

Extracting OWL Ontologies from Relational Databases Using Data Analysis and Machine Learning

Lama AL KHUZAYEM^{a,b,1} and Peter MCBRIEN^a

^a*Department of Computing, Imperial College, London SW7 2AZ, UK*

^b*Computer Science Department, Faculty of Computing and IT, King Abdulaziz University, Jeddah 21583, Saudi Arabia*

Abstract. Extracting OWL ontologies from relational databases is extremely helpful for realising the Semantic Web vision. However, most of the approaches in this context often drop many of the expressive features of OWL. This is because highly expressive axioms can not be detected from database schema alone, but instead require a combined analysis of the database schema and data. In this paper, we present an approach that transforms a relational schema to a basic OWL schema, and then enhances it with rich OWL 2 constructs using schema and data analysis techniques. We then rely on the user for the verification of these features. Furthermore, we apply machine learning algorithms to help in ranking the resulting features based on user supplied relevance scores. Testing our tool on a number of databases demonstrates that our proposed approach is feasible and effective.

Keywords. Ontology Learning, OWL 2 Ontologies, Data Analysis, Machine Learning.

1. Introduction

Constructing ontologies, which are the basic building blocks of the **Semantic Web (SW)**, from scratch is tedious, error-prone, time-consuming, and requires deep understanding of the domain [1,2,3]. Attention has therefore shifted towards the generation of ontologies from existing data. The work in this paper focuses on extracting ontologies from relational models. The issue that arises in this context is that only ontologies with limited expressivity can be obtained in an automated fashion. Although leading ontology definition languages such as the **Web Ontology Language (OWL)** [4] have great expressive capabilities, they are currently not used to their full potential. This is because richer and more expressive axioms like **PropertyChain**, **TransitiveProperty** and **SymmetricProperty** require more advanced detection techniques than just analysing the database schema. Most previous work (e.g., [1,2,5,6]) have avoided this extra task, leaving it for the users to manually detect and insert the desired axioms.

¹Corresponding Author: Lama Al Khuzayem, Department of Computing, Imperial College London, 180 Queen's Gate, London SW7 2AZ, United Kingdom; E-mail: l.al-khuzayem11@imperial.ac.uk.

The approach described in this paper, which is developed as a tool called OWLRel², takes this additional step and aims at finding more complex axioms through applying a combination of schema and data analysis techniques. This paper extends the work presented in [7], in which it was demonstrated how to automatically extract an OWL ontology from a relational schema by generating a **Bidirectional Transformation (BT)**, using the **Hypergraph Data Model (HDM)** [8], then enriching the OWL schema with more expressive axioms using schema and data analysis. Because of the subjective nature of the domain, we rely on users to validate these suggested axioms. Thus, our approach benefits from the two aspects that must be taken into consideration in this context [3]: human intervention and database content analysis.

For example, consider the Emp table shown in the database schema depicted in Figure 1. From the database schema, we can generate an ontology that contains an Emp class from the Emp table. The EID column will be transformed into a data property. The RTo? column will be transformed into an object property. The schema does not say anything about the characteristics of the object property. However, by doing an inspection on the data, one can detect that the RTo object property is an **AsymmetricProperty** (if someone reports to his/her manager, the manager does not report back to that person) and an **IrreflexiveProperty** (a person does not report back to himself). Our approach also applies schema analysis heuristics, which detects that chaining the property RTo with itself, or with its inverse, produces new properties that could possibly be useful or interesting. For instance, chaining RTo with RTo produces a property that relates employees with their second line managers, while chaining RTo with RTo⁻ produces a property that contains the employees that are colleagues of each other. These proposals would not be possible without the use of schema and data analysis.

One problem that arises in this situation is that the large number of suggested axioms (where the schema is large) might overwhelm the user. Hence, in this paper, we focus on how we can help the users by filtering or sorting the suggestions according to their relevance. Since relevance is very subjective, **Machine Learning (ML)** algorithms were used to rank the suggested axioms taking into consideration knowledge gathered from previous users over multiple sessions. Our tool is also able to continuously learn from its new users, which helps in improving the results over time.

The remainder of this paper is structured as follows. Section 2 reviews the HDM and how the relational model is presented in HDM. Section 3 describes our automatic ontology extraction approach and in Section 4, we demonstrate our ontology enrichment process. Section 5 explains how we apply machine learning to our extraction approach followed by related work in Section 6. Finally, Section 7 concludes this paper.

2. Preliminaries

An overview of the HDM [8,9,10], and how the relational model is represented in the HDM is given here. An HDM schema S is defined as a tuple $\langle Nodes, Edges, Cons, Types \rangle$ where:

- *Nodes* is a set of nodes in the graph such that each node $node: \langle \langle n, t \rangle \rangle$ is identified by its name n , and it is given an associated type $t \in Types$. A node can also be referred to by a shorthand $node: \langle \langle n \rangle \rangle$.

²automated.doc.ic.ac.uk/releases/jars/OWLRel-rel-0-1.jar

Table 1. Rules for Representing a Relational Schema as an HDM Schema ($S_r \rightarrow S_{hr}$)

Relational Construct	HDM Representation
table $\langle\langle T \rangle\rangle$	node: $\langle\langle T, \text{any} \rangle\rangle$
column $\langle\langle T, C, N, U, t \rangle\rangle$	node: $\langle\langle T:C, t \rangle\rangle$, edge: $\langle\langle \leftarrow, T, T:C \rangle\rangle$, cons: $\langle\langle \triangleright, \langle\langle T:C \rangle\rangle, \langle\langle \leftarrow, T, T:C \rangle\rangle \rangle\rangle$, cons: $\langle\langle \triangleleft, \langle\langle T \rangle\rangle, \langle\langle \leftarrow, T, T:C \rangle\rangle \rangle\rangle$
$N = \text{notnull}$	cons: $\langle\langle \triangleright, \langle\langle T \rangle\rangle, \langle\langle \leftarrow, T, T:C \rangle\rangle \rangle\rangle$
$U = \text{unique}$	cons: $\langle\langle \triangleleft, \langle\langle T:C \rangle\rangle, \langle\langle \leftarrow, T, T:C \rangle\rangle \rangle\rangle$
primary_key $\langle\langle T, C \rangle\rangle$	cons: $\langle\langle \xrightarrow{\text{id}}, \langle\langle T \rangle\rangle, \langle\langle \leftarrow, T, T:C \rangle\rangle \rangle\rangle$
primary_key $\langle\langle T, C_1, C_2 \rangle\rangle$	cons: $\langle\langle \xrightarrow{\text{id}}, \langle\langle T \rangle\rangle, \langle\langle \leftarrow, T, T:C_1 \rangle\rangle \bowtie \langle\langle \leftarrow, T, T:C_2 \rangle\rangle \rangle\rangle$
foreign_key $\langle\langle FK, T, C, T_f, C_f \rangle\rangle$	cons: $\langle\langle \subseteq, \pi_{\langle\langle T:C \rangle\rangle} \langle\langle \leftarrow, T, T:C \rangle\rangle, \pi_{\langle\langle T_f:C_f \rangle\rangle} \langle\langle \leftarrow, T_f, T_f:C_f \rangle\rangle \rangle\rangle$

- *Types* is a tuple that contains a finite set of types and a subset of the set of all possible data values consistent with this type.
- *Edges* is a set of edges in the graph such that each edge has the following scheme: edge: $\langle\langle e_1, n_1, n_2 \rangle\rangle$, where e_1 is the edge's name (can also be unnamed as “-”) and n_1 and n_2 are the two nodes that it connects.
- *Schemes* is the union of *Nodes* and *Edges*.
- *Cons* is a set of boolean-valued functions (constraints) where they form the HDM constraint language. The set of constraints used in this paper are:

- * cons: $\langle\langle \subseteq, s_1, s_2 \rangle\rangle$ is the **inclusion** constraint which states that scheme s_1 is always a subset of scheme s_2 .
- * cons: $\langle\langle \not\cap, s_1, \dots, s_n \rangle\rangle$ is the **exclusion** constraint which states that all the associate schemes are disjoint from each other.
- * cons: $\langle\langle \cup, s, s_1, \dots, s_n \rangle\rangle$ is the **union** constraint stating scheme s as the union of schemes s_1, \dots, s_n .
- * cons: $\langle\langle \triangleright, s_1, \dots, s_m, s \rangle\rangle$ is the **mandatory** constraint stating that every combination of the values that appears in schemes s_1, \dots, s_m must appear in the edge s connecting those schemes.
- * cons: $\langle\langle \triangleleft, s_1, \dots, s_m, s \rangle\rangle$ is the **unique** constraint stating that every combination of the values that appears in schemes s_1, \dots, s_m must appear no more than once in the edge s connecting those schemes.
- * cons: $\langle\langle \xrightarrow{\text{id}}, s_1, s \rangle\rangle$ is the **reflexive** constraint, stating that any value in s_1 must appear reflexively in the edge s that connects it to s_1 .

In addition to referring to schemes directly, constraints may also take joins, projections and selections of schemes as arguments.

2.1. Representing the Relational Model in the HDM

A method, summarised in Table 1, for mapping relational schemas to HDM was defined in [8] and [9]. To illustrate the method, consider the relational schema S_r , depicted in Figure 1 (where primary keys are underlined, and nullable column names are suffixed by a question mark), which represents a subset of the Northwind database³.

Each table is represented as an HDM node (illustrated in Figure 2 by a black outlined circle) with HDM type **any**. For example, the Emp table is represented by the HDM node node: $\langle\langle \text{Emp}, \text{any} \rangle\rangle$.

³<https://northwinddatabase.codeplex.com/>

Each column is represented by a node that has an HDM type based on its relational type. For example, column Emp.EID has the type INTEGER, so it is represented in the HDM as node:⟨⟨Emp:EID, int⟩⟩. Moreover, Emp.EID will be connected via an HDM edge, edge:⟨⟨-, Emp, Emp:EID⟩⟩ (illustrated in Figure 2 with a thick black line), to the node that represents the column's table. Since column values only appear with an instance of a table tuple, the edge has a mandatory constraint (illustrated with grey lines) from the column node as cons:⟨⟨▷, node:⟨⟨Emp:EID⟩⟩, edge:⟨⟨-, Emp, Emp:EID⟩⟩⟩. Furthermore, the unique constraint cons:⟨⟨◁, node:⟨⟨Emp⟩⟩, edge:⟨⟨-, Emp, Emp:EID⟩⟩⟩ ensures that each column has a single value per row.

If a column is not nullable, then it must also have a mandatory constraint. Thus, Emp.EID has cons:⟨⟨▷, node:⟨⟨Emp⟩⟩, edge:⟨⟨-, Emp, Emp:EID⟩⟩⟩. If the column is key, then we state that the column's edge is reflexive. This also applies to column Emp.EID: cons:⟨⟨^{id}→, node:⟨⟨Emp⟩⟩, edge:⟨⟨-, Emp, Emp:EID⟩⟩⟩.

Finally, foreign keys are represented as inclusion constraints. Thus, for the foreign key between Emp.RTo and Emp.EID, we create the following HDM constraint: cons:⟨⟨⊆, node:⟨⟨Emp:RTo⟩⟩, node:⟨⟨Emp:EID⟩⟩⟩.

The result of these transformations is an HDM graph (depicted in Figure 2) that is a forest of two-level trees, with subset constraints linking the leaf nodes.

Ter			Emp		ET		Reg	
TID	TDes	RID	EID	RTo?	EID	TID	RID	RDes
'20852'	'Rockville'	1	1	2	1	'06897'	1	'E'
'30346'	'Atlanta'	4	2	NULL	2	'01730'	2	'W'
'01833'	'Georgetow'	1	3	2	2	'01833'	3	'N'
'01730'	'Bedford'	1	4	2	3	'30346'	4	'S'
'06897'	'Wilton'	1	5	2	4	'20852'		
'02903'	'Providence'	1	6	5	5	'02903'		
'85014'	'Phoenix'	2	7	5	6	'85014'		

Emp(RTo) $\overset{fk}{\Rightarrow}$ Emp(EID) ET(TID) $\overset{fk}{\Rightarrow}$ Ter(TID)
 Ter(RID) $\overset{fk}{\Rightarrow}$ Reg(RID) ET(EID) $\overset{fk}{\Rightarrow}$ Emp(EID)

Figure 1. S_r , Fragment of the Northwind Relational Database Schema and Data. Entries in the Employee table are related to entries in the Territories table via ET. Each employee may optionally (indicated by a question mark) be recorded as reporting to another employee by the RTo. Each territory is in exactly one Region.

3. The Automatic Relational to OWL Transformation

Our automatic ontology extraction approach is performed via a number of steps: i) transform the relational constructs to HDM producing S_{hr} , which was described in Section 2.1, ii) perform intermodel transformations on the HDM schema producing S_{ho} and finally, iii) translate the S_{ho} to an OWL schema. The subsections below, review the final two steps which were explained in [7].

3.1. The HDM Intermodel Transformations

Our HDM intermodel transformation process aims at overcoming the fundamental differences between the relational and OWL modelling languages. We use a set of equivalence

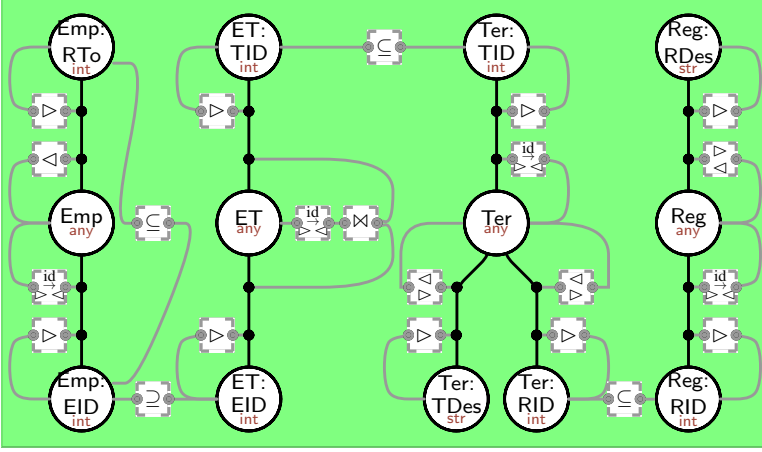


Figure 2. S_{hr} , HDM Representation of Schema S_r

both-as-view (BAV) mappings [11], presented in [9], to transform the HDM schema, S_{hr} (depicted in Figure 2) which represents the relational schema, to an equivalent (in terms of information capacity) HDM graph which represents the OWL schema.

(A) **Transform an attribute that is a foreign key, but not a primary key, into an ObjectProperty.** This can be achieved using two BAV equivalence rules: Inclusion Merge and Unique-Mandatory Redirection. For instance, using the rules shown below, we can transform the foreign key represented by the inclusion constraint between node: $\langle\langle ET:EID \rangle\rangle$ and node: $\langle\langle Emp:EID \rangle\rangle$, to become an object property represented by the edge $edge:\langle\langle -, ET, Emp \rangle\rangle$. All other foreign keys of this type are transformed similarly.

- ① $inclusion_merge(node:\langle\langle Emp:EID \rangle\rangle, edge:\langle\langle ET:EID, ET, ET:EID \rangle\rangle)$
- ② $unique_mandatory_redirection(edge:\langle\langle ET_Emp:EID, ET, Emp:EID \rangle\rangle, edge:\langle\langle ET_Emp, ET, Emp \rangle\rangle)$

(B) **Transform an attribute that is a foreign key and a primary key as a SubClassOf axiom.** This can be achieved by using the BAV equivalence rule, Identity Node Merge. The Northwind database, however, does not include an example of this particular case.

(C) **Transform a many-to-many binary relationship table to an ObjectProperty.** This can be achieved using the Identity Edge Merge rule. For instance, the rule shown below allows us to replace node: $\langle\langle ET \rangle\rangle$, edge: $\langle\langle -, ET, Emp \rangle\rangle$ and edge: $\langle\langle -, ET, Ter \rangle\rangle$ with a single edge: $\langle\langle Emp_Ter, Emp, Ter \rangle\rangle$. This newly created edge represents an object property since it connects two class nodes.

- ③ $identity_edge_merge(edge:\langle\langle ET_Emp, ET, Emp \rangle\rangle, edge:\langle\langle ET_Ter, ET, Ter \rangle\rangle)$

Figure 3, illustrates the result of these transformations. Remaining nodes and edges do not need to be transformed as they will be interpreted as classes, datatypes and data properties in the OWL model, as will be described in the next section.

3.2. Transforming the HDM Schema to an OWL Schema

We now review how HDM schemas in ‘OWL Compatible’ form, such as S_{ho} and S_{ho+} , can be translated to corresponding OWL files S_o and S_{o+} , respectively. In Table 2, we list

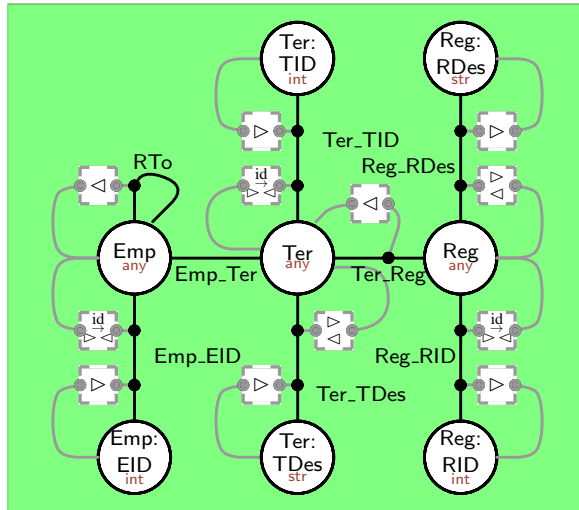


Figure 3. S_{ho} , HDM Representation of Schema S_o

some OWL 2 constructs and how we translate them into HDM. The four more important transformations are:

1. HDM nodes without a type restriction are transformed to OWL classes. Thus, the nodes $\langle\langle\text{Emp}\rangle\rangle$, $\langle\langle\text{Ter}\rangle\rangle$ and $\langle\langle\text{Reg}\rangle\rangle$ with HDM type *any* in Figure 3 are mapped to the OWL classes: *Emp*, *Ter* and *Reg*.
2. HDM nodes with HDM types are transformed to datatypes. For example, the node $\langle\langle\text{Emp_EID}\rangle\rangle$ with HDM type *int* is mapped to *xsd:integer*.
3. Edges which connect two HDM nodes, both without type restrictions, are transformed to object properties. For example, edge: $\langle\langle\text{Emp_Ter}, \text{Emp}, \text{Ter}\rangle\rangle$ is mapped to the *ObjectProperty Emp_Ter*.
4. Edges which connect nodes without type restriction, to nodes with a type restriction, are transformed to data properties. For example, the *DataProperty Emp_EID*, represents the edge: $\langle\langle\text{Emp_EID}, \text{Emp}, \text{Emp:EID}\rangle\rangle$.

Subsequently, combinations of HDM constraints are considered to represent more complex OWL 2 constructs, according to the mappings listed in Table 2. For instance, the combination of the following constraints $\text{cons}:\langle\langle\triangleleft, \text{EID}, \langle\langle\text{Emp_EID}, \text{Emp}, \text{EID}\rangle\rangle\rangle\rangle$, $\text{cons}:\langle\langle\triangleright, \text{EID}, \langle\langle\text{Emp_EID}, \text{Emp}, \text{EID}\rangle\rangle\rangle\rangle$ and $\text{cons}:\langle\langle\overset{\text{id}}{\rightarrow}, \text{EID}, \langle\langle\text{Emp_EID}, \text{Emp}, \text{EID}\rangle\rangle\rangle\rangle$, generates *HasKey(Emp, Emp_EID)*. The complete OWL schema resulting from transforming S_{ho} is listed in Figure 4.

4. HDM Schema Enhancement

This section presents an extension to the schema enhancement approach, presented in [7], in which we explain new heuristics. Our schema enhancement process applies schema analysis and data analysis heuristics on the HDM schema S_{ho} producing an enriched schema S_{ho+} .

Table 2. HDM Representations for Some OWL 2 Constructs

OWL 2 Construct	DL Syntax	HDM Representation
Class	C	node:⟨⟨C⟩⟩
SubClassOf	$C_1 \sqsubseteq C_2$	cons:⟨⟨ \sqsubseteq , C_1 , C_2 ⟩⟩
ObjectProperty	P	edge:⟨⟨P, C_1 , C_2 ⟩⟩
DataProperty	R	edge:⟨⟨R, C_1 , rdfs:Literal⟩⟩
FunctionalProperty	$\top \sqsubseteq (\leq 1 P)$	cons:⟨⟨ \triangleleft , C_1 , ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩
TransitiveProperty	$P \circ P \sqsubseteq P$	cons:⟨⟨ \sqsubseteq , $\pi_{\langle P/P_1, C_1, P/P_2, C_2 \rangle} P \triangleright \triangleleft P$, ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩
IntersectionOf	$C_1 \sqcap C_2$	node:⟨⟨NC ₁ ⟩⟩, cons:⟨⟨ \cap , C_1 , NC ₁ ⟩⟩, cons:⟨⟨ \cup , owl:Thing, C_1 , NC ₁ ⟩⟩, node:⟨⟨NC ₂ ⟩⟩, cons:⟨⟨ \cap , C_2 , NC ₂ ⟩⟩, cons:⟨⟨ \cup , owl:Thing, C_2 , NC ₂ ⟩⟩, node:⟨⟨UNC ₁ _NC ₂ ⟩⟩, cons:⟨⟨ \cup , UNC ₁ _NC ₂ , NC ₁ , NC ₂ ⟩⟩, node:⟨⟨IC ₁ _C ₂ ⟩⟩, cons:⟨⟨ \cap , IC ₁ _C ₂ , UNC ₁ _NC ₂ ⟩⟩, cons:⟨⟨ \cup , owl:Thing, IC ₁ _C ₂ , UNC ₁ _NC ₂ ⟩⟩
HasValue	$\exists P.\{a\}$	node:⟨⟨VP.a⟩⟩, edge:⟨⟨P_IE, VP.a, a⟩⟩, cons:⟨⟨ \triangleright , VP.a, ⟨⟨P_IE, VP.a, a⟩⟩⟩⟩, cons:⟨⟨ \sqsubseteq , ⟨⟨P_IE, VP.a, a⟩⟩, ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩
Cardinality	$= nP$	node:⟨⟨=nP⟩⟩, edge:⟨⟨P_IE, =nP, owl:Thing⟩⟩, cons:⟨⟨ \triangleleft^n , =nP, ⟨⟨P_IE, =nP, owl:Thing⟩⟩⟩⟩, cons:⟨⟨ \triangleright^n , =nP, ⟨⟨P_IE, =nP, owl:Thing⟩⟩⟩⟩, cons:⟨⟨ \sqsubseteq , ⟨⟨P_IE, =nP, owl:Thing⟩⟩, ⟨⟨P, C, D⟩⟩⟩⟩
MinCardinality	$\geq nP$	node:⟨⟨ $\geq nP$ ⟩⟩, edge:⟨⟨P_IE, $\geq nP$, owl:Thing⟩⟩, cons:⟨⟨ \triangleleft^n , $\geq nP$, ⟨⟨P_IE, $\geq nP$, owl:Thing⟩⟩⟩⟩, cons:⟨⟨ \sqsubseteq , ⟨⟨P_IE, $\geq nP$, owl:Thing⟩⟩, ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩
MaxCardinality	$\leq nP$	node:⟨⟨ $\leq nP$ ⟩⟩, edge:⟨⟨P_IE, $\leq nP$, owl:Thing⟩⟩, cons:⟨⟨ \triangleright^n , $\leq nP$, ⟨⟨P_IE, $\leq nP$, owl:Thing⟩⟩⟩⟩, cons:⟨⟨ \sqsubseteq , ⟨⟨P_IE, $\leq nP$, owl:Thing⟩⟩, ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩
HasKey (Key)		cons:⟨⟨ \triangleright , C_1 , ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩, cons:⟨⟨ \triangleleft , C_1 , ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩, cons:⟨⟨ \xrightarrow{id} , C_1 , ⟨⟨P, C_1 , C_2 ⟩⟩⟩⟩

$\exists RTo.T \sqsubseteq Emp$	(1)	$\top \sqsubseteq \forall Ter_TID.xsd:string$	(10)	$\top \sqsubseteq \forall Reg_RID.xsd:integer$	(19)
$\top \sqsubseteq \forall RTo.Emp$	(2)	$\top \sqsubseteq (\leq 1 Ter_TID)$	(11)	$\top \sqsubseteq (\leq 1 Reg_RID)$	(20)
$\top \sqsubseteq (\leq 1 RTo)$	(3)	$\exists Ter_TDes.T \sqsubseteq Territories$	(12)	$\exists Reg_RDes.T \sqsubseteq Reg$	(21)
$\exists Emp_EID.T \sqsubseteq Emp$	(4)	$\top \sqsubseteq \forall Ter_TDes.xsd:string$	(13)	$\top \sqsubseteq \forall Reg_RDes.xsd:string$	(22)
$\top \sqsubseteq \forall Emp_EID.xsd:integer$	(5)	$\top \sqsubseteq (\leq 1 Ter_TDes)$	(14)	$\top \sqsubseteq (\leq 1 Reg_RDes)$	(23)
$\top \sqsubseteq (\leq 1 Emp_EID)$	(6)	$\exists Ter_Reg.T \sqsubseteq Ter$	(15)	HasKey(Emp, Emp_EID)	(24)
$\exists Emp_Ter.T \sqsubseteq Emp$	(7)	$\top \sqsubseteq \forall Ter_Reg.Reg$	(16)	HasKey(Ter, Ter_TID)	(25)
$\top \sqsubseteq \forall Emp_Ter.Ter$	(8)	$\top \sqsubseteq (\leq 1 Ter_Reg)$	(17)	HasKey(Reg, Reg_RID)	(26)
$\exists Ter_TID.T \sqsubseteq Ter$	(9)	$\exists Reg_RID.T \sqsubseteq Reg$	(18)	Emp $\sqsubseteq = 1 Reg_RID$	(27)
				Ter $\sqsubseteq = 1 Ter_TDes$	(28)

Figure 4. S_o , OWL Schema (in DL syntax) Representing Schema S_r

4.1. Schema Analysis

Our **Schema Analysis (SA)** heuristic rules, presented in this section, take the form $\text{pattern} \rightsquigarrow \text{action}$, where pattern is a pattern to match against the constructs of S_{ho} , and action is either a set of BAV transformations to apply to S_{ho} , or instructions to perform further data analysis (which if positive, will add BAV transformations to S_{ho}). Below, we list some of the SA heuristics we have implemented in our prototype.

Object Property Characteristics Heuristics: In OWL, the domain and range of **Symmetric** and **Transitive** properties must match [12]. Hence, when we search for such constraints on properties, we only need to consider a property P between a class D and itself, or between D and a superclass C of D. In addition, we propose the heuristic that **Reflexive** is most likely to occur on the same type of property, hence also search for it.

Given the above, we now consider whether P is **Functional** or not. A **Functional** property cannot be **Transitive** [12], and is not useful if **Reflexive** (since it would only relate instances to themselves). Finally, **Symmetric** properties are less likely to be **Functional**, since it would restrict the instances to be in pairs. Hence, we only search for these three types of property constraints when the edge is non-functional:

SA₁: edge:⟨⟨P, D, D⟩⟩ ∧ ¬ cons:⟨⟨⊂, D, P⟩⟩ ∼
 DA_{symm}(edge:⟨⟨P, D, D⟩⟩); DA_{refl}(edge:⟨⟨P, D, D⟩⟩); DA_{tran}(edge:⟨⟨P, D, D⟩⟩)
 SA₂: edge:⟨⟨P, C, D⟩⟩ ∧ ¬ cons:⟨⟨⊂, C, P⟩⟩ ∧ cons:⟨⟨⊆, D, C⟩⟩ ∼
 DA_{symm}(edge:⟨⟨P, C, D⟩⟩); DA_{refl}(edge:⟨⟨P, C, D⟩⟩); DA_{tran}(edge:⟨⟨P, C, D⟩⟩)

Note that DA_{symm}, DA_{refl}, and DA_{tran} are data analysis heuristics that will be described in the next section.

If circumstances where a property might be **Symmetric** or **Reflexive**, we consider it sensible to also identify whether they are **Asymmetric** or **Irreflexive** properties. Here, it is possible that the property might be **Functional**, so there is no restriction on the property being **Functional** in the rules.

SA₃: edge:⟨⟨P, D, D⟩⟩ ∼
 DA_{asym}(edge:⟨⟨P, D, D⟩⟩); DA_{irre}(edge:⟨⟨P, D, D⟩⟩)
 SA₄: edge:⟨⟨P, C, D⟩⟩ ∧ cons:⟨⟨⊆, D, C⟩⟩ ∼
 DA_{asym}(edge:⟨⟨P, C, D⟩⟩); DA_{irre}(edge:⟨⟨P, C, D⟩⟩)

Again, DA_{asym} and DA_{irre} are further data analysis heuristics that will later be described.

Property Cardinality Restriction Heuristic: **MinCardinality** and **MaxCardinality** axioms can be proposed on any non-**Functional**/ non-**InverseFunctional** edge. We rely on data analysis to confirm the **MaxCardinality** of the instances of the class on this property.

SA₅: edge:⟨⟨P, C, D⟩⟩ ∧ ¬ cons:⟨⟨⊂, C, P⟩⟩ ∼
 DA_{card}(edge:⟨⟨P, C, D⟩⟩)

Class Axioms Heuristic: An **IntersectionOf** expression can be between classes or between classes and other expressions such as the **HasValue** expression. Although an **IntersectionOf** expression that contains a **HasValue** expression is complicated to search for, nevertheless, it can be found in relational databases by exploiting schema and data analysis. Here, the condition is that the schema should include a hierarchy (at least one **SubClassOf**) and the action is a data analysis heuristic.

SA₆: cons:⟨⟨⊆, D, C⟩⟩ ∼
 DA_{interHasValue}(node:⟨⟨C⟩⟩, node:⟨⟨D⟩⟩)

4.2. Data Analysis

Each **Data Analysis (DA)** heuristic rule has the form pattern : condition | probability ∼ action where pattern is a pattern to match a DA instruction from SA, condition is a query to execute against S_{ho} that must return true, and probability is a query generating a number [0, 1]. We rely on a domain expert to validate that proposed additions in action

to the ontology are correct (guided by the probabilities associated with each rule).

Transitive Heuristic: DA heuristics such as DA_{symm} , DA_{asym} , DA_{refl} , DA_{irre} and DA_{tran} all take the same general form. We illustrate the approach by presenting DA_{tran} . In the rule below, the notation P/P_1 indicates aliasing of P as P_1 .

$$DA_{\text{tran}}(\text{edge}:\langle\langle P, C, D \rangle\rangle): \text{totalT} \neq \text{notT} \mid \frac{\text{totalT} - \text{notT}}{\text{totalT}} \rightsquigarrow$$

$$\text{addCons}(\langle\langle \langle \langle P/P_1, D=P/P_2, C \rangle, \langle P, C, D \rangle \rangle, \langle P, C, D \rangle \rangle \rangle);$$

$$\text{addCons}(\langle\langle \langle \langle \pi_{(P/P_1, C, P/P_2, D)} \rangle, \langle \langle \langle P/P_1, D=P/P_2, C \rangle, \langle P, C, D \rangle \rangle \rangle, \langle P, C, D \rangle \rangle \rangle);$$

$$\text{addCons}(\langle\langle \langle \langle \subseteq, \langle \langle \pi_{(P/P_1, C, P/P_2, D)} \rangle, \langle \langle \langle P/P_1, D=P/P_2, C \rangle, \langle P, C, D \rangle \rangle \rangle \rangle, \langle P, C, D \rangle \rangle \rangle);$$

where totalT is the number of transitive instances of P calculated as:

$$\text{totalT} = |\{(x, z) \mid (x, y) \in \langle\langle P, C, D \rangle\rangle; (y, z) \in \langle\langle P, C, D \rangle\rangle\}|$$

and notT is the number of non-transitive instances of P , calculated as:

$$\text{notT} = |\{(x, z) \mid (x, y) \in \langle\langle P, C, D \rangle\rangle; (y, z) \in \langle\langle P, C, D \rangle\rangle - \langle\langle P, C, D \rangle\rangle\}|.$$

The **Transitive** heuristic, was found not to apply to the Northwind database.

Min/Max Cardinality Heuristic: The **MaxCardinality** axiom can be proposed from counting the instances of data involved in properties while a **MinCardinality** axiom of 1 can be added on any non-**Functional** / non-**InverseFunctional** property. The action below pertains to the **MaxCardinality** axiom. The BAV transformations for adding a **MinCardinality** axiom is very similar and can be found in Table 2.

$$DA_{\text{card}}(\text{edge}:\langle\langle P, C, D \rangle\rangle): N = \text{count}(\langle\langle P, C, D \rangle\rangle) \mid 1 \rightsquigarrow$$

$$\text{addNode}(\langle\langle \leq NP \rangle\rangle); \text{addEdge}(\langle\langle P_IE, \leq N P, D \rangle\rangle); \text{addCons}(\langle\langle \triangleright^N, \leq N P, \langle\langle P_IE \rangle\rangle \rangle \rangle);$$

$$\text{addCons}(\langle\langle \subseteq, P_IE, P \rangle\rangle)$$

Suppose that we count the maximum number of territories that an employee can work in is 8, therefore we can propose to add to the HDM schema a **MaxCardinality** constraint between the node $\langle\langle \text{Emp} \rangle\rangle$ and the edge $\langle\langle \text{Emp_Ter}, \text{Emp}, \text{Ter} \rangle\rangle$. The **MaxCardinality** construct can then be added to the ontology as in the DL rule: $\text{Emp} \subseteq \leq 8 \text{Emp_Ter.Ter}$.

IntersectionOf including HasValue Heuristic: The DA heuristic for finding an **IntersectionOf** that contains a **HasValue** expression is different than the previous heuristics. In here, the condition is an algorithm (presented in Algorithm 1) that returns a set of axioms.

$$DA_{\text{interHasValue}}(\text{node}:\langle\langle C \rangle\rangle, \text{node}:\langle\langle D \rangle\rangle): \text{checkIntersectionOfHasValue}(C, D) \neq \emptyset \mid 1 \rightsquigarrow$$

$$\text{addNode}(\langle\langle VP.a \rangle\rangle); \text{addEdge}(\langle\langle P_IE, VP.a, a \rangle\rangle); \text{addCons}(\langle\langle \triangleright, VP.a, \langle\langle P_IE, VP.a, a \rangle\rangle \rangle \rangle);$$

$$\text{addCons}(\langle\langle \subseteq, \langle\langle P_IE, VP.a, a \rangle\rangle, \langle\langle P, C_1, C_2 \rangle\rangle \rangle \rangle); \text{addNode}(\langle\langle NVP.a \rangle\rangle);$$

$$\text{addCons}(\langle\langle \overline{\cap}, VP.a, NVP.a \rangle\rangle); \text{addCons}(\langle\langle \cup, owl:Thing, VP.a, NVP.a \rangle\rangle);$$

$$\text{addNode}(\langle\langle NC \rangle\rangle); \text{addCons}(\langle\langle \overline{\cap}, C, NC \rangle\rangle); \text{addCons}(\langle\langle \cup, owl:Thing, C, NC \rangle\rangle);$$

$$\text{addNode}(\langle\langle UNVP.a_NC \rangle\rangle); \text{addCons}(\langle\langle \cup, UNVP.a_NC, NVP.a, NC \rangle\rangle);$$

$$\text{addNode}(\langle\langle IC.VP.a \rangle\rangle); \text{addCons}(\langle\langle \overline{\cap}, IC.VP.a, UNVP.a_NC \rangle\rangle);$$

$$\text{addCons}(\langle\langle \cup, owl:Thing, IC.VP.a, UNVP.a_NC \rangle\rangle)$$

The algorithm loops all properties in the schema and for each property queries the schema to find the y values of members of the parent class (*i.e.*, $\text{node}:\langle\langle C \rangle\rangle$) which participate in the property denoted as $PValue$. Then, it queries the schema again to find the y values of members of the child class (*i.e.*, $\text{node}:\langle\langle D \rangle\rangle$) which participate in the same property denoted as $CValue$. If there is only one $CValue$ and more than one $PValue$,

Algorithm 1. checkIntersectionOfHasValue

```

1: procedure CHECKINTERSECTION
2:   inputs: Parent, Child
3:   Axioms  $\leftarrow \emptyset$ 
4:   for  $P \in \text{SchemaProperties}$  do
5:      $P\text{Member} \leftarrow \{ \langle x, y \rangle \mid \langle x, y \rangle \in \text{edge}:\langle\langle P \rangle\rangle; \langle x \rangle \in \text{node}:\langle\langle \text{Parent} \rangle\rangle \}$ 
6:      $P\text{Value} \leftarrow \{ y \mid \langle x, y \rangle \in P\text{Member} \}$ 
7:      $C\text{Member} \leftarrow \{ \langle x, y \rangle \mid \langle x, y \rangle \in \text{edge}:\langle\langle P \rangle\rangle; \langle x \rangle \in \text{node}:\langle\langle \text{Child} \rangle\rangle \}$ 
8:      $C\text{Value} \leftarrow \{ y \mid \langle x, y \rangle \in C\text{Member} \}$ 
9:     if  $|C\text{Value}| = 1 \wedge |P\text{Value}| > 1$  then
10:       Axioms  $\leftarrow \text{Axioms} \cup (\text{Child} \equiv \text{Parent} \sqcap P.C\text{Value})$ 
11:     end if
12:   end for
13:   return Axioms
14: end procedure

```

the algorithm returns an axiom. For instance, if we envisage that the schema in Figure 1 also includes a table $\text{Mgr}(\underline{\text{MID}})$ with $\text{Mgr}(\underline{\text{MID}}) \xrightarrow{fk} \text{Emp}(\underline{\text{EID}})$, and that the table has only the data (6,7), then using this heuristic we could detect the following axiom: $\text{Mgr} \equiv \text{Emp} \sqcap \text{RTTo}.\{5\}$.

5. Applying Machine Learning

This section explains how the machine learning module has been implemented and integrated with the OWLRel tool. It is structured in such a way that it follows a data pipeline, from the way data is collected and processed to how it is used in forecasting. This encompasses the data collection step, the reconstruction of the HDM graph in an in memory representation which facilitates the feature extraction process, and finally the training of the prediction model.

The general assumption behind our approach is to find patterns that transcend the particularities of individual users, which then can be used to produce a universal score of relevance for any axiom. In order to find such patterns, the users' interaction with the OWLRel tool is recorded in log files. During the execution, in each iteration when axioms are being suggested, they are first processed through the aforementioned pipeline. Then, the models are used to predict whether an axiom will be selected as interesting or not, and this prediction takes the form of a probability that will be displayed to the user through the GUI shown in Figure 5.

5.1. Data Collection

As the tool suggests axioms to its users, their answers are stored in log files. These log files contain information such as the suggested axiom name, whether it was marked as interesting or not, and the iteration number in which it has appeared and been clicked on. Any entry from a log file would therefore have the following schema: [axiom, clickedInteresting, clickedUninteresting, iterationAppeared, iterationClicked], and example for the Northwind database being [irreflexiveproperty:⟨⟨RTTo, Emp, Emp⟩⟩, 1, 0, 0, 1]. The

axiom is suggesting that the property `RTo` is an `IrreflexiveProperty`. The axiom has been clicked as interesting by the user in the second iteration, and it was suggested in the first iteration.

5.2. Internal Graph Reconstruction

After collecting the log files, the next step is to extract features that allow interesting axioms to be distinguished from uninteresting ones. As a prerequisite to compute such features, the graph representation of the ontology is reconstructed in memory. The initial state of the graph is given by the original ontology (such as the one presented in Figure 3), and with each interaction the graph is updated, enriching the ontology, and therefore changing the features for subsequent axioms. Accepting a suggested axiom from the GUI will inevitably change the properties of the underlying edges and nodes.

For example, acknowledging that `RTo` is an `IrreflexiveProperty`, will update the list of characteristics of that particular edge. At this stage, the `Emp` class, represented as a node, has two properties: `RTo` and `Emp_Ter`, which are in return represented as edges. If a new axiom, for example a `PropertyChain`, would be suggested and accepted, it would introduce a new property *i.e.*, a new edge in the graph. These structural changes will then change the features derived from the graph. In the case of the `PropertyChain`, it will among other things, increase the number of connections a node has. As a consequence of working with a graph, features can now be derived separately for any node and edge.

5.3. Feature Engineering

The internal graph representation allows the system to extract features that best describe an axiom. The feature set of an axiom is based on information about the classes and properties that appear in that axiom. The features collected from nodes are: (1) importance, (2) number of connections, (3–5) number of inclusions in `UnionOf`, `IntersectionOf`, `DisjointWith` and (6–7) to how many other nodes it is a super class or sub class of. The importance is governed by how many times the user has selected an axiom that directly affects the node, while the number of connections is the number of edges connecting to the node. With respect to the previous example of the `IrreflexiveProperty`, the feature vector derived for the `Emp` class is $Features_{class} = [0, 2, 0, 0, 0, 0, 0]$. The only non-zero value is the number of connections, since the class has two properties. Otherwise, the class is not part of a `UnionOf`, `IntersectionOf`, `DisjointWith` or `SubClassOf`.

The features extracted for properties are: (1) the number of characteristics it contains (`SymmetricProperty`, `TransitiveProperty`, *etc.*), (2–4) the presence of `MinCardinality`, `MaxCardinality` and `HasValue`, (5) number of properties it is an `InverseOf`, (6) its importance (as defined for nodes) and (7) the iteration it appeared in. In contrast to the nodes, new properties can dynamically appear if a `PropertyChain` is being chosen from the suggested axiom list. For the same example, the feature vector for the `RTo` property is $Features_{property} = [1, 0, 0, 0, 0, 0, 0]$. Here, the only non-zero value is the number of characteristics, since the property is a `FunctionalProperty`.

When a user selects an axiom, it enriches the ontology such that the features for all other axioms referring to nodes or edges in the neighborhood are being affected. With the addition of a new property the number of edges change, and some nodes gain in importance. Therefore, the relevance is being recomputed after each enrichment and all the suggested axioms will have different scores.

As we now have the features for classes and properties, it is possible to derive the complete feature vectors for axioms. Axioms such as cardinality, property characteristics, `InverseOf` and `HasValue` always affect properties between two classes. Therefore, they can be represented by the information of the affected nodes and edges. A `PropertyChain` axiom, on the other hand, affect two sets of nodes and edges. This can be converted to the previously mentioned feature space by taking the average of the edge and node features. `UnionOf`, `IntersectionOf` and `DisjointWith` axioms do not refer to an existing property, and thus cannot be modelled together with the previous axioms. Because of this, two machine learning models have been used: one that focuses on property characteristics, `PropertyChain`, `InverseOf`, `HasValue` and cardinality, and another model that focuses only on the `UnionOf`, `IntersectionOf` and `DisjointWith` axioms. The first machine learning model receives as input the averaged features of all the nodes and properties it affects. The second model acts only on nodes, so it receives solely the averaged features of these. Additionally, each observation receives a categorical variable describing the type of axiom, along with two variables quantifying in which iteration it has appeared and when has it been clicked as interesting. Each observation gathered from the log file will therefore be processed to contain 10 features in the case of `UnionOf`, `IntersectionOf` and `DisjointWith` axioms, while 17 for the other axiom types. For an axiom affecting classes and properties the feature vector would be $Features_{axiom} = [Avg(Features_{class}), Avg(Features_{property}), itApp, itClick, axiomType]$. Finally, the feature vector given to the machine learning algorithms for the `RTo IrreflexiveProperty` has the following value: $[0, 2, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, Irreflexive]$.

5.4. Training and Validation

The nature of the collected data allows for supervised learning methods to be deployed, as the ‘correct’ label for each suggested axiom has been recorded. The machine learning models taken into consideration for this task were the Naive Bayes classifier and the Logistic Regression due to their simplicity, and the fact that they output probabilities rather than just class labels. The Logistic Regression model is defined as in equation (29) where the linear model described by the coefficients beta is applied onto the input variables x_i (e.g., the feature vector) after which it is being squashed by the logistic function to output a probability.

$$y = \frac{1}{1 + e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}} \quad (29)$$

Data has been collected from three users enriching the Wine Ontology⁴ with axioms proposed from the schema enhancement phase described in Section 4. The total amount of axioms interacted with, amounts to roughly 300. Using the Weka toolkit⁵, the machine learning algorithms were run with the default model and regularization parameters provided. The performance of the models was assessed using a standard 10-Fold Cross Validation leading to the conclusion that the Logistic Regression is the better choice, since it obtained an average accuracy of 86%. Besides of the standard, numerical evaluation method, the applicability of the models has been assessed in practice by manual

⁴<http://www.w3.org/TR/owl-guide/wine.rdf>

⁵<http://www.cs.waikato.ac.nz/ml/weka/>

inspection, showing no significant discrepancy between the predicted relevance and the users' actual preference.

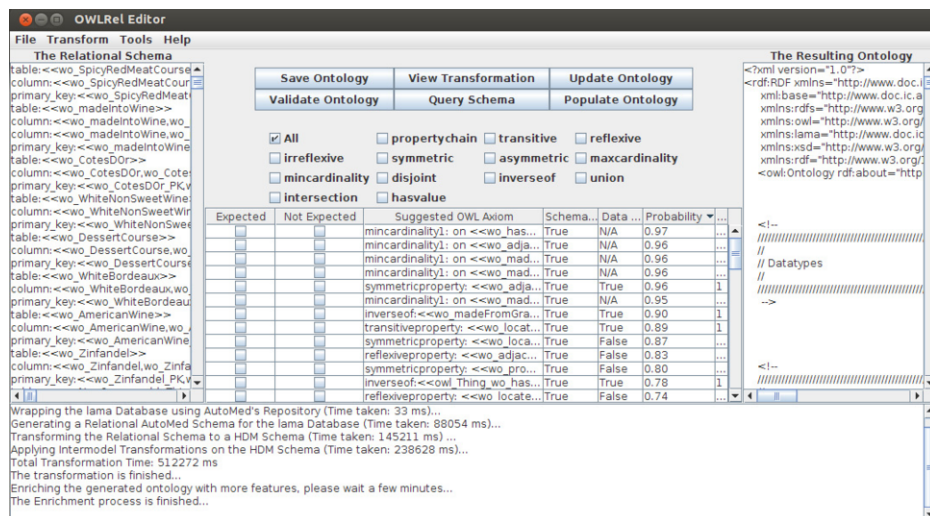


Figure 5. Ontology Enrichment using OWLRel Editor

6. Related Work

Most of the approaches that extract a domain-specific ontology from a relational database, such as: [13], [14] and [5], suffer from one or more of the following problems [15]: 1) Do not map highly expressive OWL features. 2) Miss interpret primary keys. 3) Do not properly identify class hierarchies.

The survey [15] points out that transforming a relational schema on its own is not enough, and some sort of data inspection is needed in order to solve problems 1 and 3. Problem 3 was addressed by the work of [14] and [16] in which they apply data mining techniques to detect class hierarchies. Moreover, primary keys have been interpreted in various ways in the OWL model. Some approaches translate a PK into an [InverseFunctional DataProperty](#) with [MinCardinality](#) 1, or use [MinCardinality](#) and [MaxCardinality](#) 1, or use [FunctionalProperty](#) and [Cardinality](#) 1. Only a few, like us, have considered using the [HasKey](#) construct such as [13] and [17]. Furthermore, eliciting the highly expressive constructs of OWL, such as [HasValue](#), [TransitiveProperty](#) and [SymmetricProperty](#), was only considered by [6]. However, their rules for generating these constructs are ambiguous and do not rely on the analysis of the data. In contrast to this, we have proposed an approach that analyses the schema and the data to suggest possible highly expressive OWL 2 features for a given relational database and represented primary keys using the [HasKey](#) construct, thus providing solutions for problems 1 and 2 above. Moreover, our approach applies machine learning to the extraction of OWL ontologies from relational databases, something that was not considered by any of the existing state-of-the-art approaches.

7. Conclusions and Future Work

This paper has presented an approach that extracts OWL ontologies from relational databases using a combination of data analysis and machine learning. The data analysis technique used helped in producing rich ontologies that contain more semantic information than found in the relational databases. Applying machine learning algorithms to the extraction of ontologies helped in ranking the large number of suggestions produced by the data analysis technique. Testing our approach showed that it is feasible and promising. Future work will be directed towards enhancing our heuristics and the machine learning algorithms used.

References

- [1] X. D. M. Li and S. Wang, "Learning Ontology from Relational Database," in *In proc. of Inter. Conf. on Machine Learning*, vol. 6, pp. 3410–3415, 2005.
- [2] L. Lin, Z. Xu and Y. Ding, "OWL Ontology Extraction from Relational Databases via Database Reverse Engineering," *Journal of Software*, vol. 8, no. 11, pp. 2749–2760, 2013.
- [3] B. El Idrissi, S. Baina, and K. Baina, "Automatic Generation of Ontology from Data Models: A Practical Evaluation of Existing Approaches," in *IEEE 7th Inter. Conf. on RCIS*, pp. 1–12, 2013.
- [4] W3C, "OWL 2 Web Ontology Language New Features and Rationale," June 2009. <http://www.w3.org/TR/2009/WD-owl2-new-features-20090611/>.
- [5] S. H. Tirmizi, J. Sequeda, and D. Miranker, "Translating SQL Applications to the Semantic Web," in *DEXA*, pp. 450–464, Springer, 2008.
- [6] I. Astrova, N. Korda, and A. Kalja, "Rule-Based Transformation of SQL Relational Databases to OWL Ontologies," in *Proc. of the 2nd Inter. Conf. on Metadata & Semantics Research*, 2007.
- [7] L. Al Khuzayem and P. McBrien, "OWLRel: Learning Rich Ontologies from Relational Databases," *Baltic Journal of Modern Computing*, vol. 4, no. 3, p. 466, 2016.
- [8] A. Poulouvasilis and P. McBrien, "A General Formal Framework for Schema Transformation," *DKE*, vol. 28, no. 1, pp. 47–71, 1998.
- [9] M. Boyd and P. McBrien, "Comparing and Transforming between Data Models via an Intermediate Hypergraph Data Model," *DS*, vol. IV, pp. 69–109, 2005.
- [10] A. C. Smith and P. McBrien, "Inter Model Data Exchange of Type Information via a Common Type Hierarchy," in *DISWEB*, 2006.
- [11] P. McBrien and A. Poulouvasilis, "Data Integration by Bi-Directional Schema Transformation Rules," in *Proc. of ICDE*, pp. 227–238, IEEE, 2003.
- [12] L. Lacy, *OWL: Representing Information Using the Web Ontology Language*. Victoria BC, Canada: Trafford Publishing, 2005.
- [13] J. Sequeda, M. Arenas, and D. Miranker, "On Directly Mapping Relational Databases to RDF and OWL," in *Proc. of the 21st WWW*, pp. 649–658, 2012.
- [14] I. Astrova, "Reverse Engineering of Relational Databases to Ontologies," in *The Semantic Web: Research and Applications*, pp. 327–341, Springer, 2004.
- [15] D.-E. Spanos, P. Stavrou, and N. Mitrou, "Bringing Relational Databases into the Semantic Web: A Survey," *Semantic Web*, vol. 3, no. 2, pp. 169–209, 2012.
- [16] F. Cerbah, "Mining the Content of Relational Databases to Learn Ontologies with Deeper Taxonomies," in *Inter. Conf. on WI-IAT*, vol. 1, pp. 553–557, IEEE, 2008.
- [17] E. Vysniauskas, L. Nemuraite, R. Butleris, and B. Paradauskas, "Reversible Lossless Transformation From OWL 2 Ontologies into Relational Databases," *IJITCA*, vol. 40, no. 4, pp. 293–306, 2011.