Streamlining Knowledge Graph Creation with PyRML

Andrea Giovanni Nuzzolese $^{1[0000-0003-2928-9496]}$

CNR - Institute of Cognitive Sciences and Technologies, Bologna, Italy andreagiovanni.nuzzolese@cnr.it

Abstract. Knowledge Graphs (KGs) are increasingly adopted as a foundational technology for integrating heterogeneous data in domains such as climate science, cultural heritage, and the life sciences. Declarative mapping languages like R2RML and RML have played a central role in enabling scalable and reusable KG construction, offering a transparent means of transforming structured and semi-structured data into RDF. In this paper, we present PyRML, a lightweight, Python-native library for building Knowledge Graphs through declarative mappings. PyRML supports core RML constructs and provides a programmable interface for authoring, executing, and testing mappings directly within Python environments. It integrates with popular data and semantic web libraries (e.g., Pandas and RDFlib), enabling transparent and modular workflows. By lowering the barrier to entry for KG creation and fostering reproducible, ontology-aligned data integration, PyRML bridges the gap between declarative semantics and practical KG engineering.

Keywords: RML · Declarative mappings · Knowledge Graph Creation · Python · Pandas · RDFLib

1 Introduction

Knowledge Graphs [7] (KGs) have become a foundational technology for integrating and querying heterogeneous data across domains such as climate science, cultural heritage, and life sciences. A core strength of KGs lies in their ability to make data explicit, interoperable, and semantically rich, aligning with the principles of FAIR [14] (Findable, Accessible, Interoperable, and Reusable) data.

Among the various approaches to constructing KGs, declarative mapping languages, such as R2RML¹ and RML [4, 9], have emerged as key enablers in both literature and practice. By explicitly stating the rules for transforming data from structured and semi-structured sources (e.g., relational databases, CSV, JSON, XML, SPARQL services, etc.) into RDF, declarative mappings promote separation of concerns, reusability, and cross-source interoperability. These characteristics are particularly beneficial for collaborative and long-lived data integration efforts, where mapping logic must be shared, adapted, and audited over time.

¹ http://www.w3.org/TR/r2rml/

A variety of RML-compliant engines have been developed to support the execution of declarative mappings. Examples are RMLMapper², CARML³, SDM-RDFizer [8], and Morph-KGC [1]. These tools have proven effective in translating structured data into RDF at scale and have contributed significantly to the maturation of the semantic data integration ecosystem. However, their usage often assumes familiarity with command-line interfaces, specific configuration formats, or specific programming language environments, which can present barriers to integration in modern data science workflows. In addition, features such as mapping modularity, incremental generation, unit testing, and tight coupling with ontological reasoning are still underdeveloped or inconsistently supported across tools. In this context, PyRML is conceived as a Python-native alternative that supports interactive, programmable, and transparent KG construction. It complements existing engines while focusing on usability, extensibility, and seamless integration with the Python data ecosystem. Additionally, by abstracting some of the technical complexity while maintaining expressive power, PyRML contributes to bridging the gap between declarative semantics and practical KG engineering.

The remainder of this paper is the following. Section 2 provides an overview of the related work. Section 3 details the proposed system architecture with usage examples. Section 4 describes the evaluation methodology and results. Finally, Section 5 discusses conclusions and future directions.

2 Related work

The construction of KGs from structured and semi-structured data has been extensively explored in the Semantic Web community⁴. The foundational approaches focused on the integration of semantic data. These include the use of a view-based paradigm [11], such as, Global-As-View [6] (GAV), Local-As-View [13] (LAV), and Global-Local-As-View [5] (GLAV). A view defines the relationships between heterogeneous data sources and a unified mediated schema. These paradigms, originally developed in the context of data warehouse and federated databases, provided the theoretical foundation for later declarative mapping languages used in the construction of KG. In particular, the GAV approach, where each element of the mediated schema is defined as a query over the sources, closely resembles modern RML mappings, where ontology terms are defined in terms of data source structure. In contrast, LAV and GLAV underpin more expressive approaches such as Ontology-Based Data Access (OBDA), where mappings specify how source data can satisfy arbitrary ontology queries. Although powerful, OBDA systems often rely on complex reasoning services, which can hinder scalability and accessibility for practitioners [10].

² https://github.com/RMLio/rmlmapper-java

³ https://github.com/carml/carml

⁴ Please refer to https://kg-construct.github.io/awesome-kgc-tools/ for an overview on existing tools for KG construction.

In recent years, declarative mapping languages such as the RDB to RDF Mapping Language (R2RML) and RDF Mapping Language [4,9] (RML) have become a standard mechanism for aligning raw data with RDF vocabularies and ontologies in a transparent and maintainable way. More specifically R2RML is designed to cope with the transformation of relational databases to RDF, whilst RML generalises the mapping model of R2RML to support diverse semistructured data sources. Accordingly, several tools have been developed to implement and execute RML mappings. The RMLMapper, written in Java, was among the first engines to support full RML core semantics and has been widely used for research and data publication tasks. CARML builds on the same paradigm, offering improved modularity and performance. More recently, tools like SDM-RDFizer [8] and Morph-KGC [1] have focused on scalability and performance, enabling the efficient generation of large KGs from relational databases and tabular data. These engines have been successfully applied in large-scale projects, such as iASiS⁵, which is an EU funded project to enable precision medicine approaches by utilising insights from patient data.

Despite their robustness, existing RML engines often assume specific technological stacks (e.g., Java or Docker-based deployments), and their integration with modern data science workflows—typically centred around Python—is limited. PyRML contributes to this landscape by offering a Python-native, programmable interface for RML-based KG construction. Unlike black-box engines, PyRML enables fine-grained control over mapping composition, execution, and testing, and is designed to integrate with widely used Python libraries such as Pandas and RDFlib. This positions PyRML as a complementary tool in the RML ecosystem, addressing the need for flexible, scriptable, and developer-friendly solutions in data-centric environments.

A complementary approach to data integration is presented by SPARQL Anything [2], which enables querying heterogeneous data sources directly using SPARQL, without the need for upfront data transformation into RDF. By overloading the SERVICE clause in SPARQL 1.1, SPARQL Anything allows users to access data from various formats through a uniform SPARQL interface. This approach leverages the Facade-X meta-model [3] to provide a simplified RDF representation of diverse data sources, facilitating rapid prototyping and adhoc querying. While SPARQL Anything excels in on-the-fly data access, it does not produce persistent RDF graphs, which may be a limitation for applications requiring long-term data storage and reasoning capabilities.

3 The PyRML system

3.1 Architecture

Figure 1 shows the modular architecture of PyRML that counts of four main modules. These modules are: (i) the API module, (ii) core Framework, (iii) Functions Provider, and (iv) Mapper.

⁵ http://project-iasis.eu/

4 A.G. Nuzzolese

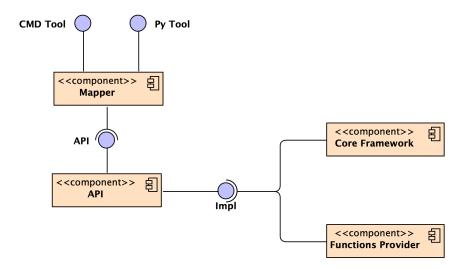


Fig. 1. The architecture of PyRML

API module. The API module provides the abstract base classes that define the core structure of the programming interface for capturing the RML model within the software platform. At the top of this structure is the TermMap abstract base class, which represents any entity in an RML mapping associated with an IRI and intended for generating RDF data from a logical table. It is worth noting that the class hierarchy defined in PyRML directly mirrors the taxonomy of classes specified in the RML and R2RML ontologies. Consequently, examples of classes derived from TermMap include SubjectMap, which specifies the mapping instructions to generate the subject of a triple from a logical table, and PredicateObjectMap, which links a PredicateMap and an ObjectMap to generate the predicate and object of a triple, respectively.

The TermMap class extends the abstract base class IdentifiedNode, implemented by RDFLib⁶, a widely used Python package for working with RDF. Hence, all instances of TermMap are valid RDF terms that can be used as the subject, predicate, or object of an RDF triple when constructing an RDFLib graph.

For example, the following code defines a TermMap through its derived class SubjectMap, and uses it as the subject of a triple to construct a graph with RDFLib.

```
from pyrml import TermMap, SubjectMap
from rdflib import Graph, Namespace, RDF

ex: Namespace = Namespace('https://foo.org/example/')
```

⁶ https://github.com/RDFLib/rdflib

```
rr: Namespace = Namespace('http://www.w3.org/ns/r2rml#')
tm: TermMap = SubjectMap(ex.SM)

g: Graph = Graph()
g.add((tm, RDF.type, rr.SubjectMap))
```

The abstract methods of TermMap include: (i) to_rdf, which converts the PyRML term into an RDFLib graph while preserving the graph structure rooted at that term; (ii) apply, which performs the mapping by using all the informations associated with a term against a given LogicalSource; and (iii) from_rdf, which is a static method that allows to instantiate a TermMap and its associated terms directly from an RDFLib graph. These three abstract methods are implemented by the base classes of TermMap that are defined in the core molude.

Core module. The core module has a twofold purpose, as reflected in its name: (i) it serves as the foundation of PyRML by implementing all the derived classes of TermMap that enable PyRML's functionality, and (ii) it addresses the core model of RML. The constructor (i.e. the __init__ method) of the class TermMap accepts a mandatory positional argument, which is the IRI of an RML term, and optional keyword arguments that can be used to associate term-specific values with an RML term, such as the values of the rr:template or rml:reference predicates in an RML mapping. The following code snippet defines RML triples map in a programmatic way.

```
from rdflib import FOAF
  ls: LogicalSource = ... // details on logical sources later
  sm: SubjectMap = SubjectMap(ex.SM,
                            template='https://foo.org/d/{ID}',
                             _classes=FOAF.Person)
  pm: PredicateMap = PredicateMap(ex.PM, constant=FOAF.name)
9
10
  om: ObjectMap = ObjectMap(ex.OM,
11
                             reference='name',
12
                             term_type=rr.Literal)
13
14
pom: PredicateObjectMap = PredicateObjectMap(ex.POM,
                                                  predicates=pm,
16
                                                 object_maps=om)
17
  tm: TripleMappings = TripleMappings(ex.Tm,
                                    logical_sources=ls,
                                    subject_maps_mp,
21
                                    predicate_object_maps=pom)
22
23
```

```
g: Graph = tm.to_rdf()
```

In the code snippet above a SubjectMap named ex:SM is created at line 5. This SubjectMap uses the template IRI https://foo.org/d/{ID}, where {ID} is replaced by the value of the ID from the input data. It also declares that each generated subject is an instance of foaf:Person, specified via the RML declaration _class=FOAF.Person.

Instead, an instance of PredicateMap is defined at line 9. This indicates that the predicate of the triple is the FOAF property foaf:name. Then, an ObjectMap named ex:OM is instantiated at line 11. The object value is taken from the name attribute in the input data (i.e. reference='name') and the term type is set to be an RDF literal (i.e. term_type=rr.Literal). At line 15 a PredicateObjectMap is created by connecting the pm predicate to the om object map. Finally, a triples mapping is instantiated at line 19. The latter represents a complete mapping rule that takes data from the logical source identified by 1s, builds a subject from sm, and predicates and objects from pom. An invocation of the method to_rdf() against the triples mapping tm is provided at line 24 to return an RDFLib graph as output. This graph is represented in the Turtle serialisation below.

```
@prefix ex: <https://foo.org/example/> .
  @prefix rr: <http://www.w3.org/ns/r2rml#>
3 @prefix rml: <http://semweb.mmlab.be/ns/rml#> .
  @prefix foaf: <http://xmlns.com/foaf/0.1/> .
  ex: Tm a rr: TriplesMapping;
        rml:logicalSource ...
        rr:subjectMap ex:SM ;
        rr:predicateObjectMap ex:POM .
  ex:SM a rr:SubjectMap ;
        rr:template 'https://foo.org/d/{ID}';
        rr:class foaf:Person .
13
14
  ex:PM a rr:PredicateMap ;
15
        rr:constant foaf:name .
16
17
  ex:OM a rr:ObjectMap ;
        rml:reference 'name';
19
        rr:termType rr:Literal .
20
  ex:POM a rr:PredicateObjectMap ;
        rr:predicateMap ex:PM ;
23
        rr:objectMap ex:OM .
```

As stated previously, any instance of TermMap can be generated directly from an RDFLib graph by invoking the from_rdf() method on the corresponding

class. For example, the triples map ex:TM can be converted to a PyRML object as shown in the code block below.

```
tm: TripleMappings = TripleMappings.from_rdf(g, parent=ex.TM)
```

The execution of the declarative mappings is enabled by the method apply that performs tranformation defined in a TermMap against a specific LogicalSource. In PyRML, a LogicalSource makes use of a Pandas DataFrame for providing a flexible and Python-native abstraction of the input data. This design decouples data access and parsing from the mapping logic, allowing users to load, clean, transform, and prepare data using familiar Pandas operations before applying declarative mappings. By working with DataFrames, PyRML integrates seamlessly into Python-based data science workflows, enabling direct manipulation of data, efficient debugging, and reuse of in-memory datasets from diverse sources (e.g., CSV, databases). This approach aligns the RML notion of logical tables with the DataFrame's tabular structure, facilitating transparent and programmable knowledge graph construction while maintaining compatibility with RML's mapping semantics. In this context, the apply method enables the vectorised application of transformation operations across a DataFrame that represents a Logical Source. Instead of processing records one by one, the apply method efficiently computes the corresponding RDF terms (e.g., subject IRIs, object literals, etc.) for all rows in a single operation, leveraging Pandas' inherent performance optimisations. This vectorised processing enhances scalability and supports the integration of complex transformation logic directly into the mapping execution pipeline.

PyRML currently supports the following data sources for instantiating a LogicalSource: (i) CSV, (ii) XML, (iii) JSON, (iv) SPARQL, (v) MySQL, (vi) SQL Server, and (vii) PostgreSQL. A LogicalSource in PyRML is built on top of a Source, which is the base abstract class used to represent various data source types. For instance, the class CSVSource extends Source to model a CSV data source. The code snippet below illustrates how to create a logical source from a CSV file.

```
from pyrml import LogicalSource, Source, CSVSource
source: Source = CSVSource(ex.CSV, 'students.csv')
source: LogicalSource(ex.LS, sources=source)
```

Functions module. Many real-world scenarios require additional data transformations beyond simple attribute retrieval or template substitution—such as string manipulation, date formatting, or value normalisation. To address these needs, the RML community introduced RML Functions⁷ (or FnO Functions), a

Thttps://kg-construct.github.io/rml-fnml/spec/docs/

mechanism to declaratively invoke functions as part of a mapping. PyRML provides a functions module to support the use of RML Functions—operations that enable data transformations within declarative mappings, following the principles of the Function Ontology (FnO). Similar to engines like RMLMapper, PyRML includes a list of built-in functions that mirror the default set provided by RMLMapper, allowing users to leverage standard transformation operations out of the box. This list is implemented in the Function module⁸.

A distinctive feature of PyRML is its support for user-defined functions. PyRML allows developers to extend the set of available functions programmatically, integrating custom transformation logic directly into the mapping execution pipeline. This is achieved by defining a standard Python function and registering it using a dedicated decorator provided by the library, called rml_function. The rml_function decorator facilitates the registration of a Python function as an RML Function by associating it with a function identifier (IRI) and parameter mappings. Its implementation follows a standard Python decorator pattern, wrapping the original function and registering it with the PyRML runtime. Once registered, the function can be invoked within RML mappings through the Function Ontology (FnO) mechanism. The following are the signature and a usage example of the rml_function decorator.

In the example above the Python function to_lower_case implements the logic to convert a string to lowercase. The @rml_function decorator registers it under the IRI corresponding to the GREL⁹ toLowerCase function from FnO. The mapping between the FnO parameter IRI and the Python argument name is specified via the value keyword. Once registered, this function can be invoked inside an RML mapping referencing its function IRI, enabling seamless integration between declarative mappings and Python-defined transformation logic.

Mapper module. The mapper module is the execution core of PyRML, responsible for transforming RML mapping definitions into RDF triples. It provides a Python-native, extensible, and efficient engine for materialising knowledge graphs from structured and semi-structured data, fully aligned with the

 $^{^{8} \ \}mathtt{https://github.com/anuzzolese/pyrml/blob/master/pyrml/functions.py}$

⁹ https://openrefine.org/docs/manual/grelfunctions

RDF Mapping Language (RML). This module enables the seamless integration of declarative knowledge graph construction into Python-based workflows. At its core is the RMLConverter class, which implements the high-level interface for executing a set of RML mappings. It supports both single-threaded and parallel execution strategies, allowing it to scale from lightweight testing to larger data transformation tasks. When processing a mapping, the converter parses the mapping document, builds an internal representation of the mapping using the classes for term maps defined in the API and Core modules, and applies each term map to the associated logical source by leveraing Pandas. When processing a mapping, the converter first renders the RML document using Jinja2¹⁰ templating, which allows users to define parameterised and reusable mappings. Template variables can be injected at runtime, enabling dynamic substitution of values such as file paths, graph names, or filter conditions, improving the flexibility and maintainability of mapping definitions. The class RMLConverter can be instantiated programmatically or used from command line through a dedicated Python script. The following is an example of the programmatic use of the class RMLConverter.

```
from pyrml import PyRML, RMLConverter

c: RMLConverter = PyRML.get_mapper()

invoke the method convert on the instance of class
    RMLConverter by:
    - using the persons.ttl RML descriptor;
    - obtaining an RDF graph as output.

g: Graph = c.convert('persons.ttl')
```

Variables for Jinja2 templating can be provided through the template_vars argument of the convert method. This argument accepts a Python dictionary, where each key corresponds to a parameter referenced in the template, and each value specifies the actual value to substitute. An example of an RML file that makes use of templating is the following.

```
1 <#Mapping> a rr:TriplesMap;
2  rml:logicalSource [
3  rml:source "{{ INPUT_CSV }}";
4  rml:referenceFormulation ql:CSV
5 ];
6  rr:subjectMap [
7  ...
8 ] .
```

¹⁰ https://jinja.palletsprojects.com/en/stable/

The template variable in the RML above is reported at line 3, i.e. {{ INPUT_CSV }} and can be set to an actual value as in the following example.

```
vars = {'INPUT_CSV': 'students.csv'}
g : Graph = c.convert('persons.ttl', template_vars=vars)
```

We note that Jinja2 templating, as implemented in PyRML, is a non-standard extension and is not part of the official RML specification or reference documentation. For more details on how templating works, please refer to the official Jinja2 documentation.

Instead, the following is a usage example of the command line tool pyrml-mapper.py that wraps the application aroung the RMLConverter.

```
python pyrml-mapper.py [-o RDF out file] [-f RDF out file] [-
m] input
```

Where:

- input is the required positional argument that specifies the RML mapping file to be used for RDF conversion;
- o filename is an optional argument that specifies the output file for saving the resulting RDF graph. If omitted, the output is written to standard output by default;
- f rdf-syntax is an optional argument to define the syntax used to serialise the RDF graph. Supported values include n3, nquads, nt, pretty-xml, trig, trix, turtle, and xml. If not specified, nt (N-Triples) is used by default;
- -m is an optional flag that enables multiprocessing to accelerate the transformation process.

3.2 Release and Availability Notes

Reusability. PyRML is released as open-source software under the Apache 2.0 license¹¹. It is well-documented, with API references, examples, and tutorials available in the GitHub repository. It is general-purpose and not tied to any specific domain, making it suitable for a wide range of use cases. The modular architecture supports extension, e.g., custom functions and mapping loaders. The documentation clearly specifies supported features, usage patterns, and known limitations, allowing users to adapt the tool confidently to their own needs.

Availability. PyRML is publicly available at:

```
- GitHub: https://github.com/anuzzolese/pyrml;
```

¹¹ https://github.com/anuzzolese/pyrml?tab=Apache-2.0-1-ov-file

- Python Package Index (PyPI): https://pypi.org/project/pyrml-lib, which
 allows to download and install PyRML directly from the official third-party
 software repository for Python via the pip command, e.g. pip install pyrml-lib;
- Documentation: included in the repository and browsable via GitHub Pages;
- Canonical citation: DOI: https://doi.org/10.5281/zenodo.15399948 released by Zenodo;
- License: Apache 2.0.

Impact and adoption. PyRML fills a critical gap in the knowledge graph construction landscape by providing a Python-native, declarative RML mapping engine that integrates seamlessly with modern data science workflows. It is of direct relevance to the Semantic Web community, particularly to researchers and practitioners engaged in FAIR data, open science, and knowledge integration initiatives. More broadly, PyRML lowers the barrier to entry for the scientific and societal adoption of Semantic Web technologies by aligning with widely adopted tools and conventions in the Python ecosystem. PyRML has already seen adoption in several EU-funded projects. It has been successfully integrated into the data transformation pipelines of HACID (Hybrid Human-AI Collective Intelligence in Open-Ended Domains), WHOW (Water Health Open Knowledge), and FOSSR (Fostering Open Science in Social Sciences and Humanities). In these projects, PyRML has been used to support transparent, modular, and reproducible workflows for transforming heterogeneous data sources into semantically rich RDF graphs aligned with domain ontologies. Its programmable and extensible architecture has proven particularly valuable in collaborative and evolving research environments. The project is actively maintained by a team of four developers and is openly available on GitHub. As of May 13th, 2025, the repository has received 37 stars, been forked 13 times, and, between April 30th and May 13th, has been cloned 30 times and viewed 155 times. 12 These indicators of adoption and engagement confirm the practical relevance of PvRML and its growing role in enabling declarative, transparent, and reproducible knowledge graph engineering across both research and applied domains.

4 Evaluation

4.1 Experimental setup

To evaluate the correctness and performance of PyRML, we designed an experimental protocol based on the official RML-Core test cases ¹³. The RML-Core test cases are a standardised suite developed by the Knowledge Graph Construction Community Group¹⁴ to evaluate the correctness and feature coverage of RML-compliant engines. Each test case defines a mapping scenario with an input data

¹² Usage statistics are available at https://github.com/anuzzolese/pyrml/graphs/ traffic.

 $^{^{13}}$ https://github.com/RMLio/rml-test-cases

¹⁴ https://www.w3.org/community/kg-construct/

source, an RML mapping file, and an expected RDF output. The suite covers key RML features such as logical sources, templates, joins, constant values, and graph maps. These tests are designed to be minimal, deterministic, and interpretable, making them ideal for verifying whether an engine behaves according to the standard. Successful execution of the core test cases provides strong evidence of RML compliance and correctness. Supported input formats in the test suite include CSV, JSON, XML, SPARQL endpoints, MySQL, SQL Server, and PostgreSQL databases. By validating against these tests, an engine demonstrates conformance to the RML specification across a wide range of source types and mapping patterns. The evaluation was carried out in two phases: feature coverage and computational performance benchmarking.

RML-Core conformance. In the first phase, we assessed the coverage of PyRML against the RML-Core specification¹⁵. This was done by systematically executing all the RML-Core test cases and verifying whether the output RDF graphs conformed to the expected results. This step ensured that PyRML adheres to the semantics and structural requirements of the RML specification. For this analysis we set up a Docker container to provide Apache Jena Fuseki¹⁶ as SPARQL endpoint, MySQL, PostgreSQL, and SQL Server. The container is available on GitHub¹⁷ and can be instantiated via docker-compose¹⁸ Then, a Python script¹⁹ was developed to assess RML-Core conformance using Python's unit testing framework, where each RML-Core test case is interpreted as a separate unit test.

Computational performance. In the second phase, we benchmarked the computational performance of PyRML and compared it to the widely used RMLMapper engine. For each test case, we executed the transformation 10 times using both engines and recorded the execution time in milliseconds. The final reported time for each engine and test case corresponds to the average execution time over the 10 runs. This procedure reduces the impact of transient system-level fluctuations and provides a more stable basis for comparison. The bash script that enabled the comparison is available on GitHub²⁰.

All benchmarks were run under identical conditions on the same hardware environment to ensure comparability, that is a 2.3 GHz Quad-Core Intel Core i7 with 32 GB of memory. The results of this evaluation are reported in Section 4.2.

4.2 Results

RML-Core conformance. Table 1 presents the coverage of PyRML against the reference RML core test suite, broken down by the type of logical data source.

 $^{^{15}\ {\}rm https://kg\text{-}construct.github.io/rml\text{-}core/spec/docs/}$

¹⁶ https://jena.apache.org/documentation/fuseki2/

¹⁷ https://github.com/anuzzolese/pyrml-testing/tree/main/docker-framework

¹⁸ Please refer to the readme provided along with the Docker container.

¹⁹ https://github.com/anuzzolese/pyrml-testing/blob/main/unittesting.py

²⁰ https://github.com/anuzzolese/pyrml-testing/blob/main/ comparative-analysis.sh

Each row corresponds to a different source type supported in RML, and the table reports three key metrics: (i) the number of test cases successfully executed by PyRML, where the generated RDF output matches the expected result; (ii) the number of test cases where PyRML did not produce the expected output, either due to missing feature support or incorrect behaviour; and (iii) the total number of test cases defined in the suite for that source type.

Source type	Passed	Failed	# of test cases	
CSV	39	0	39	
JSON	40	0	40	
777 67	0.0		2.0	

55

PostSQL

SQL Server

Table 1. PyRML coverage of test cases for RML core.

XML38 SPARQL 24 2 26 MySQL 60 55

5

60

Complete coverage (100% pass rate) is achieved for CSV, JSON, and XML. indicating robust and stable support for these commonly used structured and semi-structured formats. Among the 26 RML core test cases involving SPARQL data sources, PyRML successfully passes 24 and fails 2. The two failed cases, i.e. RMLTC0008b-SPARQL and RMLTC0009a-SPARQL, show limitations in the current implementation of join semantics when SPARQL is used as a logical source. In fact, RMLTC0008b-SPARQL tests the generation of triples that involve a referencing object map, where data from one source must be joined with another in a specific predicate object map. Similarly, RMLTC0009a-SPARQL evaluates the handling of foreign key-style relations between logical sources—an operation conceptually analogous to joins in relational databases. While PyRML supports referencing object maps for sources such as CSV and JSON, support for such joins across SPARQL requires more investigation and testing. Among the 60 RML core test cases available for each of the relational database systems, i.e. MySQL, PostgreSQL, and SQL Server, PyRML passes 55 and fails 5, which are: (i) RMLTC0009d; (ii) RMLTC0011a; (iii) RMLTC0013a; (iv) RMLTC0015a; and (v) RMLTC0016d. These failures are consistent across all three systems and are attributed to specific limitations in the current implementation of PvRML's SQL mapping layer. The test case RMLTC0009d checks the ability to handle column names that match SQL reserved keywords. PyRML currently does not implement automatic quoting or escaping of such identifiers, leading to parsing or execution errors during mapping. RMLTC0011a involves the mapping of manyto-many (M:N) relationships via custom SQL queries embedded in the logical source. RMLTC0013a tests the behaviour of referencing object maps when joined columns contain null values. In accordance with the RML specification, no triples should be generated in such cases; however, PyRML does not yet suppress triple generation in the presence of nulls, resulting in incorrect output. RMLTC0015a evaluates the correct generation of language-tagged literals based on values from a source column. Finally, RMLTC0016d tests the handling of datatype conversions, specifically boolean values. PyRML does not yet support automatic casting of source values to xsd:boolean, which leads to failures in producing semantically correct RDF literals when boolean datatypes are required. All these test cases represent clear, bounded limitations rather than architectural constraints. All identified issues are part of PyRML's ongoing development roadmap, and their resolution is planned in upcoming releases to support full RML compliance across relational database backends.

Computational performance. Figure 2 shows the results of the comparative analysis between PyRML and RMLMapper. The results are expressed in seconds, whilst error bars represent standard deviotions, which is reported among brackets.

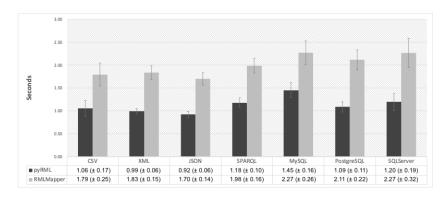


Fig. 2. Comparison of PyRML with RMLMapper respect to execution time of test cases expressed in seconds.

The results demonstrate a consistent performance advantage for PyRML across all source types. For example, on CSV sources, PyRML achieved an average execution time of 1.06 seconds, compared to 1.79 seconds for RMLMapper. Similar gains are observed for XML (0.99s vs. 1.83s), JSON (0.92s vs. 1.70s), and SPARQL (1.18s vs. 1.98s). The difference is particularly notable for relational database sources, where PyRML outperformed RMLMapper on both MySQL (1.45s vs. 2.27s) and PostgreSQL (1.09s vs. 2.11s). In addition to lower average execution times, PyRML also exhibited lower standard deviation across all source types, indicating more stable performance. For instance, in the case of CSV sources, the standard deviation for PyRML was 0.17s, compared to 0.25s for RMLMapper. This pattern is consistent for all data sources, reflecting the deterministic and efficient design of PyRML's mapping engine. These results highlight PyRML's efficiency and robustness, confirming its suitability for interactive and automated data integration pipelines, especially where low-latency transformation is required.

5 Conclusion and future work

This paper introduced PyRML, a Python-native engine for declarative knowledge graph construction based on the RML specification. Our experimental evaluation demonstrates that PyRML is both standards-compliant and computationally efficient, achieving full or near-full conformance across a wide spectrum of data sources—including structured (CSV), semi-structured (JSON, XML), and relational databases (MySQL, PostgreSQL, SQL Server), as well as SPARQL endpoints. These results validate PyRML's design, which closely mirrors the RML ontology and offers a transparent, modular, and extensible software architecture. When benchmarked against RMLMapper, one of the most widely adopted engines, PyRML exhibits consistently lower execution times and reduced performance variance, making it especially well-suited for automated or latency-sensitive knowledge graph construction pipelines. The native integration with Python's data ecosystem (e.g., pandas, RDFLib, Jinja2) further enhances its usability for data scientists and knowledge engineers working in FAIR and open science contexts. Nevertheless, the evaluation highlights several implementation-level limitations that are planned for resolution in upcoming releases. These include improved support for referencing object maps over SPARQL sources, quoting of SQL reserved keywords, handling of language-tagged literals, datatype conversions (e.g., xsd:boolean), and execution of custom SQL queries as logical sources. These are well-defined, bounded challenges that do not affect the core architecture, and their resolution is within reach in the short term. Beyond these refinements, several promising directions emerge for future work. One major area of development is the extension of PvRML's design and implementation to support RML extensions, such as those involving nested records, graph provenance, or incremental transformation. These features are increasingly relevant in contemporary use cases of knowledge graph construction, particularly in dynamic or federated data environments. Another research trajectory lies in the integration of large language models (LLMs) into mapping workflows. As recent studies in knowledge engineering suggest [12], LLMs can assist with the semi-automatic generation, validation, and explanation of mappings, potentially lowering the expertise required for building semantically rich data pipelines. We plan to explore how PyRML can serve as a backbone for such LLM-assisted workflows, offering human-in-the-loop interfaces and programmatic scaffolding for co-creating mappings that are both syntactically correct and semantically meaningful. In conclusion, PyRML stands as a robust, efficient, and extensible platform for declarative RDF generation in Python. It addresses a key gap in the tooling ecosystem and opens up new possibilities for research and practice in transparent, scalable, and intelligent knowledge graph construction.

Acknowledgments. This work has been supported by the Water Health Open knoWledge (WHOW) project co-financed by the Connecting European Facility programme of the European Union under grant agreement INEA/CEF/ICT/A2019/206322. Additional financial support to this project was provided by NextGenerationEU under NRRP Grant agreement n. MUR IR0000008 - FOSSR (CUP B83C22003950001).

References

- Arenas-Guerrero, J., Chaves-Fraga, D., Toledo, J., Pérez, M.S., Corcho, O.: Morph-KGC: Scalable knowledge graph materialization with mapping partitions. Semantic Web 15(1), 1–20 (2024). https://doi.org/10.3233/SW-223135
- 2. Asprino, L., Daga, E., Dowdy, J., Mulholland, P., Gangemi, A., Ratta, M.: Streamlining knowledge graph construction with a fa\c {c} ade: The SPARQL Anything project. arXiv preprint arXiv:2310.16700 (2023)
- 3. Daga, E., Asprino, L., Mulholland, P., Gangemi, A.: Facade-X: an opinionated approach to SPARQL anything. In: Further with Knowledge Graphs, pp. 58–73. IOS Press (2021)
- Dimou, A., Sande, M.V., Colpaert, P., Verborgh, R., Mannens, E., de Walle, R.V.: Rml: A generic language for integrated rdf mappings of heterogeneous data. In: Bizer, C., Heath, T., Auer, S., Berners-Lee, T. (eds.) LDOW. CEUR Workshop Proceedings, vol. 1184. CEUR-WS.org (2014)
- 5. Friedman, M., Levy, A., Millstein, T.: Navigational plans for data integration. In: Proceedings of the sixteenth national conference on Artificial intelligence and the eleventh Innovative applications of artificial intelligence conference innovative applications of artificial intelligence. pp. 67–73. AAAI '99/IAAI '99, American Association for Artificial Intelligence, Menlo Park, CA, USA (1999)
- Halevy, A.Y.: Answering queries using views: A survey. VLDB Journal 10(4), 270– 294 (2001)
- Hogan, A., Blomqvist, E., Cochez, M., d'amato, C., Melo, G.D., Gutierrez, C., Kirrane, S., Gayo, J.E.L., Navigli, R., Neumaier, S., Ngomo, A.C.N., Polleres, A., Rashid, S.M., Rula, A., Schmelzeisen, L., Sequeda, J., Staab, S., Zimmermann, A.: Knowledge Graphs. ACM Computing Surveys 54(4), 1–37 (May 2022). https://doi.org/10.1145/3447772, https://doi.org/doi/10.1145/3447772
- 8. Iglesias, E., Vidal, M.E., Collarana, D., Chaves-Fraga, D.: Empowering the sdm-rdfizer tool for scaling up to complex knowledge graph creation pipelines. Semantic Web 16(2), SW-243580 (2025)
- Iglesias-Molina, A., Van Assche, D., Arenas-Guerrero, J., De Meester, B., Debruyne, C., Jozashoori, S., Maria, P., Michel, F., Chaves-Fraga, D., Dimou, A.: The RML Ontology: A Community-Driven Modular Redesign After a Decade of Experience in Mapping Heterogeneous Data to RDF. In: Payne, T.R., Presutti, V., Qi, G., Poveda-Villalón, M., Stoilos, G., Hollink, L., Kaoudi, Z., Cheng, G., Li, J. (eds.) The Semantic Web ISWC 2023. pp. 152–175. Springer Nature Switzerland, Cham (2023)
- Lembo, D., Mora, J., Rosati, R., Savo, D.F., Thorstensen, E.: Mapping analysis in ontology-based data access: Algorithms and complexity. In: The Semantic Web-ISWC 2015: 14th International Semantic Web Conference, Bethlehem, PA, USA, October 11-15, 2015, Proceedings, Part I 14, pp. 217-234. Springer (2015)
- 11. Lenzerini, M.: Data integration: a theoretical perspective. In: Proceedings of the twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems. pp. 233–246. PODS '02, ACM, New York, NY, USA (2002). https://doi.org/http://doi.acm.org/10.1145/543613.543644
- 12. Lippolis, A.S., Saeedizade, M.J., Keskisärkkä, R., Zuppiroli, S., Ceriani, M., Gangemi, A., Blomqvist, E., Nuzzolese, A.G.: Ontology generation using large language models. arXiv preprint arXiv:2503.05388 (2025)
- Ullman, J.D.: Information integration using logical views. Theoretical Computer Science 239(2), 189–210 (2000). https://doi.org/DOI: 10.1016/S0304-3975(99)00219-4

Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., Appleton, G., Axton, M., Baak, A., Blomberg, N., Boiten, J.W., da Silva Santos, L.B., Bourne, P.E., Bouwman, J., Brookes, A.J., Clark, T., Crosas, M., Dillo, I., Dumon, O., Edmunds, S., Evelo, C.T., Finkers, R., Gonzalez-Beltran, A., Gray, A.J., Groth, P., Goble, C., Grethe, J.S., Heringa, J., Hoen, P.A., Hooft, R., Kuhn, T., Kok, R., Kok, J., Lusher, S.J., Martone, M.E., Mons, A., Packer, A.L., Persson, B., Rocca-Serra, P., Roos, M., van Schaik, R., Sansone, S.A., Schultes, E., Sengstag, T., Slater, T., Strawn, G., Swertz, M.A., Thompson, M., van der Lei, J., van Mulligen, E., Velterop, J., Waagmeester, A., Wittenburg, P., Wolstencroft, K., Zhao, J., Mons, B.: The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data 3(1), 160018 (Mar 2016). https://doi.org/10.1038/sdata.2016.18, https://doi.org/10.1038/sdata.2016.18