

Efficient Database Model to Represent Numerical Research Data of Material Flow Analysis.

Master Thesis - Computer Science

Albert-Ludwigs-Universität Freiburg

06 June 2018

Mohammad Mahadi Hasan



**UNI
FREIBURG**



- **Introduction**
 - Problem Statement
 - Approach

- **Background**
 - Relational Database
 - The Semantic Web
 - Triple stores

- **Construction of Relational Database**
 - MySQL architecture and design process
 - Implementation of web interface

- **Construction of RDF Triple Store**
 - Mapping relational database to RDF triples
 - Apache Jena TDB architecture

- **Database Performance Analysis**

- **Discussion**

■ Problem Statement

Usually, the research data of Industrial Ecology (IE) is stored in .mat or .csv files on a local machine which imposes challenges like:

- Lack of an efficient data storing mechanism.
- Lack of an online platform to query and retrieve data publicly.
- Data sharing and re-using terms has not been discussed for such data.

■ Approach

- Study on Relational Database (RDB)
 - MySQL database.
 - Implementation of web interface.
- Study on Semantic Web and its components
 - RDF - Data Model.
 - Apache Jena TDB triple store.

■ Relational Database (RDB)

- A collection of data sets organized and stored in relational tables.
- Each of these tables has one primary key column and shares at least one column (referred as foreign key) with another table to establish relationship.
- Data can be extracted by querying with query language like SQL.
- RDBMS is a database management system that is used to create and maintain relational databases.
- Popular RDBMS examples include Microsoft Access, SQL Server, Oracle Database, MySQL, PostgreSQL etc.

MySQL Database Architecture

- **The application layer:** client-side layer which handles the connection string, authentication and most importantly security.
- **The server layer:** the brain of the overall architecture. Any kind of query statement is executed in this layer.
- **The storage engine layer:** offers different storage engines.
- **Indexes:**
 - Indexes are used to find a data entry quickly.
 - Let us consider following example:

```
SELECT * FROM table WHERE id=1;
```

- Without the indexes, the query will go through every row and column.
- The optimization here is to add an index, for example, primary key. It runs the query only against the table indexes rather than all the column data.

The Semantic Web

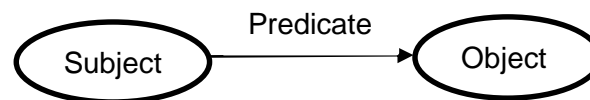
- "The Semantic Web provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries".
- Data is structured and published following semantic web standards: RDF.

RDF- The Data Model

Resource Description Framework is a "framework for representing information in the Web".

- The basic elements of RDF are triples.
- A triple is a set of three entities that form a statement in the form of *subject-predicate-object* expressions.

Harry Potter → *has author* → *J. K. Rowling*



URIs, Vocabularies & Ontologies

- Uniform Resource Identifiers (URIs) are short strings that identify resources in the web.

| | |
|--|--------|
| <code><http://dbpedia.org/resource/Berlin></code> | .. (1) |
| <code><http://www.w3.org/1999/02/22-rdf-syntax-ns#type></code> | .. (2) |
| <code><http://dbpedia.org/class/yago/CapitalsInEurope></code> | .. (3) |

subject → ***predicate*** → ***object***
(1) → **(2)** → **(3)**

This triple refers that resource (1) has a relationship (2) with the resource (3).

- The “vocabularies” on the Semantic Web is used for data integration between data sets.
- The term “ontology” is for more complex, and more formal collection of terms, where “vocabulary” is used for more basic use.

SPARQL – Query Language

- SPARQL is the standardized query language for RDF, just like SQL is the standardized query language for relational databases.
- Like SQL, SPARQL also follow the same “SELECT...FROM...WHERE...” query structure.

```
:id1 foaf:name "André Schürle"  
:id1 foaf:based_near : Dortmund  
:id2 foaf:name "Nils Petersen"  
:id2 foaf:based_near : Freiburg
```

```
SELECT ?name  
WHERE {  
  ?x foaf:name ?name .  
  ?x foaf:based_near : Freiburg .  
}
```

The predicate has a constant value of foaf:based_near and the object has a constant value of : *Freiburg* match to one of RDF triples. The result is *Nils Petersen*.

Triple Stores

- Triple Store (RDF triple store) is a specialized DBMS for RDF triples.
- RDF data can be stored in two ways:
 - in *files*, triples are stored following one of the serialization formats, or
 - in special kind of databases for triples, called *triple stores*.
- *Triple stores* have three possible architectures:
 - In-memory: stores the triples in main memory. It is fast but expensive.
 - Native Store: storage systems with own database. For example- Jena TDB, Sesame Native, Virtuoso, AllegroGraph, Oracle 11g etc.
 - Non-native Store: storage using a third-party RDBMS. For example- Jena SDB backed by MySQL database.



■ MySQL Design Process:

- Six different CSV files containing the data of country, unit specification, dataset (list of datasets), process list, stock data and flow data between processes.
- Database Schema, Tables, and Constraints
 - Database Schema => database name.
 - 6 tables.
 - 6 primary keys and 11 foreign keys; are said to be database constraints.
 - 589165 data rows.
 - Roughly 55mb of SQL dump.

Construction of Relational Database



| stocks |
|-----------------------------------|
| <code>seq_id</code> |
| <code>system_id</code> |
| <code>stock_dataset_number</code> |
| <code>process_id</code> |
| <code>ISO_code</code> |
| ... |
| ... |
| <code>Unit_id</code> |
| <code>comment</code> |

| countries |
|--------------------------------|
| <code>seq_id</code> |
| <code>ISO_code</code> |
| <code>name</code> |
| ... |
| ... |
| <code>alternative_name7</code> |

| flows |
|----------------------------------|
| <code>seq_id</code> |
| <code>system_id</code> |
| <code>flow_dataset_number</code> |
| <code>process_id_source</code> |
| <code>process_id_target</code> |
| <code>region_source</code> |
| <code>region_target</code> |
| ... |
| ... |
| <code>Unit_id</code> |
| <code>comment</code> |

| process_list |
|---------------------------|
| <code>system_id</code> |
| <code>process_id</code> |
| <code>process_name</code> |
| ... |
| ... |
| <code>YPos</code> |

| dataset |
|--------------------------------|
| <code>system_id</code> |
| <code>system_definition</code> |
| <code>dataset_name</code> |
| ... |
| ... |
| <code>comment</code> |

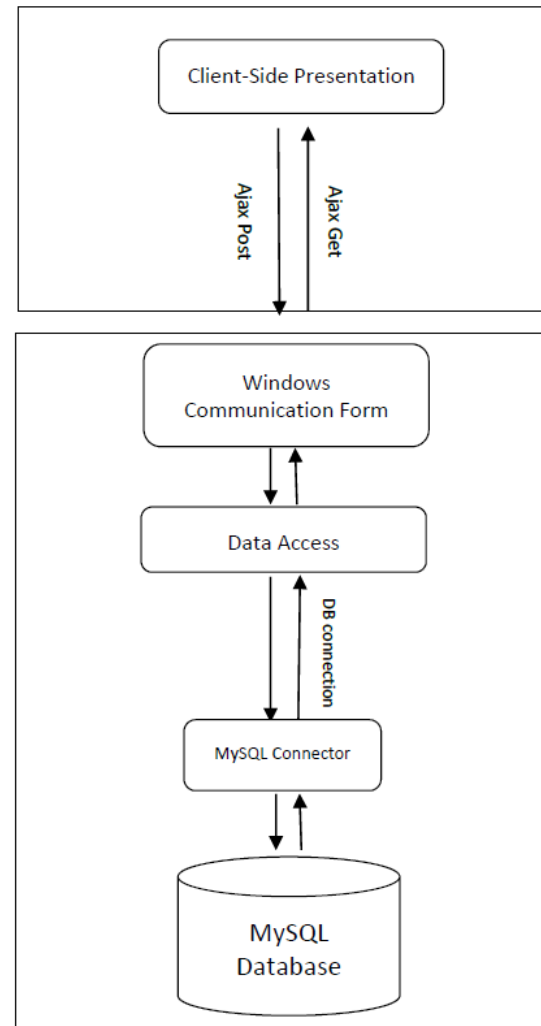
| Unit_classification |
|-------------------------|
| <code>unit_id</code> |
| <code>si_unit_id</code> |
| <code>unit_code</code> |
| ... |
| ... |
| <code>factor</code> |

Green columns = primary keys.

Purple columns = foreign keys.

Implementation of web interface

- Back-End Development
 - Mechanisms to communicate with server-side database.
- Front-End Development
 - Mechanisms to present data in client-side.



Client-Side Architecture

Server-Side Architecture

Figure: web interface architecture

Construction of Relational Database



Data Table with the catalogue of datasets:

- **Query 1:** *Select all the dataset available in the database from dataset table;*

Catalogue of Data Sets

| Select DataSet | | | | | | | | | | |
|--------------------------|-----------|-------------------------|---------------------|---------------------|--------------------|----------------------|------------------------|--------|---------------------|-------------------|
| | System ID | Dataset Name | Reference Date | Most Recent Update | Corresponding Auth | Other Author Info | Document Reference | Region | Time Period Of Anal | Indicator Element |
| <input type="checkbox"/> | 1 | Global multiregional st | 05.08.2012 00:00:00 | 05.08.2012 00:00:00 | Pauliuk, Stefan | Tao Wang, Daniel B M | 10.1016/j.resconrec.20 | World | 1700-2008 | Iron |

Export Selected Rows To CSV

Page 1 of 1

View 1 - 1 of 1

Figure: DataGrid with list of datasets available in the database

Data Table with material stocks:

- **Query 2:** Select all the stocks within the dataset (identified by the id) and filter the data with year (2007, 2008) and country (Austria, Germany, Switzerland and Italy) from stock table;

Processes with stocks

Choose Stock Process: Slag piles

Filter by Year and Region

2008 × 2007 ×
Tags: 2008, 2007

Austria × Germany × Switzerland × Italy ×
Tags: Austria, Germany, Switzerland, Italy

| stock dataset num | process name | country name | year | age cohort | indicator element | aspect of dataset | value | error type | error value 1 | error value 2 | data quality | unit id | comment |
|--------------------------|--------------|--------------|-------------|------------|-------------------|-------------------|---------------|------------|---------------|---------------|--------------|---------|---------|
| <input type="checkbox"/> | 762 | Slag piles | Austria | 2007 | slag | Stock | 14364,0962395 | | | | | kt | |
| <input type="checkbox"/> | 763 | Slag piles | Austria | 2008 | slag | Stock | 14742,1002613 | | | | | kt | |
| <input type="checkbox"/> | 5340 | Slag piles | Germany | 2007 | slag | Stock | 157162,243482 | | | | | kt | |
| <input type="checkbox"/> | 5341 | Slag piles | Germany | 2008 | slag | Stock | 159458,228723 | | | | | kt | |
| <input type="checkbox"/> | 7302 | Slag piles | Italy | 2007 | slag | Stock | 59069,7802223 | | | | | kt | |
| <input type="checkbox"/> | 7303 | Slag piles | Italy | 2008 | slag | Stock | 60640,8354917 | | | | | kt | |
| <input type="checkbox"/> | 13515 | Slag piles | Switzerland | 2007 | slag | Stock | 2070,50257304 | | | | | kt | |
| <input type="checkbox"/> | 13516 | Slag piles | Switzerland | 2008 | slag | Stock | 2136,16484003 | | | | | kt | |

Export Selected Rows To CSV Page 1 of 1 20 View 1 - 8 of 8

Figure: Data Grid with stock data between processes

Data Table with material flows:

- Query 3:** Select all the flows within the dataset (identified by the id) and filter the data with year (2005) and source process 'scrap market' and target process 'foundries' from flow table;

Flows between processes

Choose Source and Target Process

Scrap market Foundries

Filter by Year and Region

2005 x

Tags: 2005

Tags:

Select Flow Processes

| <input type="checkbox"/> | flow dataset num | source process n | target process n | source country n | target country n | year | age cohort | indicator element | value | error type | error value 1 | error value 2 | data quality | unit id | comment |
|--------------------------|------------------|------------------|------------------|------------------|------------------|------|------------|-------------------|----------------|------------|---------------|---------------|--------------|---------|---------|
| <input type="checkbox"/> | 106 | Scrap market | Foundries | Afghanistan | Afghanistan | 2005 | | steel scrap | 0 | | | | | kt/yr | |
| <input type="checkbox"/> | 215 | Scrap market | Foundries | Albania | Albania | 2005 | | steel scrap | 0 | | | | | kt/yr | |
| <input type="checkbox"/> | 324 | Scrap market | Foundries | Algeria | Algeria | 2005 | | steel scrap | 69,2790266983 | | | | | kt/yr | |
| <input type="checkbox"/> | 433 | Scrap market | Foundries | Angola | Angola | 2005 | | steel scrap | 10,25217303 | | | | | kt/yr | |
| <input type="checkbox"/> | 542 | Scrap market | Foundries | Argentina | Argentina | 2005 | | steel scrap | 331,052613857 | | | | | kt/yr | |
| <input type="checkbox"/> | 651 | Scrap market | Foundries | Australia | Australia | 2005 | | steel scrap | 473,021242061 | | | | | kt/yr | |
| <input type="checkbox"/> | 760 | Scrap market | Foundries | Austria | Austria | 2005 | | steel scrap | 411,972432819 | | | | | kt/yr | |
| <input type="checkbox"/> | 869 | Scrap market | Foundries | Bahrain | Bahrain | 2005 | | steel scrap | 0 | | | | | kt/yr | |
| <input type="checkbox"/> | 978 | Scrap market | Foundries | Bangladesh | Bangladesh | 2005 | | steel scrap | 530,380352549 | | | | | kt/yr | |
| <input type="checkbox"/> | 1087 | Scrap market | Foundries | Barbados | Barbados | 2005 | | steel scrap | 4,61843659142 | | | | | kt/yr | |
| <input type="checkbox"/> | 1196 | Scrap market | Foundries | Belgium-Luxembou | Belgium-Luxembou | 2005 | | steel scrap | 532,4552064 | | | | | kt/yr | |
| <input type="checkbox"/> | 1305 | Scrap market | Foundries | Bermuda | Bermuda | 2005 | | steel scrap | 0,952782763078 | | | | | kt/yr | |
| <input type="checkbox"/> | 1414 | Scrap market | Foundries | Bolivia | Bolivia | 2005 | | steel scrap | 18,4313739029 | | | | | kt/yr | |
| <input type="checkbox"/> | 1523 | Scrap market | Foundries | Brazil | Brazil | 2005 | | steel scrap | 1992,28516852 | | | | | kt/yr | |

Page 1 of 8

View 1 - 20 of 146

Figure: Data Grid with flow data between processes

Construction of Triple Store



- The majority of data on the current Web is stored in relational databases.
- Semantic web is useful, especially if data from different sources has to be shared or integrated.
- Therefore, it is important to introduce mapping technologies between relational database and RDF.

- **D2RQ**: A mapping mechanism to use RDB as RDF graphs (in addition, it exposes RDB to RDF triples)

D2RQ architecture

- A D2RQ Engine accesses a Non-RDF database.
- The mapping language describes the relationship between an ontology (or vocabulary) and a relational data model.
- RDF APIs can be embedded with Java applications via Jena or Sesame APIs.
- The RDF dump file can be stored in any triple store.
- D2RQ server provides an HTML view to explore the mapped database.

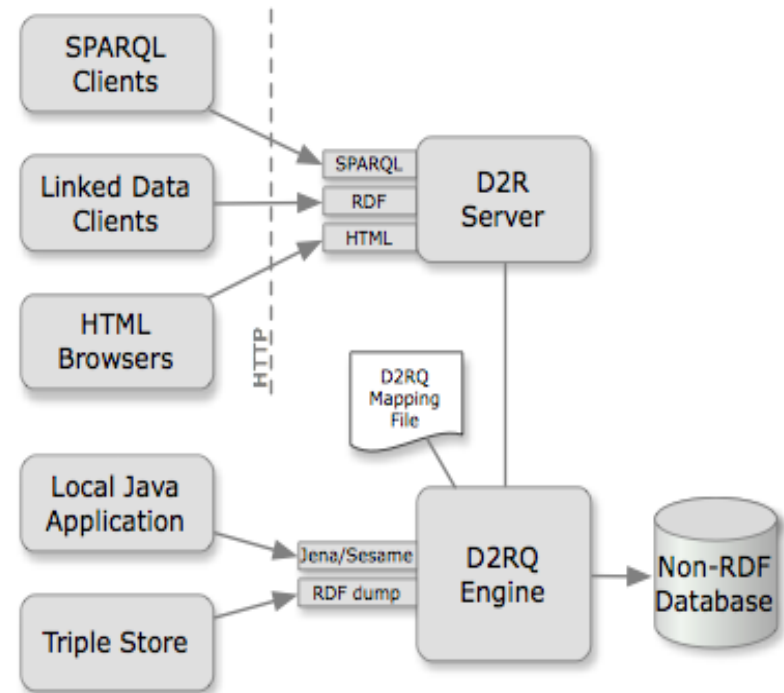


Figure: Architecture of D2RQ

D2RQ: Overview and Features

- The D2RQ platform serves different tools to offer RDF-based access to the content of relational databases.
 - The ***generate-mapping*** tool creates a D2RQ mapping file by analysing the relational database schema.

```
>generate-mapping -u username -p password -o mapping.ttl jdbc:mysql:///ief
```

- The ***dump-rdf*** tool is used to dump the contents of the entire database into a single RDF file.

```
>dump-rdf -f TURTLE -b http://localhost:2020/mapping.ttl > iefdatabase.ttl
```

- The ***D2R-server*** takes the mapping file as input and provides a web interface where RDF data can be browsed.

```
>d2r-server mapping.ttl
```

Construction of Triple Store



Database table to RDF triples with D2RQ:

<"primary key column value" "tableName_columnName" "columnValue">

| id | country | aspect_of_data_set | indicator_element | process | system_id | unit_id | value | year |
|-------|---------|--------------------|-------------------|--------------|-----------|----------|----------|------|
| 17011 | Germany | Stock | slag | Scrap market | 1 | kilo ton | 444.0907 | 1906 |



```
<"stocks17011" "stocks_country" "Germany">  
<"stocks17011" "stocks_aspect_of_dataset" "Stock">  
<"stocks17011" "stocks_indicator_element" "slag">  
<"stocks17011" "stocks_process" "Scrap market">  
<"stocks17011" "stocks_system_id" "1">  
<"stocks17011" "stocks_unit_id" "kilo ton">  
<"stocks17011" "stocks_value" "444.0907">  
<"stocks17011" "stocks_year" "1906">
```

■ Apache Jena TDB

- TDB is a component of Jena that used as native RDF storage.
- TDB is stored in a single directory in the filing system backed by a dataset.
- A complete TDB dataset consists of:
 - **The node table:** stores the representation of RDF terms. It consists of two mappings: Node to Nodeld and Nodeld to Node.
 - **Triple and Quad indexes:** used for the default graph. Triples are held as 3-tuples of Nodelds in triple indexes where quads are held as 4-tuples of Nodelds.
 - **The prefixes table:** stores index for Graph->Prefix->URI mapping. It provides the mechanism for Jena API to serialize the triples in RDF/XML or Turtle.

■ Experimental Settings:

- Two training datasets:
 - **IEF dataset:** consisting the material stock and flow data of Industrial Ecology Freiburg research group.
 - 6 tables, 589165 rows.
 - D2RQ generates roughly 10 million triples.
 - 7596007 final triples (ignoring the empty “ ” column values).
 - **ISWC dataset:** sample dataset with the information about conferences, papers, authors and topics from ISWC 2002 conference¹.
 - 9 tables, 96 rows.
 - 322 triples after mapping.

¹<http://iswc2002.semanticweb.org/>

Results Evaluation:

| SQL Queries on MySQL | | SPARQL Queries on Jena TDB | |
|----------------------|--|----------------------------|--|
| | Query Execution Time (ms) (average of 10 run) | | Query Execution Time (ms) (average of 10 run) |
| <i>Query 1</i> | 9 | <i>Query 5</i> | 145 |
| <i>Query 2</i> | 127 | <i>Query 6</i> | 1673 |
| <i>Query 3</i> | 80 | <i>Query 7</i> | 11448 |
| <i>Query 4</i> | 16 | <i>Query 8</i> | 10257 |

Table 1: Query execution time comparisons for ief dataset

Results Evaluation:

| SQL Queries on MySQL | | SPARQL Queries on Jena TDB | |
|----------------------|--|----------------------------|--|
| | Query Execution Time (ms) (average of 10 run) | | Query Execution Time (ms) (average of 10 run) |
| <i>Query 9</i> | 3 | <i>Query 13</i> | 132 |
| <i>Query 10</i> | 3 | <i>Query 14</i> | 7 |
| <i>Query 11</i> | 2 | <i>Query 15</i> | 5 |
| <i>Query 12</i> | 2 | <i>Query 16</i> | 12 |

Table 2: Query execution time comparisons for iswc dataset

Results Evaluation:

- MySQL works faster than TDB
 - MySQL storage engines are faster architecturally.
 - Use of indexing to identify records.
- TDB works slower because of enumeration
 - TDB query execution method does not finish until results are fully enumerated.
 - Enumeration process takes longer time depending on number of predicates (number of columns in case of MySQL) and triples.

```
ResultSet results = qexec.execSelect();  
long numResults = ResultSetFormatter.consume(results);
```

The variation of query execution time in TDB (when triples are available in main memory!).



- Most of the web contents are backed by RDB.
- RDB's are easy to implement; become complex with the number of tables increases.
- RDB works as “closed loop”.
- Data can not be integrated from different sources if the database schemas are not in same structure.

- Triple Stores allow data to be extended across different data sources.
- Data can be integrated with multiple data sources by adding properties.
- A small data source can be enlarged into a bigger and richer data source.

■ Future Work

- Query optimization for SQL joins.
- Allow filtering data with a range of years (1901-2000 instead of using 1901, 1902,, 2000).
- Introduce filtering based on a collection of regions (not only country wise, also continents wise like “Europe”).
- Improvement of default RDF vocabulary terms within D2RQ (use of more meaningful and publicly well-known vocabulary).
- Build a domain-specific complete ontology from the dataset and make it public.



Thank You!