

ORIGINAL ARTICLE

Journal Section

Using Semantic Web technologies in the development of Data Warehouses: A Systematic Mapping

Ricardo Gacitua^{1*} | Jose Norberto Mazon^{2†} | Ania Cravero^{1*}

¹ Department of Computer Science and Informatics, Universidad de La Frontera, Temuco, Chile

² Department of Software and Computing Systems, Universidad de Alicante, Alicante, Spain

Correspondence

Ricardo Gacitua, Department of Computer Science and Informatics, Universidad de La Frontera, Temuco, Chile
Email: ricardo.gacitua@ufroterra.cl

Present address

[†] Department of Computer Science and Informatics, Universidad de La Frontera, Temuco, Chile

Funding information

Universidad de la Frontera, Chile, Grant Number: DI15-0020 and DI17-0043

The exploration and use of Semantic Web technologies have attracted considerable attention from researchers examining data warehouse development. However, the impact of this research and the maturity level of its results are still unclear. The objective of this study is to examine recently published research articles that take into account the use of Semantic Web technologies in the data warehouse arena with the intention of summarizing their results, classifying their contributions to the field according to publication type, evaluating the maturity level of the results and identifying future research challenges. Three main conclusions were derived from this study: (a) There is a major technological gap that inhibits the wide adoption of Semantic Web technologies in the business domain, (b) There is limited evidence that the results of the analyzed studies are applicable and transferable to industrial use, (c) Interest in researching the relationship between data warehouses and Semantic Web has decreased because new paradigms, such as Linked Open Data, have attracted the interest of researchers.

KEYWORDS

Data Warehouses, Semantic Web, Systematic Mapping

1 | INTRODUCTION

To include external Web data in the traditional Data Warehouse systems (DWs), and the traditional On-Line Analytical Processing (OLAP) processes, is a promising way for broadening traditional Business Intelligence (BI) analysis Abelló et al. (2013). In the context of BI, a data warehouse is used to collect, organize and store subject-oriented, integrated, time-variant and non-volatile data (Inmon, 1992). Traditionally, a DW has been defined as a historical data repository containing data collected from a wide variety of heterogeneous sources by means of Extraction-Transformation-Loading (ETL) processes (Kimball and Ross, 2002). On the other hand, multi-dimensional processing, also called On-Line Analytical Processing (OLAP), is an approach for responding to multi-dimensional analytical queries. Research focusing on DWs and OLAP has led to the creation of important technologies for the design, administration and use of information systems in decision-making support. Part of the interest in, and success of this field, can be attributed to the demonstrated need for software and tools that help improve data analysis and administration. This is mainly due to the large quantity of information that is being accumulated by corporations as well as scientific databases.

Traditional BI tools, such as OLAP, have been successfully applied to large amounts of data coming from operational databases. However, there is a trend whereby DWs are becoming more and more dynamic, with updates occurring almost in real-time, and with the inclusion of more complex types of data (Henschen, 2015). This situation has forced traditional BI to open its gates to external data in order to encompass a more heterogeneous and open analysis scenario (Chen et al., 2012). Current research envisions that Semantic Web technologies are required for realizing the next generation of DWs (Abelló et al., 2015), as an increasing quantity of semantically annotated data is available over the Internet¹. To include Semantic Web information in a traditional OLAP analysis process is therefore a promising way to augment traditional BI analyses (Trujillo and Maté, 2012). The Semantic Web is an extension of the Web proposed by the World Wide Web Consortium². Its intended objective is to facilitate the creation of technologies that publish legible data for informatic applications, which are implemented by adding semantic metadata and ontologies to the Web (Shadbolt et al., 2006). In practical terms, the strength of the Semantic Web lies in its ability to aggregate the semantic annotations of Web-published content so that the information can be effectively retrieved and processed, either by humans or by machines, for a wide variety of tasks (Hendler, 2001). The above is achieved through the use of a diverse variety of software technologies, such as ontologies (Coral et al., 2006) and markup languages (e.g. RDF, OWL) (Saha, 2007), which allow semantic annotations to be added to resources that can either be very simple, or very complex annotations, depending on the requirements.

Although the Semantic Web and DWs, in the context of BI applications, have gone in different research directions over the last few years, some recent results show that the convergence of these two fields is not only inevitable but also beneficial for both sides (Golfarelli et al., 2004; Lather, 2012). The exploration and use of Semantic Web technology has therefore attracted attention from researchers studying DW development (Berlanga et al., 2014; Golfarelli et al., 2004). The concept of the Linked Data Web has emerged as a mechanism to make all data (Klyne and Carroll, 2004) available using the HTTP protocol, as happens with HTML documents (Bizer et al., 2009). Linked Open data has been introduced as a promising paradigm for opening up data because it facilitates data integration on the Web (Bizer et al., 2009). Similar to Data Warehousing approaches, Linked Open Data can be prepared to enable sophisticated data analysis. As one of the main problems of DWs is data integration, some researchers propose to use Data Exchange Standards like RDF, to host structured content to publish DW content as Linked Data. In the business and enterprise domain there is still a gap between conceptual approaches for modeling architectures for Linked Data, systems, data models, as well as their implementation, operationalization and execution (Abramowicz et al., 2016). Additionally, businesses need

¹<http://linkeddata.org>

²<http://www.w3.org>

to overcome the inevitable tension between the value they traditionally assign to proprietary data, and the value of opening up.

The proposals of the Semantic Web and Data Storage come from very different areas. Thus the studies that have been published have approached the integration of these topics through varied perspectives. As the development of these areas has advanced, there has been increased growth in the number of published reports and results. Hence, summarizing and providing a general analysis of the integration of these topics is quickly becoming necessary. As a result, it is difficult to know the current state of these proposals. Essentially, the impact of these research studies, as well as the maturity level of their results, are still unclear. For example, classifying the scientific contributions made on this topic is useful for researchers, as the topics that address the Semantic Web and Data Warehouses have not necessarily been developed in a related way. Indeed, each topic has its own separate conferences and journals. Providing a general summary of the progress in integrating these two areas is therefore useful in summarizing and classifying research as well as in establishing future research challenges. The Systematic Mapping Study is a secondary study method that has recently attracted considerable attention (Petersen et al., 2008, 2015) largely because it offers a specific way of systematically reviewing the literature in regards to a particular topic (Kitchenham, 2012). A systematic map structures the type of research reports that have been realized, as well as a categorization of their published results. It usually delivers a visual summary - a map - of their results. A systematic map of the literature has been recommended mainly for research areas with a lack of relevant primary studies, as is the case for the reported use of Semantic Web technologies in DW development. In recent years there has been an increase in the number of reports related to DWs and the Semantic Web that take into account methodological proposals and techniques, among others. This goes along with the fact that these reports stem from a variety of research areas, such as Databases, Artificial Intelligence and Natural Language Processing, and the fact that the relationship between the Semantic Web and DWs is a problem of current interest. Given the massive amount of existing information (i.e. Big Data) and its use on the Web, it becomes necessary to discern the type of research that is being conducted with respect as to how Semantic Web technologies are being used in DW development. A systematic map allows for the categorization of the results and for a summary to be presented graphically.

The objective of this study is to examine recently published research articles that consider the use of Semantic Web technologies in DW development in the Business Domain. We conducted a systematic mapping study in order to study the literature with the objectives of summarizing their results, evaluating the maturity level of those results and identifying challenges for future research in this field. More specifically, we focused on the following research question: *What is the state of the recent research covering the use of Semantic Web technology in DW development?* We expounded upon the general objectives by examining more specific objectives, which were:

1. Synthesizing evidence in order to propose relevant suggestions for practical applications of Semantic Web technologies in DW development.
2. Creating a classification for the published research on DWs and the Semantic Web.
3. Identifying research trends, open problems and areas for improvement within a research body that considers both DWs and the Semantic Web together.

The rest of this article is structured in the following manner: Section 2 summarizes the motivation and contribution of this paper. Section 3 offers a general methodology and the basis of the research focus. Section 4 presents the results and a discussion of the main findings of the study. In addition, a review of some industrial tools and projects within the broad topic of DWs and the Semantic Web is provided. Section 5 presents the research challenges posed by this study. Finally, section 8 offers the conclusions of the study and recommendations for future work.

2 | MOTIVATION AND CONTRIBUTION

The main contributions of the paper can be summarized as follows:

1. The identification of the main proposals that use Semantic Web technology in DW development.
2. The identification of the Semantic Web technologies that are being used.
3. A classification framework for the contributions of scientific studies.
4. A discussion of the gaps identified in the articles reviewed.
5. An outline detailing future research challenges.

In accordance with Kitchenham et al. (2011), the contributions of a systematic map (as the one presented here) are addressed by the following people:

- Beginner researchers who are initiated in the use of the Semantic Web in the development of DWs.
- Experienced researchers who may do so as qualitative reference work that saves time for subsequent studies, provides an understanding of the existing literature on specific topics, and allows them to identify the need for conducting additional research in specific areas. This is particularly important for highly multidisciplinary fields such as the intersection between DWs and the Semantic Web, where different subcommunities publish in different publication venues, and even experienced researchers are not necessarily aware of all the facets and contributions in the broader research field.
- Industrial actors who need a thorough introduction and overview of the research field of study will find it useful. A mapping study allows industry to get an overview of state-of-the-art innovations, and to identify trends and clusters of research studies that are suitable and applicable for their particular business use, aiding communication and knowledge transfer between academia and industry. Given the relative immaturity of research on the use of Semantic Web technologies in the development of DWs, and the rapidly changing technologies, the relevance for industry cannot be overstated, since, for example, standards, tools, and case studies that are of particular interest to industry are not yet sufficiently present in the literature.

3 | LITERATURE REVIEW METHOD

As a research area matures, there is often a strong increase in the number of publications and research results. Thus the need to summarize and provide general overviews of the topic of interest begins to become more relevant. The method used in this research is a systematic mapping study, which provides a methodical and objective approach for identifying the nature and extent of the empirical study data that is available to answer a particular research question (Kitchenham et al., 2010). Systematic mapping aims to organize the research undertaken rather than to answer detailed research questions (Kitchenham, 2012). The main objective of a systematic mapping study is to provide a synthesis of the research area and identify the quantity, the type of research and the available results. Indeed, it often becomes necessary to visualize the frequency of publications over time in order to determine trends. We used ideas from da Mota et al. da Mota Silveira Neto et al. (2011) in order to use an extended systematic mapping process. This extended process includes topics not covered by Petersen et al. (2008), such as: a protocol and a classification schema. Figure-1 depicts the systematic mapping process. The basic steps of the process are: definition of the protocol and the research questions, execution of a search of relevant articles, filtration of articles, search of key concepts used in summaries

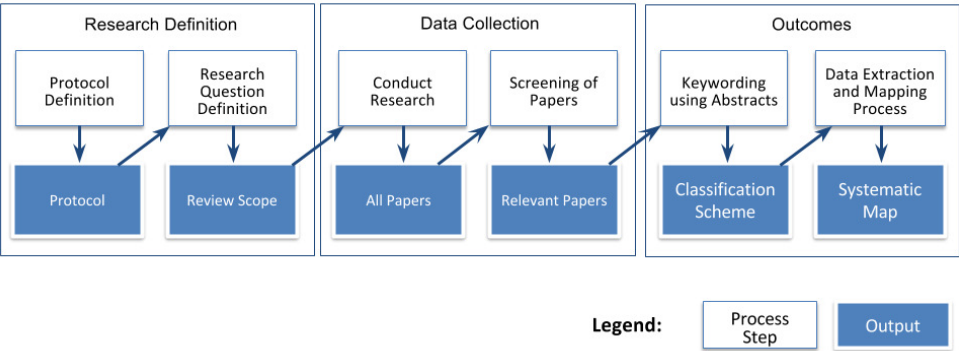


FIGURE 1 The systematic mapping process.

and, finally, the extraction of data and the processes of mapping. Each process produces output. The final output of the process is the systematic map.

The following sections describe the stages that were realized.

3.1 | First Phase: Research Definition

This section presents the first phase of the mapping study process, in which the protocol and research questions are defined.

3.1.1 | Protocol

In this study the purpose of the protocol is to manage the research objectives, and to define clearly how they can be achieved by defining the research questions and planning how the sources and selected studies will be used to respond to these questions. We adopted this task from systematic review guidelines. The first activity in this study was to develop a protocol, i.e. a strategy, defining the basic mapping study procedure.

3.1.2 | Research Questions

The objectives of systematic mapping are reflected by the *Research Questions (RQ)*. The research questions were framed by four criteria:

- **Population:** The scientific literature that shows the relationship between the Semantic Web and DWs.
- **Study design:** Describes the relationship with a focus on the study. In this case, it is defined as the use of Semantic Web technologies in DW development.
- **Intervention:** Any study that shows some degree of Semantic Web technology is used during some stage of the DW development process, such as ETL or OLAP.
- **Outcomes:** This refers to the quantity and type of related evidence regarding Semantic Web technologies and DWs, with particular interest in determining the degree of improvement on the DW development process as a result of the use of Semantic Web technologies.

In this study, restrictions regarding specific outcomes were not considered. As it is necessary to obtain a wide view of the entire area investigated based mainly on the existence of references regarding the Semantic Web and DWs, not specific Semantic Web technologies, the objective is to determine which technologies, defined as part of the Semantic Web, are referenced and used by the research community. Table-1 shows the research questions, which correspond to this study's objectives and are in accordance with the standards indicated by Petersen et al. (2008).

In order to answer the research questions data was collected from the literature. This task, in general terms, involves defining a search strategy, identifying data sources, selecting studies, and analyzing and synthesizing data.

TABLE 1 Definition of the research questions (Scope of the literature review).

ID	Questions	Objectives
RQ1	Where were the research studies published?	To determine the place of a study's publication (e.g. journals, conferences, others). This makes it possible to know where the publications are concentrated and, using that information, determine the maturity level of their results.
RQ2	What Semantic Web technologies are being used in DW development?	To determine the set of associated Semantic Web technologies that are being used in DW development. For example, ontologies (as a concept) or some of its related languages, for example, representation languages such as RDF and OWL, query languages such as SPARQL, among others. The answer to this question will make it possible to determine the most commonly used technologies.
RQ3	During which stages in the DW development process are these technologies being used?	To determine in which phases of the DW development process a particular technology is being used. For example, Requirements, ETL, or OLAP. This makes it possible to determine the area of greatest interest for researchers.
RQ4	What type of research study was conducted?	To determine what type of study was performed in order to show the use of Semantic Web technology. For example, Experiment, Case Study, Project Description, among others. This allows the published studies to be categorized.
RQ5	Is there evidence of improvement due to the use of Semantic Web technology?	To determine if the evidence shows that the presented results correspond to an improvement in the DW development process. This allows the impact of the results on DW development to be evaluated.

3.2 | Second Phase: Data Collection

In order to answer the research questions, data were collected from the research literature. This activity involved conducting research and the screening of papers.

3.2.1 | Conduct research

Conducting research involved developing a search strategy, and identifying data sources.

Search strategy The search strategy was developed through the review of the data needed to respond to each research question. Primary studies were identified using search strings on scientific databases, or searching manually through conference proceedings or publications in specialized journals. An initial set of keywords was refined after a preliminary search that retrieved many results of little relevance. Several combinations of search items were used until an appropriate set of keywords was reached. The search chain consisted of Boolean expressions composed of key words that described the stages of DW development. Table-2 shows the composition of the search string applied. Terms within a row are connected by the operation "OR", while the different parts of the search string are connected by the operator "AND" in order to improve result completeness. This results in the search string ((*"Data Sources"* OR *"ETL"* OR *"OLAP"* OR *"Analytics"* AND (*"Data Warehouse"* OR *"Data Warehousing"*) AND (*"Semantic Web"* OR *"Linked Data"*)))

TABLE 2 Table of Search Terms

No	Terms
1	Data Sources, ETL, OLAP, Analytics
2	Data Warehouse, Data Warehousing
3	Semantic Web, Linked Data

Data sources The search included important journals and related conferences with research topics related to Data Warehousing and the Semantic Web. The search was restricted to studies published between the years 2002 and 2017 in order to achieve wide coverage of the study area. This is due to the fact that several important problems in data storage creation, which were identified years ago, are still considered unresolved problems (Nguyen et al., 2005).

The initial step was conducting a search using the terms previously described in 3.2, via digital library search engines. Publications obtained from ScienceDirect, SCOPUS, IEEE Xplore, ACM Digital Library and Spring Tools were considered. The second step was the search of international review journals with articles considered relevant, which were published by IEEE, ACM, Elsevier and Springer, as they are considered high-level publication editors (McGregor, 2002). Conference proceedings were also searched. Where a conference displayed its proceedings on a published website, this was also accessed. When the proceedings were not available on a conference website, a search was conducted via the DBLP Computer Science Bibliography. The search for conference proceedings and journals produced many results that had already been obtained through the digital library search. In this case, the last results were discarded and only the first results were considered, given that these had already been included in the final list. After the search was executed for conferences and journals using digital libraries and proceedings, it became possible to identify which known publications - commonly referenced by other studies in this area, such as technical reports and theses - had not been included in the resulting final list. It was therefore decided to include these as gray literature entries. Gray literature is the term used to describe materials not commercially published or not indexed by the main databases.

3.2.2 | Screening of papers

The screening of papers involved selecting studies to analyze, and data analysis and synthesis.

Study Selection. The set of search strings was applied to the search engines, specifically in those mentioned in the previous section. The criteria for inclusion and exclusion were used to filter studies that were not relevant for answering the research questions. Inclusion criteria were used to select all the studies during the search stage. Afterwards, the criteria of exclusion were mainly applied to the titles of studies, and then to the summaries and conclusions. Regarding the inclusion criteria, the studies were only considered if they included:

- **DW development approaches that included semantic aspects.** The summary of the study explicitly mentioned the term Semantic Web or semantics in the context of DWs. From the summary, the reviewer was able to deduce that the focus of the research study was on using a Semantic Web technology at some stage of DW development.

Studies were excluded if they met the following criteria:

- **Research focus unrelated to DWs.** The article was outside the field of DWs.
- **Research focus related to DWs but insufficient information regarding the Semantic Web.** The Semantic Web was not part of the contribution to the article and its terms were only mentioned in the introductory sentence of the summary.
- **Research focus related to DWs but not related to the Semantic Web or Linked Data.** The article was within the field of DWs, but the Semantic Web or Linked Open Data was not part of the contribution to the article.
- **Duplicate studies.** When the same study was published in different articles, only the most recent was included.
- **The study had already been included from another source.**
- **The article was in a language other than English.** All studies not written in English were filtered.
- **Technical reports and theses.** Given that they neither ensure an in-depth peer review, nor are widely validated by the scientific community, technical reports and theses were excluded. In the case of theses, it was assumed that, if they offered some contributions, these originated from other publications, explaining why they were excluded from this study.

The selection of studies involved a process of analysis composed of three filters which were intended to select the most appropriate results, as the probability of retrieving inappropriate studies would have been high. Figure-2 briefly describes what was considered in each filter. Additionally, the figure presents the number of written pieces that were obtained after the application of each filter.

Data extraction The method for data extraction was designed to extract all the information needed to answer the research questions. The following information was extracted from each article: *authors; source; conference/journal; publication year; summary; a brief opinion regarding its strengths and weaknesses and the study's objectives*. It was decided that when several studies were reported in the same paper, each relevant study would be treated separately. However, this situation did not occur.

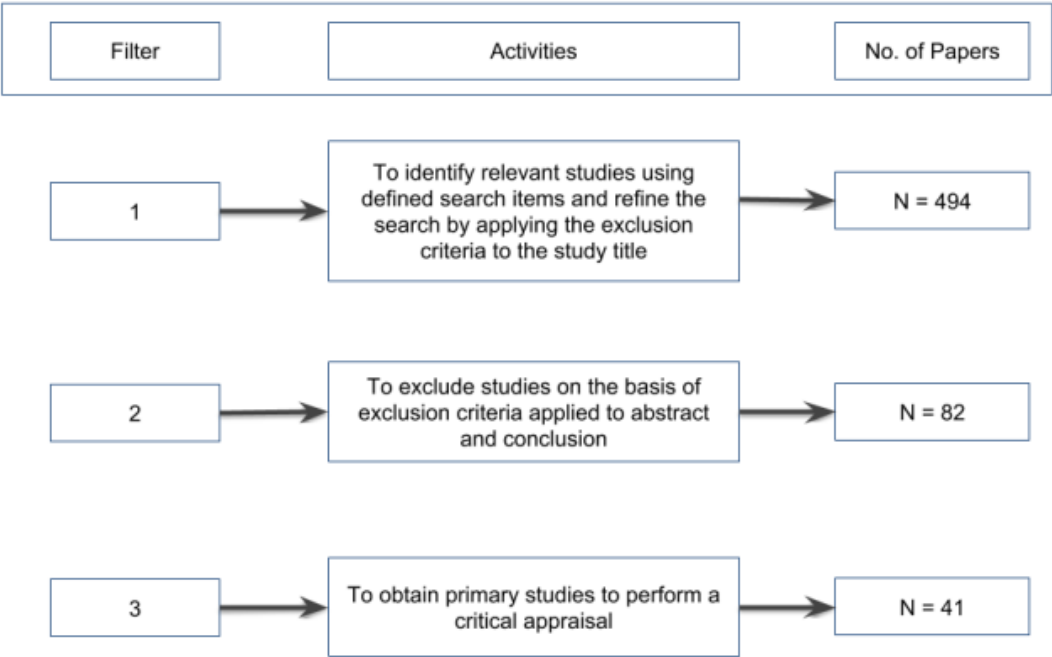


FIGURE 2 Stages of the article selection process.

3.3 | Third Phase: Outcomes

This section describes the classification framework (or schema) used and the results of data extraction. Once the framework was defined, the relevant studies were ordered according to the framework. The output of this stage is the study map, which is presented at the end of this section together with the discussion.

3.3.1 | Classification Schema

This study followed the systematic process shown in Figure-3. We used the idea proposed by Petersen et al. (2008) in order to categorize studies in facets. In our case, we defined two facets. One facet examined the type of research, and the other arranged the topic in terms of the research questions. The search for key concepts is one way of reducing the time necessary to develop the classification schema, and ensure that the schema considers all articles. This search is conducted in two steps. First, the reviewers read the article summaries. Then they look for keywords and concepts that reflect the contribution of the article, and subsequently identify the research context. Once the above has been completed a set of key words is obtained from different summaries, which are combined in order to achieve a high level of understanding with respect to the nature and contribution of the article, in addition to its reported use of Semantic Web technologies.

In this study, three aspects of interest were considered for each article: (i) the structure of the objective topic (i.e. the use of Semantic Web technologies in DW development); (ii) the names of the Semantic Web technologies used, which were derived from the list of key words, and (iii) the stage of DW development at which the said technology was

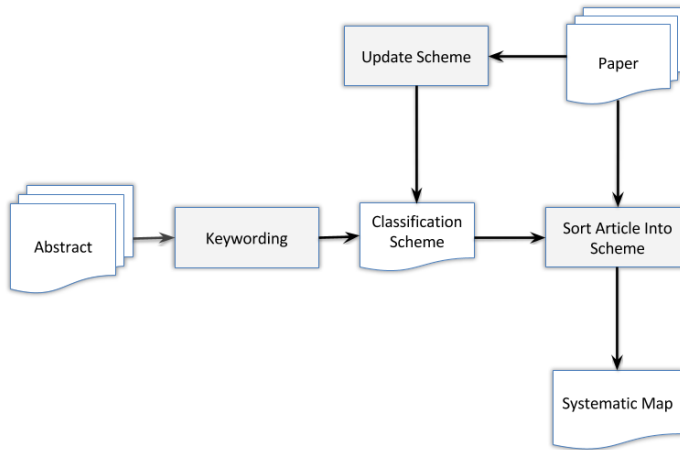


FIGURE 3 Creation of the Classification Schema.

used. Next, a search was conducted for evidence of improvement of the DW development process, which was expressed in terms of evaluations and/or validations of the results. For the first facet related to the type of research, analyzed studies were categorized according to the classification of research proposed by Wieringa et al. (2005), which was used with the intention of determining the type of research proposed. This classification identified six types of articles, defined as follows:

1. **Research Evaluation Techniques** are implemented in practice and an evaluation is conducted. This means that the way in which the technique is implemented (implementation of the solution) is presented, as are the consequences of the implementation in terms of its benefits and inconveniences (evaluation of the implementation). This also includes identifying problems within industry.
2. **Proposed Solution.** A solution is proposed for the problem. The solution could be a new technique, or the extension of an already existing one. The potential benefits and the applicability of the solution are shown through a small example or a short, but effective argument.
3. **Research Validation.** The techniques researched are recent and have still not been implemented in practice. The techniques used are, for example, experiments or, in other words, work performed in a laboratory.
4. **Philosophical Study.** These articles establish a new way of analyzing current topics of interest by structuring the area in a taxonomic way or through a conceptual framework.
5. **Opinion Articles.** These articles express the personal opinion of someone with respect to whether a certain technique is good or bad, or in regards to how it should be performed. These opinions are not based on related studies or research methodologies.
6. **Personal Experience Study.** Articles based on personal experience that explain what and how something was done in practice.

3.4 | Results

In this section, we present the elements that were found through the analysis of the selected articles with the intention of answering the research questions. The evidence obtained from the data extraction process is highlighted, as are the particularities found. Prior to the above, a classification of the collected studies analyzed is included. The final list of primary studies used in this study is shown in Table-5 and Table-6.

3.4.1 | Classification of the Studies

A classification of the obtained studies is presented below. The classification process was intended to determine relevant aspects of the set of collected studies, particularly the primary studies that were ultimately selected. First, the distribution by data sources (e.g. Journal, Conference, Book Chapters and others) is given. Next, the distribution of articles according to the place they were published is shown. Finally, the distribution of articles by publication year is presented. At this point it is worth presenting a more complete view of the distribution of publications in order to better characterize what has been published regarding the interrelationship between DWs and the Semantic Web. Significant numbers of studies were published in journals, followed by significant numbers of publications that were not submitted for peer review, (e.g. technical reports, opinion articles) though they can still be found in some Web search engines. The above shows that attempts to address the topic have been based on the opinions of experts, or on studies that require more development. Indeed, one can easily find a plethora of writings authored by consultants, technological companies, and postgraduate students³. Given what was stipulated in section 3.2.2 as exclusion criteria - that a research article must be published by high level editors (McGregor, 2002) - such studies were excluded from the final list, as is illustrated in Figure-4. This presents a complete distribution of the studies and their place of publication. All of the publications initially found, along with the results of applying the defined filters, were considered.

In terms of publication year, the analyzed studies covered a range of years beginning with the year 2000 up to the year 2017. The distribution of primary studies by publication year is illustrated in Figure-5. According to the figure, the publications analyzed first began to take into account the relationship between DWs and the Semantic Web in the year 2003. It was also evident that the area of the Semantic Web was of considerable interest to researchers in 2009. Interest in the Semantic Web began in 2007 and peaked in 2009. Since that year, interest in the Semantic Web has decreased substantially. However, interest began to increase again in 2012, and subsequently declined between 2014 and 2017. In 2017, zero publications were presented showing relations between the process of building a DW and the Semantic Web. This is likely to be due to the fact that there is a growing migration of interest from DWs to Linked Open Data. In 2016 there was some research work showing the relationship between DWs and Linked Open Data (LOD). For instance, Ravat et al. (2016). The trend of DW development proposals employing Semantic Web technologies in the Business Domain has declined. Considering that the increase in the number of publications is an indicator of interest within the scientific community regarding a particular topic, this result shows that the interest of the scientific community has probably shifted to other related areas such as: Linked Open Data (LOD) and Big Data, among others. In fact LOD have become one of the most important sources of information, allowing the enhancement of business analyses, based on warehoused data, with external data. However, DWs do not directly cooperate with LOD datasets because of the differences between data models. Moreover, in the Business domain, companies need to overcome the inevitable tension between the value they traditionally assign to proprietary data, and the value of opening up to be able to adopt LOD for sophisticated data analysis. Thus an important challenge is to determine how effective a particular solution for building a DW, based on Semantic Web technologies has been, as well as evaluating its impact on

³(e.g. <http://www.dataversity.net/down-with-the-data-warehouse-long-live-the-semantic-data-warehouse/>)

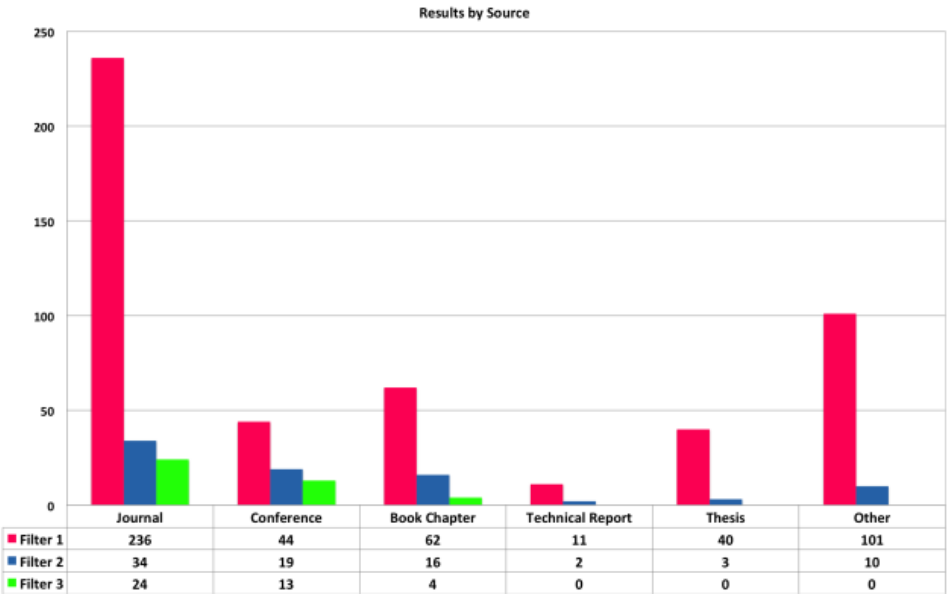


FIGURE 4 Place of Study Publication For All Filters.

the industry, if it has been implemented.

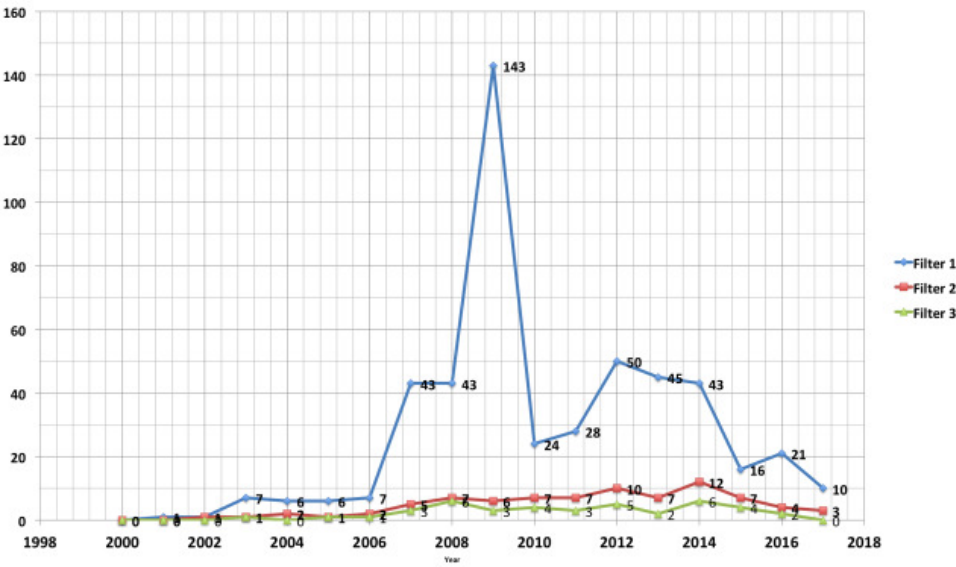


FIGURE 5 Primary source distribution by publication year.

3.4.2 | Research Questions

RQ1 - Where were the research studies published?

The 41 publications analyzed and published between the years 2002 and 2017, were distributed as follows: 24 were published in journals, 13 were published in conference proceedings, and 4 were published as book chapters.

Figure-6 presents a map of the distribution of publication place and year. The results revealed that the highest number of studies was concentrated in journals, with the year 2015 having the highest number of all.

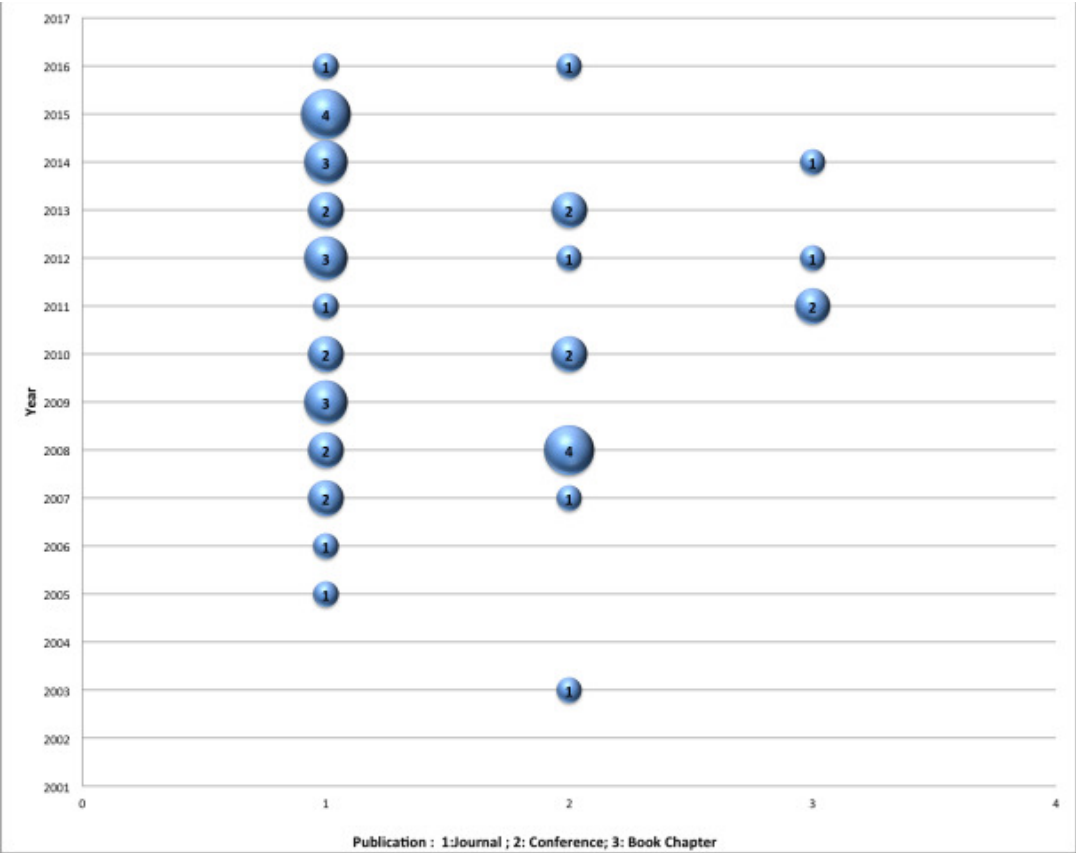


FIGURE 6 Distribution map of primary studies by publication venue and year.

RQ2 -What Semantic Web technologies are being used in DW development?

The analysis of the primary studies found two main types of Semantic Web technologies that were used in data storage development:

- (a) ontologies and vocabularies; and (b) languages for representing ontologies

In the primary studies examined, the individual and combined use of these Semantic Web technologies were reported. Some studies refers only to ontologies as a concept and others used ontologies along with a language such as RDF or OWL. The attention of most researchers was concentrated on the use of ontologies, and subsequently on the

use of the OWL together with the ontologies. The primary studies that used ontologies included: Priebe and Pernul (2003), Sell et al. (2005), Skoutas and Simitsis (2006), Skoutas and Simitsis (2007), Niemi et al. (2007), Nazri et al. (2008), Salguero et al. (2008a), Sell et al. (2008), Diamantini and Potena (2008), Salguero et al. (2008b), Spahn et al. (2008), Skoutas et al. (2009), Niinimäki and Niemi (2009), Nebot et al. (2009), Simitsis et al. (2010), Romero and Abelló (2010), Khouri and Ladjel (2010), Jiang et al. (2010), Romero et al. (2011), Bergamaschi et al. (2011), Martin et al. (2011), Nebot and Berlanga (2012), Selma et al. (2012), Sell et al. (2012), Thenmozhi and Vivekanandan (2012), Thenmozhi and Vivekanandan (2013), Ali et al. (2013), Khouri et al. (2014), Elamin and Feki (2014), Talebzadeh et al. (2014), Ramasamy and Palanivel (2014), Samuel (2014), El Sarraj et al. (2014), Di et al. (2015), Abelló et al. (2015), Lamolle et al. (2015) and Steiner et al. (2015).

The main objective of the analyzed studies was to map data through an ontology in order to manage unstructured, or slightly structured data. Some approaches assumed that an ontology is easily constructed, in the event that it was not provided in advance. However, in many cases finding an appropriate ontology for a specific domain is not a trivial task.

RQ3 - During what stages in the DW development process are these technologies being used?

The analysis of the primary studies revealed four stages of DW development: a) Data acquisition; b) Transformation and loading; c) Data integration; d) Data Access, that were mentioned in the primary studies, as illustrated in Figure-7. All of the primary studies considered at least one stage of DW development. From the studies, two approaches were identified. One focused on the automation of multi-dimensional design based on the use of ontologies. For example, Niinimäki and Niemi (2009) used ontologies to populate OLAP cubes. The other approach examined methods that focused on analyzing large amounts of Semantic Web data. For instance, Nebot and Berlanga (2012) proposed a semi-automated method for the extraction of semantic data, and for storing it in a multi-dimensional database, for analysis using traditional OLAP techniques. Figure-7 shows that the majority of primary studies are concentrated in the stage of Transformation and Loading, and in the use of ontologies, rather than in the use of some language such as RDF or OWL.

Studies that focused on the stage of Transformation and Loading and the use of ontologies included: Skoutas and Simitsis (2006), Skoutas and Simitsis (2007), Nazri et al. (2008), Salguero et al. (2008a), Sell et al. (2008), Salguero et al. (2008b), Spahn et al. (2008), Skoutas et al. (2009), Niinimäki and Niemi (2009), Simitsis et al. (2010), Jiang et al. (2010), Romero et al. (2011), Bergamaschi et al. (2011), Martin et al. (2011), Thenmozhi and Vivekanandan (2012), Ali et al. (2013), Elamin and Feki (2014), Ramasamy and Palanivel (2014), Samuel (2014), Di et al. (2015), Steiner et al. (2015), Ravat et al. (2016), and Nebot and Berlanga (2016).

The following primary studies focused on the stage of Transformation and Loading, the use of ontologies and some language, such as RDF or OWL:

Skoutas and Simitsis (2006), Skoutas and Simitsis (2007), Nazri et al. (2008), Salguero et al. (2008a), Salguero et al. (2008b), Niinimäki and Niemi (2009), Simitsis et al. (2010), Romero et al. (2011), Thenmozhi and Vivekanandan (2012), Steiner et al. (2015), Ravat et al. (2016), and Nebot and Berlanga (2016).

RQ4 - What type of research study was conducted?

The classification of research approaches proposed by Wieringa et al. (2005) was used to relate the type of research approach of a particular study to Semantic Web technologies. The greatest number of articles concentrated on "solution proposal" approaches, followed by "research evaluation." The systematic search did not find any articles that documented experiences with respect to industrial applications of Semantic Web technologies. Thus this area represents an important challenge for future research. Examples of primary studies that show the relationship between solution proposals and the use of ontologies included: Priebe and Pernul (2003), Sell et al. (2005), Skoutas and Simitsis

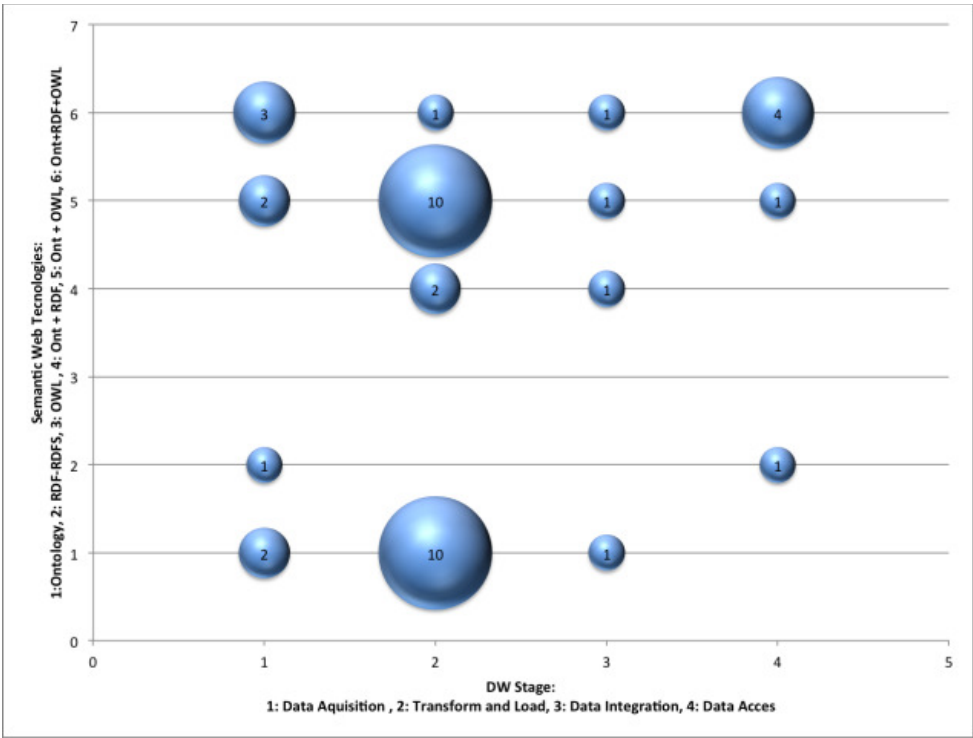


FIGURE 7 Visualization of a systematic map relating the stages of DW development to Semantic Web technologies using a bubble chart.

(2007), Niemi et al. (2007), Salguero et al. (2008a), Sell et al. (2008), Diamantini and Potena (2008), Salguero et al. (2008b), Spahn et al. (2008), Skoutas et al. (2009), Niinimäki and Niemi (2009), Simitsis et al. (2010), Khouri and Ladjel (2010), Jiang et al. (2010), Bergamaschi et al. (2011), Martin et al. (2011), Selma et al. (2012), Sell et al. (2012), Thenmozhi and Vivekanandan (2012), Thenmozhi and Vivekanandan (2013), Elamin and Feki (2014), Talebzadeh et al. (2014), Ramasamy and Palanivel (2014), Samuel (2014), El Sarraj et al. (2014) and Di et al. (2015). Primary studies that showed the relationship between solution proposals, the use of ontologies and a language included: Niemi et al. (2007), Khouri and Ladjel (2010), Selma et al. (2012), Sell et al. (2012), Thenmozhi and Vivekanandan (2013) Talebzadeh et al. (2014), Ravat et al. (2016), and Nebot and Berlanga (2016). The results of the classification are presented in Figure-8, which presents a summary of the type of article alongside the use of Semantic Web technologies.

RQ5 - Is there evidence of improvement due to the use of Semantic Web technology?

In order to analyze the evidence of improvement presented in primary studies, the classification framework proposed by Shaw (2003) was used. This framework classifies a) **the type of research results** and b) **the validation type** found in the studies. The classification by type of research results, as well as validation type, allowed the evidence of improvement resulting from the use of Semantic Web technologies to be evaluated objectively. On one side, it can be established that the type of result (ranging from the most formal to the least formal), and the validation type (ranging from a formally validated result to a simple declaration of improvements), both constitute evidence of a proposal's value. For example, if the proposal was a report in which only declarations of improvements were made, it was not considered to contain

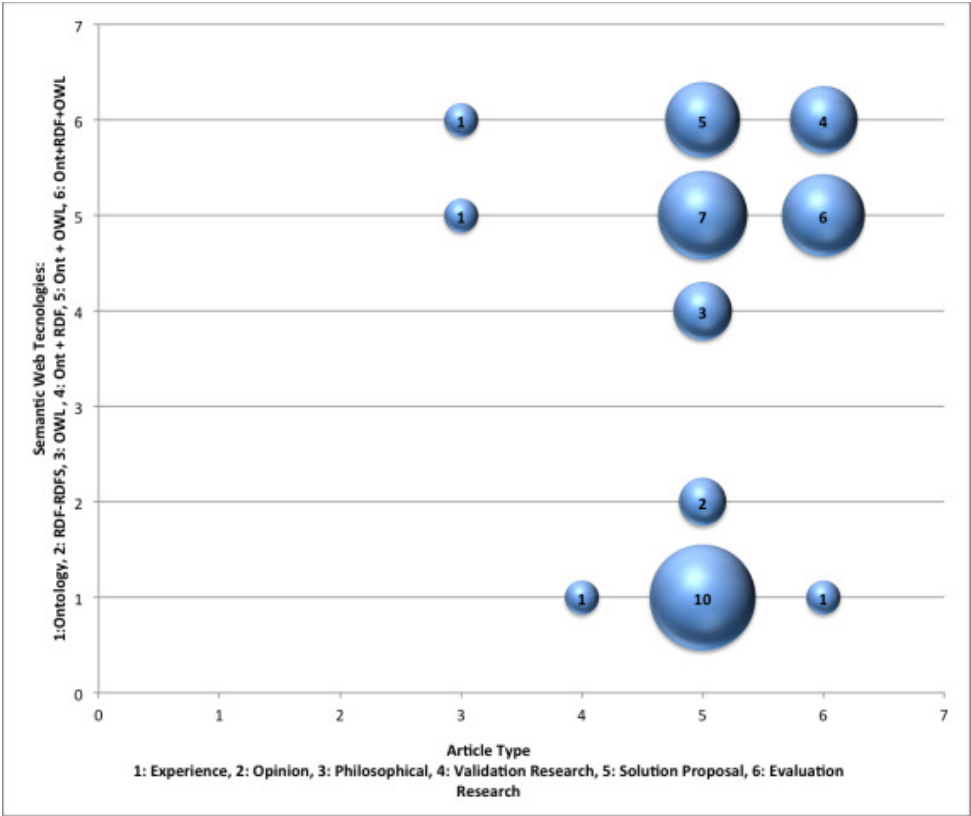


FIGURE 8 Visualization of a systematic map relating article type to the use of Semantic Web technologies using a bubble chart.

evidence of improvement according to this framework. In contrast, if the proposal presented was a rigorous analysis of a formal or empirical model, its results were considered to show valid evidence of improvement. In addition to the above, the text of each study when analyzed was searched for some explicit mention of the results indicating an improvement due to the use of Semantic Web technologies. Research results may include procedures or techniques for development or analysis; models that generalize examples, specific tools, or solutions for particular systems. Table-3 describes the type of research results that were reported in the studies analyzed, which were classified according to the framework proposed by Shaw (2003).

For the type of result, the greatest concentration of primary studies was found in the RI-6 category (Specific Solution). This is shown in Figure-9. This means that the majority of examined proposals responded to specific solutions, or implementations of a prototype that used some type of Semantic Web technology. On the other hand, the research studies also used different classes of evidence to support their results. Thus it is necessary to select a form of validation that is appropriate for the type of result and the methods used to obtain the result. For example, a formal model should be supported by rigorous testing and not by examples of use. Table-4 presents the possible evidence (validation) types found using the framework proposed by Shaw (2003). Among the primary studies analyzed, the most common validation types were RV-2 (Evaluation) and RV-4 (Example). This is shown in Figure-10. In the case of validation type

TABLE 3 Types of Research Results. Shaw Shaw (2003).

ID	Type of Results	Example
RI-1	Procedure or Technique	new or better ways of realizing a given task, such as design, implementation, maintenance, measurement, evaluation or selection of alternatives.
RI-2	Qualitative or Descriptive Model	Taxonomic structure for a problem area; architectural style, framework or design pattern; analysis of an informal domain, checklist or well-argued generalizations.
RI-3	Empirical Model	Predictive, empirical model based on observed data.
RI-4	Analytical Model	Structural model that permit formal analysis or automated manipulation.
RI-5	Notation or Tool	Software tool that implements a technique or a formal language to support a technique or model.
RI-6	Specific Solution, Prototype, Response or Judgment	Solution to a problem that shows the application of some principles - such as design, prototype or complete implementation and careful analysis of a system or its development, resulting in a specific analysis, evaluation or comparison.
RI-7	Report	Interesting observations. Golden rules. It is not sufficiently general or systematic to be considered a descriptive model.

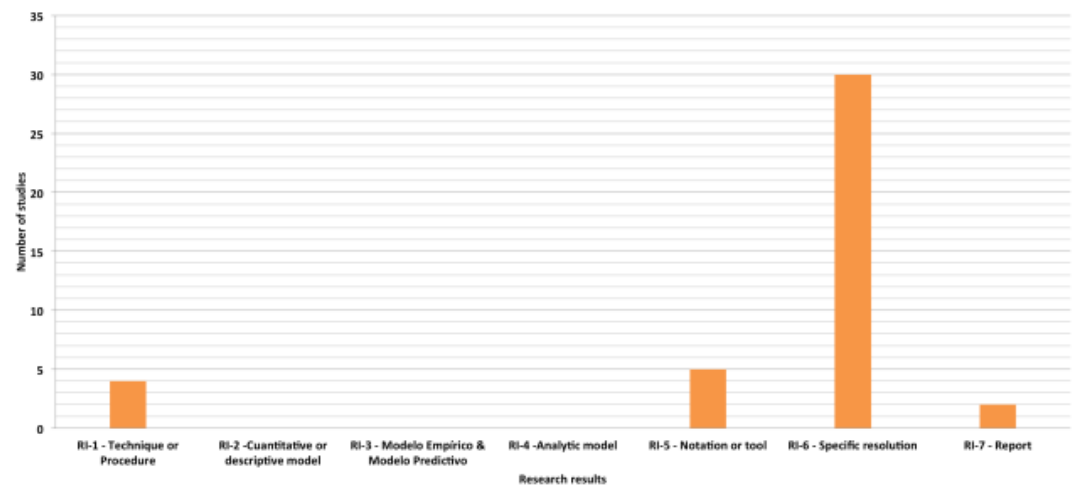


FIGURE 9 Type of research results.

RV-2 (Evaluation), the authors created a proposal and then evaluated the obtained results. For validation type RV-4 (Example), the authors presented a description of how a Web technology is used and then applied it to an example. The results of the systematic map on this topic are illustrated in Figure-11.

Regarding the systematic map of the relationship between the type of result and validation type, the majority of primary studies focused on specific proposals, or solution prototypes in which well-argued frameworks and generalizations

TABLE 4 Validation Type. Shaw (2003)

ID	Validation Type	Example
RV-1	Analysis	Provides a rigorous analysis of a formal model, empirical model or controlled experiment with rigorous testing, data in controlled situations, or experiments with abundant and significant statistical results, respectively.
RV-2	Evaluation	Includes studies of feasibility, pilot projects for descriptive, qualitative or empirical models.
RV-3	Experience	The results have been used in real examples and the evidence of their effectiveness, usefulness and accuracy is given by qualitative, empirical models, tools or techniques through narrative descriptions, statistical data or predictions that fit the current data.
RV-4	Example	Example of a given technique or procedure is used. A fragment of reality (simplified example of reality) can be convincing especially if it is accompanied by an explanation of why this simplified example retains the essence of the problem being solved. Small examples or examples found in books often fail to provide convincing validations (except for standard examples used as model problems in the area).
RV-5	Persuasion	Validation only by persuasion makes it difficult to consider an article as research study.
RV-6	Brazen assertion	No serious attempt is made to evaluate the results.

were proposed. Such studies were considered evaluations of the proposal. Studies that proposed specific solutions and evaluated their proposals included: Skoutas and Simitsis (2007), Spahn et al. (2008), Niinimäki and Niemi (2009), Romero and Abelló (2010), Jiang et al. (2010), Romero et al. (2011), Selma et al. (2012), Thenmozhi and Vivekanandan (2012), Etcheverry and Vaisman (2012), Samuel (2014), Lamolle et al. (2015), Nebot and Berlanga (2016) and Ravat et al. (2016). Studies that proposed specific solutions that were evaluated using an example included: Skoutas and Simitsis (2006), Niemi et al. (2007), Nazri et al. (2008), Salguero et al. (2008a), Sell et al. (2008), Salguero et al. (2008b), Bergamaschi et al. (2011), Martin et al. (2011), Sell et al. (2012), Thenmozhi and Vivekanandan (2013), Ali et al. (2013), Khouri et al. (2014), Elamin and Feki (2014), Talebzadeh et al. (2014), El Sarraj et al. (2014), Di et al. (2015), and Steiner et al. (2015). The analyzed studies did not show evidence of comparisons being made to other similar proposals. As a result, it is not possible to determine if the results obtained were better in comparison to the results of similar proposals.

4 | ANALYSIS OF THE RESULTS AND DISCUSSION

Of the results described in Section 3.3, the following elements stood out:

1. The most commonly used Semantic Web technologies to build DWs in the Business Domain were ontologies and markup languages, such as RDF or OWL. However, the use of ontologies, and the markup languages, has been limited to representing terms and improving model equivalence as well as specifying registry structures for databases, according to the analysis of the studies. There were no uses that showed greater complexity than the representation of conceptualizations that describe aspects of DW domain information. Such an approach would

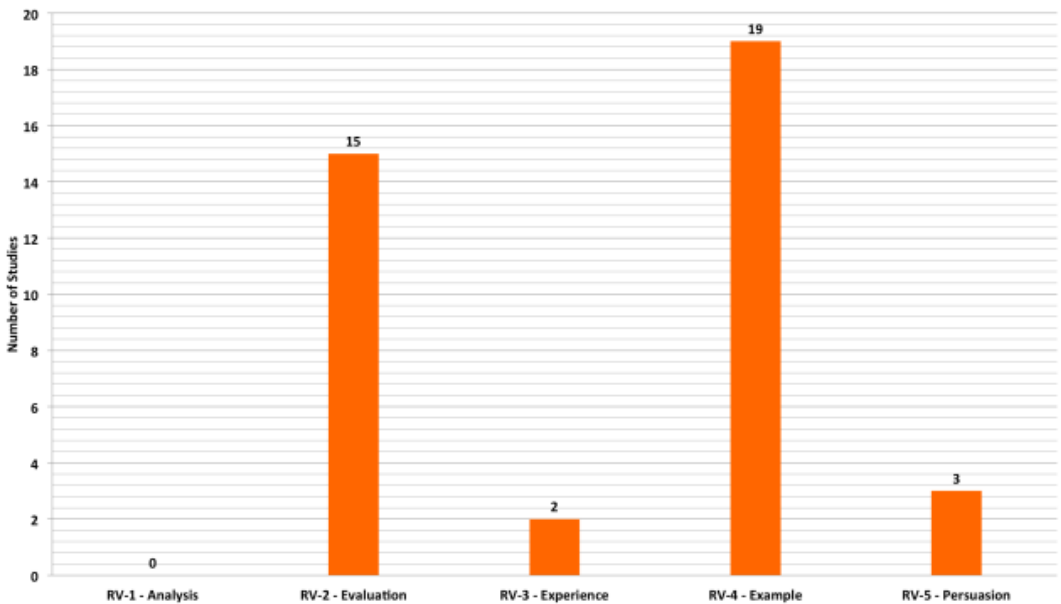


FIGURE 10 Validation Type.

enable, for example, the creation of recommendations that support the DW development process.

2. A total of 58% of the primary studies analyzed were published in journals, which indicates a higher degree of maturity among these proposals. This can be explained by the fact that both DW development and the Semantic Web are areas that have shown important advances and growth in their maturity levels over the time period examined. Thus the combination of both areas in one proposal assumes a higher maturity level.
3. Of the primary studies analyzed, there is a lack of sufficient evidence showing that the results of the studies had been transferred successfully to industry. The majority of studies analyzed proposed solutions intended to improve certain aspects of DW development. However, concrete support elements were not derived from these studies, such as standardized models or publicly available tools.
4. The results indicate an increasing lack of interest in applying Semantic Web tools to DW development. Unfortunately, there is also no evidence of mature solutions being transferred to industry. There are two possible explanations for the above: 1) The topic stopped being of interest to researchers or 2) New topics of interest are attracting the attention of researchers. We believe that the second option is the more viable, as the emergence of new topics such as: Linked Open Data, Cloud Computing, Big Data and Cognitive Computing, among others, have caused the focus of researchers to become centered on other topics of current interest, such as: security, cloud computing, efficient processing and new Web architectures, and so on. Since Linked Open Data has attracted the interest of researchers, it seems the Semantic Web community continues to develop the field (i.e. tools, technologies) by following the Linked Data approach.
5. Considering that some Semantic Web technologies were present (e.g. ontologies and markup languages) in the analyzed studies, it is clear that there is a need for greater development of the semantic aspects of DW development. The creation of metamodels, the inclusion of semantic rules that support DW development and the intensive use of natural language, are fundamental to improving aspects of DW creation. A combination of BI technologies with the

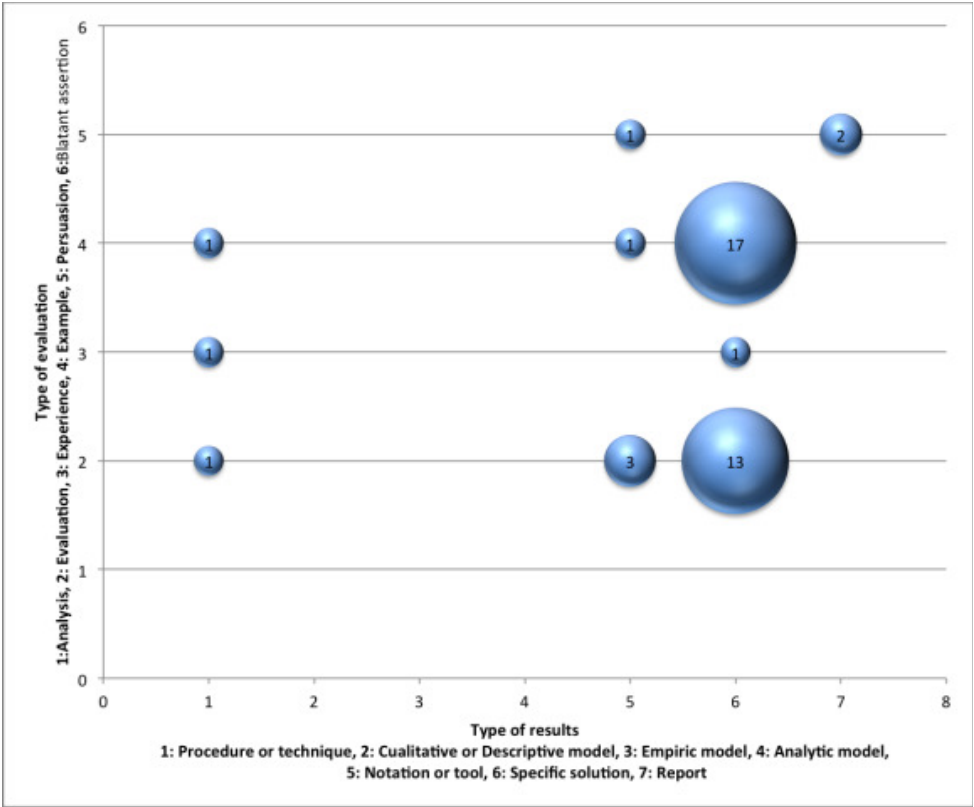


FIGURE 11 Visualization of a systematic map relating the type of result to the validation type using a bubble chart.

Semantic Web will both enhance BI analysis with web data, and allow the analyzing of Semantic Web data through BI tools. In fact, business can get benefits by using Semantic Web technologies for data analytics, data integration, and knowledge discovery. However, combining BI with the Semantic Web is not a trivial task due to the scalability, complexity and heterogeneity of Semantic Web data.

Currently, a DW is widely used as a consistent and integrated data repository in BI systems. According to Ravat et al. (2016), additional information coming from outside an organization, mostly found on the Web, should also be included in analyses to provide multiple perspectives for decision-makers, within today's highly competitive business context. Linked Open Data has been introduced as a promising paradigm, based on the Semantic Web, for opening up data which integrates published datasets, because it facilitates data integration on the Web (Bizer et al., 2009). Similar to Data Warehousing approaches, Linked Open Data can be prepared to enable sophisticated data analysis. In 2015, the report entitled: *Gartner's Hype Cycle for Enterprise Information Management, 2015*⁴ has put Linked Data into the trough of disillusionment stage. According to Garner's Hype Cycle Research Methodology the phase *Trough of Disillusionment* is one of five key phases of a technology's life cycle, which is defined as follows: "Trough of Disillusionment: Interest dies out as experiments and implementations fail to deliver. Producers of the technology shake out or fail. Investments

⁴<https://www.gartner.com/doc/3096424/hype-cycle-enterprise-information-management>

continue only if the surviving providers improve their products to the satisfaction of early adopters".⁵ Some issues that impede its adoption in the business domain are: (i) *Cost-benefits trade off*. Businesses are not clear on how to get benefits, value, and extra earnings from adopting Linked Open Data. One problem to overcome is the tension between the value they traditionally assign to proprietary data, and the value of opening up. (ii) *Reliability and Data quality*. Over the last few years, Linked Data has developed into a large number of datasets, with open access from several domains leading to the linking open data (LOD) cloud. Similar to other types of information such as structured data, Linked Data suffers from quality problems such as inconsistency, inaccuracy, being out-of-date, and incompleteness, which are frequent, and imply serious limitations to the full exploitation of such data. (iii) *High computational costs for queries*. To get information from the global Web of data, the query language SPARQL is used. The problem is that queries can require jumping from one server to another. Thus joint operations can be more costly than in the relational data base (iv) *Complexity and heterogeneity of the technologies*. Currently there are many different and complex semantic technologies such as frameworks, RDF, Turtle, Microdata, language query, amongst others. New technologies will appear, and some others will survive. The challenge is, therefore, how to use and integrate them.

In 2017, Gartner's Hype Cycle for Analytics and Business Intelligence ⁶ classified Logical Data Warehouse as a future useful technology for BI. Logical Data Warehouse is defined as an architectural layer that sits at the top the usual data warehouse (DW) store of persisting data. The logical data warehouse complements the traditional core warehouse (and its primary function of a priori data aggregation, transformation, and persistence) with functions that fetch and transform data, in real time (or near to it), thereby instantiating non-persistent data structures, as needed.

4.1 | Main findings of the study

Given the above, three important challenges can be identified regarding the use of Semantic Web technologies in DW development in the BI domain.

1. **New uses and applications based on ontologies.** Using ontologies to represent information aspects and not just terminological expressions, allows for greater capabilities to be included in applications by creating tools with greater capacities, such as advisers and intelligent assistance systems to be used in DW development.
2. **Measuring solution effectiveness and impact.** Another important challenge is in determining how effective Semantic Web-based technologies have been, and what their impact has been on industry. This is important, as the ultimate objective of research in this field should be aimed at advancing existing knowledge in order to be able to apply these advances to industry. Aspects of cost, performance, quality assurance and DW design approaches have not been extensively considered.
3. **Inclusion of semantic aspects in addressing new paradigms.** Another important challenge is to consider semantic elements in remedying current difficulties in DW development and other technological challenges. For example, the minimal use of software tools for ETL opens a space for creating new tools, or improving existing ones with semantic aspects. In the same way, the emergence of Cloud Computing and Big Data involves assuming new considerations in DW design, which could be supported with tools or Semantic Web-based technologies.

⁵<https://www.gartner.com/technology/research/methodologies/hype-cycle.jsp>

⁶<https://www.gartner.com/doc/3096424/hype-cycle-enterprise-information-management>

4.2 | Review of Industrial Projects

This section reviews some industrial tools and projects within the broad topic of DWs and the Semantic Web. First, we review briefly some existing industrial tools, and then we focus on reviewing industrial projects.

Unfortunately, research on the intersection between DWs and the Semantic Web has not moved towards industrial tools (Abelló et al., 2015). Apart from software that provides functionality to derive Linked Data from existing databases, such as Virtuoso⁷ (from OpenLink vendor) or Apache Jena⁸, there are few industrial tools that fully consider the potential of the Semantic Web together with DWs. Interestingly, most of them are OLAP tools, such as TARGIT Decision Suite (Middelfart and Pedersen, 2011)⁹ which uses extended semantics to extend the traditional multidimensional model with semantic associations between measures and dimensions; or the OpenCube Toolkit (Kalampokis et al., 2014)¹⁰ which provides facilities for analysing RDF data as multidimensional cubes. On the other hand, graph databases (such as GraphDB from OntoText¹¹) have emerged as important tools for exploring how interlinking data could be useful in developing a new generation of DWs that unleashes the potential of unstructured data, and for inferring new knowledge. **The graph paradigm differs from other database models that make it interesting for enterprises, since it naturally model real-world concepts and relations according to Bryce Merkl Sasaki¹². Furthermore, the growing prevalence of graph databases in the industrial and marketing arenas will require considering this paradigm as a driving force for storing and analyzing data. Consequently, approaches that use graph databases in the development of data warehouses, and specifically, those considering enterprise knowledge graphs, are a must. In fact, a Forrester Research survey¹³, states that “51 percent of global data and analytics technology decision makers either are implementing, have already implemented, or are upgrading their graph databases.”. According to Yu Xu¹⁴, which is the founder and CEO of TigerGraph, the world’s first native parallel graph database, graph databases (i) are simple and provide more natural data modelling, (ii) provide better, faster queries and analytics, as well as (iii) support flexibility for evolving data structures. These key benefits of graph databases might benefit the building of data warehouses.**

Finally, we have to highlight DWs for Linked Data such as CM-Well (Bennett et al., 2017) from Thomson-Reuters¹⁵ that allows a company to manage large volumes of linked data and to warehouse its knowledge graph.

Importantly, some efforts have been made in research projects, e.g. as funded by the European Union. We have conducted a search of projects on DWs and the Semantic Web in the search engine provided by CORDIS (Community Research and Development Information Service, which is the European Commission’s primary public repository and portal for disseminating information on all EU-funded research projects and their results)¹⁶. We have found the following projects:

- GeoKnow¹⁷: this project aims to make the Web a data source for geospatial knowledge, i.e., moving from Geographic Information Systems (GIS) to a Web of interlinked Geographic Knowledge Systems (GKS). GeoKnow focuses on achieving the evolution of the Web from a medium for information exchange, to a medium for (spatial) knowledge

⁷<https://virtuoso.openlinksw.com/>

⁸<https://jena.apache.org>

⁹<https://www.targit.com/en/software/decision-suite/>

¹⁰<http://opencube-toolkit.eu/>

¹¹<https://ontotext.com>

¹²<https://neo4j.com/blog/why-graph-databases-are-the-future/>

¹³Forrester Research, Forrester Vendor Landscape: Graph Databases, Yuhanna, 6 Oct. 2017

¹⁴<https://www.datanami.com/2017/11/30/look-graph-database-landscape/>

¹⁵<https://github.com/thomsonreuters/CM-Well>

¹⁶<https://cordis.europa.eu/projects>

¹⁷<http://geoknow.eu>

integration.

- SEEK¹⁸: this project aims to enrich semantically trajectory knowledge discovery, i.e., taking advantage of data collected from personal devices such as mobile phones and other location-aware devices. SEEK focuses on extracting behavioral patterns through a knowledge discovery process, where positioning data collected from mobile devices are first transformed into semantically enriched trajectory data, and then stored in a data warehouse and analysed with OLAP operations that allow the summarization of the trajectory features.
- OpenCube¹⁹: this project focuses on processing RDF data as OLAP cubes, i.e., multidimensional data represented as RDF and structured according to the RDF Data Cube ontology.

5 | RESEARCH CHALLENGES

Designing a DW is a complex task that often forces designers to acquire wide knowledge of the domain, thus requiring a high level of expertise and so becoming a prone-to-fail task. Taking into consideration the results of our systematic mapping, in the following section our findings and experience are combined with some proposed solutions for the issues that have been discovered in the systematic mapping. Importantly, based on our experience, we have detected a set of situations in which we believe that the use of some kind of Semantic Web technology will improve the design of DWs.

The main focus of research using Semantic Web technologies in the development of DWs is in the development of ETL processes, as stated in our systematic mapping, however there are other important challenges in which Semantic Web technologies may formalize the use of specific-domain knowledge for data warehousing, namely:

- Reusing expert knowledge from different domains.
- Enriching specific metadata by completing definitions and annotating their semantics.
- Enabling metadata interchange among repositories.
- Populating the designed databases from public data sources.
- Empowering data integration, and analysis.
- Automatizing reasoning on metadata.
- Validating data instances and models.

Other more elaborated proposals for using Semantic Web Technologies in the development of DWs are described in the following points:

- Reconciling information requirements and data sources. Due to the special idiosyncrasy of DWs, it should not only be the information requirements of decision makers that are considered for multidimensional design, but also a second driving force: data sources, which need to be reconciled with those information requirements. Requirements analysis for multidimensional modelling aims to elicit the information needs of decision makers in order to deploy a DW that will satisfy their expectations. Furthermore, internal and external data sources should be taken into account and matched with information requirements, since these data will populate the DW. To do so, Semantic Web technologies can be used to match formally the data sources with information requirements in the early stages of development, in order to make the not-so-obvious correspondence between information requirements and their counterparts in data sources. However, according to our systematic mapping, these topics

¹⁸<http://www.seek-project.eu>

¹⁹<http://opencube-project.eu/>

have not been investigated so far, and more research is required.

- **Incompleteness in multidimensional models.** Semantic Web technologies may enrich a multidimensional model in aspects that have not been taken into account during requirement analysis or reconciliation of data sources. For example, “public” Linked Open Data sources can be used to complete the unsupported elements within a dimension hierarchy. Moreover, taxonomies such as the Computing Classification System (CCS²⁰) can be also used as additional data sources for designing multidimensional models. CSS provides a classification for topics on computations that can be, and usually is, used for classifying computer literature. For instance, there is a taxonomy for Information Systems (H) where, e.g., Database Management (H.2) topics are classified. This taxonomy can also be used for completing aggregation hierarchies, whenever a DW requires a dimension on computer-related topics.
- **Data types of measures of DWs.** Measures of analysis described in DWs by multidimensional models in fact do not aid either designers or end-users to gather important details, such as their units or scale. Therefore, more research in the semantic description of measures in multidimensional models (addition of units, magnitudes, scales etc.) should be conducted in order to aid designers to annotate the identified measures.
- **Semantic-aware summarizability.** DWs use statistic functions to aggregate measures to provide data analysis. The application of aggregation functions depends on which sort of measure and aggregation criteria are involved according to the compatibility given by some additivity constraints. In particular, type compatibility states the compatibility of category attributes (aggregation levels), summary attributes (measures) and statistical functions (aggregation functions). Some taxonomies of measure have been proposed to deal with this compatibility, such as those in Lenz and Shoshani (1997) where measure can be classified as either a “flow”, a “stock”, or a “value-per-unit”. However, the proposed taxonomies for dealing with summarizability constraints lack a wider consideration of semantics. This means that designers do not have the necessary knowledge for working out what really is a given piece of information and thus, which constraints should be held. More research is required on this topic, to gain knowledge for a better understanding of the very nature of summarizability constraints, and of being able to understand why, when, and how to hold them.
- **Semantically-traceable models.** Mapping from multidimensional models to relational models is accomplished by structural matching. Since structure and semantics are closely related, whenever structure is manually changed, the source semantics are lost.

Automatic model transformations for the deployment of DWs are managed with the manual transformations that software engineers make, i.e., modifying the current label of a modeling element, removing some data, relating some elements, and so on. Traceability is the mechanism for propagating changes in a transformation chain, taking into account only the changes and not entire models. It solves the problem whenever automatic transformations are involved. However, modifying the label implies changing its semantics, or even worse, losing its meaning because of the misuse of semantics. This fact is easily shown in data sources where a data item may be labelled in a cryptic form (s0702) from which its semantics cannot be inferred (the sales from 2007 to 2009). Therefore more research on this topic will define novel approaches for using Semantic Web technologies, that may solve this drawback in the development of DWs by having semantically annotated multidimensional models.

6 | LIMITATIONS OF THE STUDY

Some concerns regarding the validity of our study are presented below:

²⁰<http://portal.acm.org/ccs.cfm?part=author&coll=portal&dl=GUIDE>

- *Consideration of Synonyms.* The terms used in the search string could have synonyms. As a result, it is possible that the search passed over some studies. On occasion, authors may only mention the term Semantics, without alluding to the Semantic Web, or use some particular technology of the Semantic Web, without making it explicit.
- *Quality of the Evaluation.* The quality of the selection and evaluation of the studies, as well as the weighting factor used to quantify each, could have failed to appropriately represent the importance of the selected sets. In order to mitigate this problem the quality attributes were grouped into sub-sets in order to facilitate future classifications and improve the selections.

7 | RELATED WORK

Several studies have focused on automating multidimensional design using semantic web artifacts such as existing ontologies. Some surveys have been published focusing on the above issue. For instance: Abello et al. (2015) goes on to survey the use of Semantic Web technologies for data modeling and data provisioning, including semantic data annotation and semantic-aware extract, transform, and load (ETL) processes. A characterization of DW/OLAP environments is presented, followed by an introduction to the relevant Semantic Web foundation concepts. The study describes the relationship of multidimensional (MD) models and Semantic Web technologies, including the relationship between MD models and Semantic Web formalisms. Wache et al. (2001) reviews the use of ontologies for the integration of heterogeneous information sources. Based on an in-depth evaluation of existing approaches to this problem. They discuss how ontologies are used to support the integration task. Chakraborty et al. (2017) present a survey of research in the area of data integration using ETL approaches including some semantic issues. The objective of this survey paper is to deliver a review of current approaches to data integration, various ETL tools that are already in use, linked data generation techniques, semantic-based approaches to ETL and data integration.

Although the above studies present a survey of research, they do not identify, appraise or synthesize research evidence from individual studies based on a strict protocol in order to summarize their results.

8 | CONCLUSIONS AND FUTURE WORK

The main motivation for this study was to investigate the current status of Semantic Web technology use in DW development. It was conducted using a Systematic Mapping process, a useful technique for identifying the areas where there is sufficient information to describe the current state of research, as well as those areas that require more information. Through this work, we determined the existence of two lines of study. The first focused on the automation of multi-dimensional design based on the use of ontologies, and the other examined methods that focus on analyzing large amounts of Semantic Web data. Additionally, we identified that the proposals for the use of Semantic Web technologies addressed all stages of DW development. However, some important aspects were not considered and, when they were considered, only a brief review was provided. For example, aspects such as cost and performance were barely discussed. In relation to the experiences of implementation in industry, we found that such experiences were very limited. Many of the cases discussed in the examined studies describe small projects, with results obtained in specific applications, which makes their reproduction in other contexts impractical due to a lack of details. This underscores the necessity of experimenting with more approaches that use Semantic Web technologies in industry, rather than in academia. Research on the intersection between DWs and the Semantic Web has not moved towards industrial tools. Some concerns regarding the validity of our study are presented below: (i) *Quality of the Evaluation.* The quality of the selection and evaluation of the studies, as well as the weighting factor used to quantify each, may

have failed to represent appropriately the importance of the selected sets. In order to mitigate this problem, the quality attributes were grouped into sub-sets to facilitate future classifications, and improve the selections. (i) *Consideration of Synonyms*. The terms used in the search string may have synonyms. As a result, it is possible that the search passed over some studies. On occasion, authors may only mention the term Semantics, without alluding to the Semantic Web, or use some particular technology of the Semantic Web, without making it explicit.

This study identified that the interest in Semantic Web technologies in the development of DWs in the business domain steadily decreased after the year 2012, and that the year of greatest interest was 2009. However, interest began to increase again in 2012, and subsequently declined from 2014 to 2017. This is likely to be due to the fact that there is a growing migration of interest from DWs to Linked Open Data and/or Logical Data Warehouse. This systematic map highlights some aspects that require greater research, such as aspects regarding cost, performance, quality assurance and DW design approaches that consider new technologies. Also, since the use of graph databases in the industrial and marketing arenas is increasing, it is reasonable to think that this technology might benefit the building of data warehouses in several key issues not only for data modeling but also for other complex issues, such as the consideration of knowledge graphs for evolving data warehouse design or analytic graphs for using advanced data analysis algebra. In our future work, we plan to combine the evidence identified in this study and industrial projects, in order to define hypotheses and theories that will form the basis for the design of new methods, processes and tools to integrate Semantic Web technologies into the DW development process.

9 | BEYOND SYSTEMATIC MAPPING

This section provides a description of each research article analyzed in the study. Table-5 and Table-6 present the list of analyzed studies. Table-7 to Table-13 present the summary and the contributions of each research article analyzed in the study.

ACKNOWLEDGEMENTS

This study was supported by the Universidad de La Frontera, Chile, PROY. DI15-0020.

REFERENCES

- Abelló, A., Darmont, J., Etcheverry, L., Golfarelli, M., Mazón, J.-N., Naumann, F., Pedersen, T., Rizzi, S. B., Trujillo, J., Vassiliadis, P. and Vossen, G. (2013) Fusion cubes: Towards self-service business intelligence. *Int. J. Data Warehous. Min.*, **9**, 66–88. URL: <http://dx.doi.org/10.4018/jdwm.2013040104>.
- Abelló, A., Romero, O., Bach Pedersen, T., Berlanga, R., Nebot, V., Aramburu, M. J. and Simitsis, A. (2015) Using semantic web technologies for exploratory olap: a survey. *Knowledge and Data Engineering, IEEE Transactions on*, **27**, 571–588.
- Abello, A., Romero, O., Pedersen, T. B., Berlanga, R., Nebot, V., Aramburu, M. J. and Simitsis, A. (2015) Using semantic web technologies for exploratory olap: A survey. *IEEE Transactions on Knowledge Data Engineering*, **27**, 571–588. URL: [doi.ieeecomputersociety.org/10.1109/TKDE.2014.2330822](http://ieeecomputersociety.org/10.1109/TKDE.2014.2330822).
- Abramowicz, W., Auer, S. and Heath, T. (2016) Linked data in business. *Business & Information Systems Engineering*, **58**, 323–326. URL: <https://doi.org/10.1007/s12599-016-0446-0>.
- Ali, A. A., Abdelrahman, T. A. and Mohamed, W. M. (2013) Using schema matching in data transformation for warehousing web data. *International Journal Information Technologies & Knowledge*, **7**.

- Bennett, D., Engelbrecht, J. and Landau, D. (2017) Cm-well: A data warehouse for linked data. In *Proceedings of the ISWC 2017 Posters & Demonstrations and Industry Tracks co-located with 16th International Semantic Web Conference (ISWC 2017)*, Vienna, Austria, October 23rd - to - 25th, 2017. URL: <http://ceur-ws.org/Vol-1963/paper518.pdf>.
- Bergamaschi, S., Guerra, F., Orsini, M., Sartori, C. and Vincini, M. (2011) A semantic approach to etl technologies. *Data & Knowledge Engineering*, **70**, 717–731.
- Berlanga, R., Aramburu, M. J., Llidó, D. M. and García-Moya, L. (2014) Towards a semantic data infrastructure for social business intelligence. In *New Trends in Databases and Information Systems* (eds. B. Catania, T. Cerquitelli, S. Chiusano, G. Guerrini, M. Kämpf, A. Kemper, B. Novikov, T. Palpanas, J. Pokorný and A. Vakali), 319–327. Cham: Springer International Publishing.
- Bizer, C., Heath, T. and Berners-Lee, T. (2009) Linked data - the story so far. URL: <https://eprints.soton.ac.uk/271285/>.
- Chakraborty, J., Padki, A. and Bansal, S. K. (2017) Semantic etl — state-of-the-art and open research challenges. In *2017 IEEE 11th International Conference on Semantic Computing (ICSC)*, vol. 00, 413–418. URL: doi.ieeecomputersociety.org/10.1109/ICSC.2017.94.
- Chen, H., Chiang, R. H. L. and Storey, V. C. (2012) Business intelligence and analytics: From big data to big impact. *MIS Q.*, **36**, 1165–1188. URL: <http://dl.acm.org/citation.cfm?id=2481674.2481683>.
- Coral, C., Francisco, R. and Mario, P. (2006) *Ontologies for Software Engineering and Software Technology*. Berlin, Heidelberg: Springer-Verlag.
- Di, T. F., Lefons, E. and Tangorra, F. (2015) Academic data warehouse design using a hybrid methodology. *Computer Science and Information Systems*, **12**, 135–160.
- Diamantini, C. and Potena, D. (2008) Semantic enrichment of strategic datacubes. In *Proceedings of the ACM 11th international workshop on Data warehousing and OLAP*, 81–88. ACM.
- El Sarraj, L., Espinasse, B. and Libourel, T. (2014) An ontology-driven personalization approach for data warehouse exploitation. *International Journal on Advances in Software*, **7**, 253–265.
- Elamin, E. and Feki, J. (2014) Toward and ontology based approach for data warehousing. In *The international arab conference on information technology (ACIT2014)*, 170–179.
- Etcheverry, L. and Vaisman, A. (2012) Qb4olap: a new vocabulary for olap cubes on the semantic web. *Proceedings of COLD*.
- Golfarelli, M., Rizzi, S. and Cella, I. (2004) Beyond data warehousing: What's next in business intelligence? In *Proceedings of the 7th ACM International Workshop on Data Warehousing and OLAP, DOLAP '04*, 1–6. New York, NY, USA: ACM. URL: <http://doi.acm.org/10.1145/1031763.1031765>.
- Hendler, J. (2001) Agents and the semantic web. *IEEE Intelligent Systems*, **16**, 30–37. URL: <http://dx.doi.org/10.1109/5254.920597>.
- Henschen, D. (2015) Analytics, bi, data management trends for 2015. URL: <http://www.informationweek.com/big-data/big-data-analytics/5-analytics-bi-data-management-trends-for-2015/a/d-id/1318551>.
- Inmon, W. H. (1992) *Building the Data Warehouse*. New York, NY, USA: John Wiley & Sons, Inc.
- Jiang, L., Cai, H. and Xu, B. (2010) A domain ontology approach in the etl process of data warehousing. In *e-Business Engineering (ICEBE), 2010 IEEE 7th International Conference on*, 30–35. IEEE.
- Jossen, C. and Dittrich, K. R. (2007) The process of metadata modeling in industrial data warehouse environments. In *BTW Workshops*, 16–27. Citeseer.
- Kalampokis, E., Nikolov, A., Haase, P., Cyganiak, R., Stasiewicz, A., Karamanou, A., Zotou, M., Zeginis, D., Tambouris, E. and Tarabanis, K. A. (2014) Exploiting linked data cubes with opencube toolkit. In *International Semantic Web Conference (Posters & Demos)*, vol. 1272, 137–140.

- Khoury, S., Bellatreche, L., Jean, S. and Ait-Ameur, Y. (2014) Requirements driven data warehouse design: We can go further. In *Leveraging Applications of Formal Methods, Verification and Validation. Specialized Techniques and Applications*, 588–603. Springer.
- Khoury, S. and Ladjel, B. (2010) A methodology and tool for conceptual designing a data warehouse from ontology-based sources. In *Proceedings of the ACM 13th international workshop on Data warehousing and OLAP*, 19–24. ACM.
- Kimball, R. and Ross, M. (2002) *The Data Warehouse Toolkit: The Complete Guide to Dimensional Modeling*. New York, NY, USA: John Wiley & Sons, Inc., 2nd edn.
- Kitchenham, B., Pretorius, R., Budgen, D., Pearl Brereton, O., Turner, M., Niazi, M. and Linkman, S. (2010) Systematic literature reviews in software engineering - a tertiary study. *Inf. Softw. Technol.*, **52**, 792–805. URL: <http://dx.doi.org/10.1016/j.infsof.2010.03.006>.
- Kitchenham, B. A. (2012) Systematic review in software engineering: Where we are and where we should be going. In *Proceedings of the 2Nd International Workshop on Evidential Assessment of Software Technologies, EAST '12*, 1–2. New York, NY, USA: ACM. URL: <http://doi.acm.org/10.1145/2372233.2372235>.
- Kitchenham, B. A., Budgen, D. and Brereton, O. P. (2011) Using mapping studies as the basis for further research—a participant-observer case study. *Information and Software Technology*, **53**, 638–651.
- Klyne, G. and Carroll, J. J. (2004) Resource description framework (RDF): Concepts and abstract syntax. World Wide Web Consortium, Recommendation REC-rdf-concepts-20040210.
- Lamolle, M., Menet, L. and Le Duc, C. (2015) Incremental checking of master data management model based on contextual graphs. *Enterprise Information Systems*, **9**, 681–708.
- Lather, S. (2012) *Business Intelligence and Data Warehousing*. India: Narosa Publishing House.
- Lenz, H.-J. and Shoshani, A. (1997) Summarizability in olap and statistical data bases. In *Scientific and Statistical Database Management, 1997. Proceedings., Ninth International Conference on*, 132–143. IEEE.
- Martin, A., Maladhy, D. and Venkatesan, V. P. (2011) A framework for business intelligence application using ontological classification. *arXiv preprint arXiv:1109.1088*.
- McGregor, J. D. (2002) Building reusable test assets for a product line. In *Software Reuse: Methods, Techniques, and Tools* (ed. C. Gacek), 345–346. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Middelfart, M. and Pedersen, T. B. (2011) The meta-morphing model used in targit bi suite. In *International Conference on Conceptual Modeling*, 364–370. Springer.
- da Mota Silveira Neto, P. A., Carmo Machado, I. d., McGregor, J. D., de Almeida, E. S. and de Lemos Meira, S. R. (2011) A systematic mapping study of software product lines testing. *Inf. Softw. Technol.*, **53**, 407–423. URL: <http://dx.doi.org/10.1016/j.infsof.2010.12.003>.
- Nazri, M. N. M., Noah, S. A. M. and Hamid, Z. (2008) Automatic data warehouse conceptual design. In *Information Technology, 2008. ITSIM 2008. International Symposium on*, vol. 2, 1–7. IEEE.
- Nebot, V. and Berlanga, R. (2012) Building data warehouses with semantic web data. *Decision Support Systems*, **52**, 853–868.
- (2016) Statistically-driven generation of multidimensional analytical schemas from linked data. *Knowledge-Based Systems*, **110**, 15 – 29. URL: <http://www.sciencedirect.com/science/article/pii/S0950705116302143>.
- Nebot, V., Berlanga, R., Pérez, J. M., Aramburu, M. J. and Pedersen, T. B. (2009) Multidimensional integrated ontologies: a framework for designing semantic data warehouses. In *Journal on Data Semantics XIII*, 1–36. Springer.

- Nguyen, T. M., Tjoa, A. M. and Trujillo, J. (2005) Data warehousing and knowledge discovery: A chronological view of research challenges. In *Data Warehousing and Knowledge Discovery* (eds. A. M. Tjoa and J. Trujillo), 530–535. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Niemi, T., Toivonen, S., Niinimäki, M. and Nummenmaa, J. (2007) Ontologies with semantic web/grid in data integration for olap. *International Journal on Semantic Web and Information Systems (IJSWIS)*, **3**, 25–49.
- Niinimäki, M. and Niemi, T. (2009) Journal on data semantics xiii. In *Journal on Data Semantics XIII* (eds. S. Spaccapietra, E. Zimányi and I.-Y. Song), chap. An ETL Process for OLAP Using RDF/OWL Ontologies, 97–119. Berlin, Heidelberg: Springer-Verlag. URL: <http://dl.acm.org/citation.cfm?id=2172259.2172263>.
- Petersen, K., Feldt, R., Mujtaba, S. and Mattsson, M. (2008) Systematic mapping studies in software engineering. In *Proceedings of the 12th International Conference on Evaluation and Assessment in Software Engineering, EASE'08*, 68–77. Swinton, UK, UK: British Computer Society. URL: <http://dl.acm.org/citation.cfm?id=2227115.2227123>.
- Petersen, K., Vakkalanka, S. and Kuzniarz, L. (2015) Guidelines for conducting systematic mapping studies in software engineering: An update. *Information and Software Technology*, **64**, 1 – 18. URL: <http://www.sciencedirect.com/science/article/pii/S0950584915000646>.
- Priebe, T. and Pernul, G. (2003) Towards integrative enterprise knowledge portals. In *Proceedings of the twelfth international conference on Information and knowledge management*, 216–223. ACM.
- Ramasamy, V. and Palanivel, K. (2014) Triplet dependency views of university data warehousing towards decision support system. *International Journal of Advanced Technology & Engineering Research (IJATER)*, 1–7.
- Ravat, F., Song, J. and Teste, O. (2016) Designing multidimensional cubes from warehoused data and linked open data. In *10th International IEEE Conference on Research Challenges in Information Science (RCIS 2016) co-located with the 34th French Conference INFORSID*, pp. 199–200. Grenoble, FR: INformatique des ORganisations et Systèmes d'Information et de Décision (INFORSID). URL: <http://oatao.univ-toulouse.fr/17167/>.
- Romero, O. and Abelló, A. (2010) A framework for multidimensional design of data warehouses from ontologies. *Data & Knowledge Engineering*, **69**, 1138–1157.
- Romero, O., Simitsis, A. and Abelló, A. (2011) Gem: requirement-driven generation of etl and multidimensional conceptual designs. In *Data Warehousing and Knowledge Discovery*, 80–95. Springer.
- Saha, G. K. (2007) Web ontology language (owl) and semantic web. *Ubiquity*, **2007**, 1:1–1:1. URL: <http://doi.acm.org/10.1145/1295289.1295290>.
- Salguero, A., Araque, F. and Delgado, C. (2008a) Ontology based framework for data integration. *WSEAS Transactions on Information Science and Applications*, **5**, 953–962.
- (2008b) Spatio-temporal ontology based model for data warehousing. In *Proceedings of the 7th WSEAS International Conference on Telecommunications and Informatics, TELE-INFO'08*, 125–130. Stevens Point, Wisconsin, USA: World Scientific and Engineering Academy and Society (WSEAS). URL: <http://dl.acm.org/citation.cfm?id=1404049.1404072>.
- Samuel, J. (2014) Towards a data warehouse fed with web services. In *The Semantic Web: Trends and Challenges*, 874–884. Springer.
- Sell, D., Cabral, L., Motta, E., Domingue, J. and Pacheco, R. (2005) Adding semantics to business intelligence. In *Database and Expert Systems Applications, 2005. Proceedings. Sixteenth International Workshop on*, 543–547. IEEE.
- Sell, D., da Silva, D. C., Beppler, F. D., Napoli, M., Ghisi, F. B., Pacheco, R. and Todesco, J. L. (2008) Sbi: a semantic framework to support business intelligence. In *Proceedings of the first international workshop on Ontology-supported business intelligence*, 11. ACM.

- Sell, D., da Silva, D. C., Ghisi, F. B., Todesco, J. L. and Napoli, M. (2012) *Adding Semantics to Business Intelligence: Towards a Smarter Generation of Analytical Tools*. INTECH Open Access Publisher.
- Selma, K., Ilyès, B., Ladjel, B., Eric, S., Stéphane, J. and Michael, B. (2012) Ontology-based structured web data warehouses for sustainable interoperability: requirement modeling, design methodology and tool. *Computers in Industry*, **63**, 799–812.
- Shadbolt, N., Berners-Lee, T. and Hall, W. (2006) The semantic web revisited. *IEEE Intelligent Systems*, **21**, 96–101. URL: <http://dx.doi.org/10.1109/MIS.2006.62>.
- Shaw, M. (2003) Writing good software engineering research papers: Minitutorial. In *Proceedings of the 25th International Conference on Software Engineering, ICSE '03*, 726–736. Washington, DC, USA: IEEE Computer Society. URL: <http://dl.acm.org/citation.cfm?id=776816.776925>.
- Simitsis, A., Skoutas, D. and Castellanos, M. (2010) Representation of conceptual etl designs in natural language using semantic web technology. *Data & Knowledge Engineering*, **69**, 96–115.
- Skoutas, D. and Simitsis, A. (2006) Designing etl processes using semantic web technologies. In *Proceedings of the 9th ACM international workshop on Data warehousing and OLAP*, 67–74. ACM.
- (2007) Ontology-based conceptual design of etl processes for both structured and semi-structured data. *International Journal on Semantic Web and Information Systems (IJSWIS)*, **3**, 1–24.
- Skoutas, D., Simitsis, A. and Sellis, T. (2009) Ontology-driven conceptual design of etl processes using graph transformations. In *Journal on Data Semantics XIII*, 120–146. Springer.
- Spahn, M., Kleb, J., Grimm, S. and Scheidl, S. (2008) Supporting business intelligence by providing ontology-based end-user information self-service. In *Proceedings of the First international Workshop on ontology-Supported Business intelligence*, 10. ACM.
- Steiner, D., Neumayr, B. and Schrefl, M. (2015) Judgement and analysis rules for ontology-driven comparative data analysis in data warehouses. In *Proceedings of the 11th Asia-Pacific Conference on Conceptual Modelling (APCCM 2015)*, vol. 27, 30.
- Talebzadeh, S., Seyyedi, M. A. and Salajegheh, A. (2014) Automated creating a data warehouse from unstructured semantic data. *International Journal of Computer Applications*, **88**.
- Thenmozhi, M. and Vivekanandan, K. (2012) An ontology based hybrid approach to derive multidimensional schema for data warehouse. *International Journal of Computer Applications*, **54**.
- (2013) A tool for data warehouse multidimensional schema design using ontology. *Int. J. Comput. Sci. Issues (IJCSI)*, **10**, 161–168.
- Trujillo, J. and Maté, A. (2012) *Business Intelligence: First European Summer School, eBISS 2011, Paris, France, July 3-8, 2011, Tutorial Lectures*, chap. Business Intelligence 2.0: A General Overview, 98–116. Berlin, Heidelberg: Springer Berlin Heidelberg. URL: http://dx.doi.org/10.1007/978-3-642-27358-2_5.
- Wache, H., Vögele, T., Visser, U., Stuckenschmidt, H., Schuster, G., Neumann, H. and Hübner, S. (2001) Ontology-based integration of information - a survey of existing approaches. 108–117.
- Wieringa, R., Maiden, N., Mead, N. and Rolland, C. (2005) Requirements engineering paper classification and evaluation criteria: A proposal and a discussion. *Requir. Eng.*, **11**, 102–107. URL: <http://dx.doi.org/10.1007/s00766-005-0021-6>.



RICARDO GACITÚA received his B.S in Informatics Engineering from the Universidad de Concepcion, Chile in 1996, his M.S. degree in Industrial Engineering from the Universidad de Chile, Chile in 2004 and his Ph.D. degree in Computer Science from Lancaster University, UK in 2009. From 2008 to 2011, he was a Research Assistant with the School of Computing and

Communication at Lancaster University, UK. Since 2014, he has been a Professor with the Computer Science and Informatics Department at the Universidad de La Frontera, Temuco, Chile. His research interests include requirements engineering, ontology learning, text mining and the Semantic Web. He was Chair of the 36th International Conference of the Chilean Computer Science Society (SCCC 2017), and the Latin-American Symposium on Software Engineering (2016). Mr. Gacitua's awards and honors include the ORSA Award (Overseas Research Students Awards Scheme, UK), and the British Computer Society Specialist Group on Artificial Intelligence - Outstanding Paper. He was named a Fellow of British Computer Society (FBCS) in 2010.

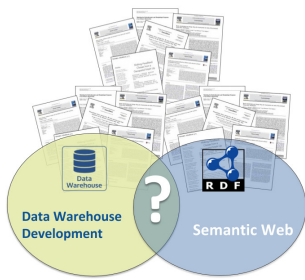


JOSE NORBERTO MAZÓN is member of the WaKe research group at the University of Alicante, Spain. His research work focuses on open data management, data integration and business intelligence within "big data" scenarios. He has published his research in international journals, such as Decision Support Systems, Information Sciences, Data & Knowledge Engineering or ACM Transaction on the Web.



ANIA CRAVERO is a Civil Engineering Engineer m. Informatics, at La Frontera University, Temuco, Chile. She obtained her Ph.D. in Computer Science and Computer Systems from Atlantic International University, USA (2010). She is an Academic in the Computer Science Department, La Frontera University. Her research interests are in Databases, Data Warehouses and Goal Alignment.

GRAPHICAL ABSTRACT



The objective of this study is to examine recently published research articles that take into account the use of Semantic Web technologies in data warehouse development with the intention of summarizing their results, classifying their contributions to the field according to publication type, evaluating the maturity level of the results and identifying future research challenges.

TABLE 5 Summary of analyzed studies.

ID	Study Title	Year	Source
P01	Towards integrative enterprise knowledge portals	2003	(Priebe and Pernul, 2003)
P02	Adding semantics to business intelligence	2005	Sell et al. (2005)
P03	Designing ETL processes using semantic web technologies	2006	Skoutas and Simitsis (2006)
P04	The process of metadata modelling in industrial Data Warehouse environments	2007	Jossen and Dittrich (2007)
P05	Ontology-based conceptual design of ETL processes for both structured and semi-structured data	2007	Skoutas and Simitsis (2007)
P06	Ontologies with semantic web/grid in data integration for olap	2007	Niemi et al. (2007)
P07	Automatic Data Warehouse Conceptual Design	2008	Nazri et al. (2008)
P08	Ontology based framework for data integration	2008	Salguero et al. (2008a)
P09	SBI: a semantic framework to support business intelligence	2008	Sell et al. (2008)
P10	Diamantini adn Potena	2008	Diamantini and Potena (2008)
P11	Spatio-temporal ontology based model for Data Warehousing	2008	Salguero et al. (2008b)
P12	Supporting business intelligence by providing ontology-based end-user information self-service	2008	Spahn et al. (2008)
P13	Ontology-driven conceptual design of ETL processes using graph transformations	2009	Skoutas et al. (2009)
P14	An ETL process for OLAP using RDF/OWL ontologies	2009	Niinimäki and Niemi (2009)
P15	Multidimensional integrated ontologies: a framework for designing semantic data warehouses	2009	Nebot et al. (2009)
P16	Representation of conceptual ETL designs in natural language using Semantic Web technology	2019	Simitsis et al. (2010)
P17	A framework for multidimensional design of data warehouses from ontologies Romero and Abelló (2010)	2010	Simitsis et al. (2010)
P18	A methodology and tool for conceptual designing a data warehouse from ontology-based sources	2010	Khourri and Ladjel (2010)
P19	A domain ontology approach in the ETL process of data warehousing	2010	Jiang et al. (2010)
20	GEM: requirement-driven generation of ETL and multidimensional conceptual designs	2011	Romero et al. (2011)
P21	A semantic approach to ETL technologies	2011	Bergamaschi et al. (2011)
P22	A framework for business intelligence application using ontological classification	2011	Martin et al. (2011)
P23	Building data warehouses with semantic web data	2012	Nebot and Berlanga (2012)

TABLE 6 Continuation of the Table Resumen.

ID	Study Title	Year	Source
P24	Ontology-based structured web data warehouses for sustainable interoperability: requirement modeling, design methodology and tool	2012	Selma et al. (2012)
P25	Adding Semantics to Business Intelligence: Towards a Smarter Generation of Analytic Tools	2012	Sell et al. (2012)
P26	An ontology based hybrid approach to derive multidimensional schema for data warehouse	2012	Thenmozhi and Vivekanandan (2012)
P27	QB4OLAP: a new vocabulary for OLAP cubes on the semantic web	2012	Etcheverry and Vaisman (2012)
P28	A tool for data warehouse multidimensional schema design using ontology	2013	Thenmozhi and Vivekanandan (2013)
P29	Using schema matching in data transformation for warehousing web data	2013	Ali et al. (2013)
P30	Requirements driven data warehouse design: We can go further	2014	Khoury et al. (2014)
P31	Toward an ontology based approach for data warehousing	2014	Elamin and Feki (2014)
P32	Automated Creating a Data Warehouse from Unstructured Semantic Data	2014	Talebzadeh et al. (2014)
P33	Triplet Dependency Views of University Data Warehousing Towards Decision Support System	2014	Ramasamy and Palanivel (2014)
P34	Towards a data warehouse fed with web services	2014	Samuel (2014)
P35	An Ontology-Driven Personalization Approach for Data Warehouse Exploitation	2014	El Sarraj et al. (2014)
P36	Academic data warehouse design using a hybrid methodology	2015	Di et al. (2015)
P37	Using semantic web technologies for exploratory OLAP: a survey	2015	Abelló et al. (2015)x
P38	Incremental checking of Master Data Management model based on contextual graphs	2015	Lamolle et al. (2015)
P39	Judgment and analysis rules for ontology-driven comparative data analysis in data warehouses	2015	Steiner et al. (2015)
P40	Statistically-driven generation of multidimensional analytical schemas from linked data	2016	Nebot and Berlanga (2016)
P41	Designing Multidimensional Cubes from Warehoused Data and Linked Open Data	2016	Ravat et al. (2016)

TABLE 7 Summary and contributions of each analyzed study.

ID	Summary	Contribution
P01	This paper discusses integration aspects within enterprise knowledge portals and presents an OLAP-based approach for communicating the user context (revealing the user's information need) among portlets, utilizing Semantic Web technologies.	A prototype to evaluate the approach, demonstrating how OLAP can be integrated with information retrieval.
P02	This paper proposes an approach, which aims at integrating business semantics into analytical tools by providing semantic descriptions of exploratory functionalities and available services..	An architecture for business intelligence, which uses Semantic Web technologies based on IRS-III. In addition, a prototype tool based on this architecture.
P03	This paper argues that ontologies constitute a suitable model for the purpose of establishing the appropriate mappings between the attributes of the data sources and the attributes of the data warehouse tables, which is critical in specifying the required transformations in an ETL workflow. The authors show how the usage of ontologies can enable a high degree of automation regarding the construction of an ETL design.	There are several contributions of this paper: (1) construction of a domain vocabulary; (2) annotation of the data stores; (3) generation of a domain ontology; and (iv) generation of conceptual ETL design.
P04	This paper discusses the process of identifying metadata model requirements, defining a new metadata model and finally implementing it in a metadata schema. The paper describes the implementation of the metadata model based on the metadata standards Resource Description Framework (RDF) and RDF Schema (RDFS).	A new metadata model using RDF and RDFS for a data warehouse environment that fulfills the needs of an industrial setting. In addition, the implementation of the model in a real-world scenario.
P05	This paper proposes an ontology-based approach to facilitate the conceptual design of the integration layer of a data warehouse. The proposed approach is based on the use of Semantic Web technologies to semantically annotate the data sources and the data warehouse, so that mappings between them can be inferred, thereby resolving the issue of heterogeneity in data integration..	A graph-based representation, which is used as a conceptual model for the datastores. In addition, a suitable application ontology is created and used to annotate the data stores.
P06	This paper presents an approach for data source integration in OLAP scenarios, based on RDF and other Semantic Web technologies.	An OWL/RDF ontology for OLAP data sources and OLAP cubes. In addition, they provide a prototype for testing the approach.

TABLE 8 Continuation of the Table Summary (I)

ID	Summary	Contribution
P07	This paper presents a methodology to develop a data warehouse (based on multidimensional modeling) by integrating all three development approaches such as supply-driven, goal-driven and demand-driven.	An ontology embedded with specific knowledge domain. The ontology aims to identify the semantics of the domain. In addition, a prototype is developed.
P08	This paper describes a framework which encompasses the entire data integration process. They deal with the problem of assisting in the process of designing information systems based on DW and the inclusion of new data sources, since it may involve re-designing the scheme of DW.	An ontology-based model to support the data integration process.
P09	This paper presents a framework called SBI (Semantic Business Intelligence) in which applied ontologies are used for the description of business rules and concepts in order to support semantic-analytical functionalities that extend traditional OLAP operations..	An analytical tool, called Extracta, which relies on SBI ontologies and modules.
P10	This paper proposes a novel model for semantic annotation of a data warehouse schema that takes into account domain ontologies as well as a mathematical ontology.	An ontology, which describes mathematical formulas underlying elements of the data cube schema, including the semantics of operators as the basis for OLAP analysis.
P11	This paper describes a spatio-temporal extension of an ontology language which facilitates the generation of the scheme of a data warehouse as well as the design of the processes which extract, transform and load (ETL) the data from the sources in the data warehouse according to the temporal characteristics of the data sources.	An ontology-based model to deal with the data sources schemes integration problem.
P12	This paper presents an ontology-based architecture and end-user tool, enabling easy data access and query creation for business users. The approach is based on a semantic middleware integrating data from heterogeneous information systems and providing a comprehensible data model in the form of a business level ontology (BO).	An ontology-based architecture and an end-user tool: Semantic Query Designer (SQD).
P13	This paper proposes a customizable and extensible ontology-driven approach for the conceptual design of ETL processes.	A graph-based representation for the source and target data stores, and a method for devising flows of ETL operations by means of graph transformations.

TABLE 9 Continuation of the Table Summary (II)

ID	Summary	Contribution
P14	This paper presents a method for on-demand construction of OLAP cubes based on relational algebra. The authors combine Semantic Web and OLAP, resulting in a method of creating a fully functional tool for data analysis.	A method and an ontology-based tool that works as a user interface for the system, from design to actual analysis.
P15	This paper proposes the Semantic Data Warehouse to be a repository of ontologies and semantically annotated data resources. In addition, an ontology-driven framework to design multidimensional analysis models for Semantic Data Warehouses is proposed.	An ontology-driven framework to design multidimensional analysis models for Semantic Data Warehouses.
P16	This paper provides a method for the representation of a conceptual ETL design as a narrative, which is the most natural means of communication and does not require particular technical skills or familiarity with any specific model.	A flexible and customizable template-based mechanism for the representation of the ETL design as a narrative.
P17	This paper presents a user-centered approach to support the end-user requirements elicitation and the data warehouse multidimensional design tasks. This proposal is based on a reengineering process that derives the multidimensional schema from a conceptual formalization of the domain.	A user-centered approach to support the end-user requirements elicitation and the data warehouse multidimensional design tasks.
P18	This paper proposes a conceptual design methodology of data warehouses from data residing in various ontology-based databases that takes into account sources and decision maker requirements.	A methodology and a case tool, called S2RWC, supporting this proposal.
P19	This paper introduces an ontology model to integrate semantic heterogeneous data through ETL process for data warehouse development. A domain ontology is introduced into ETL process of finding the data sources, defining the rules of data transformation, and eliminating the heterogeneity.	An ontology-based method to eliminate data heterogeneity so as to support the development of data warehouses.

TABLE 10 Continuation of the Table Summary (III).

ID	Summary	Contribution
20	This paper presents a system called GEM (Generating ETL and Multidimensional designs). GEM starts with a set of source data stores and business requirements (e.g., business queries, service level agreements (SLAs)) and based on these, it produces a multidimensional design for the target data stores, along with a set of ETL operations required for the population of the target data warehouse.	A system that facilitates the production of ETL and multidimensional designs, starting from a set of business requirements and source data stores. Additionally, some novel algorithms finding and validating an ontology subset as a multidimensional schema, and identifying ETL operators at the same time.
P21	This paper provides support for ETL by proposing a tool that: (1) allows the semi-automatic definition of inter-attribute semantic mappings, by identifying the parts of the data source schemas which are related to the data warehouse schema, thus supporting the extraction process; and (2) groups the attribute values semantically related thus defining a transformation function for populating with homogeneous values the data warehouse.	A ETL tool that couples and extends the functionalities of two previously developed systems: the MOMIS integration system and the RELEVANT data analysis system.
P22	This paper proposes to consider the Web as an information repository. Since retrieving specific information from the web is challenging, an ontological model is developed to capture specific information by using web semantics. From the ontology model, the relations between the data are mined using decision trees.	An ontology model and an architecture for business intelligence.
P23	This paper presents an approach for efficiently analyzing and exploring large amounts of semantic data by combining the inference power from the annotation semantics with the analysis capabilities provided by OLAP-style aggregations, navigation, and reporting.	A semi-automatic method to dynamically build multidimensional models from semantic data guided by the user requirements. In addition, two novel algorithms are presented to dynamically create dimension hierarchies complying with OLAP properties from the taxonomic relations of domain ontologies for multidimensional models generated with the previous method.
P24	This paper presents a methodology for designing data warehousing applications from various sources by using ontologies. Each source has its local ontology referencing (and specialize/extend) a global one.	A new classification of data warehouse design methods showing their convergence to ontology-base conceptualisation. In addition, a methodology for designing data warehouses from ontology-based databases and users' requirements supported by a case tool.

TABLE 11 Continuation of the Table Summary (IV).

ID	Summary	Contribution
P25	This paper presents an integrated Semantic Business Intelligence (SBI) architecture for analytical tools. The architecture is supported by business semantics that, in turn, are applied to contextualize the organizations' resources (i.e. logic, data sources and services).	An ontology for business intelligence, which supplies the business terminology used to enable data sources annotation. In addition, a business intelligence architecture, which is supported by a tool.
P26	This paper proposes a framework that uses an ontology for the design of multidimensional models. The framework uses a hybrid approach where the reconciliation of requirements and data source are done at the early stage of design. An ontology reasoning is adopted in order to automatically derive multidimensional elements such as facts and dimensions.	A comprehensive framework using a hybrid methodology to derive a multidimensional model from multiple ontology sources based on requirements. In addition, a set of ontology matching algorithms to map requirements with sources.
P27	This paper proposes a new vocabulary, denoted QB4OLAP, which extends QB vocabulary (RDF data cube vocabulary from W3) to fully support OLAP models and operators. This provides algorithms that build the QB4OLAP structures needed to analyze observations already published using QB, and vice versa.	A new vocabulary, which extends QB to fully support OLAP models and operators. In addition, some algorithms that build the structures that allow performing both kinds of analysis, and show compatibility between QB and QB4OLA.
P28	This paper presents an ontology-driven tool which helps to automatically derive the conceptual model and logical model for the data warehouse from a data source and business requirements.	Main contributions are: (1) representation of requirements formally using data warehousing requirement ontology; (2) automatically deriving multidimensional elements present in the data source ontology; (3) ontology-based data warehouse schema design tool development.
P29	This paper proposes a semi-automatic approach to support the schema mapping transformation process. This approach is based on the use a XML Schema representation of Web data and the existing warehouse schema.	A warehousing Web data extraction framework. In addition, some rules to transform Web data into a data warehouse.
P30	This paper proposes an explicit requirements engineering phase in data warehouse development, which is goal-oriented.	A goal-oriented approach to develop data warehouses.

TABLE 12 Continuation of the Table Summary (V).

ID	Summary	Contribution
P31	This paper presents an ontology- based hybrid approach for data warehouse design consisting in three main steps: (1) building multidimensional models based on business requirements; (2) building multidimensional models from an operational data source; and (3) matching the two sets of schemas to produce a data warehouse model that agrees to user and data source simultaneously.	A comparative study involving some research works tackling the data warehouse design process. In addition, an ontology- based hybrid approach for the data warehouse design and a framework to produce a data warehouse model.
P32	This paper proposes a solution for automated creating a data warehouse from unstructured semantic data using ontology. An algorithm for creating a data warehouse from semantic data by using ontology is presented.	An algorithm can create the desired data warehouse automatically and store the available data from an ontology OWL or RDF.
P33	This paper presents three-tier perspectives of the educational data warehousing to abstract ontologies from structural metadata of relational databases which are the resultant of ETL processing as well as generating reports from the data re-fined from the ETL processing.	An ontology-based approach to build an educational data warehouse.
P34	This paper presents a prototype named DaWeS (Data warehouse fed with Web Services) and explores how ETL using the mediation approach benefits this trade-off for enterprises with complex data warehousing requirements. The goal of this study is to investigate and devise an approach to address the trade-off between scalability and adaptability in large scale integration with numerous ever-evolving web services.	A prototype to use DaWeS in the business analytics subject. In addition, some optimization heuristics.
P35	This paper presents an Ontology-driven Personalization System (OPS) based on three connected ontologies: domain ontology, data warehouse ontology and resources ontology. OPS return a set of personalized resources search based on users' domain and his recurring interests.	An architecture of the "ontology-driven personalization system" and the use-cases supported by this system. In addition, a methodology used to develop the knowledge base of the personalization system.
P36	This paper presents the architecture of a business intelligence system for academic organizations. Then, It is illustrated the design process of the data warehouse devoted to the analysis of the main factors affecting the importance and the quality level of every university, such as the evaluation of the research and the education. The design process is based on a hybrid methodology that is largely automatic and relies on an ontological approach for the integration of the different data sources.	A hybrid methodology. In addition the design of an academic data warehouse and a research data mart.

TABLE 13 Continuation of the Table Summary (VI).

ID	Summary	Contribution
P37	The paper first presents a characterization of OLAP environments, followed by an introduction to the relevant Semantic Web foundation concepts. Then, it describes the relationship of multidimensional models and Semantic Web technologies. Next, the paper goes on to survey the use of Semantic Web technologies for data modeling and data provisioning, including semantic data annotation and semantic-aware ETL processes. Finally, all the findings are discussed and a number of directions for future research are outlined.	A survey to categorize how Semantic Web technologies have been applied to solve the new requirements of exploratory OLAP systems. In addition, future challenges are identified.
P38	An incremental validation method based on model driven engineering is used to develop a Master Data Management system represented by XML Schema models. To do so, authors define an abstraction layer using UML class diagrams. The validation method aims to minimise the model errors and to optimise the process of model checking.	A validation model developed in ArgoUML IDE.
P39	This paper presents a conceptual modelling of judgement and analysis rules together with their organisation and multidimensional contextualisation in rule families, explain different rule evaluation strategies, and briefly report on the implementation of the approach in order to detect multidimensional elements and OLAP operations.	A metamodel for the representation of comparative multidimensional ontologies, including comparative scores, comparative concepts, and comparative cubes. In addition, a modelling of judgement and analysis rules, their organisation in rule families, their specialisation along subsumption hierarchies of comparative concepts, and their contextualised evaluation according to different rule evaluation strategies.
P40	This paper presents an automatic approach to generate useful multidimensional model analytical patterns from Linked Data sources. Multidimensional patterns are based on the semantics of the data.	Mappings from multidimensional models to a statistical layer on top of Linked Data sources. In addition, a statistical framework to discover multidimensional patterns in an automatic way.
P41	This paper provides a new multidimensional model, named Unified Cube, which offers a generic representation for both warehoused data and Linked Open Data at the conceptual level. A two-stage process is proposed to build a Unified Cube according to decision-makers' needs. As a first step, schemas published with specific modeling languages are transformed into a common conceptual representation. The second step is to associate together related data to form a Unified Cube containing all useful information about an analysis subject.	A generic modelling solution, named Unified Cube, for both warehoused data and Linked Open Data. In addition, a process to build a Unified Curve. Moreover, a high-level declarative language is provided to enable non-expert users to define the relevance between data according to their analysis needs.