# Supporting a Data-Driven World through Data Integration and Data Cleaning

# Mourad Ouzzani

# Agenda

- Why is this an important problem?

- Data Civilizer - An end-to-end system

- Overview of some key components

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

http://visit.crowdflower.com/rs/416-ZBE-142/images/CrowdFlower_DataScienceReport_2016.pdf

QCRI
معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute
جامعة حمد بن خليفة
HAMAD BIN KHALIFA UNIVERSITY
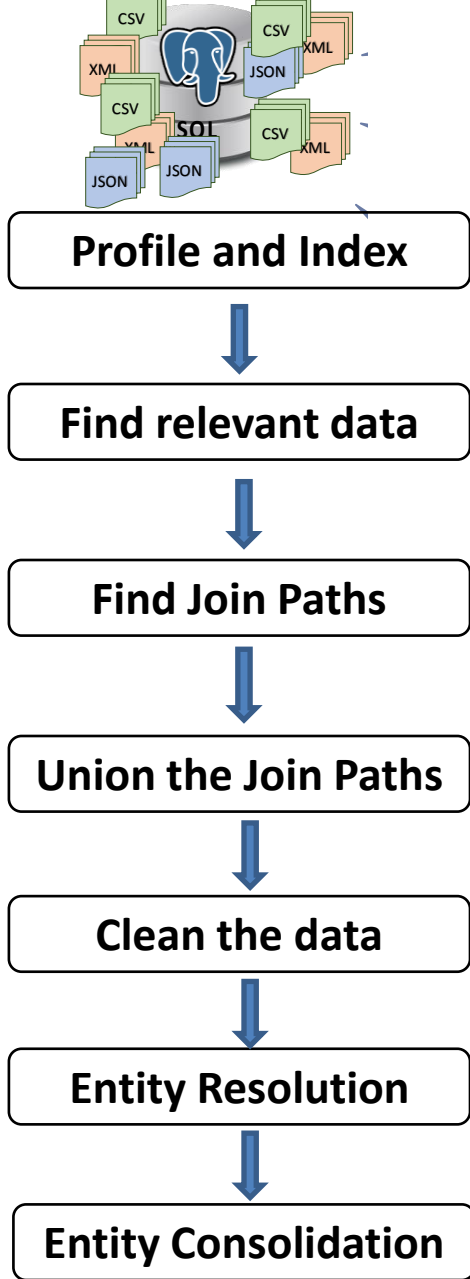
What's the least enjoyable part of data science?

- Mark Schreiber (Merck) reports that his data scientists spend 98% of their time, i.e. 39 hours/week, in grunt work and only 1 hour/week doing the job for which they were hired
- For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights (The New York Times https://www.nytimes.com/2014/08/18/technology/for-big-data-scientists-hurdle-to-insights-is-janitor-work.html)
- Nobody reports less than 80% grunt work

http://visit.crowdflower.com/rs/416-ZBE-142/images/CrowdFlower_DataScienceReport_2016.pdf

QCRI
معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

جامعة حمد بن خليفة
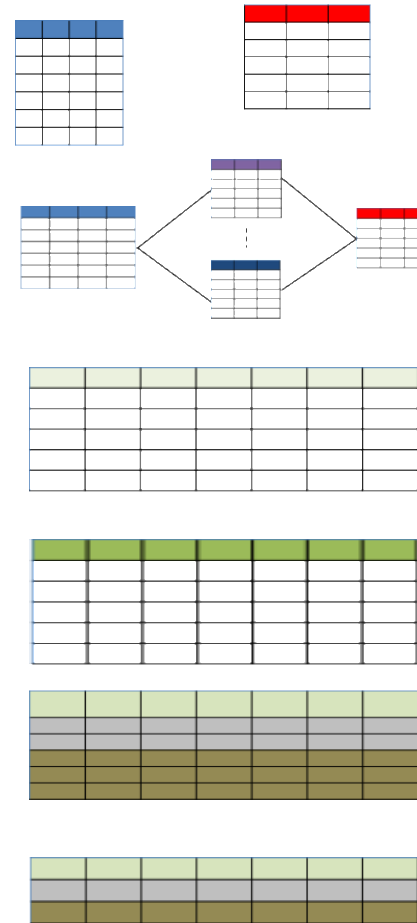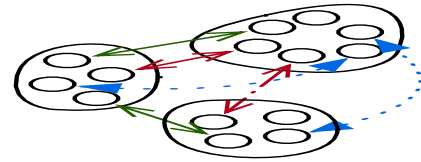HAMAD BIN KHALIFA UNIVERSITY

# We're building Data Civilizer to help ...

✔ **discover** data of interest from large numbers of data sets;

✔ **link** and **enrich** relevant data sets;

✔ **deduplicate** and **consolidate** the data;

✔ **clean** the data; and

✔ **iterate** through these tasks using a workflow system.

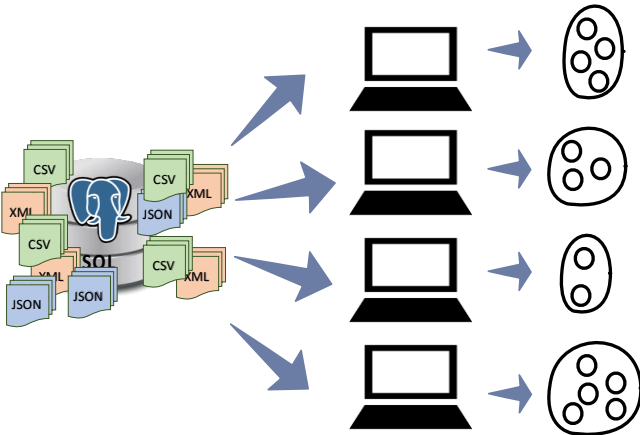Algorithms do the grunt work (80% of the pain) while data scientists can do what *they* are good at

QCRI
معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

جامعة حمد بن خليفة
HAMAD BIN KHALIFA UNIVERSITY

# Data Discovery



**Profiler**
Create Summaries

**Graph Builder**
Connect Summaries

**SRQL Query Processing**
Find relevant data

Edge and
Hyperedge
Indexes

Distributed
architecture to scale
data summarization

Scalable all-pairs
comparison of multiple
data types

Concise in-memory
indexes for interactive
query answering

# Entity Resolution using Deep Learning



A turn key solution using distributed representation (DR) and deep learning (DL)

- Tuples → high dimensional vectors where (semantically) similar tuples have a high (cosine) similarity
- Using pre-trained DR dictionaries (e.g., GloVe which is trained on a corpus of 840B tokens) → no need for manual feature engineering
- Much less training data
- Competitive or superior results wrt prior state-of-the-art methods
- Locality Sensitive Hashing-based blocking
  - automated and semantic blocking based on the entire tuple
  - no need for blocking functions from domain experts

QCRI
معهد قطر لبحوث الحوسبة
Qatar Computing Research Institute

جامعة حمد بن خليفة
HAMAD BIN KHALIFA UNIVERSITY

# Entity Consolidation



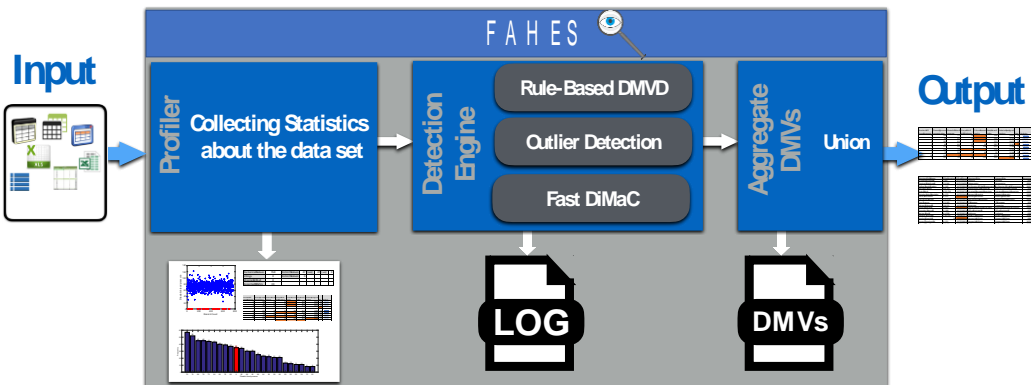From clusters of duplicate records to Golden Records



Cluster duplicates, detect matchings and group them, and ask a human

# Detecting Disguised Missing Values

| Source | Table Name | Column Name | DMV |
|---|---|---|---|
| | Pima Indians Diabetes | Diastolic Blood Pressurs | 0 |
| UCI Machine Learning Repository | adult | workclass | ? |
| | | education | Some-college |
| U.S. Food and Drug Administration | Adverse Event Reporting System (AERS) | EVENT_DT | 20010101, 20030101 |
| data.gov | Alleghency County WIC Vendor Location | Ref_ID | -1 |
| data.gov | Graduation Outcomes - School Level - Classes of 2005 - 2011 - SWD | Advanced Regents Num | s, - |
| data.gov.uk | Accidents 2015 | Junction Control | -1 |

DMV in different databases



- Rules to detect DMVs with special patterns, e.g., strings with repeated substrings

- Outlier detection algorithms

- A fast algorithm for detecting DMVs following a missing at random model

# The Civilizer Studio – Gluing Things Together

# Next Steps …

- Open-source release (ver 0.1)

- Get our technology in as many users' hands as possible

- Run tutorials in Spring 2018

شكراً

أسئلة؟

Thank You

Questions?