

Enhancing Virtual Ontology Based Access over Tabular Data with Morph-CSV

David Chaves-Fraga^{a,*}, Edna Ruckhaus^a, Freddy Priyatna^a, Maria-Esther Vidal^b and Oscar Corcho^a

^a *Ontology Engineering Group, Universidad Politécnica de Madrid, Spain*

E-mails: dchaves@fi.upm.es, eruckhaus@fi.upm.es, fpriyatna@fi.upm.es, ocorcho@fi.upm.es

^b *TIB - Leibniz Information Centre for Science and Technology and L3S Leibniz University of Hannover, Germany*

E-mail: maria.vidal@tib.eu

Abstract. Ontology-Based Data Access (OBDA) has traditionally focused on providing a unified view of heterogeneous datasets (e.g., relational databases, CSV and JSON files), either by materializing integrated data into RDF or by performing on-the-fly querying via SPARQL query translation. In the specific case of tabular datasets represented as several CSV or Excel files, query translation approaches have been applied by considering each source as a single table that can be loaded into a relational database management system (RDBMS). Nevertheless, constraints over these tables are not represented (e.g., referential integrity among sources, datatypes, or data integrity); thus, neither consistency among attributes nor indexes over tables are enforced. As a consequence, efficiency of the SPARQL-to-SQL translation process may be affected, as well as the completeness of the answers produced during the evaluation of the generated SQL query. Our work is focused on applying implicit constraints on the OBDA query translation process over tabular data. We propose Morph-CSV, a framework for querying tabular data that exploits information from typical OBDA inputs (e.g., mappings, queries) to enforce constraints and can be used together with any SPARQL-to-SQL OBDA engine. Morph-CSV relies on both a constraint component and a set of constraint operators. For a given set of constraints, the operators are applied to each type of constraint with the aim of enhancing query completeness and performance. We evaluate Morph-CSV in several domains: e-commerce with the BSBM benchmark; transportation with a benchmark using GTFS dataset from the Madrid subway; and biology with a use case extracted from the Bio2RDF project. We compare and report the performance of two SPARQL-to-SQL OBDA engines, without and with the incorporation of Morph-CSV. The observed results suggest that Morph-CSV is able to speed up the total query execution time by up two orders of magnitude, while it is able to produce all the query answers.

Keywords: Knowledge Graphs, Tabular Data, Mapping Languages, Constraints

1. Introduction

Guided by Open Data principles, governments and private organizations are regularly publishing wide amounts of public data in open data portals. For example, almost a million of datasets are available in the European Open Data Portal (EODP)¹, and many of them are available in tabular formats (e.g., CSV, Excel), as observed in Table 1. Both the simplicity of a tabular representation and the variety of tools to manage a table (e.g., Excel, Calc) have influenced in the popularity of tabular formats to represent open data.

Although extensively utilized, tabular representations imposed various data management challenges to advanced users (e.g., developers, data scientists). The lack of a unified way to query tabular data, something that is available in other formats (e.g., RDB, JSON, XML), hinders the integration of sources, especially those having datatype inconsistencies. Moreover, data may not be normalized, and information about relationships or column names are not always descriptive or homogeneous. Hence, data consumers are usually forced to apply ad-hoc or manual data wrangling processes to consume data via open data portals.

*Corresponding author. E-mail: dchaves@fi.upm.es.

¹<https://www.europeandataportal.eu>

Following Linked Data [1] and FAIR initiatives [2]², data providers are encouraged to make data available in an RDF-based representation following the 5-star linked data principles³. The Ontology-Based Data Access (OBDA) [3] paradigm facilitates the transformation of heterogeneous data into RDF. An OBDA corresponds to a data integration system (DIS) [4] over heterogeneous data sources. A DIS unified schema is defined in terms of ontologies, while mapping rules establish correspondence between the unified schema concepts and the DIS data sources. An OBDA can be materialized or virtual. In a materialized OBDA, the integration of the DIS data sources is physically represented in RDF [3]. Contrary, in a virtual OBDA, data integration is performed on the fly during query processing; DIS mapping rules are used to rewrite SPARQL queries into queries against the DIS data sources [5, 6]. Features like functions in mappings [7, 8] and metadata [9], (i.e., annotations) are usually used in materialized OBDA to overcome the aforementioned challenges of tabular data.

Traditional virtual OBDA approaches, usually, rely on loading the tabular data into SQL-based systems^{4,5} (e.g., MySQL, Apache Drill, Spark SQL, Presto) to perform the query translation techniques. However, the correctness and optimization of these techniques are supported by the main assumption about the existence of constraints over the source data (i.e., a good physical design of the relational database instance). Their absence during a virtual OBDA process over tabular data directly impacts over completeness and performance of these techniques. Completeness is affected because of heterogeneity issues in data sources (e.g., datatype CSV columns are simply treated as string-type SQL columns). Furthermore, performance is impacted because indexes are not created based on basic relational constraints, i.e., primary and foreign key constraints are not defined in the schema. As a consequence, query translation optimization techniques that normally exploit indexes (e.g., [6, 10]) do not produce the expected results.

OBDA annotations such as the W3C recommendation to annotate tabular data, CSVW [9] and some extensions of standard mapping rules (e.g., RML+FnO [7]) are commonly used to describe constraints over an

²<https://www.go-fair.org/fair-principles/>

³<https://5stardata.info/en/>

⁴<https://github.com/oeg-upm/morph-rdb/wiki/Usage#csv-files>

⁵<https://github.com/ontop/ontop/wiki/MappingDesignTips#database-tips>

Table 1

Most commonly used formats and percentage over the total number of datasets to expose data in mature EU open data portals in October 2019. Each dataset may be shared in different formats.

Data Portal	1st Format	2nd Format	3rd Format
Spain	CSV (50%)	XLS (35%)	JSON (33%)
Norway	CSV (77%)	GEOJSON (17%)	JSON (14%)
Italy	CSV (76%)	JSON (35%)	XML (25%)
Croatia	XLS (63%)	CSV (40%)	HTML (33%)

OBDA tabular dataset. For example, we can standardize a column indicating its format, define integrity constraints or declare datatypes. The majority of OBDA query translation engines [6, 11] do not include this information. Those engines that have partially included the constraints (e.g., Squerall [12] parses RML+FnO mapping rules) are not fully documented; i.e., there is no explanation of how these constraints are taken into account. The definition of a workflow that includes the exploitation of these tabular annotations during a virtual OBDA process will ensure correct and optimized SPARQL-to-SQL translations.

Problem and Proposed Solution: We address the limitations of current OBDA query translation techniques over tabular data. Our goals are (i) define a framework that includes the application of a set of constraints over tabular data and (ii) define a set of operators that apply each type of constraint in order to improve query completeness and performance. We propose a set of new steps to be aligned with the current OBDA workflow. Further, we implement Morph-CSV, and evaluate its behavior in comparison with previous approaches.

Contributions: Our main contributions are as follows:

1. Definition of the concept of Virtual Tabular Dataset (VTD) composed by a tabular dataset and its corresponding OBDA annotations, as well as its alignment with the current definition and assumptions of the OBDA framework [13].
2. Morph-CSV, a framework that implements a constraint-based OBDA workflow for tabular datasets; it receives a VTD and a SPARQL query as inputs and outputs an OBDA instance. Morph-CSV performs the following steps: (i) generation of the constraints based on information on the VTD; (ii) selection of sources and attributes needed to answer the query; (iii) pre-processing of the selected sources applying some of the constraints; and (iv) physical implementation of the corresponding RDB instance and associated schema, ensuring effectiveness of the SPARQL-to-SQL translations and optimizations.

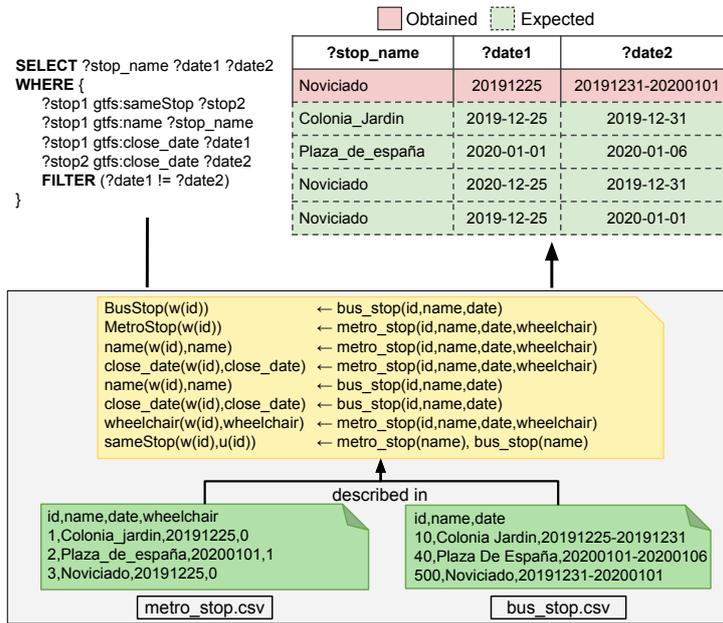


Fig. 1. **Motivating Example.** SPARQL query evaluation over two tabular data files in the transport domain through a common OBDA approach. It loads the files as single tables in an SQL-based system and uses the mapping rules for query translation. The number of results differs with respect to the expected results due the heterogeneity of the raw data. Additionally, query performance may be affected by the join condition between the two tables, the absence of indexes and the loading of columns that are not needed to answer the input query (wheelchair).

3. Evaluation of Morph-CSV over two open source engines: Morph-RDB [6] and Ontop [5]; two benchmarks (BSBM [14] and GTFS-Madrid-Bench⁶), and a real-world testbed from the Bio2RDF project [15] are used in the study.

The rest of the paper is structured as follows: Section 2 motivates the problem of OBDA query translation over tabular data with an example in the transport domain. Section 3 describes the identified challenges for querying and integrating tabular data, and current proposals of OBDA annotations for tabular data that address these challenges. Section 4 presents Morph-CSV, an approach for enhancing OBDA query translation over tabular data through the application on the fly of a set of constraints. Section 5 reports the results of our empirical study together with a general discussion in Section 6. We present the related work in Section 7 and our conclusions and future work in Section 8.

2. Motivating Example

Consider the de-facto standard for publishing open data in the transport domain, GTFS⁷. This model provides information such as *schedules*, *stops* and *routes* using 15 different inter-related CSV files called GTFS feed. Each feed usually specifies the information of one type of transportation mode (e.g., metro, train, and tram). Linking these feeds based on their stops enables route planners to offer multi-modal routes, a route that can be travelled using various types of transportation. The GTFS feeds from the metro and the buses of the city of Madrid have several stops and stations in common; they are created by different transport authorities, and the names of their stops are defined in different manners. Figure 1 depicts a SPARQL query asking for bus and metro stops with the same name, and information related to their closing dates during holidays. Since GTFS uses temporal identifiers for its resources, links have to be established joining stop names. However, as it is usual in open datasets, stop names do not follow a standard structure (e.g., “Colonia Jardin” in *bus_stops.csv* and “Colonia_jardin” in

⁶<https://github.com/oeg-upm/gtfs-bench>

⁷<https://developers.google.com/transit/gtfs/reference/>

metro_stops.csv). A similar issue is present in closing dates, where there are multi-valued cells and their format is not the standard one (e.g., yyyy-MM-dd). Following the approach commonly employed by typical OBDA engines, the two files would be loaded into an SQL-system and treated as single tables. The obtained result set only contains one answer where the stop names in the two data sources are identical (“Noviciado”). However, the expected result set should include more answers by performing an improved join among the stop names of the bus and metro, through the normalization of multi-valued date columns.

The manual and ad-hoc preparation of a tabular dataset for a virtual OBDA process is usually the most time-consuming and less reproducible task. Exploiting available standard OBDA annotations allows its generalization and automatization, as well as ensuring the effectiveness of query completeness and performance of SPARQL-to-SQL techniques, complying with OBDA assumptions.

3. Ontology Based Data Access over Tabular Data

This section describes a set of challenges demanded to be addressed whenever tabular data is queried in a virtual OBDA framework. Further, we describe relevant OBDA proposals for annotating tabular datasets and their alignment with the identified challenges.

3.1. Querying Challenges under virtual OBDA

There are specific challenges on querying tabular datasets using an OBDA approach that have not been tackled by existing techniques. We will describe those challenges and explain how they may have a negative effect in terms of completeness and performance of query-translation approaches:

- **Selection (S):** Existing frameworks load all of the files that are specified as sources in the OBDA mapping rules into a SQL database before executing the query-translation process. This step has to be repeated whenever a SPARQL query is evaluated to ensure up-to-date results, resulting in unnecessary longer loading time, affecting, thus, OBDA performance.
- **Normalization (N):** Tabular data formats do not provide restrictions on how to structure data. As a result, cells may contain multiple values, and one file may represent multiple entities. Having non-

normalized tables may affect the completeness of the query. When a tabular source with multiple-valued cells is loaded into an RDB table, the cell’s value is interpreted by the RDBMS as an atomic value, reducing, thus, completeness for queries that filter or “join” on the corresponding column. Representing several entities in a single file may lead to duplicate answers, and in turn, decrease query answering performance.

- **Heterogeneity (H):** Tabular data normally contain values that need to be transformed before query evaluation (e.g., column default values or normalization of date formats). Since there may be different formats for the same datatype or default values may have not been included in the dataset, query completeness can be affected.
- **Lightweight Schema (LS):** Most of the tabular data only provide minimal information about their underlying schema in the form of column names in the header, if at all present. Also, although there is implicit information on keys and relationships among sources, there is no way to specify primary key or foreign key constraints. The same can be said on indexes and datatypes. The existence of this type of information is assumed [13] in an OBDA approach for performing optimizations in query evaluation techniques. Therefore, the lack of this information affects the performance of OBDA engines.

Although some of the aforementioned challenges are not only specific to tabular datasets and are proposed in several data integration approaches [16–18] there are two main reasons why it is important to address these problems in this context: first, as we reflect in Section 1, the number of tabular datasets available in the web of data is enormous and still growing and these challenges were not taken into account in previous OBDA proposals; second, although there are declarative proposals to handle these issues in the state of the art like CSV on the Web [9] for metadata annotations, or mapping languages that include transformation functions to deal with heterogeneity (e.g., RML+FnO [7] or R2RML-F [19]), there is not yet a proposal that exploits the information from these inputs including their application in the form of constraints into a common OBDA workflow.

3.2. OBDA annotations for Tabular Data

R2RML [20] is a W3C Recommendation for describing transformation rules from RDB to RDF and

Table 2

Properties of CSVW and RML+FnO that can be used to address the challenges of dealing with tabular data in a virtual OBDA approach

Challenges	Relevant Properties
Describe the corresponding concept (LS)	rr:class
Describe the corresponding property (LS)	rr:predicateMap
Add header to a CSV file (H)	csvw:rowTitles
Column datatype (LS)	csvw:datatype
Constraining values (H)	csvw:minimum, csvw:maximum
Specify the format of a column (H)	csvw:format
Specify a join (H)	rr:refObjectMap, csvw:foreignKeys
Transform value (H)	fnml:functionValue
Support for multiple values in one cell (N)	csvw:separator
Primary key (N)	csvw:primaryKey
Default for missing values (H)	csvw:default
Specify NULL values (H)	csvw:null
Specify NOT NULL constraint (LS)	csvw:required
Specify columns to be transformed (H)	rr:reference, rr:template

a widely used mapping language in virtual OBDA approaches. RML [21] extends R2RML; it provides support to a variety of data formats, e.g., XML, CSV, and JSON. Both languages provide basic transformation functions to concatenate strings, which are especially useful for generating URIs from columns/fields of the dataset. Recently, RML has been integrated with the Function Ontology (FnO) [22] to support other types of transformations. Additionally, for tabular data, CSVW metadata [9] is a W3C Recommendation to describe tabular datasets. Although there are other proposals in the state of the art to deal with some of the aforementioned challenges [8, 19], Morph-CSV relies on these two proposals because they cover the identified challenges. Additionally, this election is supported by the fact that CSVW is a recommendation from the W3C and RML+FnO (in addition of being an extended version of a W3C recommendation) has been previously applied in other projects [7, 12] and is widely used by several materialization engines, e.g., RMLMapper⁸, CARML⁹ and RocketRML [23]. Finally, relevant benefits of these annotations is that both of them are defined in a declarative manner. Thus, the maintainability, the readability, and the understanding of the virtual OBDA approach is improved and independent from any specific programming language.

We now describe relevant properties of RML+FnO and CSVW, summarized in Table 2, which are useful to deal with the challenges identified:

- **Metadata.** The property `csvw:rowTitles` can be used to specify column names in case the first row is not used to specify them.
- **Transformation functions.** String concatenation functions are supported by both CSVW (`csvw:aboutUrl`, `csvw:valueUrl`) and the RML property (`rr:template`). In addition, more complex functions can be declaratively specified using RML+FnO, specifically, with the `fnml:functionValue` property. Finally, two special cases of transformation functions in the context of OBDA are related to how default values and NULL representations have to be generated in the RDB instance. These two cases can be handled by CSVW properties: `csvw:defaultValue` and `csvw:null`.
- **Domain Constraints.** CSVW allows for the specification of the datatype (`csvw:datatype` property) and format (`csvw:format` property) of tabular columns. CSVW also provides a couple of properties (e.g., `csvw:minimum` or `csvw:maximum`) to specify the range of numerical columns and a property `csvw:required` to specify the NOT NULL constraint over the column of a table.
- **Integrity Constraints.** In CSVW the property `csvw:primaryKey` can be used to declare explicitly the primary key of a table. As for the foreign key, the use of RML's `rr:joinCondition` can be seen as an indication that the parent column used over this rule could be a foreign key. CSVW provides an explicit way to declare

⁸<https://github.com/RMLio/rmlmapper-java>

⁹<https://github.com/carmil/carmil/>

whether a column is a foreign key, using the `csvw:foreignKeys` property.

- **Normalization.** The property `csvw:separator` from CSVW indicates the character used to separate multiple values in the cells of a CSV column, what is relevant when a CSV file is in 1NF. Multiple RML TriplesMap using the same data source can be used as an indication that the source contains multiple concepts (2NF).

4. The Morph-CSV Framework

The formal framework presented in [13] defines an OBDA specification as a tuple $P = \langle O, S, M \rangle$ where O is an ontology, S is the source schema, and M a set of mappings. Additionally, an OBDA instance is defined as a tuple $PI = \langle P, D \rangle$ where P is an OBDA specification and D is a data instance conforming to S . In a virtual OBDA framework, queries are posed over a conceptual layer and then translated to queries over the data layer using information in the mappings. There is a set of assumptions over the framework that support the possibility of doing query translation and ensuring semantic preservation in the process, together with the application of optimization techniques proposed in the state of the art. To motivate our proposal, we have to establish what are the main assumptions made in previous proposals and their impact when data is represented in tabular form.

4.1. OBDA assumptions

Analyzing the definition of OBDA of [13] and its extension for NoSQL databases defined in [24] we identified a set of assumptions made over the framework and their impact when the dataset is tabular:

- There is a native query language QL for D . For a tabular dataset, there is not a native query language for querying this format, which generates an important difference with other common formats for exposing raw data on the web such as JSON and XML as they include ways to query them (JSONPath, XPath). This is the main issue that needs to be solved in order to query tabular datasets in a virtual OBDA context and has a direct impact over the rest of assumptions, that have been solved in a naive manner.
- S typically includes a set of domain and integrity constraints. In the case of querying a tab-

ular dataset $D_{tabular}$, S is defined using column names extracted from $D_{tabular}$ and it does not include any type of constraint (neither domain nor integrity constraints). This has a negative impact not only in terms of query execution time but also over query result completeness as there will be queries that cannot be executed due to the lack of explicit domain constraints.

- D is an RDB instance or is a NoSQL database instance, equipped with an RDB wrapper that is able to provide a relational view over S and D . In the context of a tabular dataset $D_{tabular}$, $D=R_{wrapper}(D_{tabular})$ where $R_{wrapper}$ is a relational database wrapper that satisfies S .

4.2. From Virtual Tabular Dataset to OBDA instance

Based on the previous OBDA assumptions, we define the concepts and functions to address the problem of querying a tabular dataset in OBDA.

Definition 1. A virtual tabular dataset is defined as a tuple $VTD = \langle D_{tabular}, O, M, MD \rangle$ where $D_{tabular}$ is a tabular dataset that is composed of a set of data sources, defined as $\mathcal{D}_{tabular} = \{s_1, \dots, s_n\}$ and where each s_i is a tabular relation defined over the domains of the attributes $Att(s_i) = \{A_{i1}, \dots, A_{im}\}^{10}$, where m is the number of attributes of s_i . O is an ontology, and M is a set of global as view mappings between O and $schema(D_{tabular})^{11}$. MD is a set of metadata tabular (domain) annotations, where for each s_i there exists a set $\{(A_{i1}, Type(A_{i1})), \dots, (A_{im}, Type(A_{im}))\}$ in MD .

Given a VTD , we define the function $\theta(VTD) = PI$ where PI is an OBDA instance $PI = \langle P, D \rangle$ where $D=R_{wrapper}(D_{tabular})$ and $P = \langle O, S, M \rangle$ is an OBDA definition where S does not contain any type of constraint. We extend the function $\theta(VTD)$ with the aim of enhancing the virtual OBDA baseline approach over tabular data. We define $\theta^{++}(VTD)=PI$ as a function that extracts a set of constraints from M and MD and then applies them over $D_{tabular}$ to obtain PI . More in detail, the function can be expressed as $\theta^{++}(VTD)=\gamma(D_{tabular}, O, M, \psi(M, MD))$ where the function $\psi(M, MD) = C$ extracts a set of constraints from OBDA annotations for tabular data. Then, $\gamma(D_{tabular}, O, M, C)$ applies the constraints C

¹⁰A relation is defined as the subset of the Cartesian product of the domains of the attributes.

¹¹The set of the attributes of each tabular relation in $D_{tabular}$, i.e., $schema(D_{tabular}) = \{Att(s_1), \dots, Att(s_n)\}$

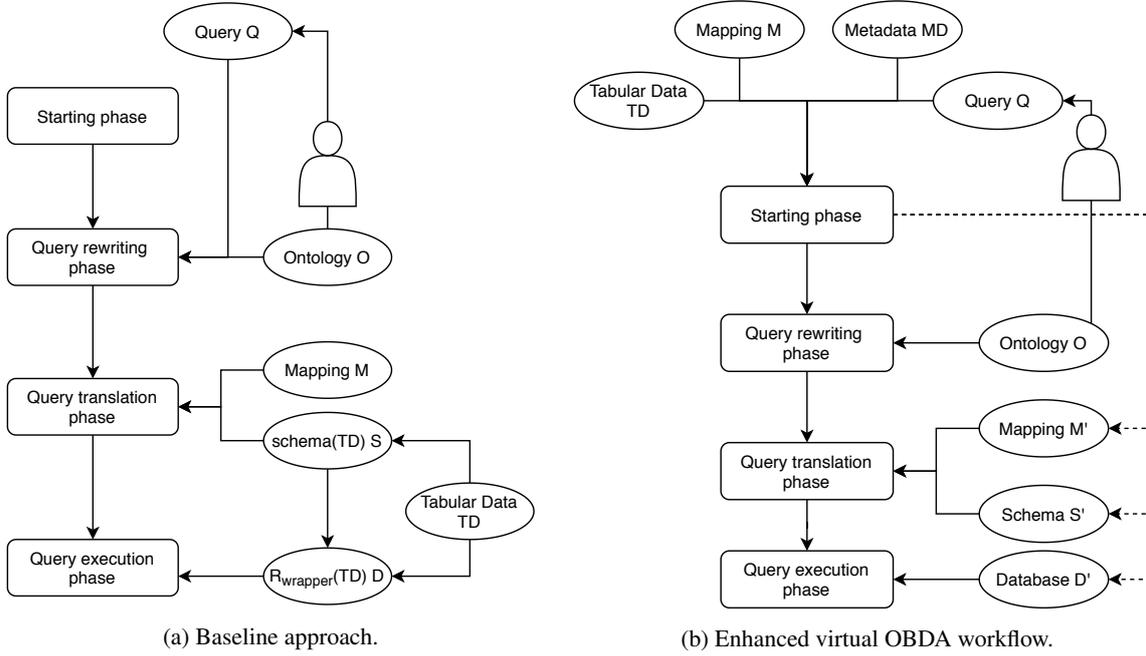


Fig. 2. **Virtual OBDA for tabular data approaches.** The baseline approach creates the schema and relational database instance extracting file and columns names from the tabular dataset. The proposed workflow exploits the information from the mapping rules and metadata to extract a set of constraints and applying them over the tabular data to generate the schema and the relational database instance.

over $D_{tabular}$ to create a relational database schema S' and its corresponding instance D' . In summary, the final output is an OBDA instance $PI' = \langle P', D' \rangle$, where D' is a relational database instance that is compliant with the main assumptions of the OBDA framework and $P' = \langle O, S', M \rangle$ where S' contains a set of domain and integrity constraints (see Figure 2b).

Constraints are conjunctive rules specified for tabular data that restrict the valid data in one or more tables. C is a set of constraints, where each constraint c is a logical statement that expresses the condition that needs to be satisfied by the data in order to be valid. Each constraint is applied through a function.

Example 1. CSVW allows expressing a primary key constraint for a table. The function $\psi(M, MD) = C$ generates the corresponding constraints in the form of a function $primaryKey(t, a)$ that applies this constraint to a source t and a set of columns a , and generates a primary key in the output schema.

Given an OBDA instance $PI = \langle P, D \rangle$, we define the function $eval(Q, PI)$, that retrieves a SPARQL answer set that is the result of the translation of Q from SPARQL to SQL using the mapping rules M defined in P , and then evaluating the query directly over D .

4.3. Problem statement and solution

Based on the preliminaries and assumptions made over the OBDA framework, we now define the problem that we address in this paper and Morph-CSV, our proposed solution.

Problem statement: Given a VTD , the problem of OBDA query translation over tabular data is defined as the problem of explicitly enforcing implicit constraints C extracted from mapping rules M and metadata MD on a tabular dataset $D_{tabular}$, such that:

- The number of results obtained in the evaluation of the SPARQL query Q over the function $eval(Q, \theta^{++}(VTD))$ is equal or greater than the number of results in the evaluation of the same query Q over the function $eval(Q, \theta(VTD))$, i.e., $\#answers(eval(Q, \theta^{++}(VTD))) \geq \#answers(eval(Q, \theta(VTD)))$.
- The total execution time of evaluating a SPARQL query Q over $eval(Q, \theta^{++}(VTD))$ is decreased compared to the evaluation of the same SPARQL query Q over the function $eval(Q, \theta(VTD))$, i.e., $time(eval(Q, \theta^{++}(VTD))) \leq time(eval(Q, \theta(VTD)))$.

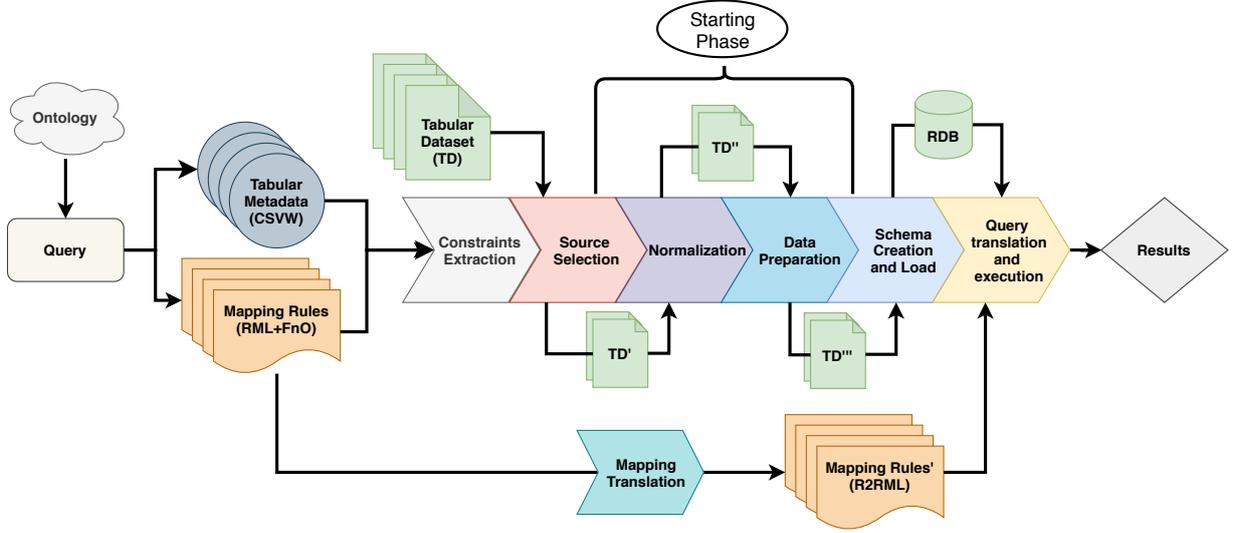


Fig. 3. **The Morph-CSV Framework.** Morph-CSV extends the starting phase of a typical OBDA system including a set of steps for dealing with the identified tabular data querying challenges. The framework first, extracts the constraints from mappings and tabular metadata and then, implements them in a set of operators that are ran before executing the query translation and query execution phases. The mapping rules are translated accordingly to the modified tabular dataset to allow its access by the underlying OBDA engine.

Proposed solution: We propose Morph-CSV, an alternative to the traditional OBDA workflow for query translation when the input is a tabular dataset (see Figure 2b and Appendix A). Morph-CSV relies on the function $eval(Q, \theta^{++}(VTD, \psi(M, MD)))$, to apply the tabular dataset constrains. Thus, Morph-CSV extends a typical OBDA workflow by including a set of steps for a maintainable extraction and efficient application of constraints. The workflow proposal is as follows:

- **Constraint Extraction:** the evaluation of the function $\psi(M, MD)$ produces as output the set of constraints C ; it exploits the information defined in the annotations of M and MD , i.e., the set of metadata tabular annotations and mapping rules, respectively. At implementation level they are expressed as CSVW specifications and RML+FnO mapping rules.
- **Source Selection:** in this step the sources required to evaluate the SPARQL query Q are selected. The required data sources correspond to the set of sources in the result of unfolding [3] Q according to the mapping rules in M .
- **Normalization:** metadata and mapping rules are utilized to extract functional dependencies between the attributes of the data sources. The algorithm by Beeri et al. [25] is followed to transform tabular data sources into tabular relations that meet 3 Normal Form (3NF).

- **Data Preparation:** application of the transformation functions based on the extracted domain constraints and on a set of optimization techniques that adapt the ideas proposed in [26] to a virtual OBDA environment.
- **Schema Creation and Load:** creation of the schema and loading the data into the database instance applying a set of rules for index creation.
- **Query Translation and Execution:** the evaluation of the query Q is delegated to any OBDA SPARQL-to-SQL engine.

We show the workflow of Morph-CSV in Figure 3 with the inputs and outputs of each step.

4.4. Steps performed in the Morph-CSV framework

We describe in detail the steps proposed in Morph-CSV together with an example extracted from the benchmark for virtual knowledge graph access, Madrid-GTFS-Bench, using the query shown in Figure 4a, the GTFS feed from the the Madrid metro as source data, and the corresponding RML+FnO mapping rules and CSVW annotations¹².

Constraint Extraction

The first step performed by Morph-CSV is the extraction of the constraints that are applied to improve

¹²Resources at: <https://github.com/oeg-upm/gtfs-bench>

Table 3
Summary of constraints, corresponding functions and OBDA annotations applied by Morph-CSV

Step	Constraint/Improvement	Rule/Annotation	Function	Challenge
Extraction	Reduce search space	SSG from Query	select_annotations	Selection
		Mapping Rules	select_sources	
Data Normalization	2NF	csvw:separator	split	Normalization
	3NF	TriplesMap with same source	cut	
Data Preparation	Standarization	csvw:null, csvw:default csvw:format, etc.	sub	Heterogeneity
		fnml:functionValue	create	
	Duplicates	-	duplicates	
Schema Creation and Load	Primary Key	csvw:primaryKey	primaryKey	Lightweight Schema
	Foreign Key	csvw:foreignKey	foreignKey	
	Data Type	csvw:datatype	datatype	
	Index	selectivity on mapping join conditions	index	

query execution and completeness. Morph-CSV benefits from having declarative and standard approaches to generalize this step: CSVW [9] for the metadata; and RML+FnO [7] for mapping rules and specific transformation functions. Thus, maintainability, understandability and readability of this process are improved in comparison with ad-hoc pre-processing approaches.

Most of the constraints such as PK-FK relations, datatypes or NULL values are explicitly declared in the metadata of the sources. However, there are a set of implicit constraints such as the conditions for the normalization of sources and the creation of indexes, that require complex rules to extract them and that are explained in detail in the corresponding steps. The summary of the constraints, associated functions, and properties used from OBDA annotations to extract them, are shown in Table 3.

Source Selection

The second step is to select the relevant sources to answer the input query. The baseline approach delegates this step to the RDBMS: it loads all the sources of the dataset in the RDB instance because it does not have information about what sources are going to be queried. This has a negative impact in the total execution time of a query. Taking the input mapping rules, Morph-CSV performs query unfolding, and pushes down source selection by executing the function $select(Q, M)$, divided into two main steps. First, Morph-CSV performs an operation to select only the relevant annotations for answering the input query,

$select_annotations(Q, M)$. It first creates the set of star shaped groups $SSG_1 \dots SSG_n$ of the query [27] (triple patterns with the same subject)¹³. Then, for each SSG_i and $rr:TriplesMap TM_j$ defined in M , the engine selects TM_j when the predicates in SSG_i are contained in the set of $rr:PredicateObjectMap$ (POMs) defined in TM_j . Finally, for each selected $rr:TriplesMap TM_j$, Morph-CSV only selects the POMs according to the predicates defined in the SSG_i , hence, removing from each TM_j irrelevant rules for the input query. Using these mapping rules M' , only relevant metadata annotations are also selected MD' . The obtained mapping rules M' and annotations MD' by this step substitute the original ones in VTD . An example of this step is shown in Figure 4, where the input query asks for trips, their route type, routes names and corresponding time frequencies. Morph-CSV first creates the SSGs, 3 in this case, and using the predicates of each SSG, the $rr:TriplesMap$ are selected from the general GTFS mapping document, discarding the rest of the rules. Then, it only selects the necessary POMs for evaluating the query such as $gtfs:startTime$, $gtfs:shortName$ and $gtfs:routeType$ (Figure 4b).

Second, Morph-CSV runs $select_sources(M)$, where it projects, from the input $D_{tabular}$, the sources and columns that are referenced in M , hence, relevant sources for the input query. The output of this func-

¹³As usual in these approaches, we assume bounded predicates in the triple patterns

```

1 SELECT ?trip ?routeName ?routeType ?startTime ?endTime ?code WHERE {
2   {
3     ?trip a gfts:Trip .
4     ?trip gfts:route ?route .
5     ?frequency a gfts:Frequency .
6     ?frequency gfts:startTime ?startTime .
7     ?frequency gfts:endTime ?endTime .
8     ?frequency gfts:trip ?trip .
9   }
10  {
11    ?route a gfts:Route .
12    ?route gfts:shortName ?routeName .
13    ?route gfts:routeType ?routeType .
14    ?routeType gfts:routeTypeCode ?code .
15  }
16 }

```

(a) Input SPARQL query.

```

11 frequencies:
12 sources:
13   - [frequencies.csv-csv]
14 s: mbench:freq/$(trip_id)-$(start_time)
15 po:
16   - [a, gfts:Frequency]
17   - [gfts:startTime, $(start_time)]
18   - [gfts:endTime, $(end_time)]
19   - [gfts:headSecs, $(headway_secs)]
20   - [gfts:exactTimes, $(exact_times)]
21   - p: gfts:trip
22 o:
23   - mapping: trips
24   condition:
25     function: equal
26     parameters:
27       - [str1, $(trip_id)]
28       - [str2, $(trip_id)]
29
30 trips:
31 sources:
32   - [trips.csv-csv]
33 s: mbench:trips/$(trip_id)
34 po:
35   - [a, gfts:Trip]
36   - [gfts:headsign, $(trip_headsign)]
37   - [gfts:shortName, $(trip_short_name)]
38   - [gfts:direction, $(direction_id)]
39   - [gfts:block, $(block_id)]
40   - p: gfts:route
41 o:
42   - mapping: routes
43   condition:
44     function: equal
45     parameters:
46       - [str1, $(route_id)]
47       - [str2, $(route_id)]
48
49 routes:
50 sources:
51   - [routes.csv-csv]
52 s: mbench:routes/$(route_id)
53 po:
54   - [a, gfts:Route]
55   - [gfts:shortName, $(route_short_name)]
56   - [gfts:longName, $(route_long_name)]
57   - [dct:description, $(route_desc)]
58   - [gfts:routeUrl, $(route_url)-iri]
59   - [gfts:color, $(route_color)]
60   - [gfts:textColor, $(route_text_color)]
61   - p: gfts:agency
62 o:
63   - mapping: agency
64   condition:
65     function: equal
66     parameters:
67       - [str1, $(agency_id)]
68       - [str2, $(agency_id)]
69
70 p: gfts:RouteType
71 o:
72   - mapping: route-type
73   condition:
74     function: equal
75     parameters:
76       - [str1, $(route_type)]
77       - [str2, $(route_type)]
78
79 route-type:
80 sources:
81   - [routes.csv-csv]
82 s: CONCAT(gfts:TRANS, $(route_type))
83 po:
84   - [a, gfts:RouteType]
85   - [gfts:routeTypeCode, $(route_code)]

```

(b) Mapping rules selection.

Fig. 4. Selection of Mapping Rules. Based on the SPARQL query relevant rules are selected (in bold), the rest are discarded.

tion generates a set of new tabular sources $s_1 \dots s_n$ that substitute the original $D_{tabular}$ of VTD . Following the previous example, Figure 5 shows the selection of the relevant columns of source *routes.csv*, where Morph-CSV has the original source as input (Figure 5a), and discards the unnecessary columns of the source based on the mapping rules, obtaining as output the source with the relevant columns for evaluating the input query (Figure 5b). Note that in this step, unnecessary sources from the input GTFS feed such as *agency.csv* and *stops.csv* are also discarded.

Normalization

There are two functions for performing data normalization. The first one is the treatment of multi-values in a column. In this case, Morph-CSV performs the function $split(A_{ij}, sep)$ where A_{ij} is the multi-valued column of source s_j and sep is the character defined in the CSVW metadata using the `csvw:separator` property. The output is a modified VTD with a new source s_t containing the sepa-

```

1 route_id,agency_id,route_short_name,route_long_name,route_desc,route_type,route_code,route_url,route_color
2 4_1,CRTM,1,Pinar de Chamartin-Valdecarros,,1401,crtm.es/metro/4_1,2DBEF0
3 4_2,CRTM,2,Las Rosas-Cuatro Caminos,,1401,crtm.es/metro/4_2,ED1C24
4 4_3,CRTM,3,Villaverde Alto-Moncloa,,1401,crtm.es/metro/4_3,FFD000
5 4_4,CRTM,4,Pinar de Chamartin-Argüelles,,1401,crtm.es/metro/4_4,B65518
6 5_C1,CRTM,C1,P.Pio-AeropuertoT4,,2,109,http://www.crtm.es/cercanias/5_1,4FB0E5,FFFFFF
7 5_C2,CRTM,C2,Guadalajara-Chamartin,,2,109,http://www.crtm.es/cercanias/5_2,008B45,FFFFFF
8 5_C3,CRTM,C3,Aranjuez-Escorial,,2,109,http://www.crtm.es/cercanias/5_3,9F2E66,FFFFFF
9 5_C4,CRTM,C4,Parla-Colmenar Viejo,,2,109,http://www.crtm.es/cercanias/5_4,005AA3,FFFFFF

```

(a) Original routes.csv input source.

```

1 route_id,route_short_name,route_type,route_code
2 4_1,Pinar de Chamartin-Valdecarros,1,401
3 4_2,Las Rosas-Cuatro Caminos,1,401
4 4_3,Villaverde Alto-Moncloa,1,401
5 4_4,Pinar de Chamartin-Argüelles,1,401
6 5_C1,P.Pio-AeropuertoT4,2,109
7 5_C2,Guadalajara-Chamartin,2,109
8 5_C3,Aranjuez-Escorial,2,109
9 5_C4,Parla-Colmenar Viejo,2,109

```

(b) Output of routes.csv source.

Fig. 5. Source Selection. Based on the selection of the rules, only `route_id` and `trip_id` columns are selected, discarding the rest fields.

rated values, and a updated mapping document M with a new `rr:TriplesMap` TM_t generated for the new source s_t and a `rr:joinCondition` between the `rr:TriplesMap` of s_j , TM_j and TM_t . The application of this function is known as the normalization step for second normal form (2NF) [28].

The second function is the treatment of multiple entities in the same source. Morph-CSV takes the mapping rules and executes the function $cut(\mathcal{M}, \mathcal{D}_{tabular})$. This function analyzes mapping rules \mathcal{M} , and performs a 3NF [28] normalization step over $D_{tabular}$ when there are two sets of mapping rules (TM_j and TM_i) that have the same source, and the intersection of their columns in the rules only contains the join condition references. Following a similar approach as in 2NF, the output is a modified VTD with a set of new sources $s_1 \dots s_n$, each one with the corresponding columns of each entity. For example, in Figure 6 we show the 3NF normalization of the *routes.csv* file, that generates an auxiliary source for the `rr:TriplesMap` with the `gfts:RouteType` entity data (Figure 6), removing that information for the *routes.csv*. In several data integration approaches, normalization steps are not taken into account in order to improve query execution (reducing the number of joins among sources). However, in the case of RDF, where each entity of a class has a unique URI (subject), joins cannot be reduced (see input mapping of Figure 4b). This means that taking into account normalization steps in an OBDA context not only helps to improve query completeness, but also helps to improve performance. Additionally, normalization is also

essential for allowing Morph-CSV to efficiently run data preparation steps, as we show in the next step.

```
route_id,route_short_name,route_type
4_1,Pinar de Chamartín-Valdecarros,1
4_2,Las Rosas-Cuatro Caminos,1
4_3,Villaverde Alto-Moncloa,1
4_4,Pinar de Chamartín-Argüelles,1
5_C1,P. Pío-AeropuertoT4,2
5_C2,Guadalajara-Chamartín,2
5_C3,Aranjuez-Escorial,2
5_C4,Parla-Colmenar Viejo,2
```

(a) Routes.csv after 3NF normalization step.

```
route_type,route_code
1,401
1,401
1,401
1,401
2,109
2,109
2,109
2,109
```

(b) Route_type.csv file generated with Morph-CSV.

Fig. 6. **Normalization.** 3NF Normalization step over the *routes.csv* file generating other file with the data for `gtfs:RouteType` class.

Data preparation

In this step, Morph-CSV addresses the challenge of *Heterogeneity* and executes three different functions: *duplicates*, *sub* and *create*. First, Morph-CSV removes all duplicates in the raw data, not only the original ones, but also other duplicates that can appear during the normalization step (see Figure 6b). It applies the ideas described in [26], performing $duplicates(s_j)$ where s_j is a source in $D_{tabular}$. As it has already been demonstrated in [26], this step not only has a high impact on the behavior of these engines, but in this case, it also reduces the number of operations performed by Morph-CSV *sub* and *create*, as they are defined as deterministic functions. The first one is defined as $sub(exp(A_{ij}), val)$ where $exp(A_{ij})$ is a boolean function over column A_{ij} of source s_j that when true, the value of A_{ij} is substituted by val . There are multiple substitution functions that Morph-CSV executes such as default values, null values and date formats. The second function creates a new column in a specific source s_j . It is defined as

$create(c(A_{nj}, \dots, A_{mj}))$, where $c(A_{nj}, \dots, A_{mj})$ is the application of a set of transformation functions over the columns A_{nj}, \dots, A_{mj} in source s_j . This function is used to push down the application of ad-hoc transformation functions, usually defined inside the mapping rules [7, 8], thus, avoiding the incorporation of them inside the SQL translated query. In Figure 7 we show the *route_type.csv* file after the execution of this step. First, Morph-CSV removes the duplicates of the file obtaining as output a file with only two rows. Then, it executes the transformation function defined in the mapping rules and creates a new column in the file, generating the desired value for the subject of the class according to the LinkedGTFSS ontology, “Subway”. Additionally, the engine substitutes the definition of the transformation functions in the mapping rules by a reference to the created column. In this manner, Morph-CSV efficiently performs the *sub* and *create* functions directly over the raw data and together with the normalization step. Thus, the number of joins in the input query is reduced.

```
route_type, route_code, route_type_fn
1, 401, Subway
2, 109, Train

route-type:
sources:
- [routes_types.csv~csv]
s: gtfs:$(route_type_fn)
po:
- [gtfs:routeTypeCode,$(route_code)]
```

Fig. 7. Data preparation of *route-types.csv* file.

Schema Creation and Load

The final step before translating and executing the query is the creation of an SQL schema applying the rest of the identified constraints, and loading the selected tabular data sources. Besides the typical integrity constraints that can be extracted from CSVW annotations (PK/FK), Morph-CSV implements a rule for creating indexes in the RDB instance in order to optimize the execution of query joins. In tabular datasets, it is common that the join conditions defined in the mapping rules are based on columns that are not part of PK-FK relations; thus, they are not indexed and OBDA optimizations do not have the desired effect. To address this problem, Morph-CSV gets the `rr:child` and `rr:parent` references of the mapping rules and calculates their selectivity on the fly. Then, taking this selectivity into account Morph-CSV decides to create, or not, an index over these columns. Figure 8 shows the RDB schema generated by Morph-CSV for the input query in Figure 4a, with the applied domain and integrity constraints.

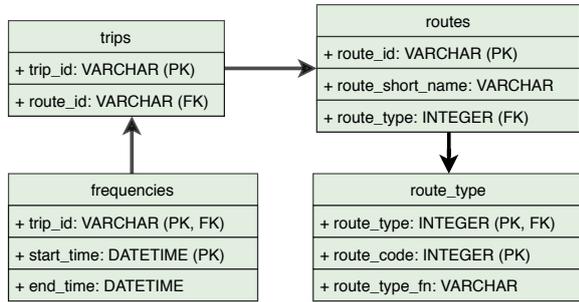


Fig. 8. **Generated schema.** The schema generated by Morph-CSV extracting domain and integrity constraints from the annotations and based on the identified sources selected from the input query.

There are two main points that make the contributions of Morph-CSV relevant: (i) it incorporates the steps to the standard OBDA workflow without modifying the rest of the steps, hence, it can also benefit from optimizations in other steps of the workflow like query rewriting (reasoning) [29] or query translation (SPARQL-to-SQL) [6], and (ii) the reliance of the approach on declarative and standard annotations for OBDA allows generalizing the proposed steps, usually solved in an ad-hoc manner, not only automatizing the process but also improving its maintainability, understandability and readability.

5. Evaluation

This section reports on the results of the empirical evaluation conducted to test the effect of respecting constraints, on the fly, during OBDA query translation over tabular data. Our aim is to answer the following research questions: **RQ1**) What is the effect of combining different types of constraints over a tabular dataset? **RQ2**) What is the impact of the constraints when the tabular dataset size increases? **RQ3**) What is the effect of different levels of data heterogeneity in the extraction and application of constraints? To answer these questions, we have performed three evaluations in different domains: e-commerce, transportation, and biology. Our first evaluation is in the e-commerce domain, in which we used the Berlin SPARQL Benchmark (BSBM) [14]. Our second evaluation is in the transportation domain in which we used the GTFS-Madrid-Bench. GTFS-Madrid-Bench benchmark focuses on measuring the performance of ontology based data access for heterogeneous data sources, based on the publicly-released public transportation data in GTFS format. One of the resources

provided by GTFS-Madrid-Bench is a tabular dataset together with its corresponding mappings and annotations. Finally, our third evaluation is in the domain of biological data, in which we extend one of our previous proposals [30] for the generation of an OBDA layer over Bio2RDF tabular datasets. Appendix B presents the features of the queries together with the constraints and number of sources used by Morph-CSV. In all of the evaluations the common configurations are:

Engines. The baselines of our study are two open source OBDA engines: Ontop^{14,15} v3.0.1 and Morph-RDB v3.9.15¹⁶. To evaluate the baseline approach, we manually generate the relational database schemas of each benchmark without any kind of constraints, and measure the load and query execution times. In order to measure the impact of the additional steps proposed by Morph-CSV^{17,18}, we integrate our solution on top of the two OBDA engines. To ensure the reproducibility of the experiments, we also provide all of the resources in a docker image. In order to test the number of answers, we use the gold standards provided by both benchmarks in RDF, loaded in a Virtuoso triple store. **Metrics.** We measure the loading time of each query and the total query execution time (including the steps proposed by Morph-CSV or baseline when it corresponds), and the number of answers obtained (see Appendix C). Each query was executed 5 times with a timeout of 2 hours in cold mode, that means that the corresponding database is generated each time a query is going to be evaluated in order to ensure up to date number of answers. The experiments were run in an Intel(R) Xeon(R) equipped with a CPU E5-2603 v3 @ 1.60GHz 20 cores, 64GB memory and with the O.S. Ubuntu 16.04LTS.

5.1. BSBM

The Berlin SPARQL Benchmark [14] is one the most popular benchmarks in the Semantic Web field that not only tests the performance of RDF triple stores, but also tests approaches that perform SPARQL-to-SQL query translations providing a RDB instance. It is the chosen benchmark to test the capabilities of many state-of-the-art OBDA engines [5, 6, 12].

¹⁴<https://github.com/ontop/ontop>

¹⁵We modified the default configuration of Ontop extending the maximum used memory from 512Mg to 8Gb

¹⁶<https://github.com/oeg-upm/morph-rdb>

¹⁷<https://doi.org/10.5281/zenodo.3731941>

¹⁸<https://github.com/oeg-upm/morph-csv>

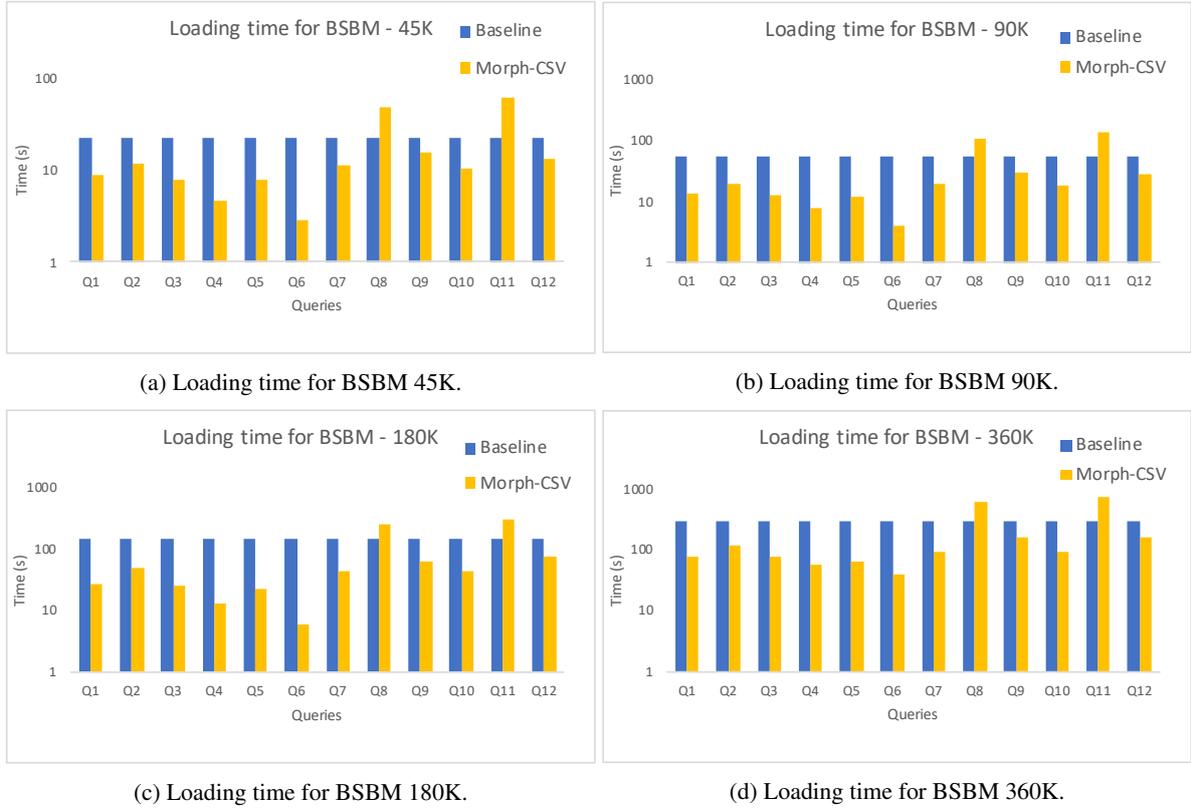


Fig. 9. **Loading Time of Tabular Datasets in BSBM.** Loading time in seconds of the tabular datasets from the BSBM benchmark with number of products 45K, 90K, 180K and 360K. The baseline approach (blue columns) is constant for each dataset and query, while Morph-CSV (orange columns) depends on the query and number of constraints to be applied over the selected sources.

Datasets, annotations and queries. In order to test our proposal we decided to adapt BSBM, extracting the tabular data sources in CSV format from the SQL generated instances. Additionally, we create the corresponding mapping rules in RML and the metadata following the CSVW specification. We measure the loading time of the two proposals (baseline and Morph-CSV) for each query in the benchmark. Since the focus Morph-CSV is not the improvement of the support of SPARQL operators in the query translation process, we only select the queries of the benchmark that include supported operators and operations by each engine. This means that Morph-RDB will be evaluated over the queries Q1, Q2, Q3, Q4, Q5, Q6, Q7, Q8, Q9, Q10 and Q12 and Ontop will be evaluated over Q1, Q3, Q4, Q5 and Q10, both of them using the corresponding R2RML mapping document. For the baseline approach we manually create the RDB schema without any constraints.

5.1.1. BSBM Results

Loading Time. The results of the load time for each query and dataset size are shown in Figure 9. The main difference between the two methods is that while the loading time for the baseline approach is constant for each size, Morph-CSV loading time depends on several input parameters such as the query and the number and type of constraints. It could be understandable that the application of a set of constraints over the raw data in order to improve query performance and completeness would have a negative impact in the loading time of our proposal. This happens in queries Q8 and Q11, where the number of sources and the application of the constraints (mainly integrity constraints) impact negatively on the loading time of the data in the RDB instance in comparison with the baseline approach. However, in the rest of the queries, the Morph-CSV step is focused on the selection of constraints, sources and columns, and exploiting the information in query and mapping rules improves the loading time for each query in comparison with the baseline loading time. This means that, although the engine is including



Fig. 10. **Query execution Time of Tabular Datasets in BSBM with Morph-RDB.** Execution time in seconds of the tabular datasets from the BSBM benchmark with scale values 45K, 90K, 180K and 360K. The baseline Morph-RDB approach (blue columns) is compared with the combination of Morph-CSV together with Morph-RDB (orange columns). Red marks on the top of the columns mean a timeout (72000 seconds) in query execution.

a set of additional steps during the starting phase of an OBDA system, the application of these steps only over the data that is required to answer the query, has a positive impact in the total query execution time. Additionally, we can observe that Morph-CSV is able to process, apply the different constraints and generate the corresponding instance of the RDB for any query.

Evaluation Time with Morph-RDB. The query execution time using Morph-RDB as the back-end OBDA engine is shown in Figure 10. The first remarkable observation can be seen in query Q5. Although this query contains operators supported by Morph-RDB, the engine reports an error when evaluating the query over the database generated by the baseline approach, because it is not able to evaluate the arithmetic expressions in the FILTER clauses. On the contrary, the datatype of each column in the database generated by Morph-CSV is properly defined, making it possible for Morph-RDB to evaluate the query without any problem and obtaining the expected results. Another remarkable difference is in query Q2, which contains

a large number of joins, Morph-RDB reports a timeout error for 180K and 360K with the database generated by the baseline approach. However, it is still able to evaluate this query in reasonable time over the databases generated by Morph-CSV. The effect of the application of integrity constraints in the generation of the RDB instance can also be seen in most of the queries (i.e., Q1, Q2, Q3, Q6, Q9, Q10) reducing considerably the query execution time in the database generated by Morph-CSV in comparison with the baseline approach. There are cases (i.e., Q4, Q7, Q12) where the amount of data to retrieve is large, minimizing the effect of the optimizations. Finally, there are cases where optimizations over the indexes cannot be applied (e.g. asking for all the properties of a class). We observe this behavior in Q8, although the difference between the two approaches is not very relevant and is maintained across the datasets.

Evaluation Time with Ontop. The query execution time using Ontop as the back-end OBDA engine is shown in Figure 11. Like Morph-RDB, Ontop needs

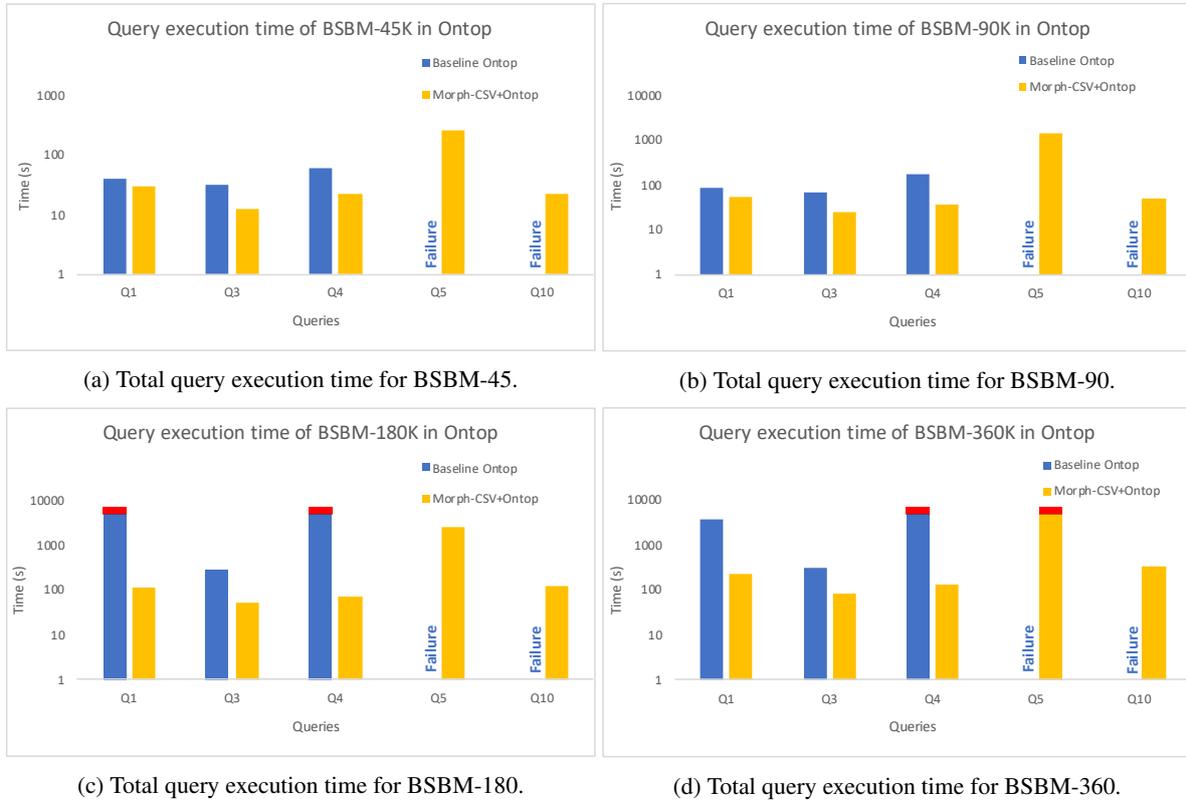


Fig. 11. **Query execution Time of Tabular Datasets in BSBM with Ontop.** Execution time in seconds of the tabular datasets from the BSBM benchmark with scale values 45K, 90K, 180K and 360K. The baseline Ontop approach (blue columns) is compared with the combination of Morph-CSV together with Ontop (orange columns). Red marks on the top of the columns mean a timeout (72000 seconds) in query execution.

the Morph-CSV generated databases to be able to evaluate Q5 due to the arithmetic expressions of its FILTER operators. Additionally, it also fails in Q10 because it cannot process a FILTER with a date value. In the rest of the queries (Q1, Q3, Q4) we can see that the query evaluation time in Ontop with Morph-CSV is lower than the query evaluation time over the baseline database. Note that in larger databases (180K and 360K), Q1 and Q4 can only be evaluated over the databases generated by Morph-CSV.

As mentioned in the Ontop repository page¹⁹, integrity constraints are essential for the correct behavior of the engine. Although it is out of the scope of this paper, we observe in our experiments that the main reason why Ontop is only able to answer half of the queries in this benchmark, is related to some issues about maintaining the desirable properties [31] when translating R2RML mapping rules to its own mappings, called OBDA. The engine also fails to evalu-

ate queries with OPTIONAL clauses when there are NULL values in the answers, as they acknowledged, it is possible that this support has not been implemented in the engine [32].

5.2. GTFS-Madrid-Bench

The GTFS-Madrid Benchmark²⁰ consists of an ontology, an initial dataset of the metro system of Madrid following the GTFS model, a set of mappings in several specifications, a set of queries according to the ontology that cover relevant features of the SPARQL query language, and a data scaler based on a state of the art proposal [33].

Datasets, annotations and queries. We select the tabular sources of this benchmark (i.e., the CSV files) and we scale up the original data in several instances (scale factors 10, 100 and 1000). Each generated dataset is

¹⁹<https://github.com/ontop/ontop/wiki/MappingDesignTips>

²⁰Paper under review. Resources available at: <https://github.com/oe-g-upm/gtfs-bench>

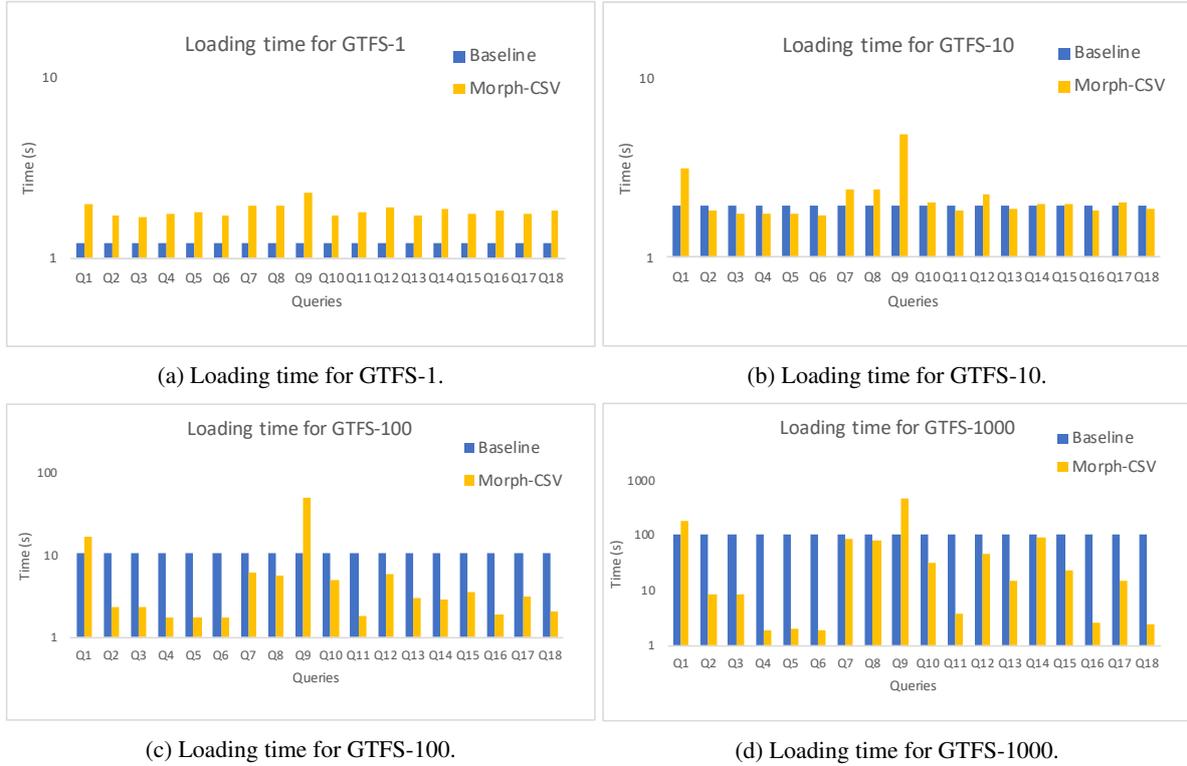


Fig. 12. **Loading Time of Tabular Datasets in GTFS.** Loading time in seconds of the tabular datasets from the Madrid-GTFS-Bench with scale values 1, 10, 100 and 1000. The baseline approach (blue columns) is constant for each dataset and query while Morph-CSV (orange columns) depends on the query and number of constraints to be applied over the selected sources.

denoted by GTFS- S where S is the scale factor. The resources of the benchmark already include the necessary mapping rules and tabular metadata. Like our previous evaluation with BSBM benchmark, we only select the queries with operators that are supported by each engine: Morph-RDB will be evaluated using queries Q1, Q2, Q4, Q6, Q7, Q9, Q12, Q13, Q14, Q17 and Ontop will be evaluated using queries Q1, Q2, Q3, Q4, Q5, Q7, Q9, Q13, Q14, Q17. The description and features of each query are also available online²¹.

5.2.1. Madrid-GTFS-Bench Results

Loading Time. The loading time of the GTFS-Madrid-Bench queries is shown in Figure 12. For GTFS-1 the baseline approach clearly has better performance than Morph-CSV. However, when the size of the datasets increases, the positive effects of applying constraints become more apparent. For most of the queries, the loading time needed by Morph-CSV is lower in comparison to the loading time in the base-

line approach. Additionally, similarly to BSBM, there are a set of queries where the application of integrity constraints has a negative impact on the loading time (queries Q1 and Q9).

Evaluation Time with Morph-RDB. The query execution time with Morph-RDB as the back-end OBDA engine is shown in Figure 13. Analyzing the results, we generally observe that the incorporation of Morph-CSV in the workflow of OBDA enhances query performance. With respect to the results of each query, we can observe that on the one hand the behavior of the engine over simple queries (Q1, Q2, Q7, Q12 and Q17) is similar. This is understandable as the selected data sources needed to answer the query do not include the application of several constraints (e.g. there are no joins in the query). On the other hand, in the case of complex queries such as Q4, Q6, Q9, Q13 and Q14, where several tabular sources are needed to answer the queries, the application of constraints has a better impact in comparison to the the baseline approach. For example, in the case of query Q9, Morph-RDB is not able to evaluate the query over the 10th scale database

²¹<https://github.com/oeg-upm/gtfs-bench/tree/master/queries>

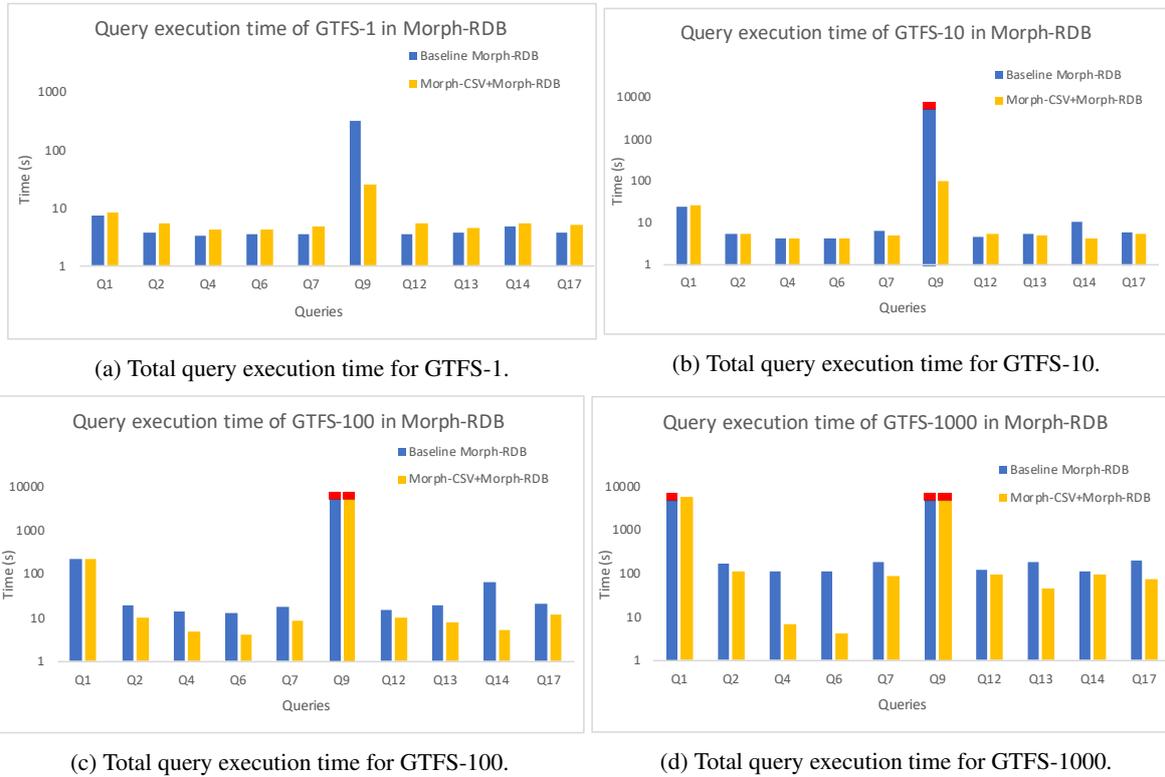


Fig. 13. **Query execution Time of Tabular Datasets in GTFS with Morph-RDB.** Execution time in seconds of the tabular datasets from the Madrid-GTFS-Bench with scale values 1, 10, 100 and 1000. The baseline Morph-RDB approach (blue columns) is compared with the combination of Morph-CSV together with Morph-RDB (orange columns). Red marks on the top of the columns mean a timeout (72000 seconds) in query execution.

generated by the baseline approach, while in the case of the database generated by Morph-CSV, the query can be answered in reasonable time. If we analyze the results obtained, we can observe that for small datasets (GTFS-1), the cost of applying the proposed steps of Morph-CSV impacts total execution time. However, when the size of the dataset increases, the baseline approach is impacted due to the fact that it has to load all of the input data sources in the RDB before executing the query, low performance is reported for GTFS-100 and GTFS-1000, including timeout in some queries of the latter. Thanks to the application of the constraints and to the source selection step, for Morph-CSV together with Morph-RDB, the return of the results of the queries has a high performance most of the time. In the cases where Morph-CSV reports a timeout (e.g., Q1 in GTFS-1000); it is because the extremely high number of obtained results cannot be handle by Morph-RDB.

Evaluation Time with Ontop. The experimental evaluation of the query execution in Ontop as the backend OBDA engine is shown in Figure 14. This engine

is more strict with datatypes in the RDB in comparison with Morph-RDB, and it is why Q2, Q5, Q7 and Q9 produce a failure in the execution over the databases generated by the baseline approach. All these queries have a FILTER clause on a specific datatype (e.g., date, integer, etc) and Ontop proceeds to check the domain constraints before executing the queries. Morph-CSV solves this problem by exploiting the annotations from the metadata and defines the correct datatypes of each column before evaluating the query. For the queries that can be answered by both approaches (Q1, Q3, Q4, Q13, Q14, Q17), the absence of integrity constraints has a negative impact in Ontop, resulting in lower execution time over the databases generated by Morph-CSV. Finally, in the case where Ontop is not able to evaluate the query under the defined threshold, we report it as a time-out.

5.3. Use Case: The Bio2RDF project

Bio2RDF is one of the most popular projects that integrates and publishes biomedical datasets as Linked

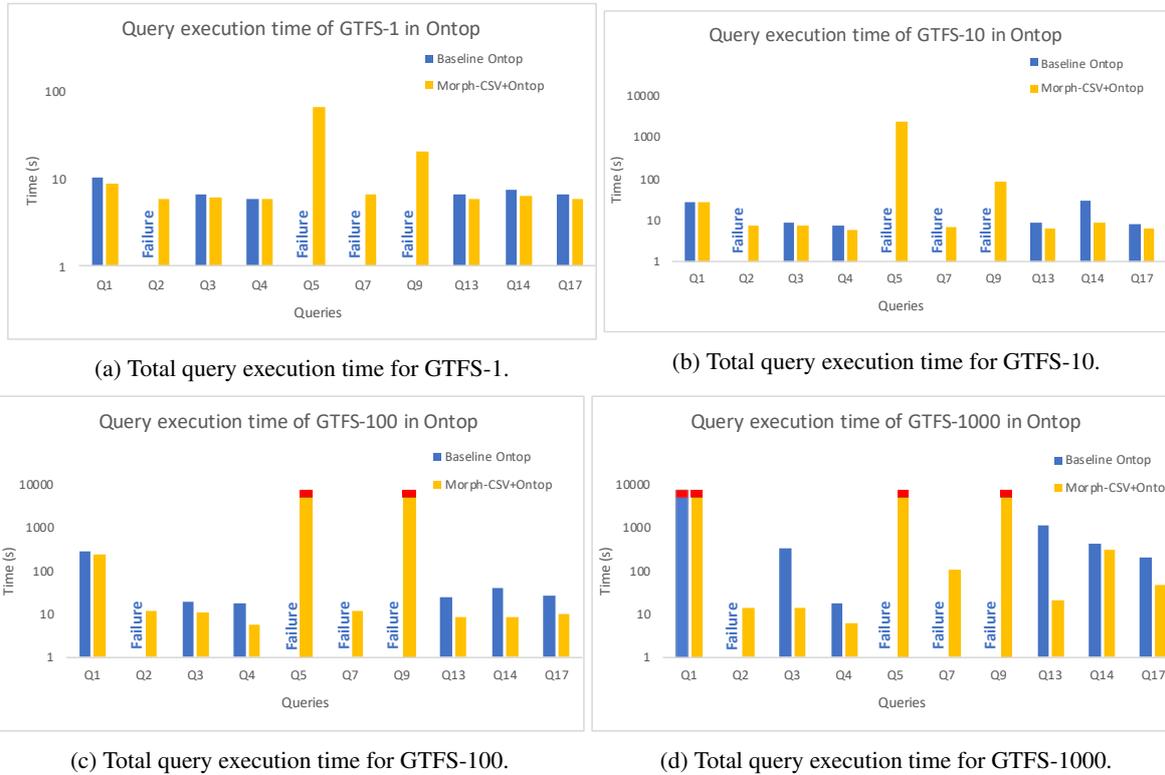


Fig. 14. **Query execution Time of Tabular Datasets in GTFS with Ontop.** Execution time in seconds of the tabular datasets from the Madrid-GTFS-Bench with scale values 1, 10, 100 and 1000. The baseline Ontop approach (blue columns) is compared with the combination of Morph-CSV together with Ontop (orange columns). Red marks on the top of the columns mean a timeout (72000 seconds) in query execution.

Data [15]. Its community has actively contributed to the generation of those datasets using ad-hoc programming scripts, such as PHP. In our previous work [30] we proposed an alternative way of generating the datasets using a set of declarative mapping rules to improve the maintainability, readability and understanding of the procedure. In comparison with the other benchmarks where the focus of the evaluation was the improvement of the query evaluation time, this real use case contains multiple heterogeneity challenges that, for example, enforce the application of ad-hoc transformation functions (i.e., mappings in the form of RML+FnO). Thus, with this use case we want to demonstrate the benefits of exploiting declarative annotations (metadata and mappings) over the raw data in order to improve query completeness and the need of incorporating the proposed steps for executing queries over real world data sources.

Dataset, annotations, and queries. Tabular datasets in CSV or Excel formats cover over 35% of the total datasets in the Bio2RDF project [30]. In order to test the capabilities of Morph-CSV, we select a sub-

set of the tabular datasets guaranteeing that they cover all of the identified challenges. Additionally, as far as we are aware, there is no standard benchmark over the Bio2RDF project; we also propose a set of SPARQL queries in order to exploit the selected data. Their main features are shown in Appendix B).

5.3.1. Bio2RDF Results

The results obtained for query evaluation in Bio2RDF are shown in Figure 15 with Morph-RDB as backend engine and in Figure 16 with Ontop. First, we can observe that there are no results for the baseline approach, this means it was not possible to create an RDB schema and load the input data. The main reasons are the heterogeneity problems of a real use case that do not exist in the previous evaluations. GTFS and BSBM have well formed and standard source data models. Problems such as the absence of column names, multiple formats of same datatype in different files (numbers, dates) and the use of delimiters inside the column data, make it impossible to generate the baseline approach without a manual and ad-hoc preprocessing step. However, exploiting declarative an-

1 notations, Morph-CSV is able to apply the proposed
2 workflow to this dataset, and successfully answer the
3 proposed queries with both back-end OBDA engines.
4 More in detail, we observe that due to the size of the
5 input datasets and the number of constraints to be ap-
6 plied, most of the total evaluation time of each query is
7 spent in the loading process. Contrary, query execution
8 is benefited by this previous step obtaining the results
9 in reasonable time for all of the queries.

12 6. Discussion of Experimental Results

14 We have run an experimental evaluation to analyze
15 what are the effects on the use of declarative annota-
16 tions to extract and apply constraints to enhance vir-
17 tual OBDA approaches. We have tested our approach
18 over three different cases: (i) a well known benchmark
19 (BSBM) from the e-commerce domain; (ii) a bench-
20 mark focused on a virtual OBDA approach for the
21 transport domain; and (iii) a real use case from the bio-
22 logical domain. We describe the main conclusions and
23 findings based on the results obtained:

- 25 – **Query complexity:** Clear benefits are obtained
26 from being able to analyze and take advantage of
27 the information provided by the input query, be-
28 fore translating and running it. It allows to only
29 select sources and constraints that are going to be
30 useful for answering the query, avoiding carrying
31 out additional and unnecessary functions over the
32 raw data. Together with the mapping rules, the
33 queries are essential to make relevant decisions
34 during the on-the-fly physical design of the RDB
35 instance (e.g., integrity constraints).
- 36 – **Data size:** The total query evaluation time is be-
37 ing impact from how the engine manages the
38 input dataset and the application of constraints.
39 The delegation of these operations to the RDBMS
40 system after loading the full dataset may not be
41 efficient enough. Morph-CSV pushes down the
42 source selection and the application of domain
43 constraints over the raw data. Although it incor-
44 porates a set of additional steps in comparison
45 with the baseline, the benefits in the query exe-
46 cution time by the SPARQL-to-SQL engine are
47 already demonstrated, enhancing the total execu-
48 tion time of the queries in most of the cases.
- 49 – **Declarative annotations:** The use of declara-
50 tive and standard mapping rules and metadata
51 makes it possible the generalization of the pro-

1 posal, avoiding ad-hoc and manual steps. It also
2 incorporates a set of important benefits for the
3 process such as the improvement of its maintain-
4 ability, readability, and understandability.

- 5 – **Querying raw data in OBDA:** Most of the data
6 shared on the web is currently raw data in well
7 known formats such as CSV, JSON, and XML.
8 Semantic Web and more specifically, OBDA
9 technologies, play a key role in starting to see the
10 web as an integrated database that can be queried.
11 With this approach, we demonstrate that query-
12 ing tabular data is: i) neither a trivial nor an easy
13 task that can be delegated to naïve querying ap-
14 proaches and ii) optimizations and improvements
15 can still be proposed taking advantage and ex-
16 ploiting current annotation proposals to not only
17 enhance performance but also completeness.

20 7. Related Work

22 In this section, we first refer to previous works in
23 data integration systems that precede the OBDA ap-
24 proach. Then, we refer to the general techniques used
25 in systems that handle raw data. Next, we describe cur-
26 rent Ontology Based Data Integration (OBDI) systems
27 that handle tabular data. Finally, we describe existing
28 tabular annotation languages and the use of transfor-
29 mation functions in mappings.

30 The most relevant concept that predates the OBDA
31 data integration approach is that of mediator [34],
32 defined in the early 90's by Wiederhold. In the pro-
33 posed architecture for information systems, mediators
34 form a middle layer that makes user applications in-
35 dependent of the data resources. The idea is to trans-
36 form heterogeneous data sources into a common data
37 model, which can then be processed and integrated.
38 Classical examples of systems that implemented the
39 original mediator architecture were TSIMMIS [35],
40 Information Manifold [36], and GARLIC [37]. The
41 problem of inconsistent formats is not new, and in
42 general mediators may convert attributes of several
43 sources into a common format. The TSIMMIS [35]
44 architecture includes a Constraints Manager compo-
45 nent which handles integrity constraints across dif-
46 ferent sources. The constraints manager supports the
47 definition of the interfaces that a source supports for
48 the constraint, e.g. a trigger, the specification of the
49 desired constraint, and the specification of the strat-
50 egy for enforcing the constraint or for detecting vi-
51 olations. Information Manifold [36] is an integration

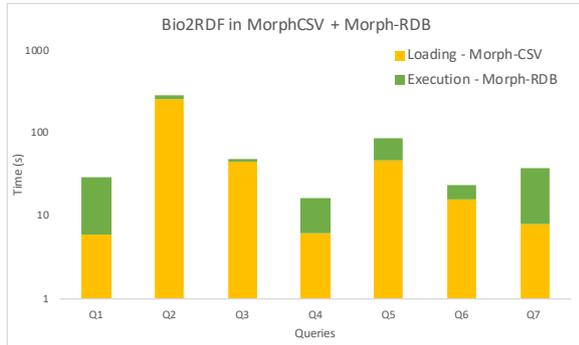


Fig. 15. **Query execution Time of Tabular Datasets in Bio2RDF with Morph-RDB.** Execution time in seconds of the tabular datasets from Bio2RDF. The loading time applying a set of constraints (yellow) and query execution with Morph-RDB (green).

system for heterogeneous sources on the Web. It uses source content and capabilities descriptions in order to prune the space of sources that are accessed to answer a query. Garlic [37] is a system that provides an integrated view over legacy data sources. Each source or repository has its own data model, schema, programming interface, and query capability. Each Garlic object has an interface and may have several implementations, corresponding to different data sources. The system uses these implementations to optimize and execute a query. Both these systems do not handle domain constraints nor constraints across sources.

The work presented in [38] provides a toolkit for the generation of wrappers for web-accessible heterogeneous sources (may be represented as HTML tables) through the description of their capabilities. It provides a specification language to define the capability for each source, and generates a wrapper according to this specification. It also provides a graphical interface for specifying domains of input attributes and built-in operators to manipulate the data that is extracted. Similarly to this work, the Morph-CSV framework takes into account the specification of domain constraints and transformation functions, but using established standards for tabular annotations and mapping function definitions.

Throughout the years these ideas have evolved from the use of description logics [39] to the use of ontologies as a common model for data access [5], what is called *Ontology-Based Data Access*. Most of the works proposed under this framework are focused on providing access to relational databases [5, 6, 40] and optimizations on the SPARQL-to-SQL translation process. In this context, the term *constraint* has been used in [41], where the authors defined two new properties

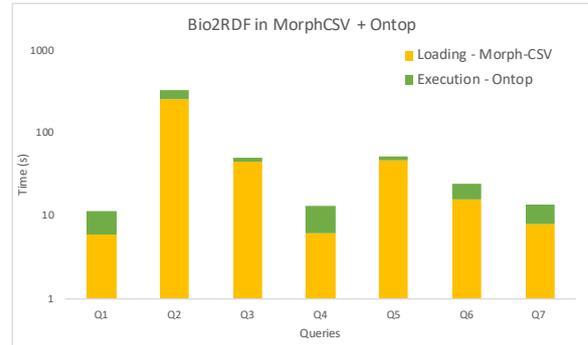


Fig. 16. **Query execution Time of Tabular Datasets in Bio2RDF with Ontop.** Execution time in seconds of the tabular datasets from Bio2RDF. The loading time applying a set of constraints (yellow) and query execution with Ontop (green).

extending the concept of OBDA instance. They propose a set of optimizations during SPARQL-to-SQL translation with techniques that take into account these constraints. However, the main assumptions made over the OBDA framework (e.g, the data source is an RDB o has an RDB wrapper, or the schema contains a set of constraints) are maintained. There are other works such as [24, 42] that apply the OBDA framework over document-based databases, i.e., MongoDB. Morph-CSV follows an OBDA approach including the exploitation of additional information from mappings, tabular metadata and queries for tabular datasets.

Related to our work are those approaches that allow querying directly information stored in flat files [43], Drill²², NoDB [44]. These systems provide a layer where “raw” data is queried, the data is adaptively loaded and stored, and then the query is executed using an assortment of strategies. Although these systems evaluate queries on raw tabular data and may exploit information encoded in the query, they do not make use of annotations or any sort of description of the data as Morph-CSV does.

Current OBDA open source systems that take tabular data as input are Ontario [11] and Squerall [12]. Ontario is a federated query processing approach for heterogeneous data sources. In its source selection step, Ontario uses source descriptions named RDF Molecule Templates [45] which keep information on the sources. The system handles tabular data among other formats, and implements a virtualization approach of query answering techniques for efficient execution. Similarly, Squerall is also an OBDA system

²²<https://drill.apache.org/>

1 that takes as its inputs data and mappings and uses a
2 middleware to aggregate the intermediate results in a
3 distributed manner. Although the aforementioned sys-
4 tems evaluate queries against raw tabular data, they do
5 not exploit the constraints declared in annotations or
6 mapping rules.

7 CSV on the Web (CSVW)²³ is a W3C proposal
8 for the definition of metadata on CSV files such as
9 datatypes, valid values, data transformations, and pri-
10 mary and foreign key constraints. A related W3C pro-
11 posal²⁴ defines a procedure and rules for the genera-
12 tion of RDF from tabular data and a few implementa-
13 tions that refer to this proposal are already available.
14 The CSV2RDF tool is presented in [46], the authors
15 define algorithms to transform CSV data into RDF
16 using CSVW metadata annotations, and their experi-
17 mental study uses datasets from the CSVW Implemen-
18 tation Report²⁵. Another tool, COW: Converter for
19 CSV on the Web²⁶ allows the conversion of datasets
20 in CSV format and uses a JSON schema expressed
21 in an extended version of the CSVW standard. Both
22 are focused on RDF materialization. To the best of
23 our knowledge, no existing tool exploits information
24 in CSVW annotations for querying tabular data in an
25 OBDA context.

26 Another area related to our work is the definition
27 and application of data transformation functions. An
28 approach independent of a specific implementation
29 context is described in [47]. It enables the descrip-
30 tion, publication and exploration of functions and in-
31 stantiation of associated implementations. The pro-
32 posed model is the Function Ontology and the publi-
33 cation method follows the Linked Data principles. Pre-
34 vious works related to this topic focus on develop-
35 ing ad-hoc and programmed functions. For example,
36 R2RML-F [19] allows using functions in the value of
37 the `rr:objectMap` property, so as to modify the
38 value of the table columns from a relational database.
39 KR2RML [48], used in Karma, extends R2RML by
40 adding transformation functions in order to deal with
41 nested values. OpenRefine enables such transforma-
42 tions with the usage of GREL functions, which can be
43 used in its RDF extension. Morph-CSV uses the exten-
44 sion of RML together with the Function Ontology [7]
45 that allows to incorporate ad-hoc transformation func-
46 tions over the data sources in a declarative manner.

47
48
49 ²³<https://www.w3.org/TR/tabular-data-primer/>

²⁴<https://www.w3.org/TR/csv2rdf/>

²⁵<https://w3c.github.io/csvw/tests/reports/index.html>

²⁶<https://csvw-converter.readthedocs.io/en/latest/>

8. Conclusions and Future Work

1
2
3 In this paper, we have presented an extension of the
4 common OBDA specification to address the problem
5 of query translation over tabular data. We describe and
6 evaluate Morph-CSV, a framework that exploits the in-
7 formation of mapping rules and metadata OBDA an-
8 notations to extract and apply a set of relevant con-
9 straints. It pushes down the application of these ele-
10 ments directly over the raw data in order to improve
11 query evaluation and query completeness. One of the
12 main contributions of this proposal is that it can be
13 used together with any OBDA framework. From the
14 set of experiments that we have performed with two
15 existing state-of-the-art OBDA engines (Morph-RDB
16 and Ontop), we can see that the use of those engines
17 inside the Morph-CSV framework brings several posi-
18 tive impacts: more queries can be answered and less
19 time is needed to answer most queries.

20 The definition, application and optimization of new
21 functions and constraints to address other challenges
22 for querying tabular data is one of the main lines
23 for future work [30]. We also want to study the per-
24 formance of the proposed workflow over OBDA dis-
25 tributed query systems such as the ones proposed
26 in [11, 12]. The results obtained can also be useful
27 to machine learning approaches for identifying when
28 the application of the integrity constraints is needed
29 or not, as we observe that there are special cases that
30 it can have a negative impact. We will also study the
31 challenges for querying other data formats (e.g., XML,
32 JSON) in an OBDA context and extend our approach
33 to incorporate them. We also want to remark the im-
34 portance of having standard and shared methods and
35 vocabularies to publish metadata of raw data on the
36 web, available for tabular data but not for tree data for-
37 mats such as XML and JSON. Finally, we will adapt
38 this proposal for a materialization process and study its
39 effects comparing it with previous proposals.

References

- 41
42
43
44
45 [1] C. Bizer, T. Heath and T. Berners-Lee, Linked data: The story
46 so far, in: *Semantic services, interoperability and web applica-*
47 *tions: emerging concepts*, IGI Global, 2011, pp. 205–227.
48 [2] M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Ap-
49 pleton, M. Axton, A. Baak, N. Blomberg, J.-W. Boiten,
50 L.B. da Silva Santos, P.E. Bourne et al., The FAIR Guiding
51 Principles for scientific data management and stewardship, *Sci-*
entific data **3** (2016).

- [3] A. Poggi, D. Lembo, D. Calvanese, G. De Giacomo, M. Lenzerini and R. Rosati, Linking data to ontologies, in: *Journal on data semantics X*, Springer, 2008, pp. 133–173.
- [4] M. Lenzerini, Data Integration: A Theoretical Perspective, in: *Proceedings of the Twenty-first ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 3-5, Madison, Wisconsin, USA, 2002*, pp. 233–246.
- [5] D. Calvanese, B. Cogrel, S. Komla-Ebri, R. Kontchakov, D. Lanti, M. Rezk, M. Rodriguez-Muro and G. Xiao, Ontop: Answering SPARQL queries over relational databases, *Semantic Web* 8(3) (2017), 471–487.
- [6] F. Priyatna, O. Corcho and J. Sequeda, Formalisation and experiences of R2RML-based SPARQL to SQL query translation using morph, in: *Proc. of WWW*, ACM, 2014, pp. 479–490.
- [7] B. De Meester, W. Maroy, A. Dimou, R. Verborgh and E. Mannens, Declarative data transformations for Linked Data generation: the case of DBpedia, in: *ESWC*, Springer, 2017, pp. 33–48.
- [8] A.C. Junior, C. Debruyne, R. Brennan and D. O’Sullivan, FunUL: a method to incorporate functions into uplift mapping languages, in: *Proceedings of the 18th International Conference on Information Integration and Web-based Applications and Services*, ACM, 2016, pp. 267–275.
- [9] J. Tennison, G. Kellogg and I. Herman, Model for tabular data and metadata on the web. W3C recommendation, *World Wide Web Consortium (W3C)* (2015).
- [10] M. Rodriguez-Muro and M. Rezk, Efficient SPARQL-to-SQL with R2RML mappings, *Web Semantics* 33 (2015), 141–169.
- [11] K.M. Endris, P.D. Rohde, M.-E. Vidal and S. Auer, Ontario: Federated Query Processing Against a Semantic Data Lake, in: *International Conference on Database and Expert Systems Applications*, Springer, 2019, pp. 379–395.
- [12] M.N. Mami, D. Graux, S. Scerri, H. Jabeen, S. Auer and J. Lehmann, Squerall: virtual ontology-based access to heterogeneous and large data sources, in: *International Semantic Web Conference*, Springer, 2019, pp. 229–245.
- [13] G. Xiao, D. Calvanese, R. Kontchakov, D. Lembo, A. Poggi, R. Rosati and M. Zakharyashev, Ontology-Based Data Access: A Survey, in: *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018*, 2018.
- [14] C. Bizer and A. Schultz, The berlin sparql benchmark, *International Journal on Semantic Web and Information Systems (IJSWIS)* 5(2) (2009), 1–24.
- [15] F. Belleau, M.-A. Nolin, N. Tourigny, P. Rigault and J. Morissette, Bio2RDF: towards a mashup to build bioinformatics knowledge systems, *Journal of biomedical informatics* 41(5) (2008), 706–716.
- [16] B. Golshan, A. Halevy, G. Mihaila and W.-C. Tan, Data integration: After the teenage years, in: *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, 2017, pp. 101–106.
- [17] A. Halevy, A. Rajaraman and J. Ordille, Data integration: The teenage years, in: *Proceedings of the 32nd international conference on Very large data bases*, 2006, pp. 9–16.
- [18] A. Doan, A. Halevy and Z. Ives, *Principles of data integration*, Elsevier, 2012.
- [19] C. Debruyne and D. O’Sullivan, R2RML-F: Towards Sharing and Executing Domain Logic in R2RML Mappings, in: *LDOW@ WWW*, 2016.
- [20] S. Das, S. Sundara and R. Cyganiak, R2RML: RDB to RDF Mapping Language, W3C Recommendation 27 September 2012, www.w3.org/TR/r2rml (2012).
- [21] A. Dimou, M. Vander Sande, P. Colpaert, R. Verborgh, E. Mannens and R. Van de Walle, RML: A Generic Language for Integrated RDF Mappings of Heterogeneous Data, in: *LDOW*, 2014.
- [22] B. De Meester, A. Dimou, R. Verborgh and E. Mannens, An ontology to semantically declare and describe functions, in: *ISWC*, Springer, 2016, pp. 46–49.
- [23] U. Şimşek, E. Kärle and D. Fensel, RocketRML-A NodeJS implementation of a use-case specific RML mapper, in: *Proceeding of the First International Workshop on Knowledge Graph Building*, 2019.
- [24] E. Botoeva, D. Calvanese, B. Cogrel, J. Corman and G. Xiao, Ontology-based data access—Beyond relational sources, *Intelligenza Artificiale* 13(1) (2019), 21–36.
- [25] C. Beeri, P.A. Bernstein and N. Goodman, A Sophisticate’s Introduction to Database Normalization Theory, in: *VLDB*, 1978.
- [26] S. Jozashoori and M.-E. Vidal, MapSDI: A Scaled-Up Semantic Data Integration Framework for Knowledge Graph Creation, in: *OTM Confederated International Conferences" On the Move to Meaningful Internet Systems"*, Springer, 2019, pp. 58–75.
- [27] M.-E. Vidal, E. Ruckhaus, T. Lampo, A. Martínez, J. Sierra and A. Polleres, Efficiently joining group patterns in SPARQL queries, in: *Extended Semantic Web Conference*, Springer, 2010, pp. 228–242.
- [28] E.F. Codd, Extending the database relational model to capture more meaning, *ACM Transactions on Database Systems (TODS)* 4(4) (1979), 397–434.
- [29] J. Mora, R. Rosati and O. Corcho, kyrie2: Query Rewriting under Extensional Constraints in $\{\text{ELHIO}\}$, in: *International Semantic Web Conference*, Springer, 2014, pp. 568–583.
- [30] A. Iglesias-Molina, D. Chaves-Fraga, F. Priyatna and O. Corcho, Enhancing the Maintainability of the Bio2RDF Project Using Declarative Mappings, in: *Proceedings of the 12th International Conference on Semantic Web Applications and Tools for Healthcare and Life Sciences*, 2019.
- [31] O. Corcho, F. Priyatna and D. Chaves-Fraga, Towards a New Generation of Ontology Based Data Access, *Semantic Web Journal* 11(1) (2020), 153–160.
- [32] G. Xiao, R. Kontchakov, B. Cogrel, D. Calvanese and E. Botoeva, Efficient handling of SPARQL optional for OBDA, in: *International Semantic Web Conference*, Springer, 2018, pp. 354–373.
- [33] D. Lanti, M. Rezk, M. Slusnys, G. Xiao and D. Calvanese, The NPD benchmark for OBDA systems, *CEUR Electronic Workshop Proceedings*, 2014.
- [34] G. Wiederhold, Mediators in the architecture of future information systems, *Computer* 25(3) (1992), 38–49.
- [35] S. Chawathe, H. Garcia-Molina, J. Hammer, K. Ireland, Y. Papakonstantinou, J. Ullman and J. Widom, The TSIMMIS Project: Integration of Heterogenous Information Sources, in: *Information Processing Society of Japan (IPSJ 1994)*, 1994.
- [36] A.L.A. Rajaraman, J. Ordille et al., Querying heterogeneous information sources using source descriptions, in: *Proc. of VLDB*, 1996.

- [37] M.T. Roth and P.M. Schwarz, Don't Scrap It, Wrap It! A Wrapper Architecture for Legacy Data Sources., in: *VLDB*, Vol. 97, 1997, pp. 25–29.
- [38] M.-E. Vidal, L. Bright, J.-r. Gruser and L. Raschid, A Wrapper Generation toolkit to specify and construct Wrappers for Web Accessible Data Sources (WebSources), *Computer Science and Engineering* **14** (1999).
- [39] D. Calvanese, G. De Giacomo and M. Lenzerini, *Description Logics for Information Integration*, in: *Computational Logic: Logic Programming and Beyond: Essays in Honour of Robert A. Kowalski Part II*, A.C. Kakas and F. Sadri, eds, Springer Berlin Heidelberg, 2002, pp. 41–60.
- [40] J.F. Sequeda and D.P. Miranker, Ultrawrap: SPARQL execution on relational data, *Journal of Web Semantics* **22** (2013), 19–39.
- [41] D. Hovland, D. Lanti, M. Rezk and G. Xiao, OBDA constraints for effective query answering, in: *International Symposium on Rules and Rule Markup Languages for the Semantic Web*, Springer, 2016, pp. 269–286.
- [42] F. Michel, L. Djiméno, C.F. Zucker and J. Montagnat, Translation of Relational and Non-Relational Databases into RDF with xR2RML, in: *11th International Conference on Web Information Systems and Technologies (WEBIST'15)*, 2015, pp. 443–454.
- [43] S. Idreos, I. Alagiannis, R. Johnson and A. Ailamaki, Here are my data files. here are my queries. where are my results?, in: *Proceedings of 5th Biennial Conference on Innovative Data Systems Research*, 2011.
- [44] I. Alagiannis, R. Borovica, M. Branco, S. Idreos and A. Ailamaki, NoDB: efficient query execution on raw data files, in: *Proc. ACM SIGMOD*, ACM, 2012, pp. 241–252.
- [45] K.M. Endris, M. Galkin, I. Lytra, M.N. Mami, M.-E. Vidal and S. Auer, Querying interlinked data by bridging RDF molecule templates, in: *Transactions on Large-Scale Data-and Knowledge-Centered Systems XXXIX*, Springer, 2018, pp. 1–42.
- [46] S.M.H. Mahmud, M. Hossin, H. Jahan, S. Noori and M. Hossain, CSV2RDF: Generating RDF Data From CSV File Using Semantic Web Technologies, *Journal of Theoretical and Applied Information Technology* **96** (2018).
- [47] B. De Meester, T. Seymoens, A. Dimou and R. Verborgh, Implementation-independent function reuse, *Future Generation Computer Systems* (2019).
- [48] J. Slepicka, C. Yin, P.A. Szekely and C.A. Knoblock, KR2RML: An Alternative Interpretation of R2RML for Heterogenous Sources, in: *COLD*, 2015.
- [49] D. Lanti, G. Xiao and D. Calvanese, VIG: Data scaling for OBDA benchmarks, *Semantic Web* **10**(2) (2019), 413–433.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

Appendix A. Morph-CSV algorithm

The Morph-CSV algorithm exploiting the mapping rules and metadata to enhance virtual ontology based data access for tabular datasets.

Algorithm 1: Morph-CSV algorithm

```

Result: SPARQL query result set
M ← mapping_rules;
MD ← metadata;
Q ← query;
D ← tabular_dataset;
SSG ← ∅;
for tp ← 0 to Q.getTP().size() do
  | SSG.add(tp);
end
for i ← 0 to SSG.size() do
  | p ← SSG.getPredicates(i);
  | for j ← 0 to M.getTM().size() do
  | | if p.isContainedIn(M.getTM(j)) then
  | | | M' ← M.getTM(j);
  | | | MD' ← getMD(M.getTM(j));
  | | end
  | end
  | M ← M';
  | MD ← MD';
end
for i ← 0 to M.getTM().size() do
  | path ← M.getTM(i).getSource();
  | ref ← M.getTM(i).getReferences();
  | ts ← D.get(path);
  | D'.add(TS.project(ref));
end
D ← D';
for i ← 0 to D.size() do
  | path ← D[i].getPath();
  | norm_2NF(D[i], MD.getMetadata(path));
  | norm_3NF(D[i], M);
  | duplicates(D[i]);
  | substitute(D[i], MD.getMetadata(path));
  | create(D[i], M.getDeclarativeFunctionFragment(path));
end
S ← schema(D, M, MD);
D' ← load(D, S);
M' ← translate(M);
PI = (O, M', S, D');
return run_query(Q, PI);

```

Appendix B. Query Features

Table 4

Query features of the evaluation of Morph-CSV. Domain constraints are described based on the function performed by Morph-CSV and reflect the number of the columns where that functions has been applied. Improvement functions (duplicates, source selection) are always applied.

Query	Query characteristics	Constraints		# Sources
		Integrity	Domain	
Madrid-GTFS-Bench				
Q1	4 TP	-	3 DataType, 4 Sub	1
Q2	5 TP, 2 OPT, 1 Filter	1 INDEX	3 DataType, 5 Sub	1
Q3	5 TP, 3 OPT, 1 Filter	1 INDEX	4 DataType, 5 Sub	1
Q4	9 TP, 1 Join, 4 OPT	2 PK, 1 FK	7 Sub	2
Q5	5 TP, 2 Join, 1 Filter	2 PK	2 DataType, 2 Sub	2
Q6	3 TP, 1 Join, 1 Filter	2 PK, 1 FK	-	2
Q7	15 TP, 5 Join, 5 OPT, 1 Filter	6 PK, 5 FK	3 DataType, 8 Sub	6
Q8	14 TP, 4 Join, 3 OPT	6 PK, 5 FK	3 DataType, 8 Sub	6
Q9	7 TP, 5 Join, 1 OPT, 1 Filter	5 PK, 3 FK	2 DataType, 3 Sub	5
Q10	4 TP, 1 Join, 1 Filter	2 PK, 1 FK	2 Sub	2
Q11	10 TP, 3 Join, 3 Filter (1 not exists)	3 PK, 2 FK	2 DataType, 2 Sub	3
Q12	10 TP, 3 Joins	4 PK, 3 FK	1 DataType, 4 Sub	4
Q13	6 TP, 1 Join, 1 OPT	1 PK, 1 FK	1 DataType, 3 Sub	1
Q14	8 TP, 3 Join, 1 OPT	4 PK, 3 FK	1 DataType, 3 Sub	3
Q15	3 TP, 1 Filter	1 PK, 1 FK	4 DataType, 11 Sub	1
Q16	8 TP, 3 Join, 2 Filter	4 PK, 2 FK	2 DataType, 2 Sub	3
Q17	9 TP, 2 Join	3 PK, 2 FK	1 DataType, 4 Sub	3
Q18	8 TP, 1 Union, 3 Join	4 PK, 3 FK	1 DataType, 3 Sub	4
Bio2RDF				
Q1	4 TP	-	3 Sub	1
Q2	4 TP, 1 Join, 1 Filter	1 PK, 1 INDEX	7 Sub	2
Q3	4 TP, 1 Join	1 PK, 3 INDEX	5 Sub	3
Q4	4 TP, 1 Join	1 PK, 1 INDEX	7 Sub	2
Q5	5 TP, 1 Join	1 PK, 2 INDEX	6 Sub	2
Q6	4 TP	-	2 Sub	1
Q7	6 TP, 1 Join, 2 Filter	1 PK	1 DataType, 4 Sub, 1 Create	1
BSBM				
Q1	5 TP, 3 Join, 1 Filter	3 PK, 2FK	7 DataType, 1 Sub	3
Q2	15 TP, 3 Join, 3 OPT	4 PK, 3 FK	10 DataType, 12 Sub	4
Q3	7 TP, 3 Join, 2 Filter, 1 OPT	3 PK, 2FK	8 DataType, 3 Sub	3
Q4	12 TP, 1 Union, 6 Join, 2 Filter	3 PK, 2FK	2 DataType, 4 Sub	2
Q5	7 TP, 2 Join, 2 Filter	2 PK, 1FK	6 DataType, 3 Sub	2
Q6	2 TP, 1 Filter	-	1 Sub	1
Q7	14 TP, 5 Join, 1 Filter, 2 OPT	5 PK, 4 FK	11 DataType, 2 Sub	5
Q8	10 TP, 2 Join, 4 OPT	3 PK, 2 FK	8 DataType, 8 Sub	3
Q9	DESCRIBE, 1 TP	-	-	1
Q10	7 TP, 3 Join, 2 Filter	3 PK, 3 FK	7 DataType, 2 Sub	3
Q11	2 TP, 1 Union	11 PK, 11 FK	29 DataType, 53 Sub	11
Q12	CONSTRUCT, 9 TP, 2 Join	3 PK, 2 FK	6 DataType, 7 Sub	3

Appendix C. Query Completeness

Table 5

Query completeness over multiple sizes of a GTFS dataset (the number indicates the scale factor: 1, 10, 100 and 100). The absence of a value means that the OBDA engine does not support the features of the SPARQL query or timeout.

Engines/Queries	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q9	Q12	Q13	Q14	Q17	Total
GTFS-1													
Virtuoso	58540	765	765	13	28	1	2	151439	6	734	2364	855	156972
Morph-RDB	58540	765	-	13	-	1	2	151439	6	734	2364	855	156179
Morph-CSV & Morph-RDB	58540	765	-	13	-	1	2	151439	6	734	2364	855	156179
Ontop	58540	-	765	13	-	-	-	-	-	734	2364	855	4731
Morph-CSV & Ontop	58540	765	765	13	28	-	2	151439	-	734	2364	855	156965
GTFS-10													
Virtuoso	353660	6312	4207	130	350	1	67	718317	130	2650	23640	8550	764354
Morph-RDB	353660	6312	-	130	-	1	67	-	130	2650	23640	8550	41480
Morph-CSV & Morph-RDB	353660	6312	-	130	-	1	67	718317	130	2650	23640	8550	759797
Ontop	353660	-	4207	130	-	-	-	-	-	2650	23640	8550	39177
Morph-CSV & Ontop	353660	6312	4207	130	350	-	67	718317	-	2650	23640	8550	764223
GTFS-100													
Virtuoso	3536600	63100	42067	1300	3500	1	67	7183874	1300	26500	236400	85500	7643609
Morph-RDB	3536600	63100	-	1300	-	1	67	-	1300	26500	236400	85500	414168
Morph-CSV & Morph-RDB	3536600	63100	-	1300	-	1	67	-	1300	26500	236400	85500	414168
Ontop	3536600	-	42067	1300	-	-	-	-	-	26500	236400	85500	391767
Morph-CSV & Ontop	3536600	63100	42067	1300	-	-	67	-	-	26500	236400	85500	454934
GTFS-1000													
Virtuoso	35366000	1261368	420667	13000	35000	1	69	19077083	13000	420666	2364000	855000	24459854
Morph-RDB	-	1261368	-	13000	-	1	69	-	13000	420666	2364000	855000	4927104
Morph-CSV & Morph-RDB	35366000	1261368	-	13000	-	1	69	-	13000	420666	2364000	855000	4927104
Ontop	-	-	420667	13000	-	-	-	-	-	420666	2364000	855000	4073333
Morph-CSV & Ontop	-	1261368	420667	13000	-	-	69	-	-	420666	2364000	855000	5334770

Table 6

Query completeness over multiple sizes of a BSBM dataset. The absence of a value means that the OBDA engine does not support the features of the SPARQL query or timeout.

Engines/Queries	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	Q12	Total
45K												
Virtuoso	10	19672	10	10	5	3	580691	20	450000	10	900000	1950431
Morph-RDB	10	19672	10	10	-	3	580691	20	450000	10	900000	1950426
Morph-CSV & Morph-RDB	10	19672	10	10	5	3	580691	20	450000	10	900000	1950431
Ontop	10	-	10	10	-	-	-	-	-	-	-	30
Morph-CSV & Ontop	10	-	10	10	5	-	-	-	-	10	-	45
90K												
Virtuoso	10	38665	10	10	5	5	1161448	20	900000	10	1800000	3900183
Morph-RDB	10	38665	10	10	-	5	1161448	20	900000	10	1800000	3900178
Morph-CSV & Morph-RDB	10	38665	10	10	5	5	1161448	20	900000	10	1800000	3900183
Ontop	10	-	10	10	-	-	-	-	-	-	-	30
Morph-CSV & Ontop	10	-	10	10	5	-	-	-	-	10	-	45
180K												
Virtuoso	10	69434	10	10	5	9	2168792	20	1800000	10	3600000	7638300
Morph-RDB	10	-	10	10	-	9	2168792	20	1800000	10	3600000	7568861
Morph-CSV & Morph-RDB	10	69434	10	10	-	9	2168792	20	1800000	10	3600000	7638295
Ontop	10	-	10	10	-	-	-	-	-	-	-	30
Morph-CSV & Ontop	10	-	10	10	5	-	-	-	-	10	-	45
360K												
Virtuoso	10	137359	10	10	5	18	4337584	20	3600000	10	7200000	15275026
Morph-RDB	10	-	10	10	-	18	-	20	3600000	10	-	3600078
Morph-CSV & Morph-RDB	10	137359	10	10	-	18	-	20	3600000	10	-	3737437
Ontop	10	-	10	10	-	-	-	-	-	-	-	30
Morph-CSV & Ontop	10	-	10	10	5	-	-	-	-	10	-	45

Table 7

Query completeness over of Bio2RDF tabular dataset.

Engines/Queries	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Total
Morph-CSV + Morph-RDB	1000	1190181	10	102594	200	28224	>10000	>1422209
Morph-CSV + Ontop	1000	1190181	10	102594	200	28224	13481	1335690