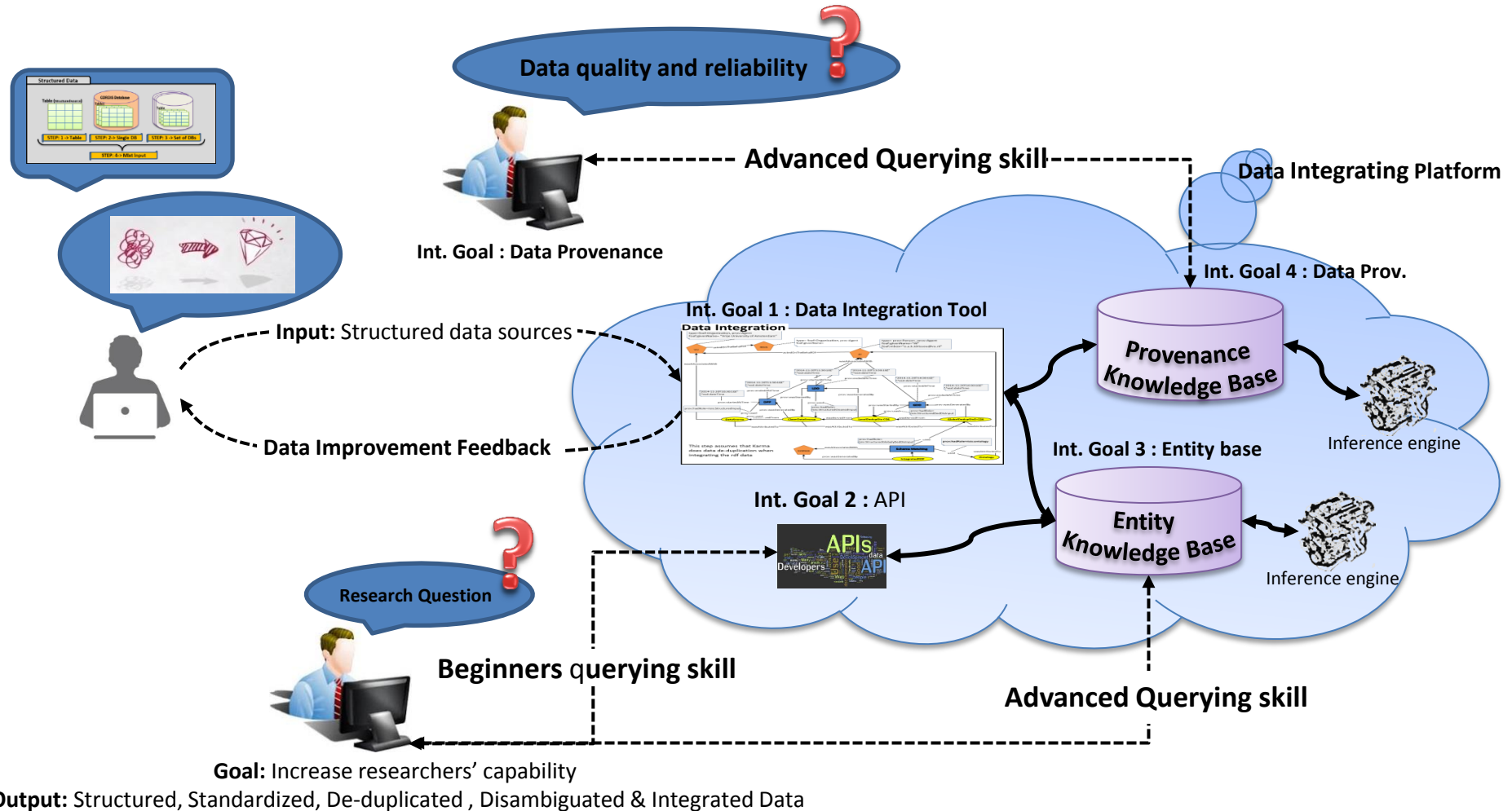# Entity-Centric Data Integration & Structured Data

Al Koudous Idrissou

January 26, 2015

# Overview

- Goal
- Problem
- Terminologies
- Data integration Tasks
- Tools towards data integration
- Plan

# Goal

**Increase researchers' capability** by creating public single instance **Entity Store** and **Data Provenance Store,** and populate them with an **Automated Data Integrating Tool** using **structured data sources** as input. Among others, the Data Integration Platform ensures a **cleaned, structured, normalized, de-duplicated** and **disambiguated** Entity Knowledge-base and a mean to assess its data's **Quality** and **Reliability** through **data provenance**.

Data quality and reliability

Advanced Querying skill

Data Integrating Platform

Int. Goal : Data Provenance

Int. Goal 4 : Data Prov.

**Input:** Structured data sources

Int. Goal 1 : Data Integration Tool

Provenance Knowledge Base

Inference engine

Int. Goal 3 : Entity base

Data Improvement Feedback

Int. Goal 2 : API

APIs

Entity Knowledge Base

Inference engine

Research Question

Beginners querying skill

Advanced Querying skill

**Goal:** Increase researchers' capability

**Output:** Structured, Standardized, De-duplicated , Disambiguated & Integrated Data

# Data in the real world

- Ideally, given 2 data sets about the same thing it should be relatively simple to combine them.

- It is not the case in real world.
  - Real data is **messy**
  - Real data is **inconsistent**
  - Real data contains **ambiguity**
  - Real data is not **normalized**
  - Real data often has no **uniquely and globally identifiable entities**

- **Difficult to link** the available data between different sources.

- Big limitation in gathering the **right information** in an acceptable time interval.

- Successfully gathered data often presents a lack of meta data about data's **origin**, **context** and **quality**

# Structured Messy and Inconsistent Data

# What is an **Entity**?

Person

Amsterdam

Capital of Kingdom of the Netherlands

Amsterdam is the capital city and most populous city of the Kingdom of the Netherlands. Its status as the Dutch capital is mandated by the Constitution of the Netherlands though it is not the seat of the Dutch government, which is The Hague. Wikipedia

**Area:** 219 km²
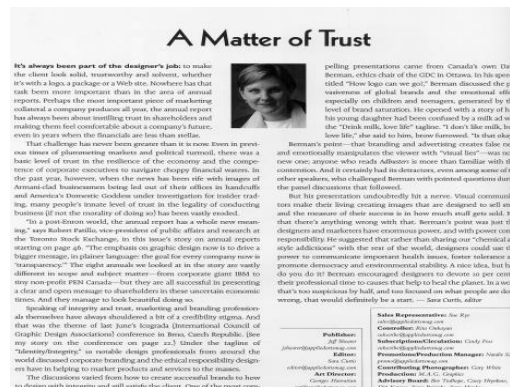**Weather:** 9°C, Wind SE at 11 km/h, 81% Humidity
**Local time:** Thursday 9:33 AM
**Province:** North Holland
**Population:** 779,808 (2011) UNdata

City

Organization

## A Matter of Trust

Article

Products

**Real World Things**
**OR**
**Anything**

W3C Provenance Ontology defines entity as a physical, digital, conceptual, or other kind of **thing with some fixed aspects**; **entities may be real or imaginary**
.

# What is **Data Provenance**?

**Def**: "**Provenance** is to **electronic data** what a record of ownership is to a work of art." [Moreau et Al, 2007]

**W3C** **Goal**: Provenance captures information about entities, activities and people involved in producing a piece of data or thing which can be used to form assessments about it quality, Reliability or trustworthiness

## Purpose
- **Understand** how data was collected
- Determine **ownership** and **wrights** over an object
- Making judgment about information to determine whether to **trust it**.
- **Verify** the **process** and steps used to obtain a result given a set of **requirements**
- **Reproduce** how something was generated

# Prov Data flow

DataSource ←—wasDerivedFrom— CleanDataSource ←—wasDerivedFrom— LocalDedupDis-CDS ←—wasDerivedFrom— GlobalDedupDisD-CDS ↑ wasDerivedFrom IntegratedRDF

# Prov Process flow

"2014-11-20T11:30:10Z"^^xsd:dateTime

"2014-11-20T13:30:10Z"^^xsd:dateTime

prov:startedAtTime

prov:endeddAtTime

"2014-11-20T11:30:10Z"^^xsd:dateTime

"2014-11-20T14:30:10Z"^^xsd:dateTime

"2014-11-20T16:30:10Z"^^xsd:dateTime

prov:endeddAtTime

LDD

prov:startedAtTime

prov:endeddAtTime

prov:wasGeneratedBy

"2014-11-20T10:30:10Z"^^xsd:dateTime

prov:startedAtTime

DC

prov:wasStartedBy

prov:hadRole=risis:StructuredCleanedInput

GDD

prov:hadRole=risis:StructuredInput

prov:used prov:wasGeneratedBy

prov:used

prov:wasStartedBy

prov:used

prov:hadRole=risis:StructuredDedDisInput

prov:wasGeneratedBy

prov:used

DataSource ←—wasDerivedFrom— CleanDataSource —wasDerivedFrom→ LocalDedupDis-CDS ←—wasDerivedFrom— GlobalDedupDisD-CDS

prov:hadRole=risis:StructuredGlobalyDedDisInput

used

prov:hadRole=risis:ontology

Schema Matching

IntegratedRDF

used

Ontology

# What is a **Knowledge-base System ?**

❖ **Data** raw signals" [ex: . . . - - - . . .]

❖ **Information** attaches **meaning** to data [Ex: S O S]

❖ **Knowledge** attaches **purpose** and **competence** to information
[Ex: emergency alert -> start rescue operation ] **->** potential to generate action

**Fact:** Alice is an animated character

A **knowledge base** (**KB**) is a technology used to sore complex information (**facts about the world**) used by a computer system.

➔ While a **knowledge-base** stores facts about the world, an **Inference engine** reasons about facts to infer new facts.

A **knowledge-base system** consists a knowledge-base and an inference engine

# Why **Data Integration?**

Need for seamlessly getting a complete and accurate view on the same thing.

**Thing**: Professors in the Netherland?
**Resources**: VU – UVA – INHOLLAND ….

**Problem**

Based on the same requirements, 2 databases designed by 2 different persons will look significantly different (Schema & representation)

**Employee Database [VU]**

**FullTimeEmp**(ssn, empID, firstName, middleName, lastName)
**Hire**(empID, hireDate, recruiter)
**TemEmplofees**(ssn, hireStart, hireEnd, name, hourlyRate)

**Semantic heterogeneity**
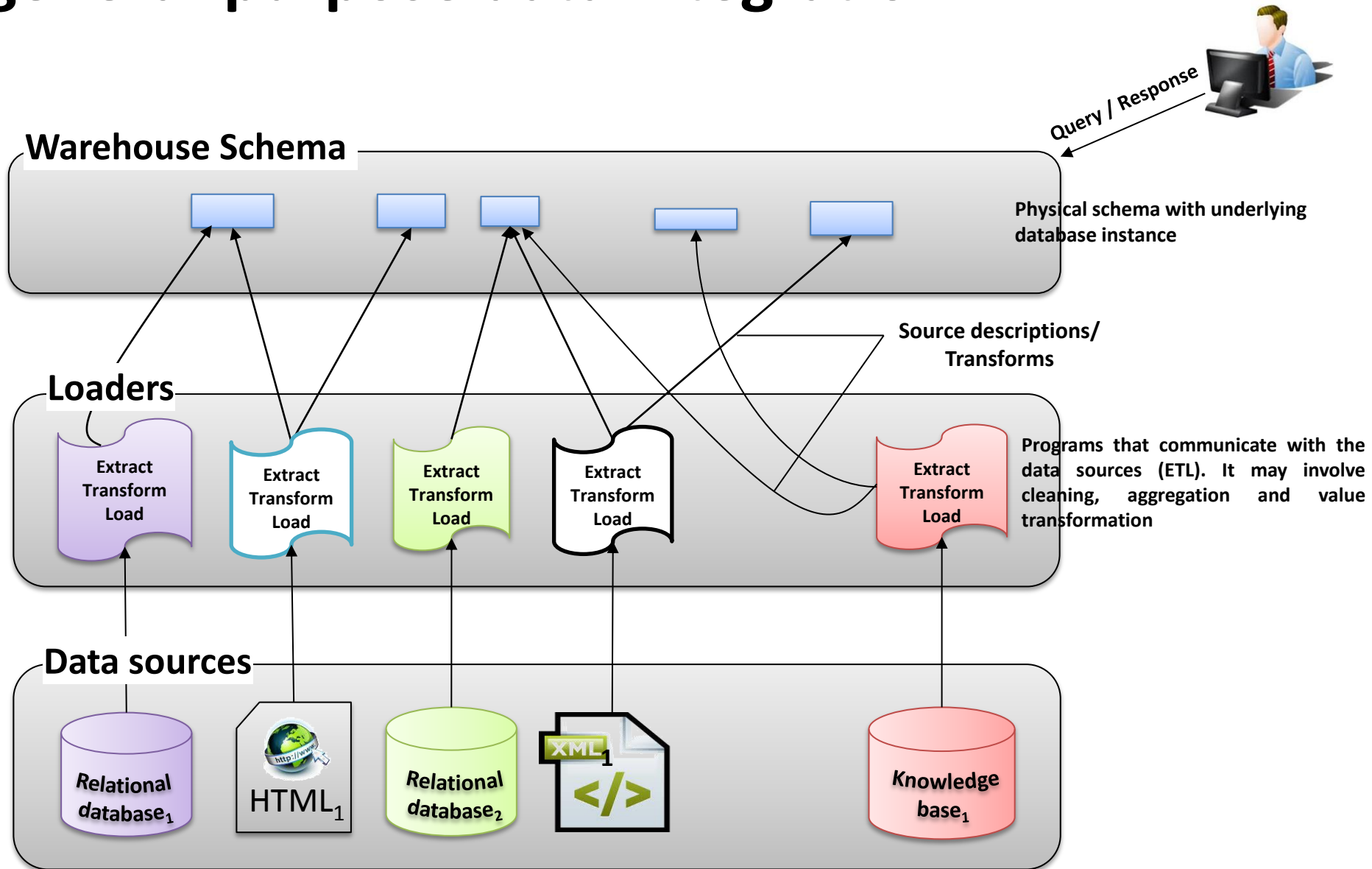
**Employee Database [UVA]**

**Emp**(ID, firstNameMiddleInitial, lastName, salary)
**Hire**(ID, hireDate, recruiter)

# Basic Architecture of a Physical general-purpose data integration

**Query / Response**

**Warehouse Schema**

Physical schema with underlying database instance

Source descriptions/ Transforms

**Loaders**

Extract Transform Load

Extract Transform Load

Extract Transform Load

Extract Transform Load

Extract Transform Load

Programs that communicate with the data sources (ETL). It may involve cleaning, aggregation and value transformation

**Data sources**

Relational database$_1$

HTML$_1$

Relational database$_2$

XML$_1$

Knowledge base$_1$

# Basic Architecture of a Virtual general-purpose data integration

# Tasks toward data Integration?

- **Input**: Structured data about Entity
- Data pre-processing
- Schema matching and Mapping
- Data Deduplication
- Data Disambiguation
- Data reconciliation
- Data Consistency
- Data Integrity
- **Output**: RDF

# Data Pre-processing

## Removing Unwanted Characters and Tokens

\FR?D?RIC JOLIOT-CURIE\" NATIONAL RESEARCH INSTITUTE FOR RADIOBIOLOGY AND RADIOHYGIENE"

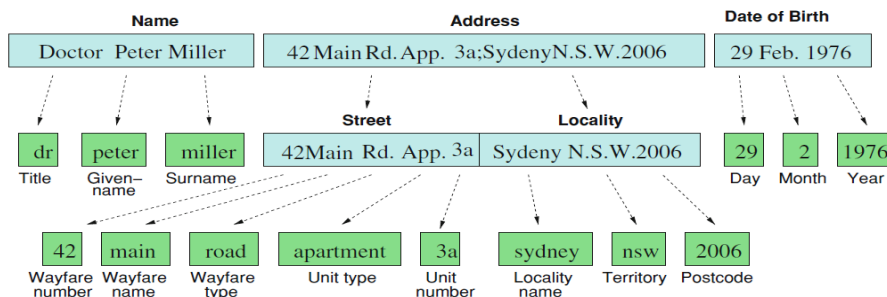\fr?d?ric joliot-curie\" national research institute for radiobiology and radiohygiene"

2.-0 LCA CONSULTANTS APS ⟶ 2.-0 lca consultants aps

anti-discriminatory
→ antidiscriminatory
→ anti discriminatory

3d Scanners Ltd
→ 3d scanners
→ 3d scanners limited
→ 3d Ssanners limited

U.S.A.
→ usa
→ united states of america

résumé vs. resume | tuebingen vs. tübingen | pena vs. peña.

## Standardisation , Tokenisation & Segmentation into Output Fields

| Name | Address | Date of Birth |
|---|---|---|
| Doctor  Peter Miller | 42 Main Rd. App.  3a;SydenyN.S.W.2006 | 29 Feb. 1976 |

| | | | Street | Locality | | | |
|---|---|---|---|---|---|---|---|
| dr | peter | miller | 42Main Rd. App. 3a | Sydeny N.S.W.2006 | 29 | 2 | 1976 |
| Title | Given–name | Surname | | | Day | Month | Year |

| 42 | main | road | apartment | 3a | sydney | nsw | 2006 |
|---|---|---|---|---|---|---|---|
| Wayfare number | Wayfare name | Wayfare type | Unit type | Unit number | Locality name | Territory | Postcode |

2/3/91

OR

3/2/91

2/3/91   2-3-91

February 3rd, 1991

3/2/91   3-2-91

Mars 2sd, 1991

## Verification

# Schema Matching and Mapping

**DVD-VENDOR**

> **Movies**(id, title, year)
> **Products**(mid, releaseDate, releaseCompany, basePrice, rating, saleLocID)
> **Locatios**(lid, name, classification, price)

**AGGREGATOR]**

> **Items**(name, releaseInfo, classification, price

**Schema Matching** identifies **correspondence** between elements of two schemas

> **Movies.**title ≈ **Items.**name
>
> **Movies.**year ≈ **Items.**year
>
> **Products.**rating ≈ **Items.**classification
>
> **Items.**price ≈ **Products.**basePrice * (1 + **Location.**taxRate)

**Schema Mapping** describes how to convert a **source schema** into a **target schema**

# Record Deduplication - Data Inconsistency Data Integrity

**Data Deduplication** is a technique for storing only one copy of repeating data (identical data)

| ID | Title | First | Last | AddressLine | City | Postcode | Telephne | DOB |
|----|-------|---------|--------|----------------|------|----------|-----------|----------|
| 1 | Miss | Catrina | Trewin | 123 Sample Road | Town | ABC 123 | 123456789 | 06/15/75 |
| 2 | Miss | Catrina | Trewin | 123 Sample Road | Town | ABC 123 | **123456780** | 06/15/75 |
| 3 | | Catrina | Trewin | 123 Sample Road | Town | ABC 123 | 123456789 | **06/15/76** |
| 4 | Miss | **C** | Trewin | 123 Sample Road | | ABC 123 | | 06/15/75 |

## Advantage

- Improve **storage** utilization
- Reduces the amount of bytes (data) that must be **transferred**
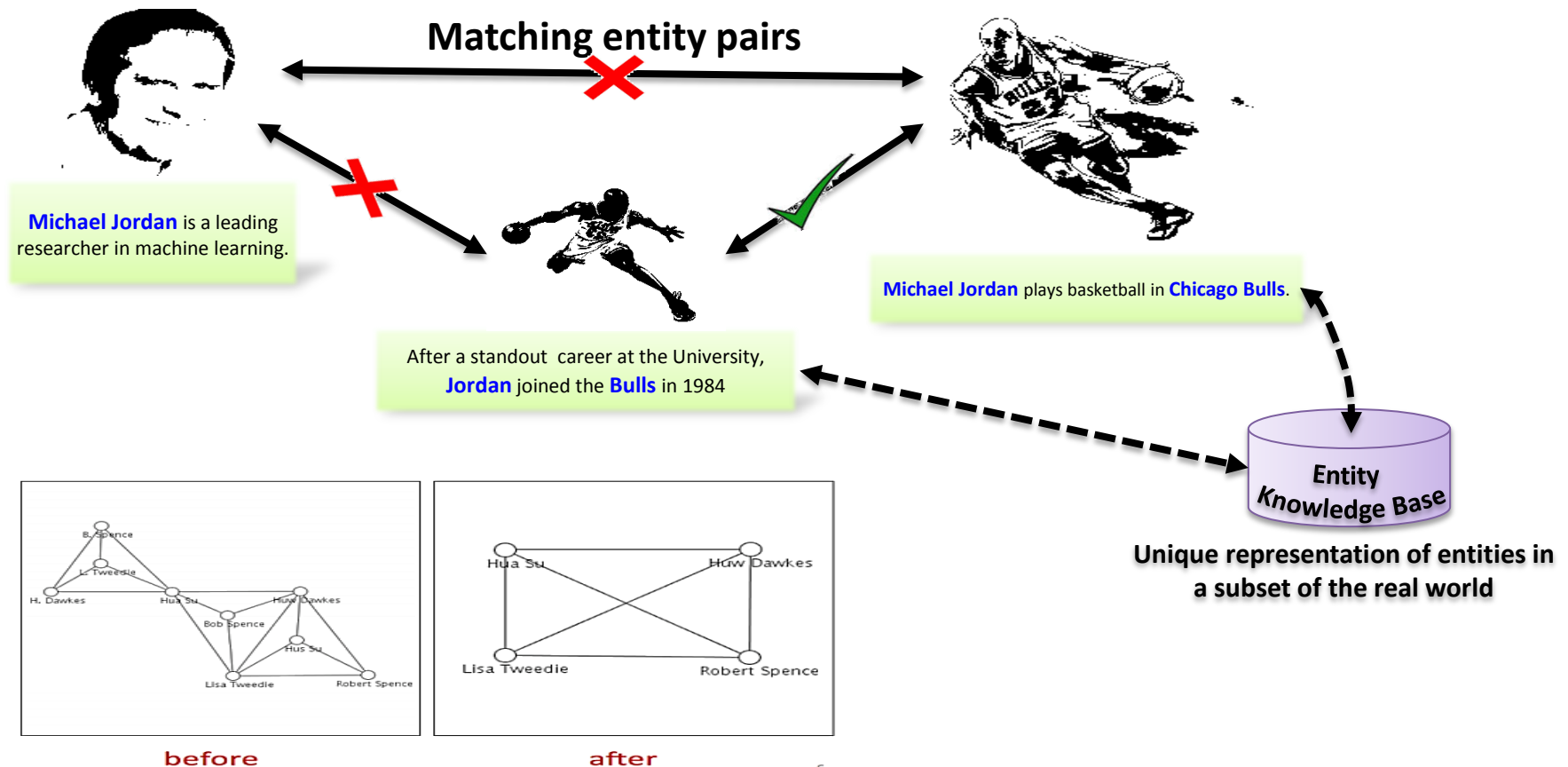- Provides accuracy for **statistical analyses**

## Drawbacks

- **Data Integrity** (potential loss of data)
- Computational **resource intensive**

# Entity Disambiguation

**Entity Disambiguation** is the process of removing **uncertainty** or **confusion** introduced by a named entity suggesting multiple interpretations. In other words, Entity Disambiguation is the practice of **identifying and linking a confusing named entity** (textual form) **to its true and unique representation** (real world) **within a knowledge base** which, in turn, provides a single semantic interpretation.



**Matching entity pairs**

**Michael Jordan** is a leading researcher in machine learning.

After a standout career at the University, **Jordan** joined the **Bulls** in 1984

**Michael Jordan** plays basketball in **Chicago Bulls**.

**Entity Knowledge Base**

**Unique representation of entities in a subset of the real world**

before

after

# Data Integration Tasks vs. Tools Overview

| | Pre-processing | Consistency Check | S. Matching | D. Integrity | Deduplication | Disambiguation | Reconciliation | Prov. |
|---|---|---|---|---|---|---|---|---|
| **Duke** | | | | | ■ | | ■ | |
| **Karma** | | | ■ | | | | | ■ |
| **Open Refine** | ■ | ■ | | ■ | ■ | ■ | ■ | |
| **D-Dupe** | | | | | ■ | | | |
| **FOXPSL** | | | | | ■ | ■ | ■ | |

# Vision Automated Data Integration Platform for Researchers

## Plan

- ➢ Create an Entity Knowledge-base
- ➢ Guidelines for Technical Harmonization
- ➢ Tools Implementation and Evaluation
- ➢ What to reuse? what to improve?
- ➢ RISIS Integration Tool
- ➢ Populate the Knowledge-base

# Reference

- Anhai Doan, Alon Halevy and Zachry Ives. Principles of data integration. Morgan Kaufmann Publishers. 2012 Elsevier. ISBN: 978-12-416044-6
- Peter Christen. Data Matching - Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection. Springer-Verlag Berlin Heidelberg 2012. ISBN 978-3-642-31164-2
- Philip A. Bernstein, Jayant Madhavan and Erhard Rahm. Generic Schema Matching, Ten Years Later. PVLDB 4(11): 695-701 (2011)