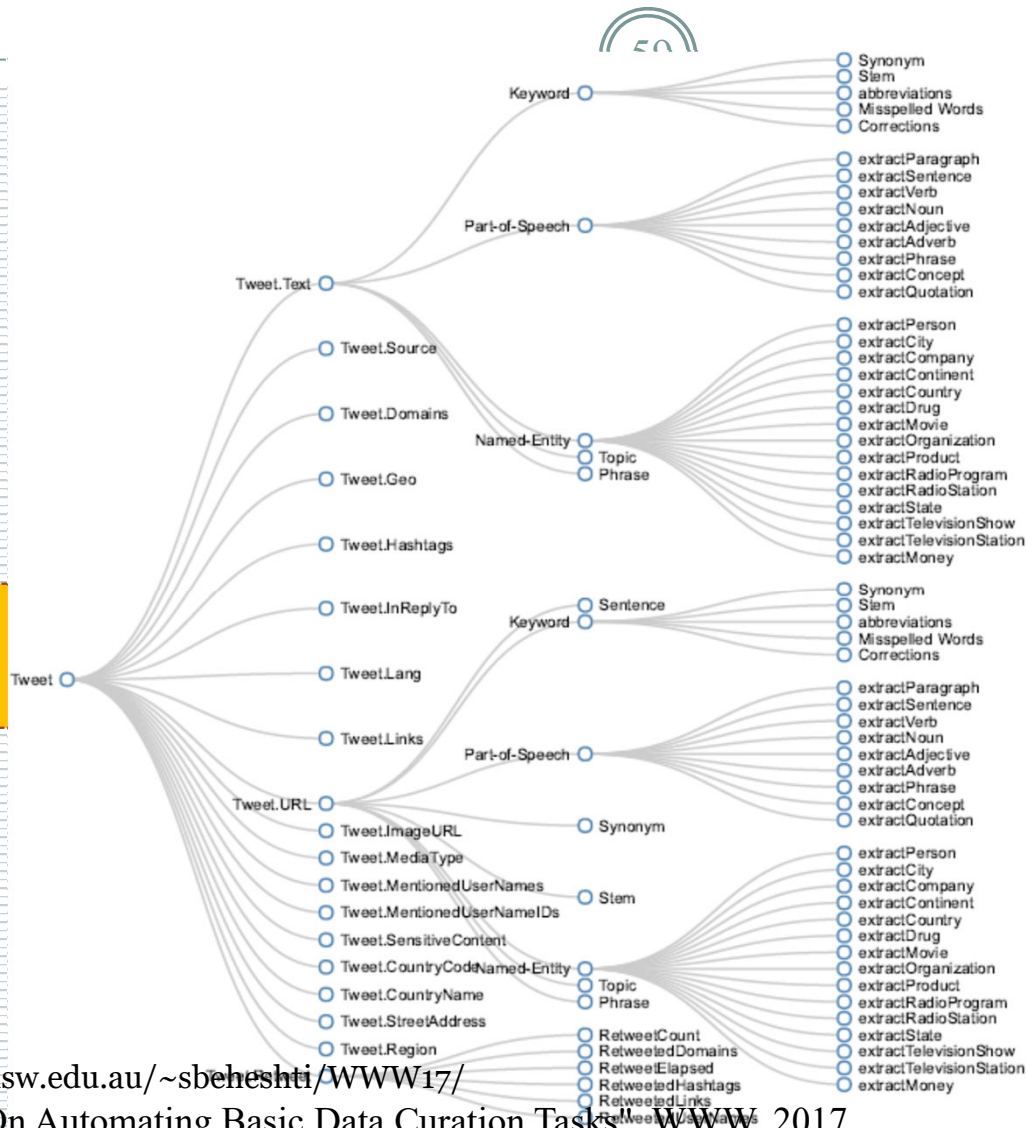


Beheshti et al., "DataSynapse: A Social Data Curation Foundry". Distributed and Parallel Databases (DAPD) Journal, 2018

Example

# Extraction – Featurized Item

Raw Tweet



Contextualized Tweet

<http://www.cse.unsw.edu.au/~sbeheshti/WWW17/>

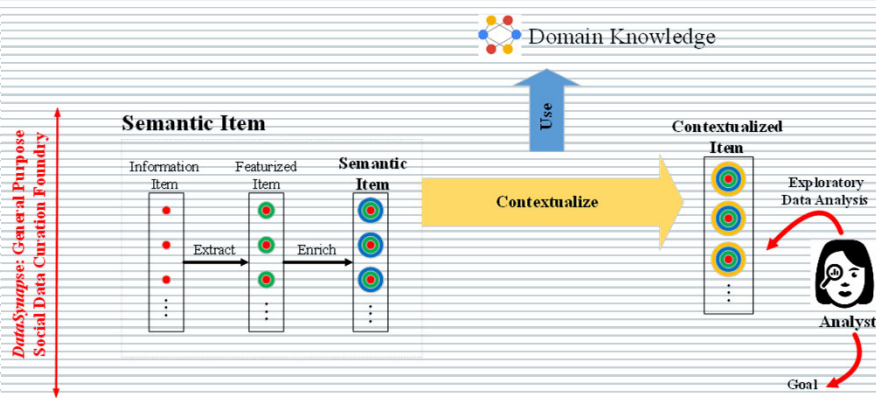
Beheshti et al., "On Automating Basic Data Curation Tasks", WWW, 2017

Copyright © DataAnalyticsResearchGroup @MQ

<https://data-science-group.github.io/>

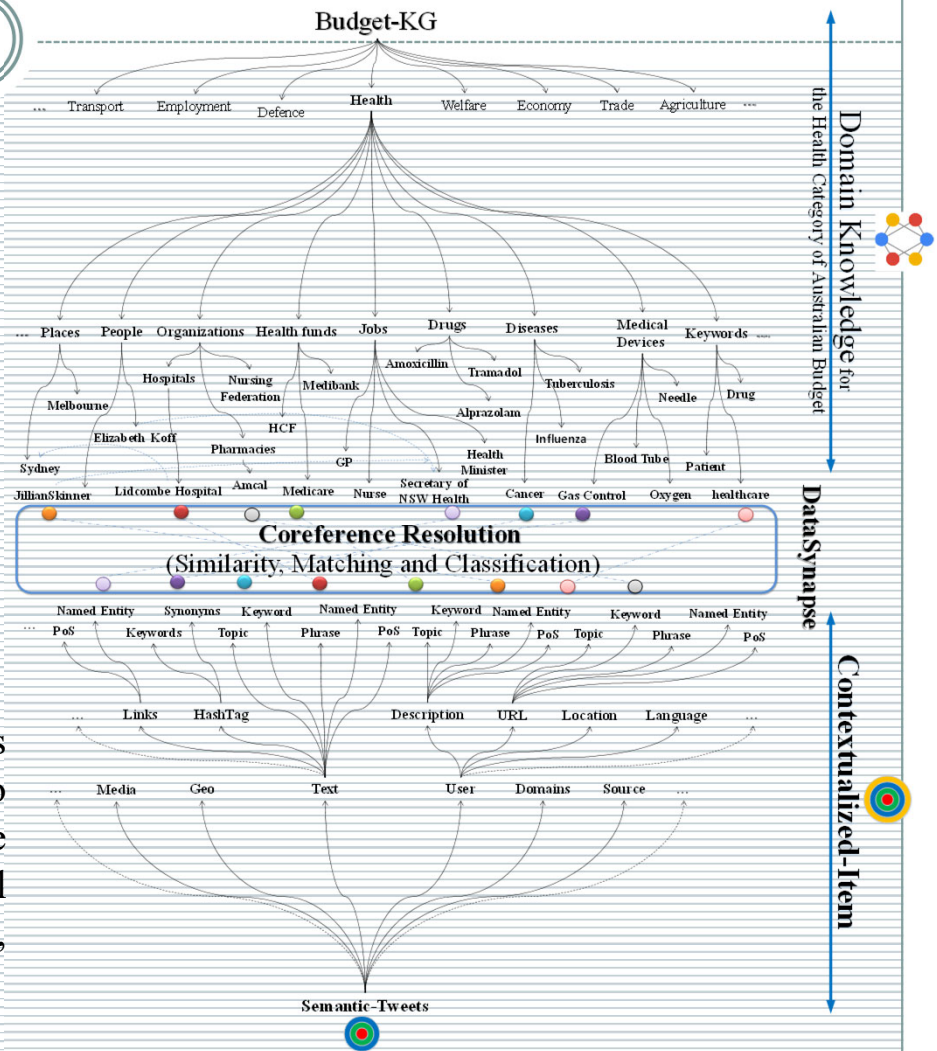
# Linking – Contextualized Item Item

60



## Motivating Example:

A typical scenario for analyzing Urban Social Issues from Twitter as it relates to the Government Budget, to highlight how DataSynapse significantly improves the quality of extracted knowledge compared to the classical curation pipeline (in the absence of feature extraction, enrichment and domain-linking contextualization).



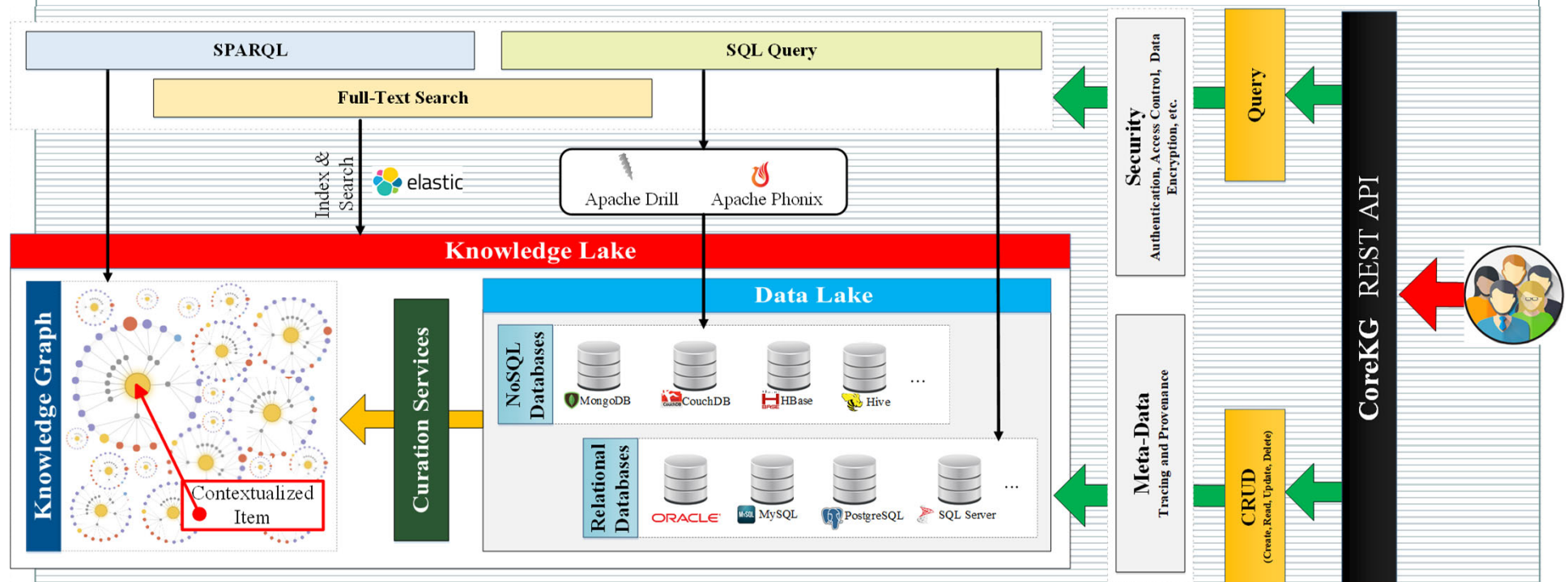
Beheshti et al., "DataSynapse: A Social Data Curation Foundry". Distributed and Parallel Databases (DAPD) Journal, 2018

# Knowledge Lake

# Knowledge Lake

62

A **Knowledge Lake**, i.e. a contextualized Data Lake, is a centralized repository containing virtually inexhaustible amounts of both data and contextualized data that is readily made available to perform analytical activities.



Beheshti et al., **CoreKG: a Knowledge Lake Service (VLDB'18)**, <https://github.com/uns-w-cse-soc/CoreKG>

# Motivating Scenario

# Motivating Scenario

64

## Police Investigation for Missing Person

### **BPM 2018:**

- "iProcess: Enabling IoT Platforms in Data-Driven Knowledge-Intensive Processes"

### **ICSOC 2018:**

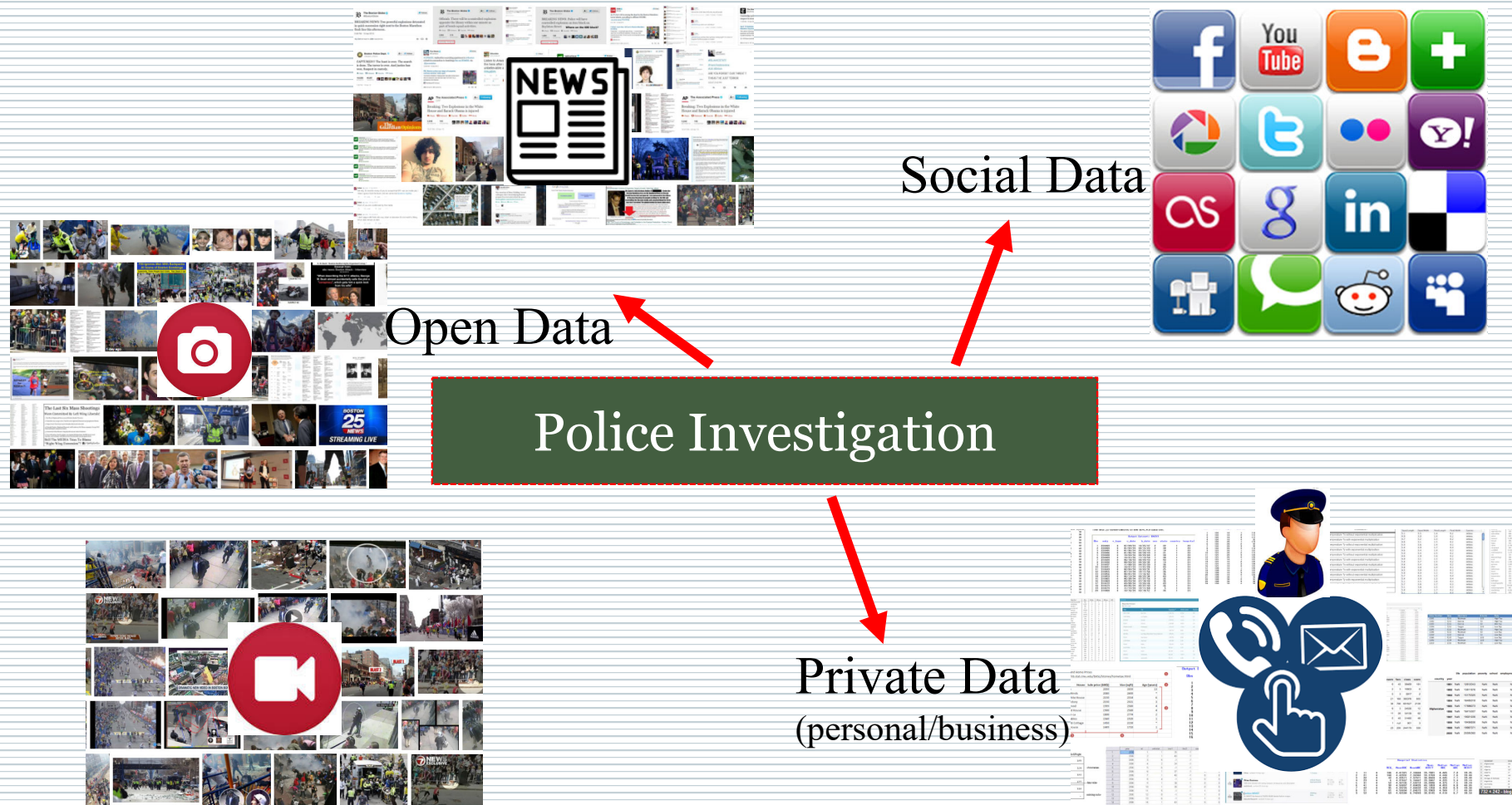
- "iCOP: IoT-enabled Policing Processes"
- "iSheets: A Spreadsheet-based Machine Learning Development Platform for Data-driven Process Analytics".





# Data-Driven Knowledge-Intensive Processes

66



Beheshti et al. "ProcessAtlas: A scalable and extensible platform for business process analytics", Software: Practice and Experience, 2018

Copyright © DataAnalyticsResearchGroup @MQ

<https://data-science-group.github.io/>

# Enabling IoT Platforms in Data-Driven Knowledge-Intensive Processes

67

Motivating Scenario:  
**Missing Person !**

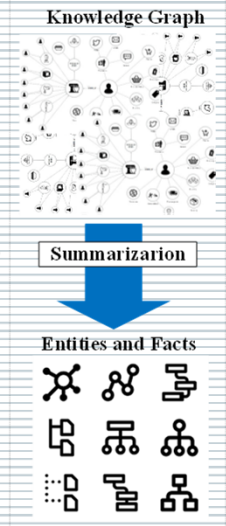
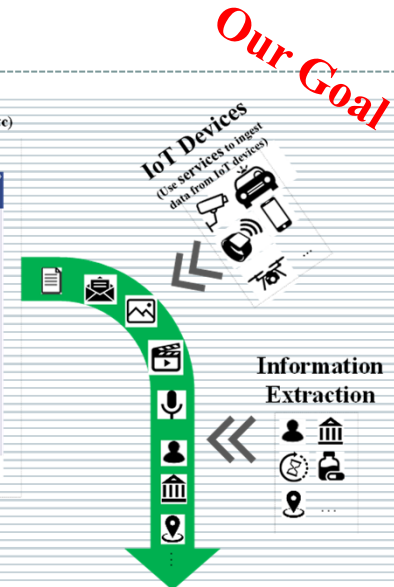
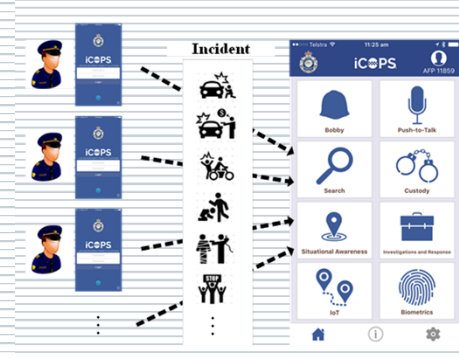
**In Australia, more than 38,000 people are reported missing each year.**

<https://missingpersons.gov.au/view-all-profiles>

**In USA, on any given day, there are as many as 100,000 active missing person's cases.**

<https://nij.gov>

**Information Collection**  
(Use iCoPs to take a photo, record the interview with witnesses, etc)

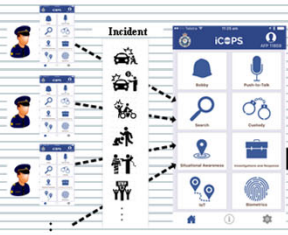


# Solution: iProcess

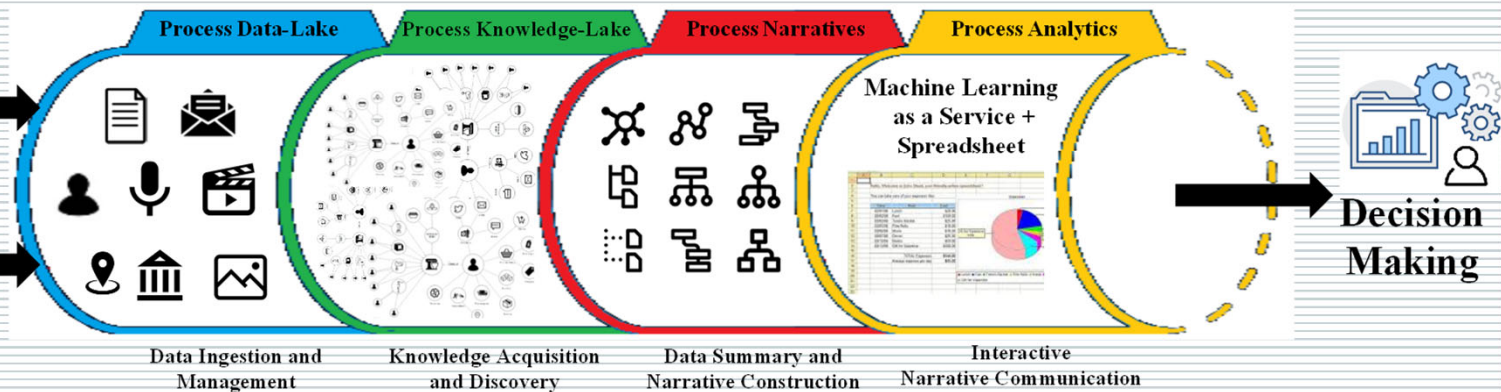
68

**iProcess:** is a scalable and extensible IoT-Enabled Process Data Analytics Pipeline to enable analysts ingest data from IoT devices, extract knowledge from this data and link them to process execution data.

**Information Collection**  
(take a photo, record the interview with witnesses, get the location, etc)



**IoT Devices**

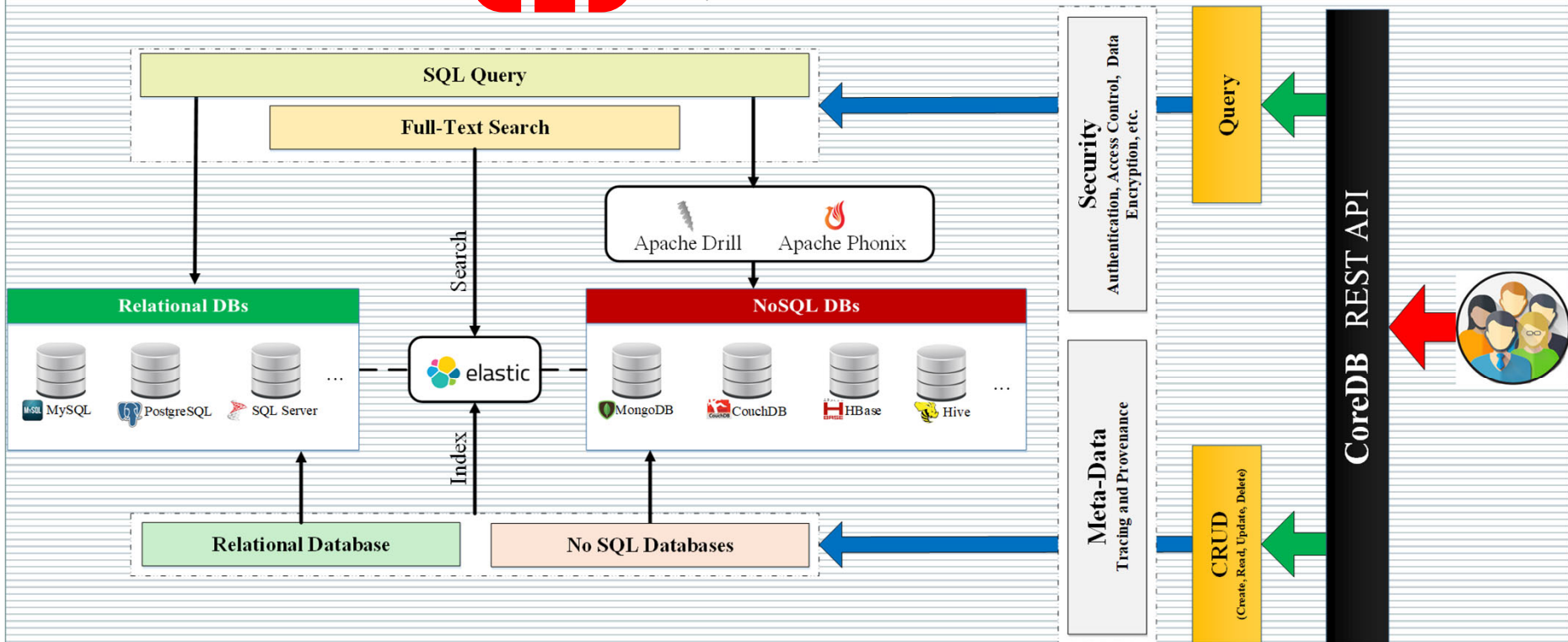
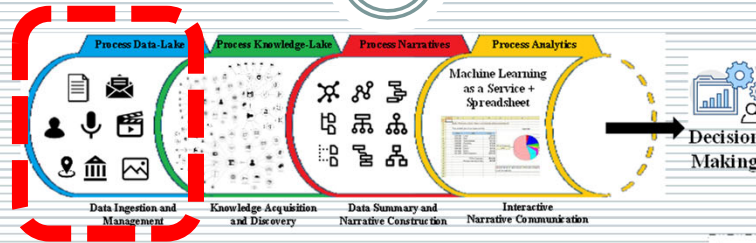


Beheshti et al., "iProcess: Enabling IoT Platforms in Data-Driven Knowledge-Intensive Processes", 16th conference on Business Process Management (BPM), Sydney, Australia, 2018

# Solution: iProcess

69

*Data Lake as a Service*



Beheshti, Benatallah, et al., **CoreDB: a Data Lake Service (CIKM'17)** <https://github.com/uns-w-cse-soc/CoreDB>

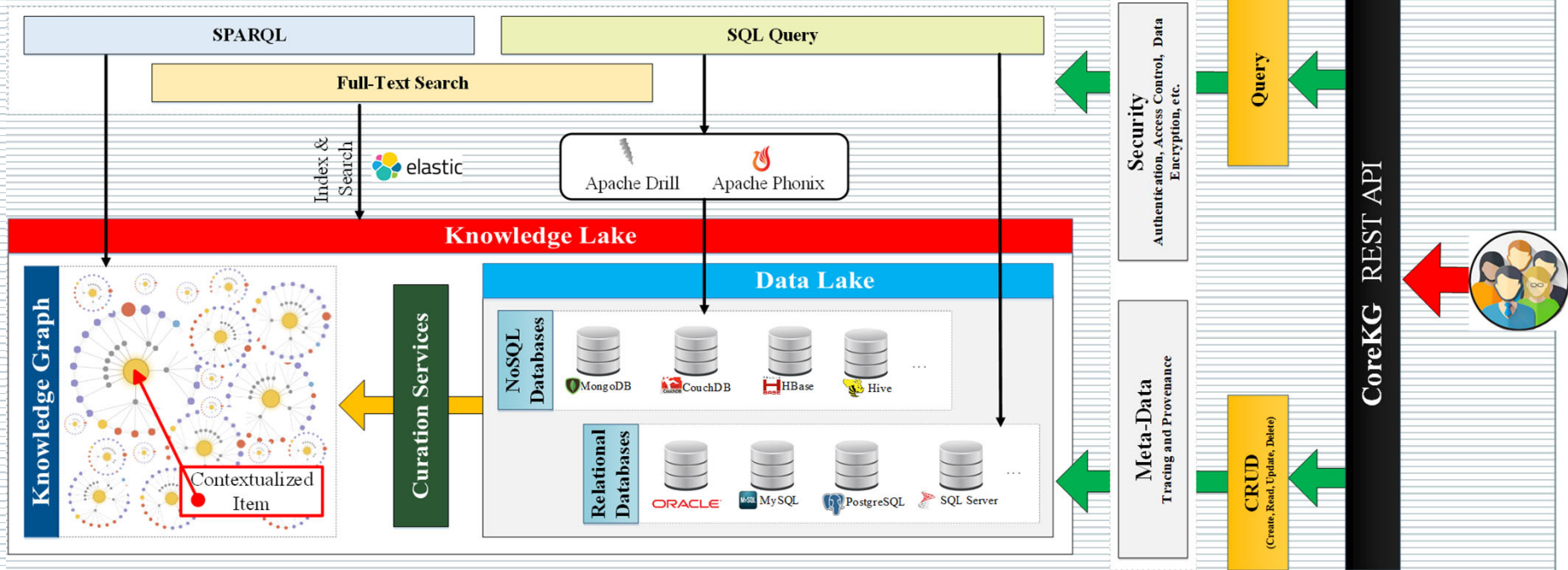
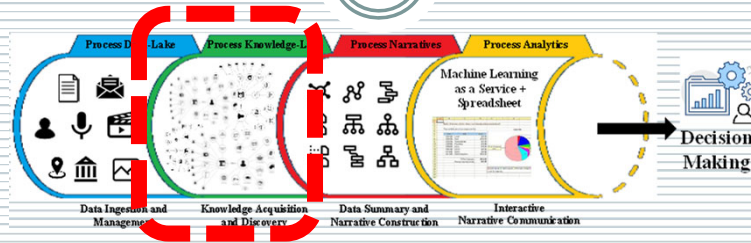
Copyright © DataAnalyticsResearchGroup @MQ

<https://data-science-group.github.io/>

# Solution: iProcess

70

**Knowledge Lake as a Service**



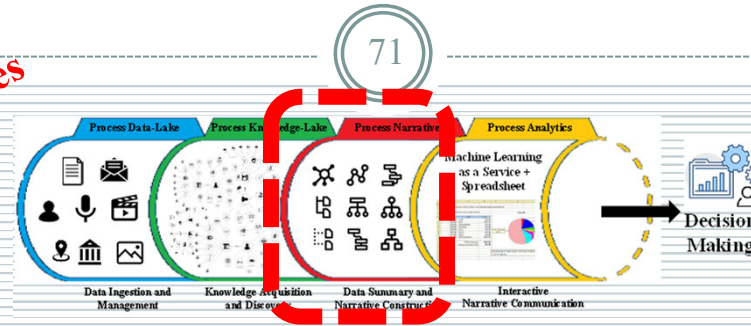
Beheshti, Benatallah, et al., **CoreKG: a Knowledge Lake Service (VLDB'18)**

<https://github.com/uns-w-cse-soc/CoreKG>

Beheshti, Benatallah, et al., **DataSynapse: A Social Data Curation Foundry (DAPD Journal, 2018)**

# Solution: iProcess

Summarizing the process data and constructing process narratives



*Process OLAP*  
*Process Cubes*  
*Dimensions*  
*Cells*  
*Measures*  
*Operations*

**OLAP**, is an approach to answering multi-dimensional analytical queries swiftly.

**OLAP**  
Online Analytical Processing



## **Problem:**

- extension of existing OLAP techniques to analysis of graphs is not straightforward.
- key business insights remain hidden in the interactions among objects.

## **Solution:**

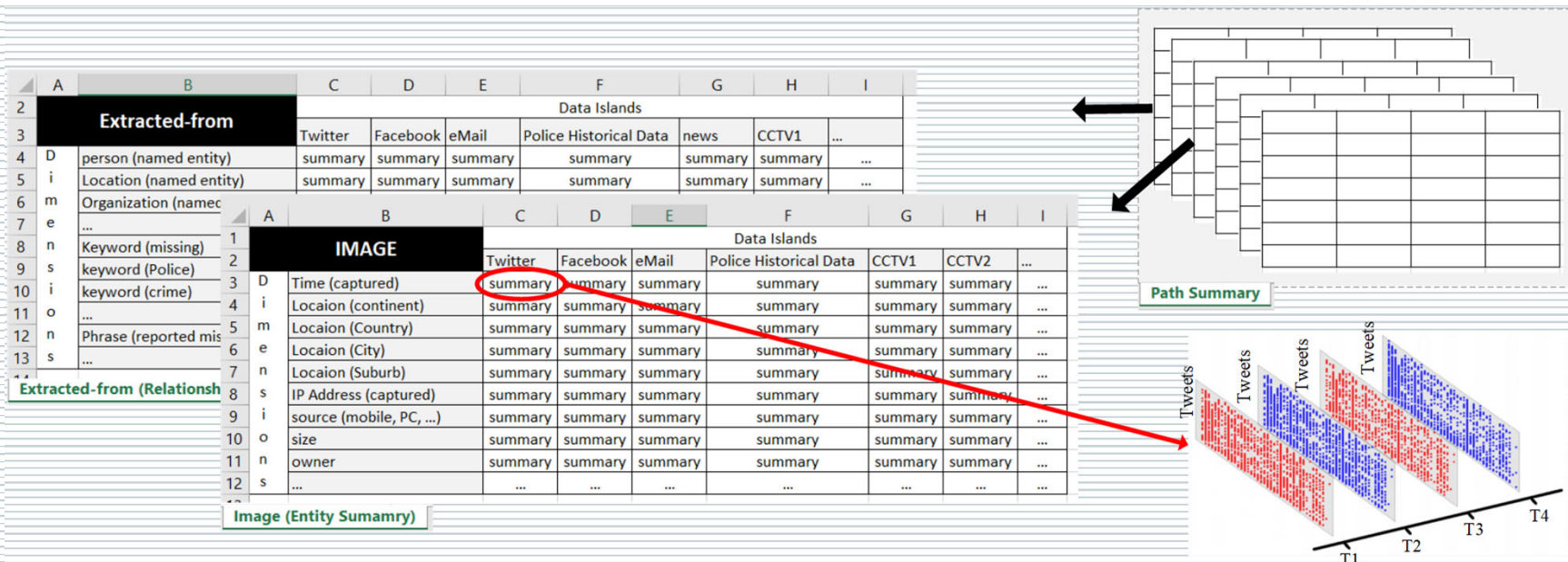
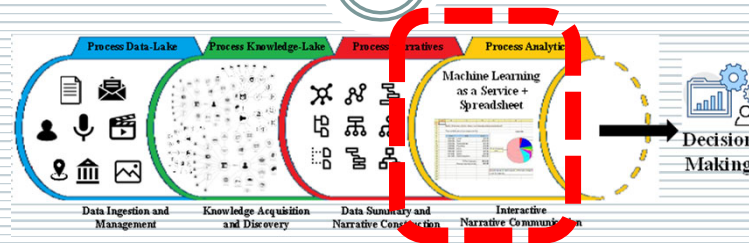
- On-Line Analytical Processing on Graphs

Beheshti et al., "Scalable Graph-based OLAP Analytics over Process Execution Data", Distributed and Parallel Databases (DAPD) Journal, 34(3), 379-423, 2016

# Solution: iProcess

72

**ML as a Service**



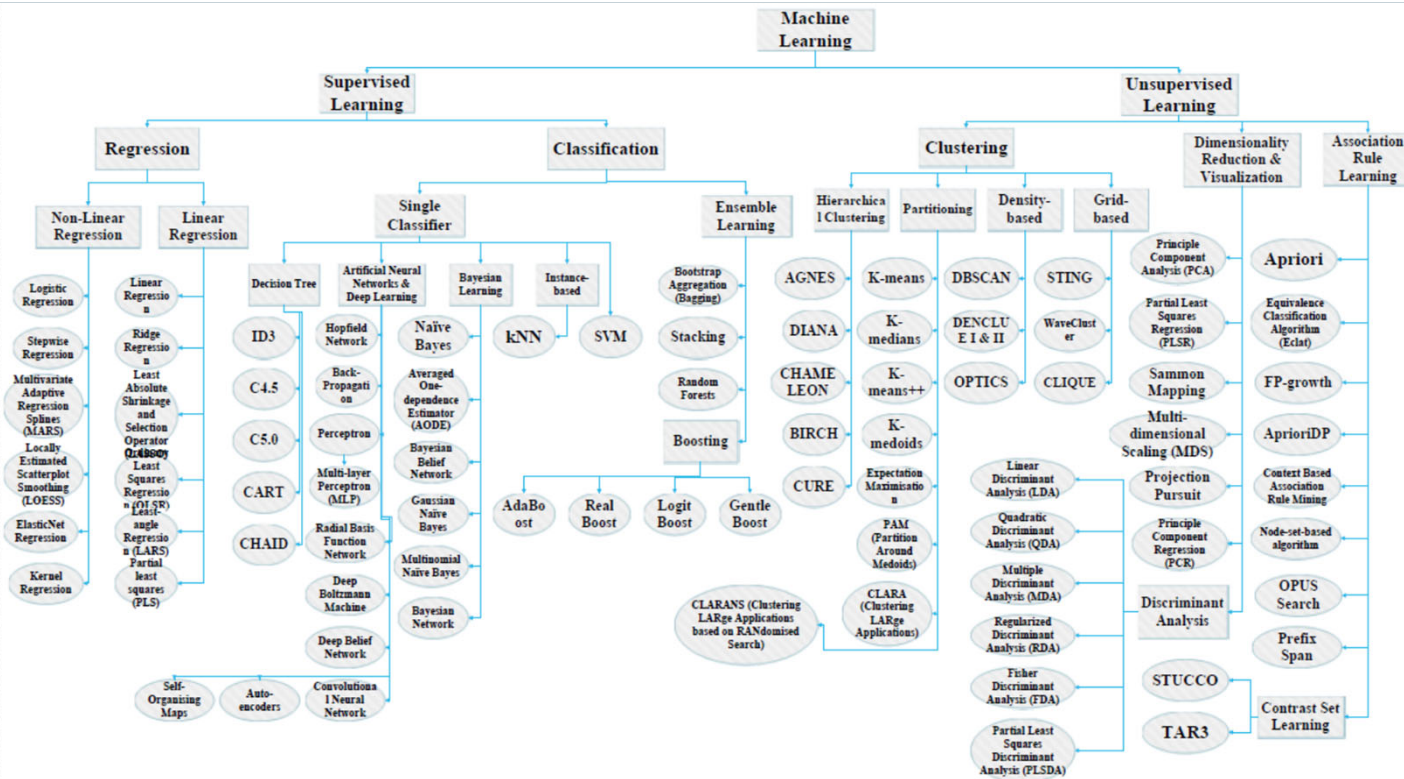
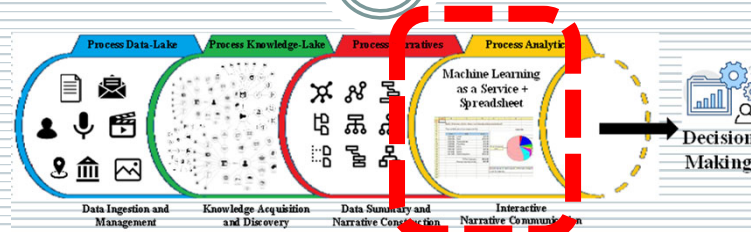
Amouzegar, Beheshti et al., " iSheets: A Spreadsheet-based Machine Learning Development Platform for Data-driven Process Analytics", ICSOC, 2018.



# Solution: iProcess

73

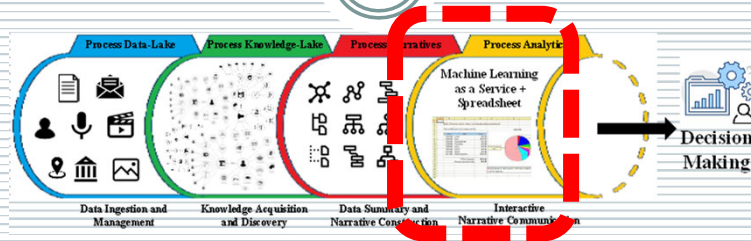
**ML as a Service**



# Solution: iProcess

74

iSheets



localhost MLAAS localhost:8080/

**Formula Bar:** group("location","time") Apply

**Dataset**

	2016-04-01	2016-04-02	2016-04-03	2016-04-04	2016-04-05	2016-04-06	2016-04-07	2016-04-08	2016-04-09	2016-04-10
Brisbane	N/A	N/A	N/A	<a href="#">Tweets</a>	<a href="#">Tweets</a>	<a href="#">Tweets</a>	N/A	N/A	N/A	N/A
Canberra	N/A	<a href="#">Tweets</a>	N/A	N/A	N/A	N/A	N/A	<a href="#">Tweets</a>	N/A	N/A
Melbourne	<a href="#">Tweets</a>	N/A	N/A	<a href="#">Tweets</a>	N/A	<a href="#">Tweets</a>	N/A	N/A	N/A	<a href="#">Tweets</a>
Perth	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	<a href="#">Tweets</a>	N/A
Sydney	<a href="#">Tweets</a>	<a href="#">Tweets</a>	<a href="#">Tweets</a>	N/A	<a href="#">Tweets</a>	<a href="#">Tweets</a>	<a href="#">Tweets</a>	<a href="#">Tweets</a>	<a href="#">Tweets</a>	N/A

classify.bytopic() Apply

	Sport	Politics	Business	Education	Entertainment
Sydney-2016-04-05	N/A	N/A	<a href="#">Tweets</a>	N/A	N/A
Sydney-2016-04-06	N/A	N/A	<a href="#">Tweets</a>	N/A	N/A
Sydney-2016-04-07	N/A	N/A	<a href="#">Tweets</a>	N/A	<a href="#">Tweets</a>

classify.bySentiment() Apply

	Positive	Negative
Sydney-2016-04-07-Sc	<a href="#">Tweets</a>	<a href="#">Tweets</a>

Twitter Group

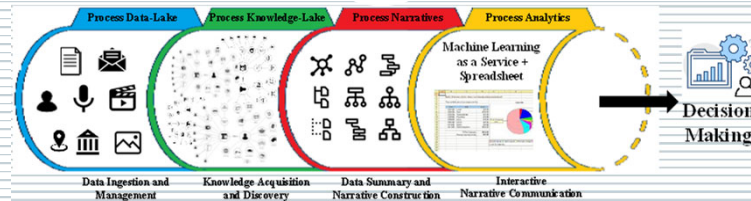
**ML as a Service**

- Classification
  - SVM
  - kNN
  - Logistic Regression
  - C5.0
  - MLP
  - Trained:Topic Classifier
  - Trained:Sentiment Classifier
- + Clustering
- + Association Learning
- Operations
  - Select
  - Group
  - Sort
  - Partition
  - Addition

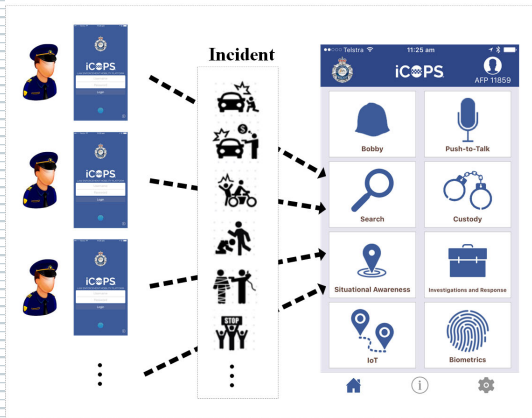
# Solution: iProcess

75

**Motivating Scenario  
Missing People!**



## Information Collection (Use iCoPs to take a photo, record the interview with witnesses, etc)



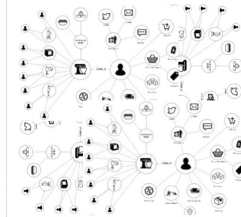
## IoT Devices (Use services to ingest data from IoT devices)



## Information Extraction

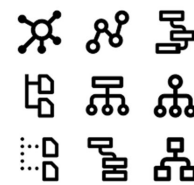


## Knowledge Graph



## Summarization

## Entities and Facts



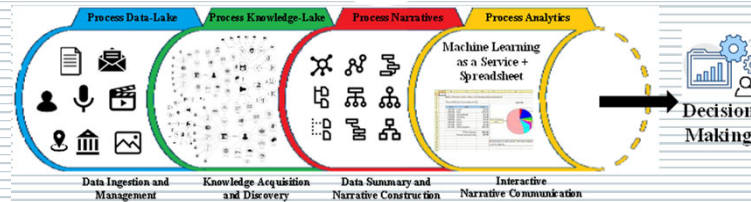
## ICSOC 2018:

- "iCOP: IoT-enabled Policing Processes"
- "iSheets: A Spreadsheet-based Machine Learning Development Platform for Data-driven Process Analytics"

# Solution: iProcess

76

**Motivating Scenario  
Missing People!**

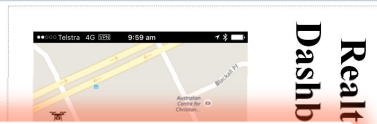


**Information Collection**  
(Use iCoPs to take a photo, record the interview with witnesses, etc)

**IoT Devices**  
(Use services to ingest)

**Information Extraction**

**Knowledge Graph**

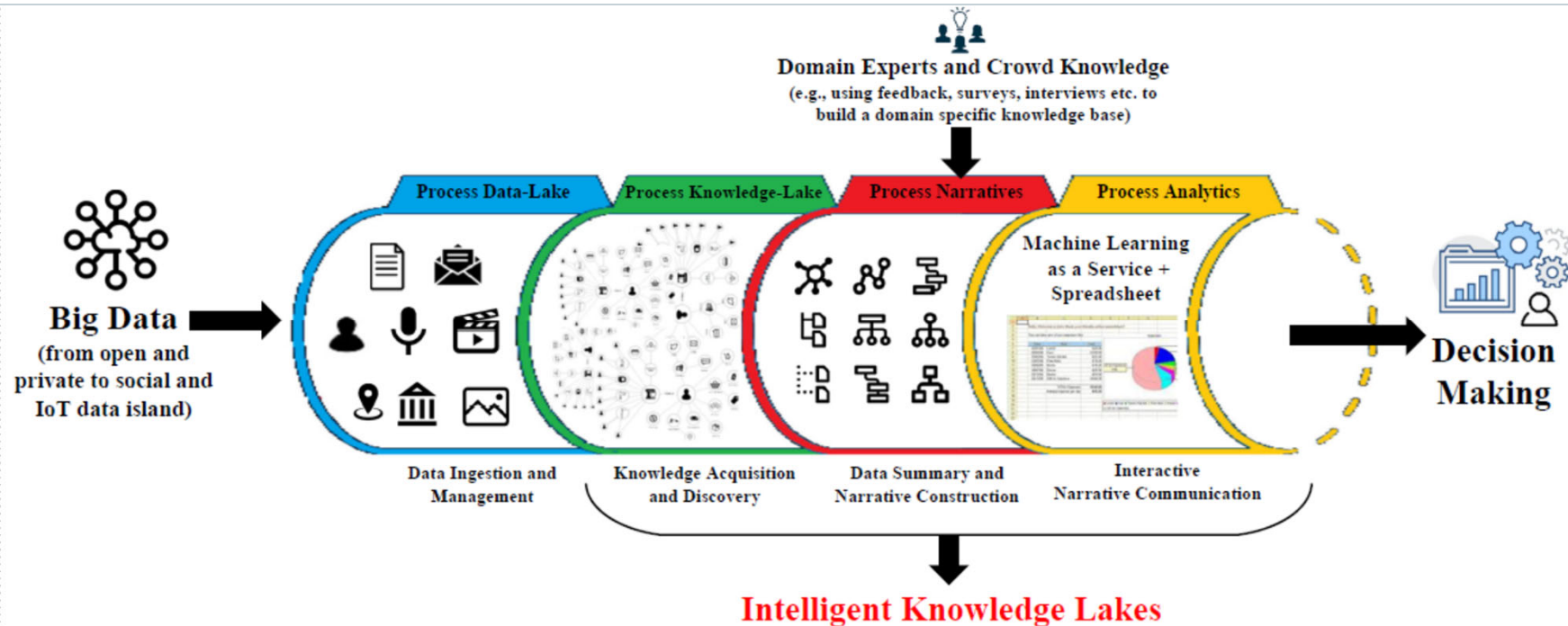


**ICSOC 2018:**

- "iCOP: IoT-enabled Policing Processes"
- "iSheets: A Spreadsheet-based Machine Learning Development Platform for Data-driven Process Analytics"

# Intelligent Knowledge Lakes !

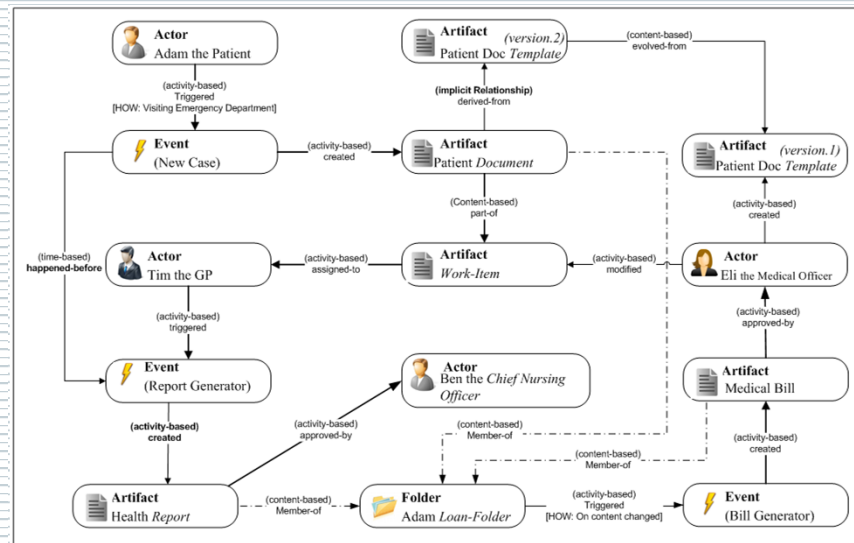
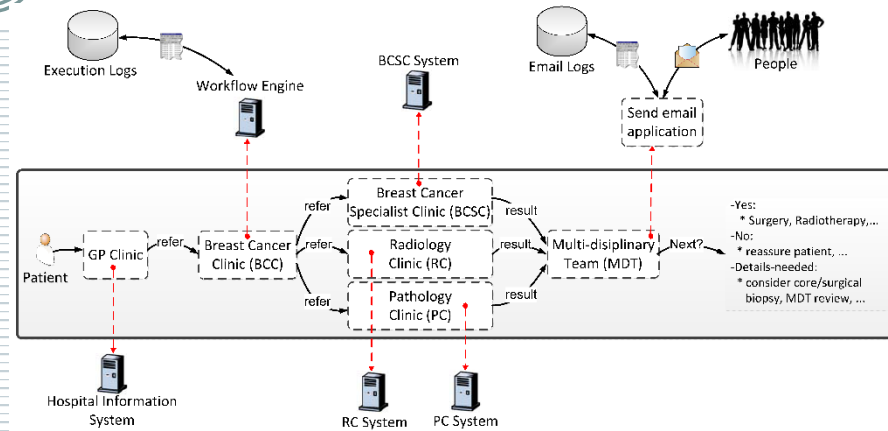
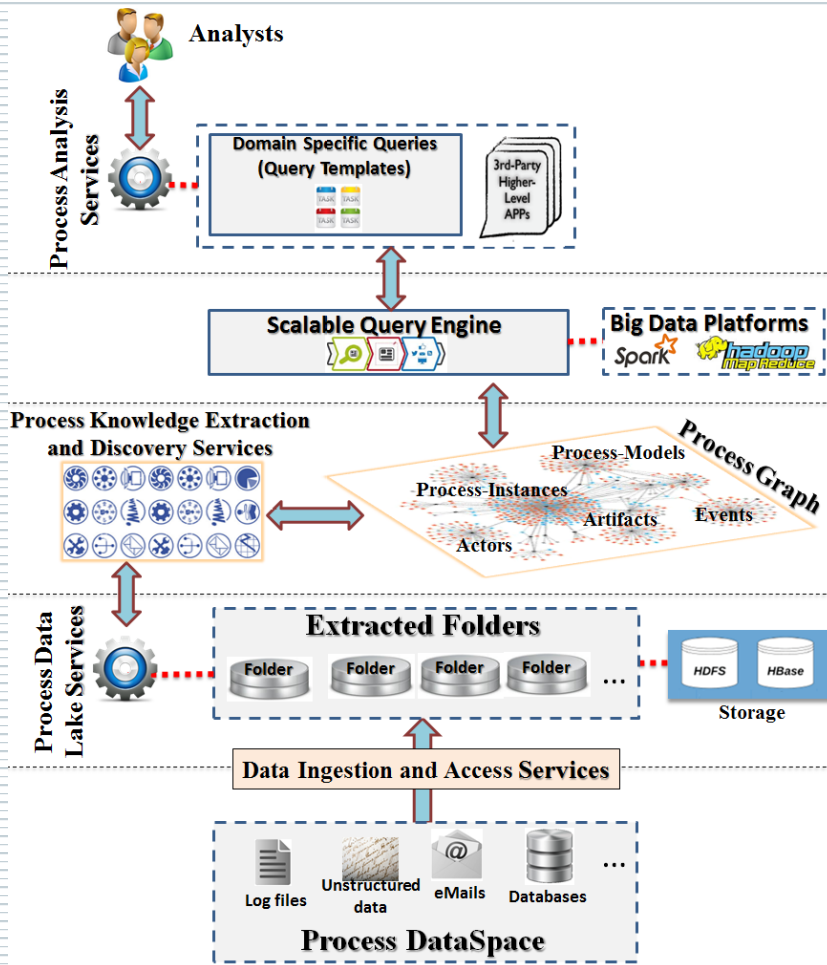
77



## Other Scenarios

# Health

79



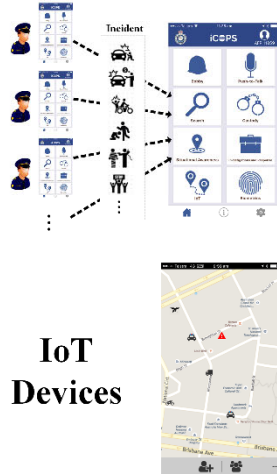
Beheshti et al., "Galaxy: A Platform for Explorative Analysis of Open Data Sources", EDBT,

# IoT

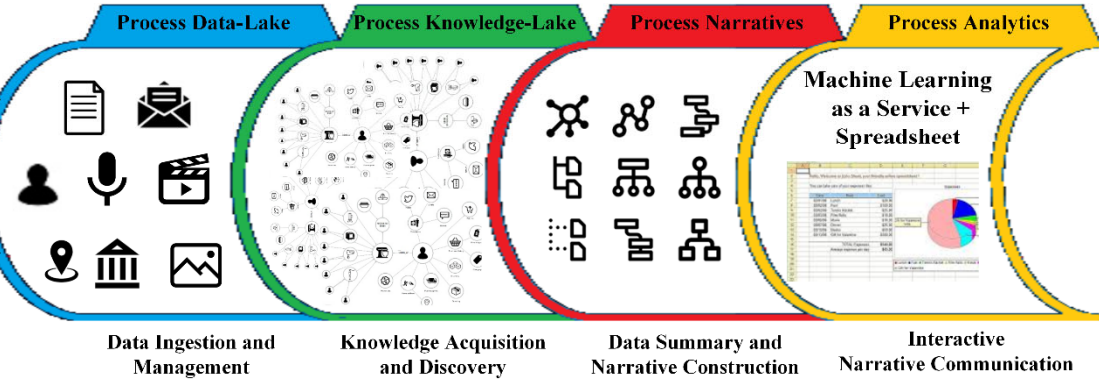
80

## Information Collection

(take a photo, record the interview with witnesses, get the location, etc)



IoT Devices



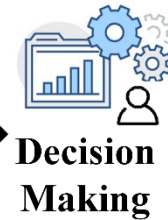
Data Ingestion and Management

Knowledge Acquisition and Discovery

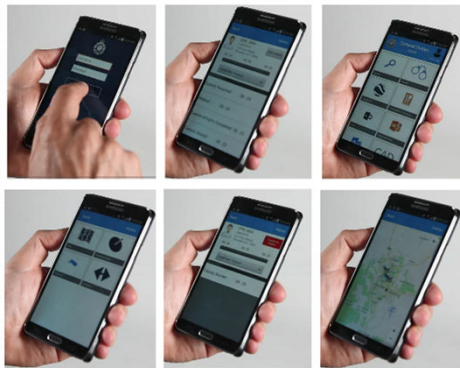
Data Summary and Narrative Construction

Interactive Narrative Communication

Machine Learning as a Service + Spreadsheet



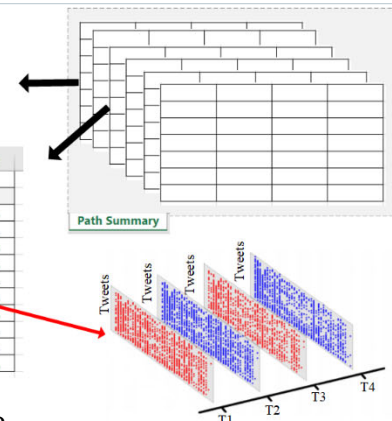
Decision Making



Extracted-from		Data Islands						
		Twitter	Facebook	eMail	Police Historical Data	news	CCTV1	...
D	person (named entity)	summary	summary	summary	summary	summary	summary	...
i	Location (named entity)	summary	summary	summary	summary	summary	summary	...
m	Organization (name)	summary	summary	summary	summary	summary	summary	...
e	...	summary	summary	summary	summary	summary	summary	...
n	Keyword (missing)	summary	summary	summary	summary	summary	summary	...
s	keyword (Police)	summary	summary	summary	summary	summary	summary	...
i	keyword (crime)	summary	summary	summary	summary	summary	summary	...
n	Phrase (reported mis)	summary	summary	summary	summary	summary	summary	...
s	...	summary	summary	summary	summary	summary	summary	...
n	Location (Suburb)	summary	summary	summary	summary	summary	summary	...
s	IP Address (captured)	summary	summary	summary	summary	summary	summary	...
i	source (mobile, PC, ...)	summary	summary	summary	summary	summary	summary	...
o	size	summary	summary	summary	summary	summary	summary	...
n	owner	summary	summary	summary	summary	summary	summary	...
s	...	summary	summary	summary	summary	summary	summary	...

Extracted-from (Relations)		Data Islands						
		Twitter	Facebook	eMail	Police Historical Data	CCTV1	CCTV2	...
D	Time (captured)	summary	summary	summary	summary	summary	summary	...
i	Location (continent)	summary	summary	summary	summary	summary	summary	...
m	Location (Country)	summary	summary	summary	summary	summary	summary	...
e	Location (City)	summary	summary	summary	summary	summary	summary	...
n	Location (Suburb)	summary	summary	summary	summary	summary	summary	...
s	IP Address (captured)	summary	summary	summary	summary	summary	summary	...
i	source (mobile, PC, ...)	summary	summary	summary	summary	summary	summary	...
o	size	summary	summary	summary	summary	summary	summary	...
n	owner	summary	summary	summary	summary	summary	summary	...
s	...	summary	summary	summary	summary	summary	summary	...



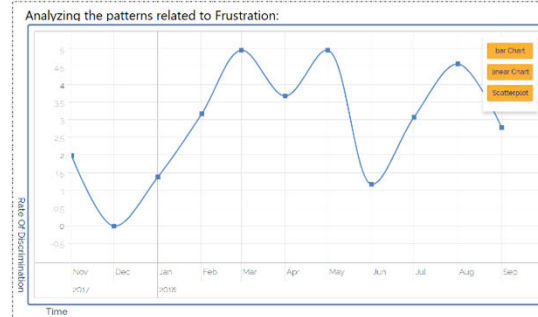
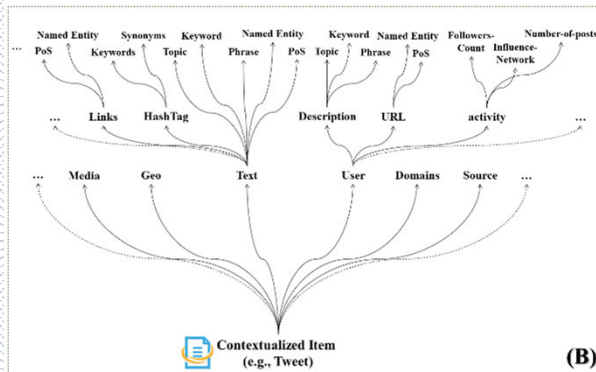
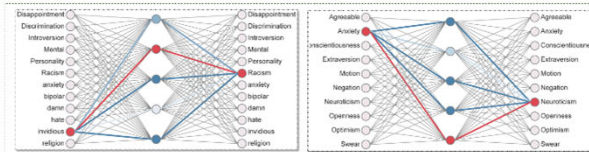
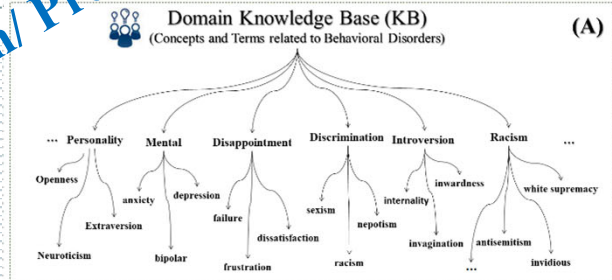
Beheshti et al., "iProcess: Enabling IoT Platforms in Data-Driven Knowledge-Intensive Processes", BPM, 2018



# AI-Enabled Policing (Applications)

Crime Detection/Prevention

81



**(D) Behavioral Analytics**

CONTENT ANALYSIS	CONTEXT ANALYSIS	ACTIVITY ANALYSIS
Content	Context	Activity
Sentiment	Mental Health	Number of followers
Extract Keywords	Frustration	Number of friends
Extract Named Entities	Discrimination	Pageview count
Extract Part of Speech	Educational level	Audience size
Extract Synonym	Introversion	Replies count
Extract Stem	Political ideology	Likes count
	Origin	Mention count
	Terrorism	Activity rate
	Racism	Klout score
	Negative ideas about Western Society	Profile rank
	Positive ideas about religion	Outreach score

Beheshti et al., "personality2vec: Enabling the Analysis of Behavioral Disorders in Social Networks", 13th ACM International WSDM Conference (WSDM), Houston, Texas, USA, 2020. (ERA Rank: A\*)

# AI-enabled Processes (AIP) Research Centre

82

MACQUARIE University

## AI-enabled Processes Research Centre

Home Industry Streams Research Programs Projects People Contact

### AI-enabled Processes (AIP) Research Centre

Business processes, i.e., set of coordinated tasks and activities carried out manually/automatically to achieve a business objective or goal, are central to the operation of public and private enterprises. Modern processes are often extremely complex, data-driven and knowledge-intensive. In such processes, it is not sufficient to focus on data storage/analysis; and the knowledge workers will need to collect, understand and relate the big data (from open, private, social and IoT data islands) to process analysis.

Today, the advancement in Artificial Intelligence (AI) and Data Science has the potential to transform business processes in fundamental ways; by assisting knowledge workers in communicating analysis findings, supporting evidences and to make decisions. The core of the idea for AI-enabled Processes (AIP) Research Centre is to advance the scientific understanding of AI-enabled processes and to assist organizations identifying novel applications of AI and Data Science: from process automation, to cognitive assistants and smart entities. Our research covers the full spectrum of topics related to AI and Processes, and to deriving knowledge from process related data: theory, algorithms, applications and software infrastructure.

Long-term objectives of understanding AI-enabled Processes are bold and ambitious, and we know that making significant progress in this field can't be done in isolation. That's why we have two complementary components in the AIP Research Centre:

- Industry Streams:** Our industry streams (e.g., AI-enabled Policing, AI-enabled Banking and AI-enabled teaching and learning) are led by our industry partners such as Australia Federal Police, TATA Consultancy Services, ITIC and more.
- Research Programs:** Our Research Programs (e.g., Process Automation, Data Curation, Cognitive Technology, Smart Entities, IoT-enabled Business Processes and Storytelling with Business Data) are led by research labs in top universities such as Macquarie University (Sydney, Australia), UNSW Sydney (Sydney, Australia), and University of Wollongong (Wollongong, Australia).

#### Industry Streams

- AI-enabled Policing
- AI-enabled Banking
- AI-enabled Education
- AI-enabled Industry
- AI-enabled Health
- AI-enabled Agriculture
- AI-enabled Transport
- AI-enabled Marketing
- [See All Streams](#)

#### Research Programs

- AI-enabled Process Automation
- Cognitive Assistants for Knowledge Intensive Processes
- Data Curation for Data-Driven Processes
- Smart Entities
- IoT-enabled Business Processes
- Storytelling with Business Data
- Cognitive Analytics
- [See All Programs](#)

#### Recent Projects

**Total Funding Since 2019: \$2,072,834.00**

- "AI-enabled Banking", Linkage: Tata Consultancy Services (TCS) and Macquarie University, 2019-2023
- "Intelligence-led Teaching and Learning", Linkage: ITIC Training and Resourcing and Macquarie University, 2019-2023
- [See All Grants](#)

**Web:** <https://aip-research-center.github.io/>

**email:** [aip@mq.edu.au](mailto:aip@mq.edu.au)

**Address:** Level 3, Becton-Dickinson (BD) Building, 4 Research Park Drive, Macquarie University, Sydney, Australia.

# Thank You !

83

