# Big Data Integration for Data-Centric AI

itaDATA2022

**Sonia Bergamaschi, Domenico Beneventano, Giovanni Simonini, Luca Gagliardelli, Adeel Aslam, Giulio De Sabbata, Luca Zecchini**

DBGroup @ Department of Engineering "Enzo Ferrari"

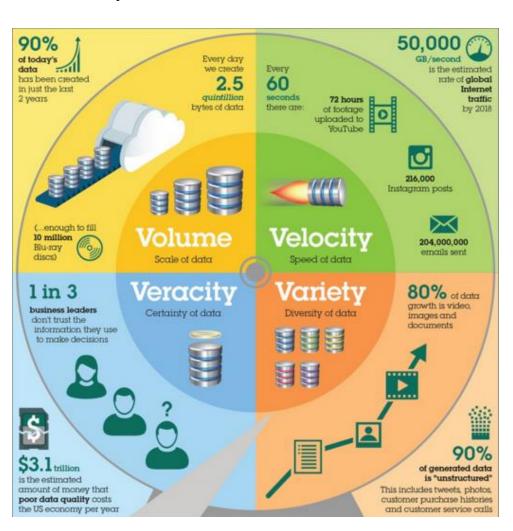sonia.bergamaschi@unimore.it

www.dbgroup.unimore.it

**DBGroup**

DB Group @ unimore

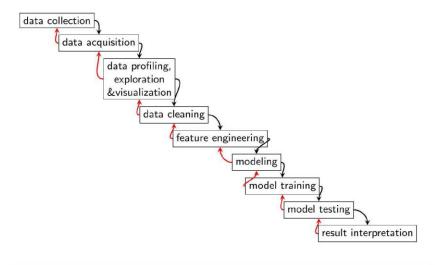In today's world, we need to deal with a huge amount of data…





Challenging for traditional paradigms…

- Big Data Management
- Big Data Science
- Big Data Integration

DB Group @ unimore

**Ingredients:**
50g statistics
120g linear algebra
200g programming
1kg visualisation
300g software
engineering

**Additional skills:**
creativity
out of the box thinking
grit
team spirit

Artificial Intelligence/ Machine Learning

Data Management

Data Mining

Application Domain

© istock.com sasilsolutions

data collection
data acquisition
data profiling, exploration &visualization
data cleaning
feature engineering
modeling
model training
model testing
result interpretation

This is **at the same time** a process model **and** a dataflow.

# Big Data Integration + Big Data Analysis
## (BI, Statistics, Data Mining, Math + Data-Centric AI)

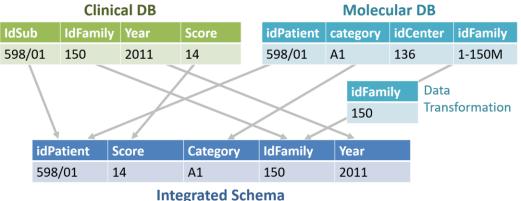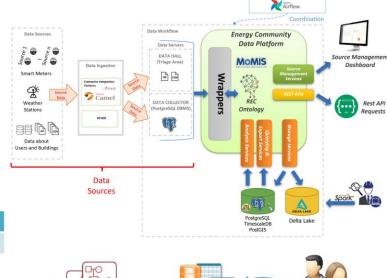| Model-Centric | Data-Centric |
|---|---|
| - Collects as much data as possible<br>- Iteratively **improves the model** to deal with the noise in the data | - Holds the model fixed<br>- Iteratively **improves the quality of the data** to obtain good results |

Andrew Ng

**MoMIS**
Mediator envirOnment for Multiple Information Sources

**DataRiver**
open source data management

DB Group @ unimore

Data Integration is the process of consolidating data from a **set of heterogeneous data sources** into a **single uniform dataset** or view on the data.
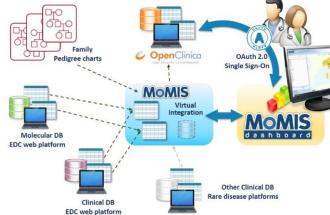
**MOMIS** is an open-source (Big) Data Integration system able to aggregate data from heterogeneous and distributed sources in a semi-automatic way, exploiting the **semantic relationships** existing in the data sources.



**Clinical DB**

| IdSub | IdFamily | Year | Score |
|---|---|---|---|
| 598/01 | 150 | 2011 | 14 |

**Molecular DB**

| idPatient | category | idCenter | idFamily |
|---|---|---|---|
| 598/01 | A1 | 136 | 1-150M |

| idFamily | Data Transformation |
|---|---|
| 150 | |

**Integrated Schema**

| idPatient | Score | Category | IdFamily | Year |
|---|---|---|---|---|
| 598/01 | 14 | A1 | 150 | 2011 |

[1] S. Bergamaschi, S. Castano, M. Vincini: *Semantic Integration of Semistructured and Structured Data Sources*. SIGMOD Rec. 28(1): 54-59 (1999)
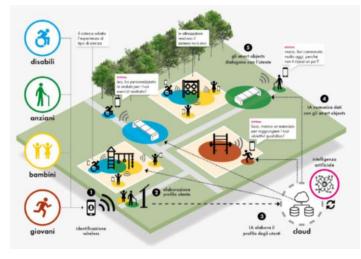
DB Group @ unimore



Creating contexts of smart inclusive parks to promote the adoption of correct active lifestyles and good health for all age groups.
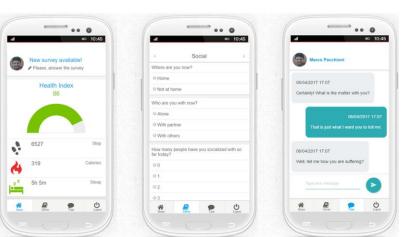
OSO (Outdoor Smart Objects): furnishings and tools (IoMT) able to recognize the user and dynamically adapt their performance.

User activities monitored through AI techniques to provide a **personalized feedback** promoting physical well-being and health.
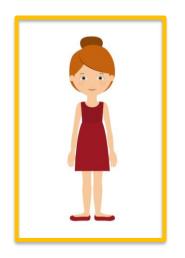
# Entity Resolution (ER) a.k.a. Record Linkage (RL)

**Data Discovery**
- Source Collection
- Source Analysis
- Data Patterns

**Data Validation**
- Correctness
- Completeness
- Data Quality Constraints

**Data Cleaning**
- Missing Values
- **Entity Resolution**
- Data Transformation

**Semantic Enrichment**
- Semantic Annotation
- Standard Vocabularies
- Ontologies

**Data Integration**
- Schema Generation
- Schema Mapping

Given one or more data sources, Entity Resolution (ER) is the task of identifying the **records (entity profiles)** that refer to the **same real-world object (entity)**.

## Data Source A

| | Name | Surname | Address | Sex |
|---|---|---|---|---|
| r1 | Mary-Ann | White | West Main Street 29, 12068, Fonda, NY, New York | F |
| r2 | Thomas J. | Franklin | 50 Liverpool Street, London | M |

## Data Source B

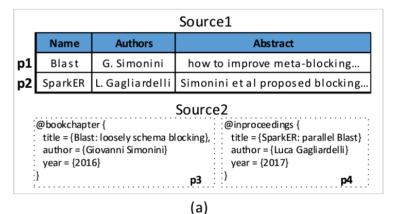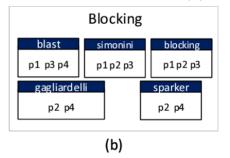| | Name | Residence | Age | Gender | |
|---|---|---|---|---|---|
| Franklin, Tom | London (UK) | 25 | Male | | r3 |
| Withe, Mary Ann | New York (USA) | 29 | Female | | r4 |

7

ER is a challenging task:

- Hard to scale (quadratic complexity) → Blocking functions
- Hard to detect duplicates → Matching functions
- Even harder with privacy constraints → PPRL

**Dirty Data**
(duplicate records)

**Clean Data**
(representative records)



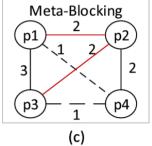| Blocking | Block Cleaning | Entity Matching | Entity Clustering | Data Fusion |
|---|---|---|---|---|
| Cluster together similar records (**blocking function**) | Refine blocks to increase precision without affecting recall | Compare the candidate pairs of records (**matching function**) | Partition the retained records into real-world entities | Obtain from each cluster a single clean record representative of the entity (**conflict resolution function**) |

DB Group @ unimore

8

Reducing the number of candidate pairs by discarding obvious non-matches.
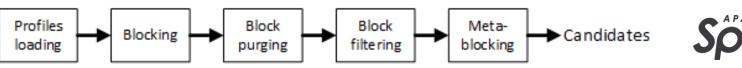


(a)

(b)

(c)

[1] G. Simonini, S. Bergamaschi, H. V. Jagadish: *BLAST: a Loosely Schema-aware Meta-blocking Approach for Entity Resolution*. PVLDB 9(12): 1173-1184 (2016)

[2] G. Simonini, L. Gagliardelli, S. Bergamaschi, H. V. Jagadish: *Scaling entity resolution: A loosely schema-aware approach*. Inf. Syst. 83: 145-165 (2019)

[3] L. Gagliardelli, G. Simonini, D. Beneventano, S. Bergamaschi: *SparkER: Scaling Entity Resolution in Spark*. EDBT, 602-605 (2019)

[4] L. Gagliardelli, G. Papadakis, G. Simonini, S. Bergamaschi, T. Palpanas: *Generalized Supervised Meta-blocking*. PVLDB 15(9): 1902-1910 (2022)
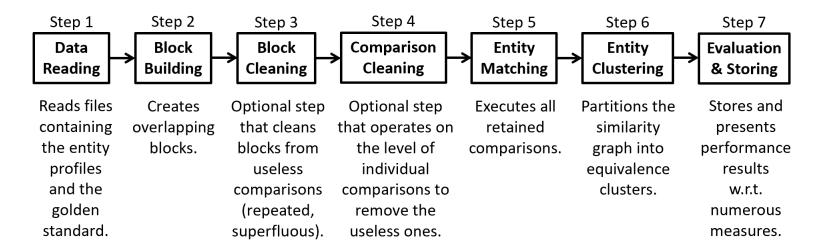
**SparkER** is an ER framework for Apache Spark. Its source code is available on GitHub and its complete and detailed documentation can be found on Read the Docs.

DB Group @ unimore

JedAI implements the whole **schema-agnostic, end-to-end workflow** for ER and can be used as an open-source library, as a desktop application, or as a workbench.

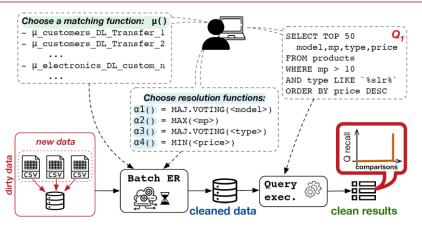| Step 1 | Step 2 | Step 3 | Step 4 | Step 5 | Step 6 | Step 7 |
|---|---|---|---|---|---|---|
| **Data Reading** | **Block Building** | **Block Cleaning** | **Comparison Cleaning** | **Entity Matching** | **Entity Clustering** | **Evaluation & Storing** |
| Reads files containing the entity profiles and the golden standard. | Creates overlapping blocks. | Optional step that cleans blocks from useless comparisons (repeated, superfluous). | Optional step that operates on the level of individual comparisons to remove the useless ones. | Executes all retained comparisons. | Partitions the similarity graph into equivalence clusters. | Stores and presents performance results w.r.t. numerous measures. |

[1] G. Papadakis, G. Mandilaras, L. Gagliardelli, G. Simonini, E. Thanos, G. Giannakopoulos, S. Bergamaschi, T. Palpanas, M. Koubarakis: *Three-dimensional Entity Resolution with JedAI*. Inf. Syst. 93: 101565 (2020)

- Project website: http://jedai.scify.org
- GitHub repository: https://github.com/scify/JedAIToolkit

# Beyond Traditional Batch ER: Progressive Approaches



Choose a matching function: μ()
- μ_customers_DL_Transfer_1
- μ_customers_DL_Transfer_2
  ...
- μ_electronics_DL_custom_n
  ...

```
SELECT TOP 50               Q₁
  model,mp,type,price
FROM products
WHERE mp > 10
AND type LIKE `%slr%`
ORDER BY price DESC
```

Choose resolution functions:
α1() = MAJ.VOTING(<model>)
α2() = MAX(<mp>)
α3() = MAJ.VOTING(<type>)
α4() = MIN(<price>)

**Useless comparisons (produce entities that will surely not appear in the result of the query)**
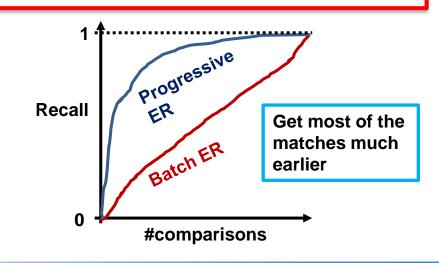
Our time, resources and affordable costs (e.g., pay-as-you-go in cloud) are often **limited**

**PROGRESSIVE ER: Maximize** the number of retrieved **matches** in a limited amount of time (driven by **matching likelihood)**
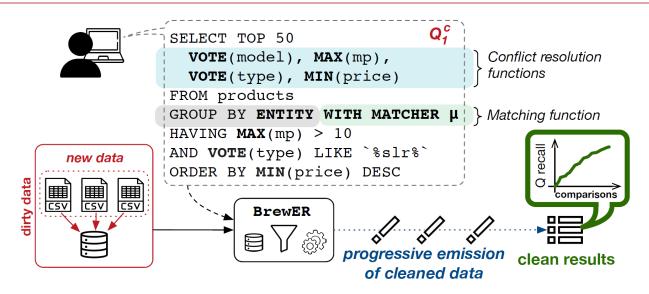
[1] G. Simonini, G. Papadakis, T. Palpanas, S. Bergamaschi: *Schema-agnostic Progressive Entity Resolution*. ICDE: 53-64 (2018)

**Get most of the matches much earlier**

11

DB Group @ unimore



```
SELECT TOP 50                        Q₁ᶜ
    VOTE(model), MAX(mp),
    VOTE(type), MIN(price)
FROM products
GROUP BY ENTITY WITH MATCHER μ
HAVING MAX(mp) > 10
AND VOTE(type) LIKE `%slr%`
ORDER BY MIN(price) DESC
```

Conflict resolution functions

Matching function

new data

dirty data

CSV CSV CSV

BrewER

progressive emission of cleaned data

clean results

**Agnostic** approach to blocking and matching functions
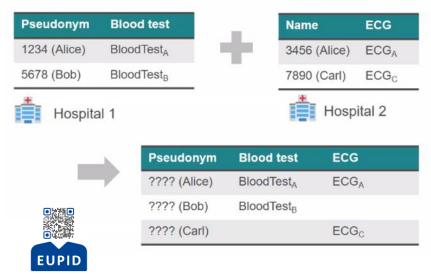
**Clean queries on dirty data**

❑ **QUERY-DRIVEN**: ER only on the portion of dataset useful to answer the query (according to the HAVING clauses)

❑ **PROGRESSIVE**: return the entities in the result <u>in the right order</u> as soon as they are obtained (according to the ORDER BY clause)
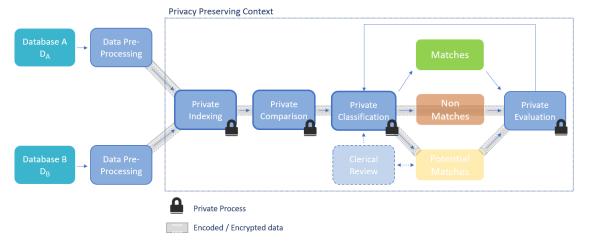
[1] G. Simonini, L. Zecchini, S. Bergamaschi, F. Naumann: *Entity Resolution On-Demand*. PVLDB 15(7): 1506-1518 (2022)

VLDB 2022

# Privacy-Preserving Record Linkage (PPRL)

Identifying and linking the **records (profiles)** that refer to the **same real-world object (entity),** across several data sources held by different parties, in a manner that prevents both the computation and the output of the computation from revealing (to any **internal parties** involved in the process and **external adversaries**) any private sensitive information about the entities represented in the data.



**Privacy Preserving Temporal Record Linkage Project:** design and implement in prototype form PPTRL techniques as part of the "Recidivism Data Mart and Criminal Data Warehouse" project, funded by the CRUI Foundation, which is part of the broader "PNRR Justice Plan Data Lake".

# THANK YOU FOR THE ATTENTION!