

# A Scalable Data Transformation Platform

– *The Example of Data Anonymization*

Xiufeng Liu, *Postdoc.*

xiuli@dtu.dk

Department of Engineering Management, DTU

26 May 2015



# Table of Contents

- 1 Introduction
- 2 System Architecture
- 3 An Example – Data Anonymization

# Introduction

- ▶ In data warehousing, up to 80% of the time is spent on data pre-processing [1], including data extraction, transformation and loading;
- ▶ It is challenging for big data transformation and analytics;
- ▶ Diverse data warehousing tools in the market incorporate big data technologies for handling scalable data sets, but most of them lack of flexibility;

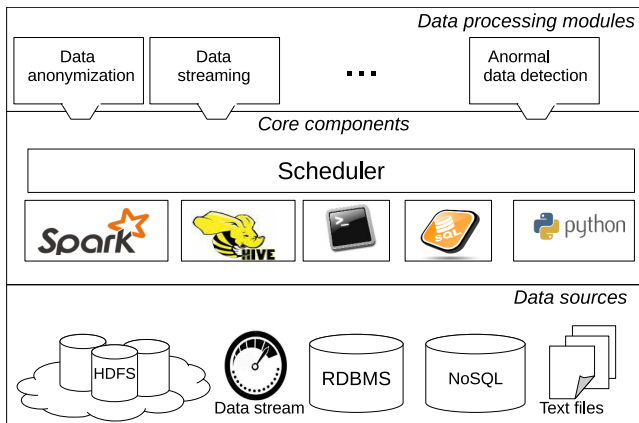
# Introduction

- ▶ In data warehousing, up to 80% of the time is spent on data pre-processing [1], including data extraction, transformation and loading;
- ▶ It is challenging for big data transformation and analytics;
- ▶ Diverse data warehousing tools in the market incorporate big data technologies for handling scalable data sets, but most of them lack of flexibility;

*Therefore, we intend to implement a scalable platform for big data transformation and analytics; and for research and production.*

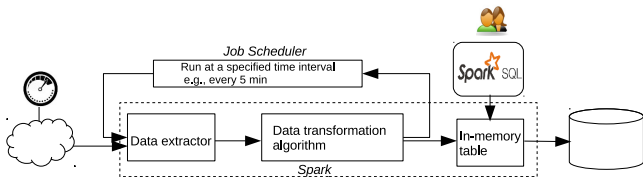
# System Architecture

- ▶ The building blocks of the system



# Realtime & Batch Processing

- ▶ Realtime stream processing
  - ▶ Spark Streaming (*in-memory based distributed computing framework*)



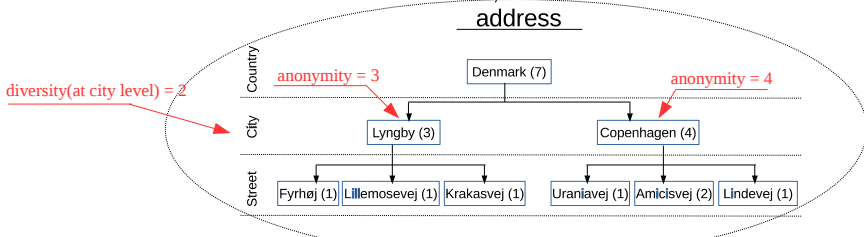
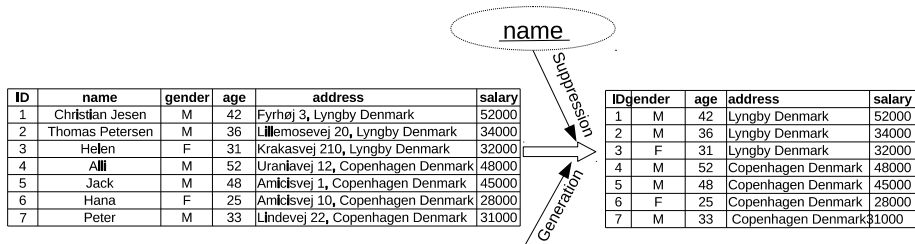
- ▶ Batch processing
  - ▶ Hive (*built on top of Hadoop MapReduce framework*)
- ▶ Job scheduling strategies
  - ▶ Deterministic – *For jobs not running in a cluster*
  - ▶ Undeterministic – *For jobs running in a cluster*

# A Transformation Example – Data Anonymization (Cont.)

- ▶ Anonymization Process:
  - ▶ *Data Source* → *Anonymize* → *Data Publishing & Sharing*
- ▶ Anonymization for two types of data:
  - ▶ Smart meter data:
    - ▶ Separate storing meter data from social-economic data, and use foreign-key referencing between tables
    - ▶ Aggregate data if possible
  - ▶ Social-economic data, e.g., customer information
    - ▶ Anonymize methods: *generalization* and *suppression*
    - ▶ Metric: *k-anonymity* and *l-diversity*
    - ▶ Balancing between *anonymity level* and *information loss*

# A Transformation Example – Data Anonymization (Cont.)

- Anonymize quasi-attribute values, i.e., *name* and *address*. If the requirement of anonymizing *address* attribute values is  $k\text{-anonymity} \geq 2$  and  $l\text{-diversity} \geq 2$ , then generalize to *city* level will suffice.





# A Transformation Example – Data Anonymization

- ▶ **Use Case:** We have a remote server containing customer data. We want to extract, anonymize, and publish the data.
- ▶ **The implementation:**
  1. Extract the raw customer data from remote server to staging database at local server, e.g., by SCP;
  2. Define the anonymization rules for suppressing/generalizing quasi-attribute values;
  3. Verify the results (by the online analytics tool)
  4. Load the anonymized data into database for publishing.
  5. Housekeeping, i.e., delete the data from staging database, and send email for notification, etc.

# References



C. Thomsen and T. B. Pedersen. pygrametl: A powerful programming framework for Extract-Transform-Load Programmers. In Proc. of DOLAP, pp. 49-56, 2011.