



Big Data Architectures

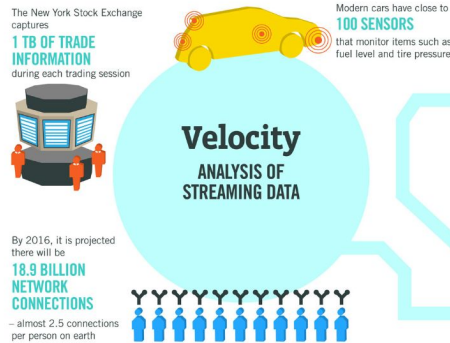
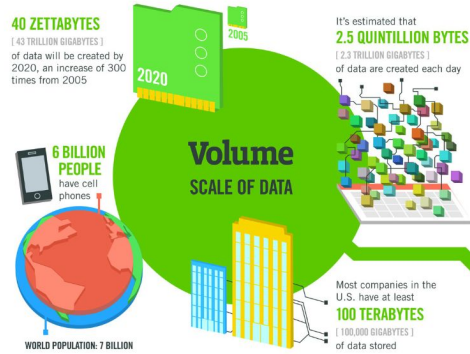
Implementing Data Value Chains at large scale

Summary

We are going to

- recap. Big Data concepts
- introduce a Use Case for BD architectures
- analyze the logic of BD architectures
- introduce the main ETLing concepts
- practice with a Use Case example

The 4 V's of Big Data



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015
4.4 MILLION IT JOBS
will be created globally to support big data, with 1.9 million in the United States

As of 2011, the global size of data in healthcare was estimated to be **150 EXABYTES**
[161 BILLION GIGABYTES]

30 BILLION PIECES OF CONTENT
are shared on Facebook every month

Variety
DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be **420 MILLION WEARABLE, WIRELESS HEALTH MONITORS**

4 BILLION+ HOURS OF VIDEO
are watched on YouTube each month

400 MILLION TWEETS
are sent per day by about 200 million monthly active users

1 IN 3 BUSINESS LEADERS
don't trust the information they use to make decisions

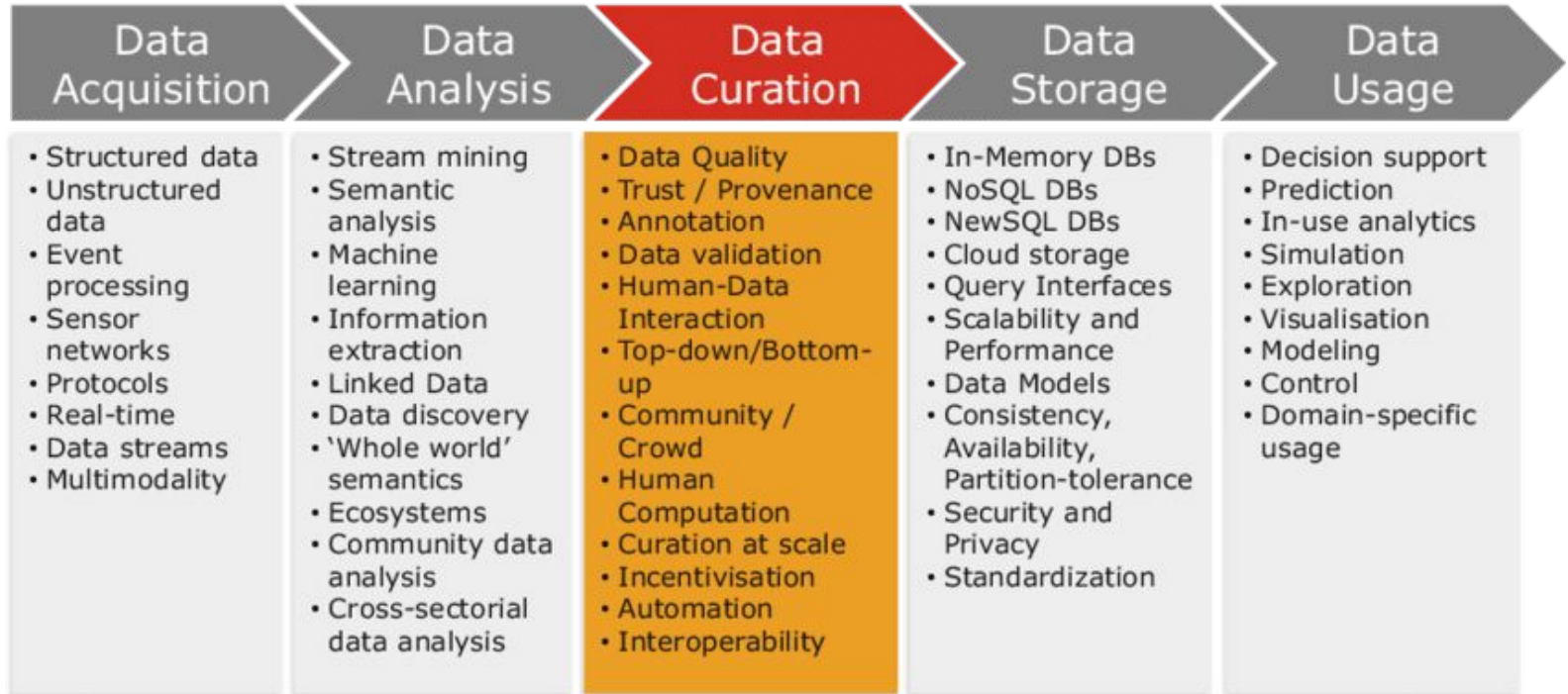
27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate

Veracity
UNCERTAINTY OF DATA

Poor data quality costs the US economy around **\$3.1 TRILLION A YEAR**

Big Data Value Chain



Use Case

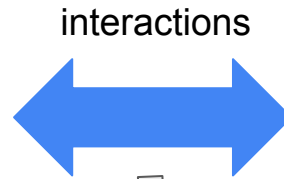
Intelligent Transport

- Global view of the transport scenario in a city
- Gather data from all the perspectives
 - Sensor data
 - Citizen interactions
 - Transport Modalities
 - City's Infrastructures
 - etc.

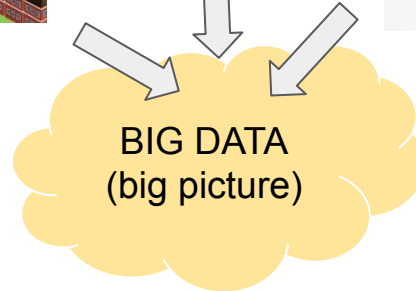
Smart Cities for Public Transportation



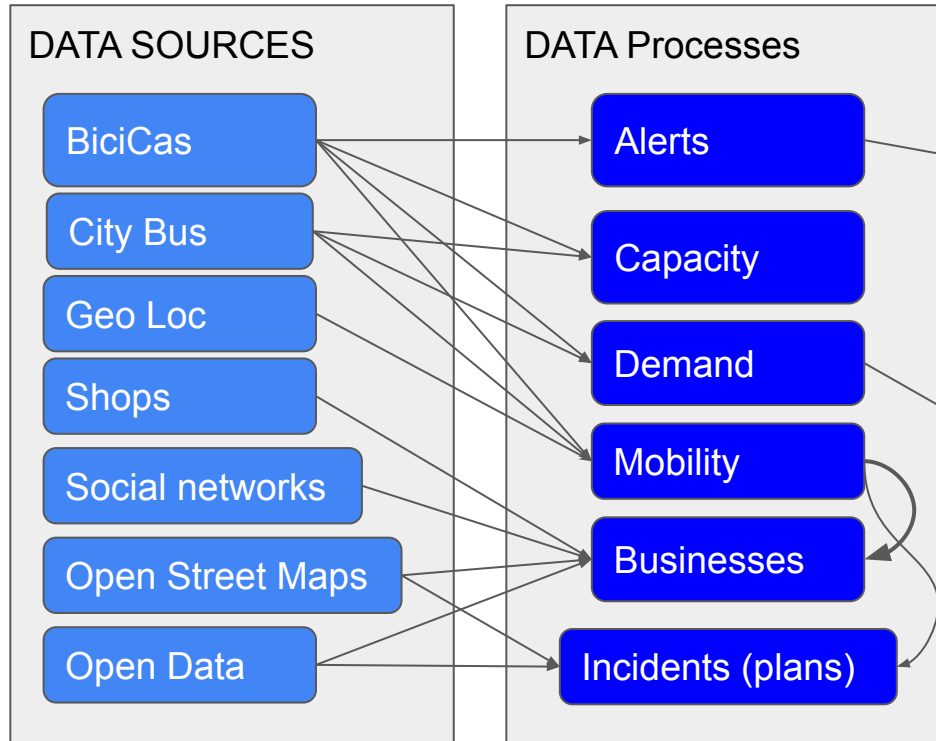
TRANSPORTS



CITIZENS



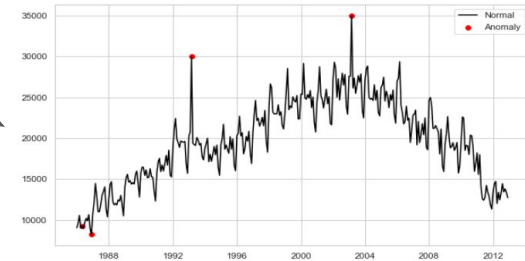
Data Chain Values



Visualizations and insights



Real-time alerts and reporting



Rent-a-bike datasets (PayNoPain)



bicycles

benches
bench_status
benches_anchors

benches_messages
benches_notifications

anchors

loans

incidents

putbacks

bans
users
subscriptions
payments



Datasets/tables

anchors(uuid, bench_uuid, number)

bans(uuid, use_uuid, date_start, date_expire, date_created, reason)

benches_anchors(bench_uuid, anchor_uuid, anchor_number, bicycle_number, date_created)

benches(uuid, name, latitude, longitude, date_created)

benches_notifications(uuid, bench_uuid, notification, priority, date_created)

bench_messages(uuid, message, date_start, date_expire, date_created)

bicycles(uuid, number, date_created, status)

incidents(uuid, error_key, use_uuid, anchor_uuid, status, date_created, bike_uuid)

loans(uuid, use_uuid, bicycle_uuid, anchor_uuid, date_created, type_access)

putbacks(uuid, bicycle_uuid, anchor_uuid, date_created)

users(change_uuid, uuid, ~~document~~, ~~name~~, ~~surname~~, ~~street~~, city, postcode, birthday, ~~phone~~, ~~email~~, gender, date_created, date_change, deleted)

bench_status(bench_uuid, date_lastseen, ip, queue, version, number_loans)

loan_historical(uuid, loan_uuid, use_uuid, bicycle_uuid, loan_anchor_uuid, loan_date_created, putback_uuid, putback_anchor_uuid, putback_date_created, type_access, app_loan)

Exercise: data chain values

HISTORICAL

Incidentes por demografía, geografía

Rutas más frecuentadas + temporalidad

Histórico de ocupación de bancadas con respecto a la demanda

Simulación de movimiento de recursos (modelo generativo probabilístico)

Estudiar el “estado” de calidad de las bancadas y las bicicletas (kms que lleva la bicicleta)

Recorrido de una bicicleta.

Popularidad de las bicicletas.

Bicicletas con problemas técnicos (se devuelven en la misma bancada con un tiempo corto entre loan y putback).

Mapa por código postal de demanda.

STREAMING

Incidentes: cálculo real de bicicletas operativas

Situación actual + modelo predictivo de demanda -> movilidad de recursos

Avisos dirigidos

Modelo predictivo de posible “relleno” de bicicletas en bancadas vacías

Recomendación de bancadas cercanas y tiempo previsto (localización del usuario)

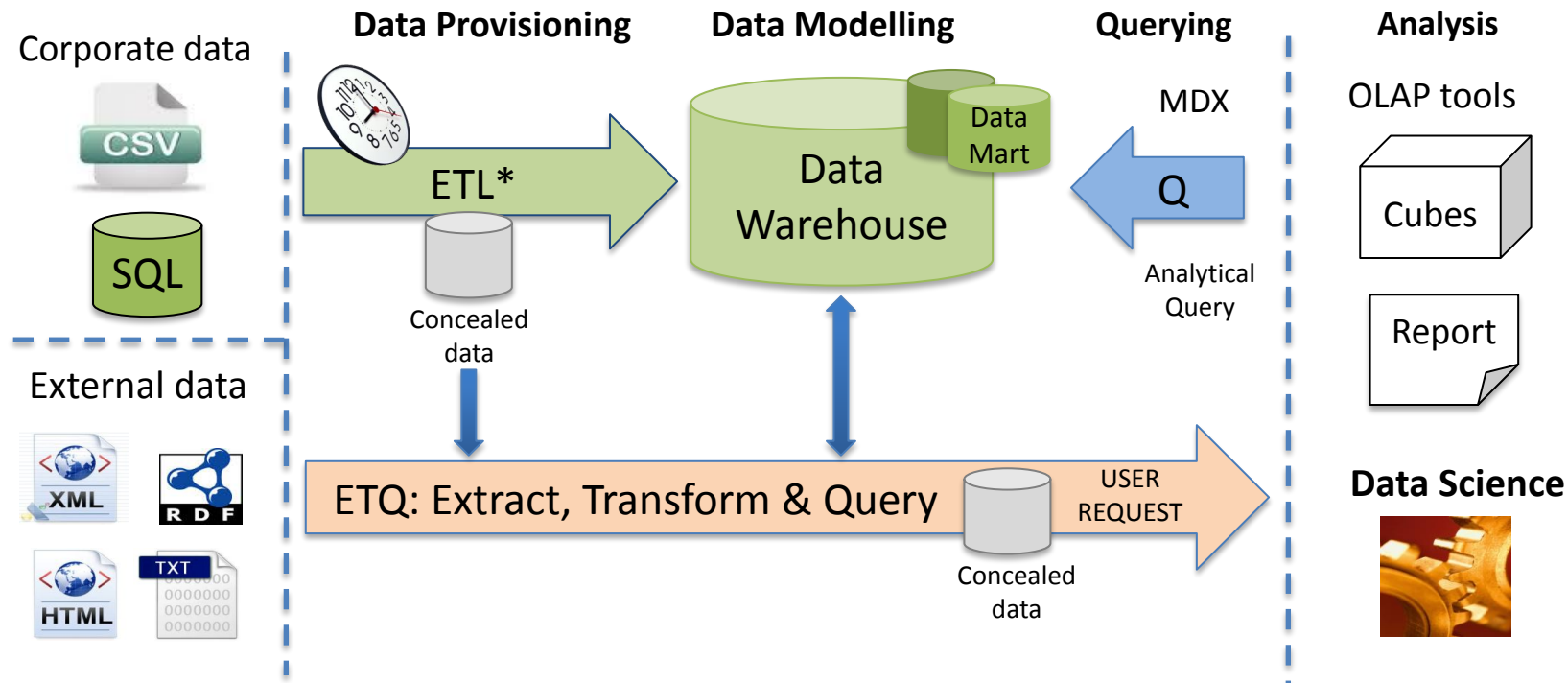
Alerta de bancada llena.

IMPLEMENTING DATA VALUE CHAINS

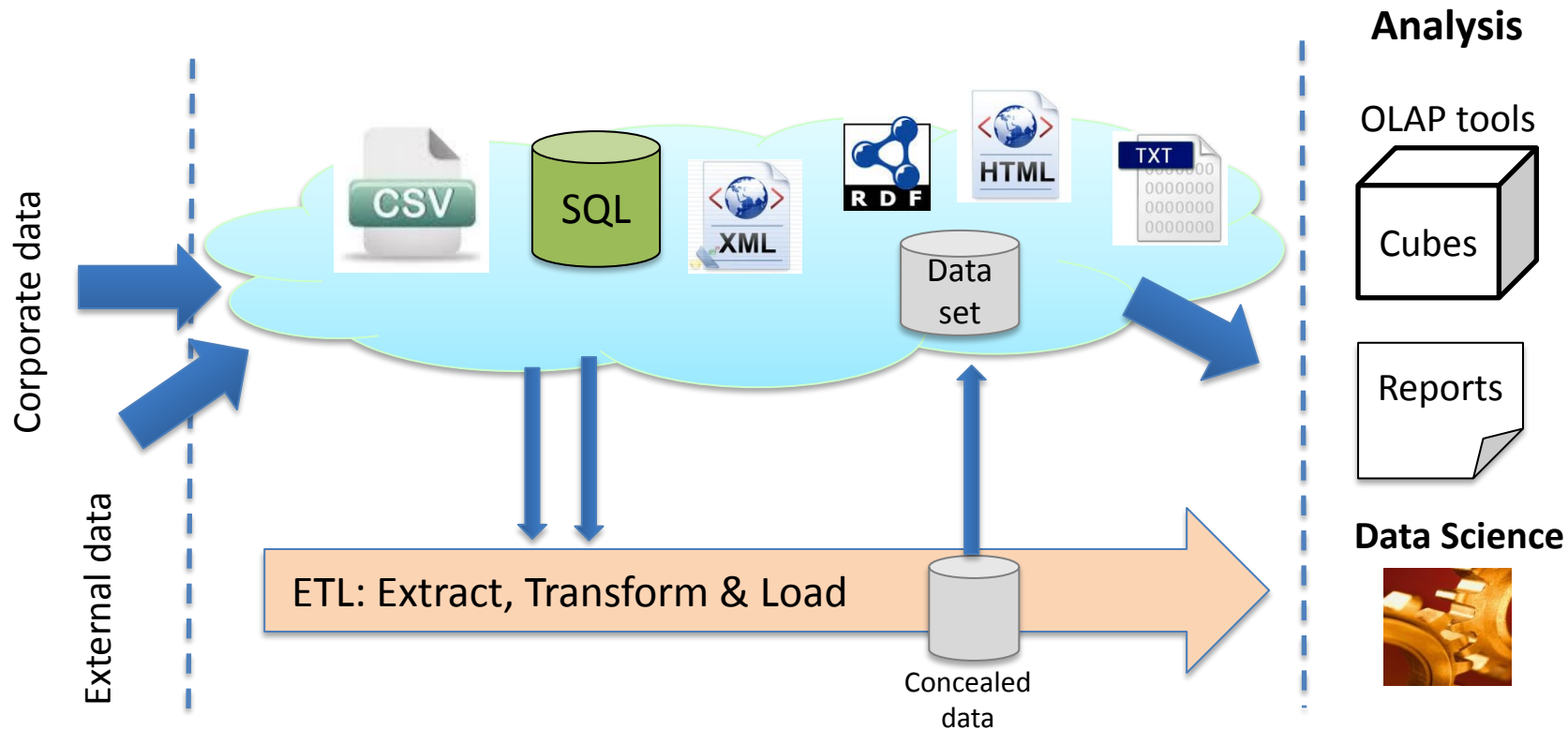
Big Data Architectures

- Data warehousing
- Data lakes
- Kappa architecture
- Lambda architecture

DATA WAREHOUSING

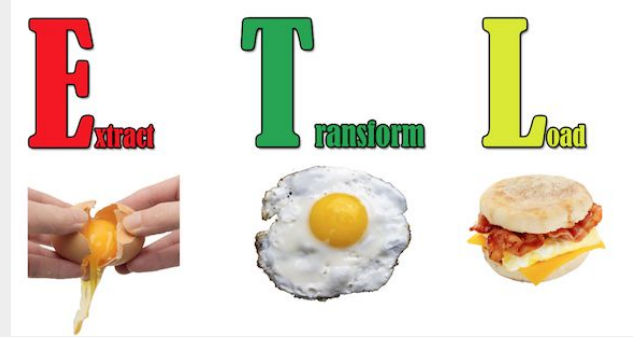


DATA LAKES



ETL/Q

Basic operations for capturing
and transforming data



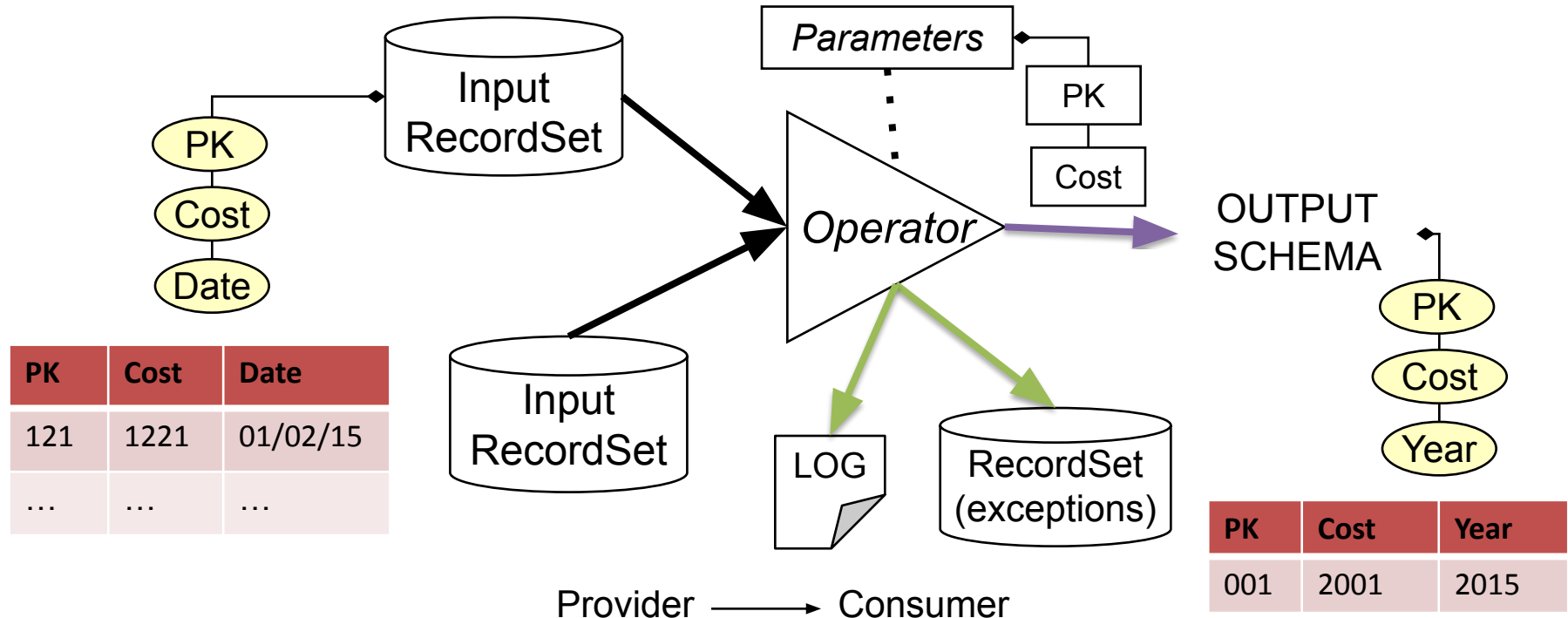
An ETL is a workflow of
transformations on source data,
which is aimed at **integrating** and
concealing data for analytical
purposes

Improving the quality of data

- **Data Cleaning:**
 - Detect and fix omitted or missing values
 - Remove noisy data (no relevant for analysis)
 - Detect incoherent data
- **Data integration:**
 - Duplicate removal
 - Gathering data of the same entities
 - Homogenization of terms and keys
- **Data transformation:**
 - Change of formats, scaling and normalization
- **Data reduction:**
 - Selection of attributes and aggregation of data

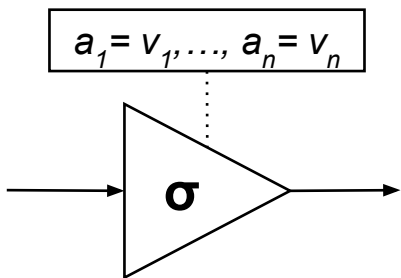
ETL: Workflows of transformation operators

Vassiliadis et al. Information Systems 30 (2005)

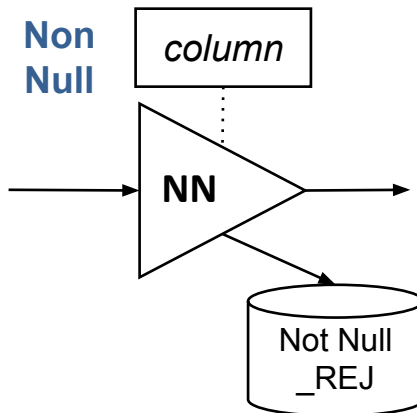


Some examples of operators

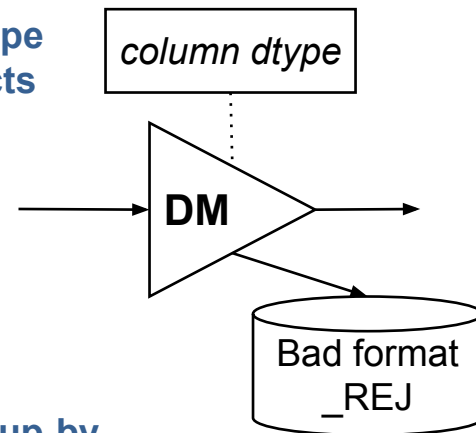
Row
selection



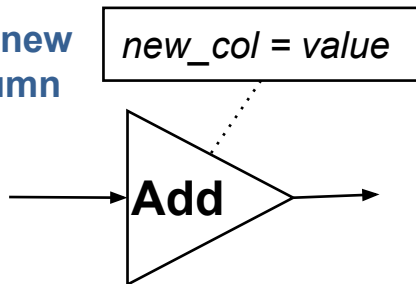
Non
Null



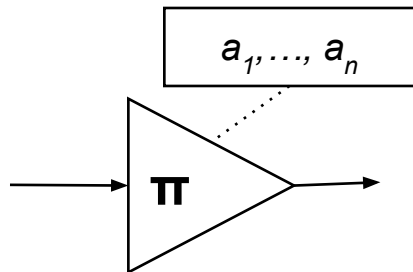
Data type
conflicts



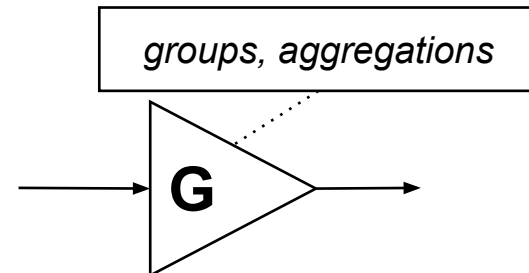
Add new
column



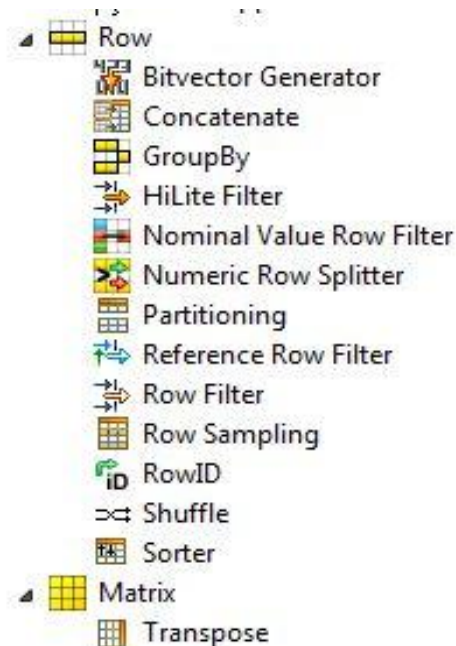
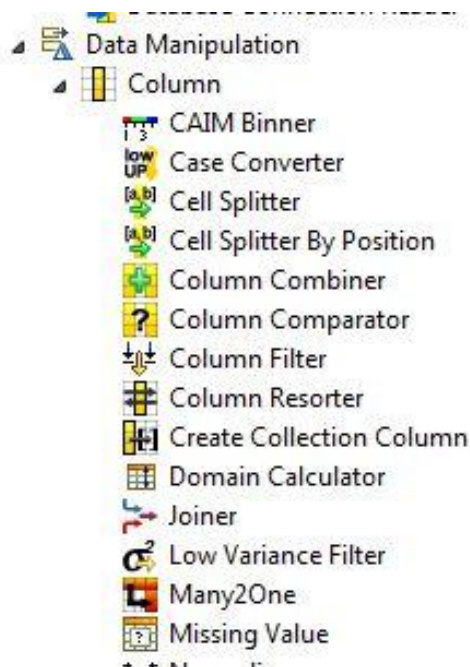
Project



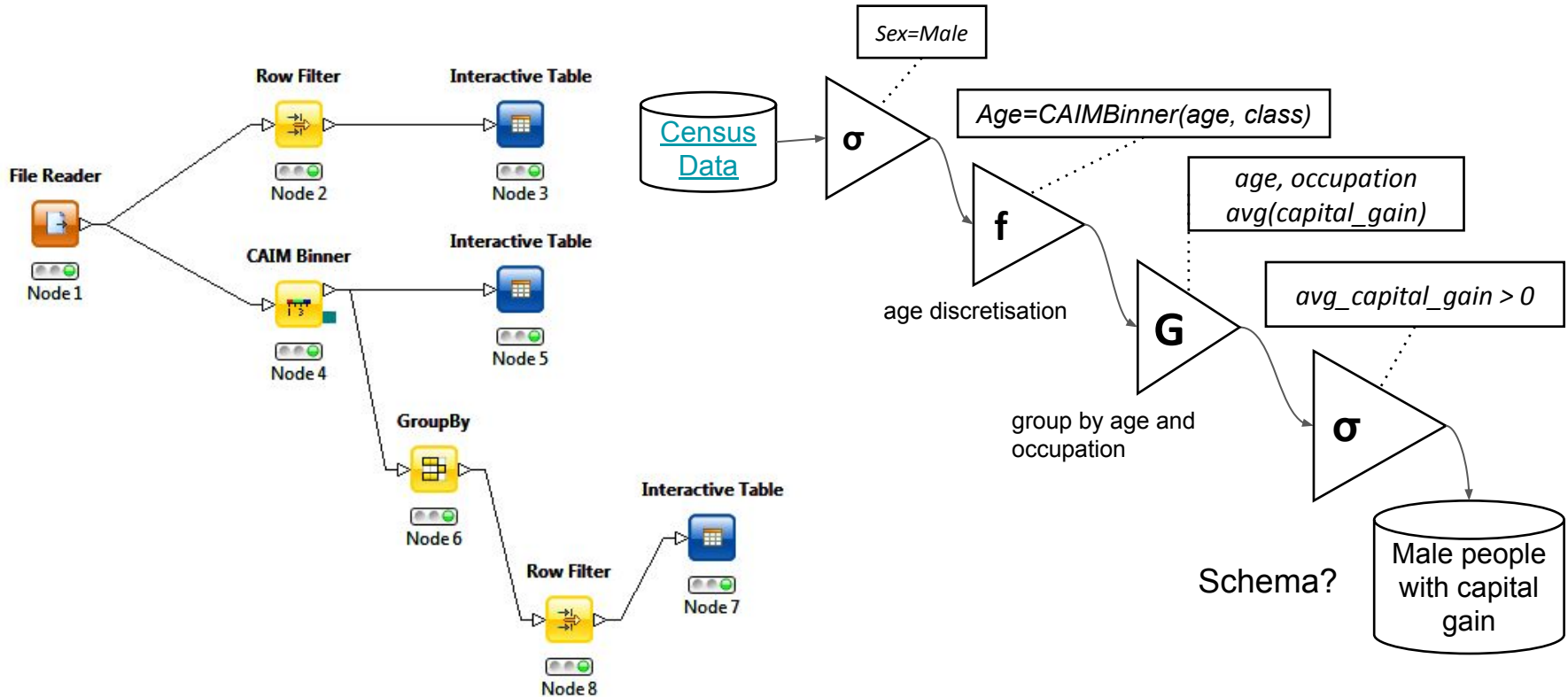
Group-by



Tools for ETLing: KNIME



Tools for ETLing: KNIME



Exercise: Pandas for ETLing

SELECTION

DUPLICATE REMOVAL

MISSING VALUES

GROUP-BY

JOIN

DIFFERENCE (anti-join)

UNION of recordsets

PROJECT columns

CONCAT columns

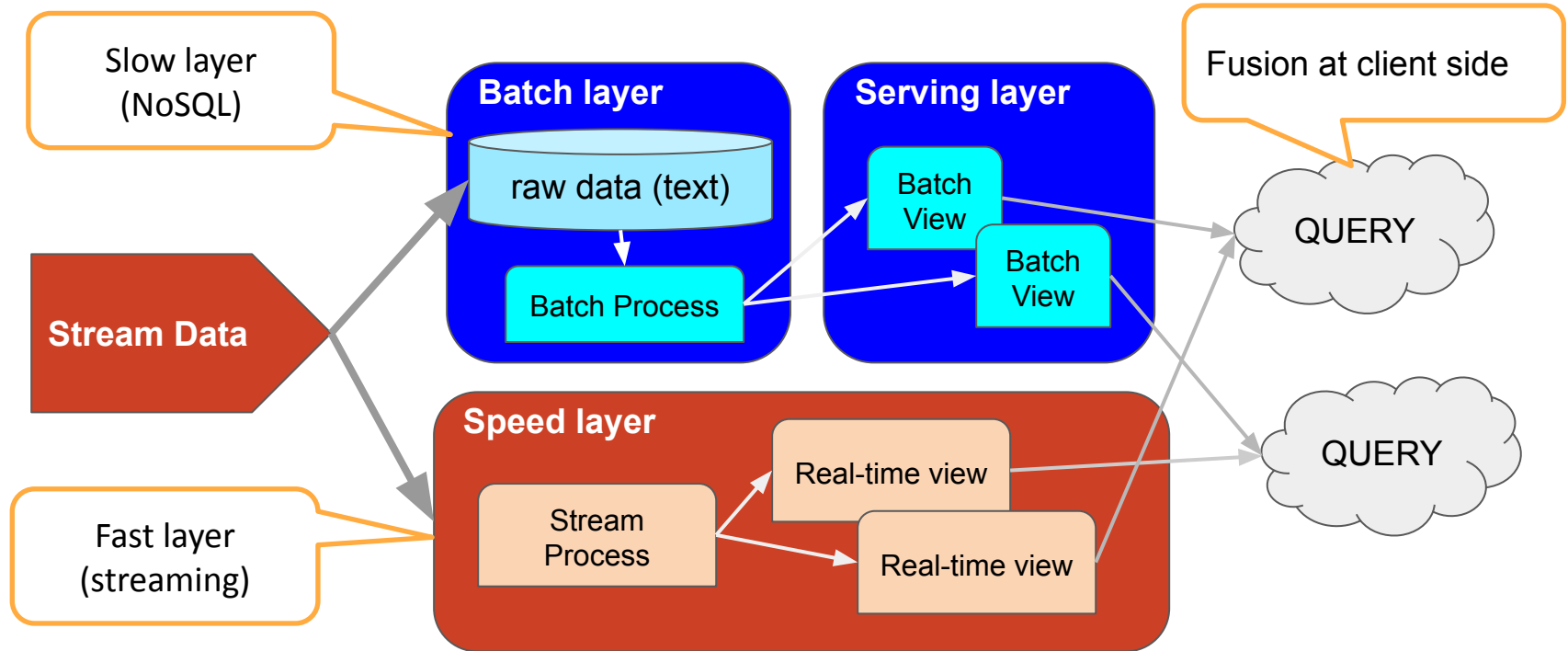
SPLIT a column

APPLY a function to a column (or several)

DISTRIBUTION (reverse of join)

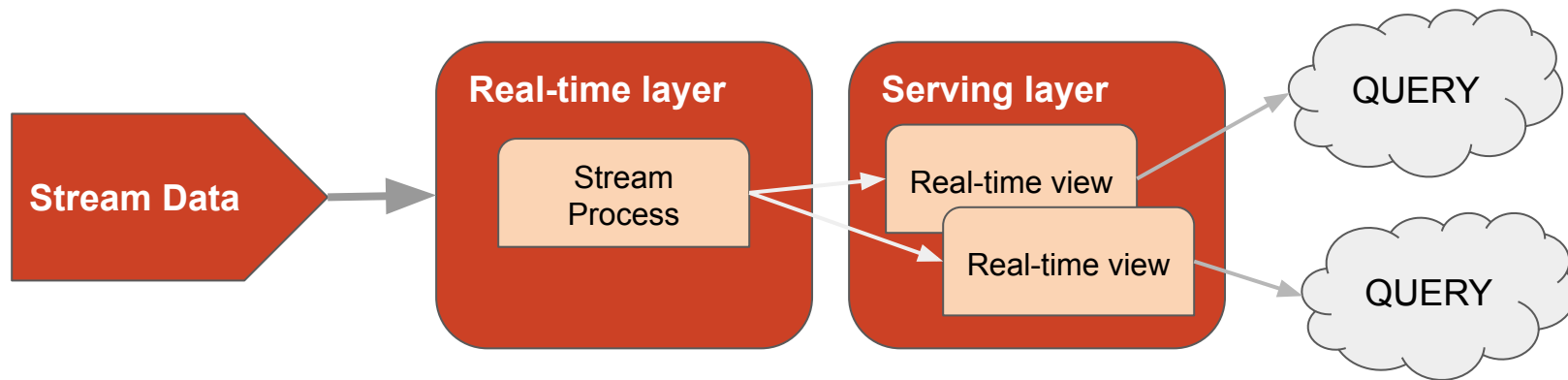
(complete with Pandas methods)

Lambda Architecture (Streaming data)



Kappa Architecture

- Use of “logs” (append-only files) to store data.
- The “batch” layer of the Lambda architecture is removed.
- Greatly simplifies code maintenance.
- Great processing speed.



Lambda or Kappa?

