

# Incremental Knowledge Graph Construction from Heterogeneous Data Sources

Semantic Web  
Vol. 17(2) 1–33  
© The Author(s) 2026  
Article reuse guidelines:  
[sagepub.com/journals-permissions](https://sagepub.com/journals-permissions)  
DOI: 10.1177/22104968251412270  
[journals.sagepub.com/home/swj](https://journals.sagepub.com/home/swj)



Dylan Van Assche<sup>1</sup> , Julián Andrés Rojas<sup>1</sup> , Ben De Meester<sup>1</sup>  and Pieter Colpaert<sup>1</sup> 

## Abstract

Sharing datasets that change (through creates, updates, deletes) poses challenges to data consumers, including reconciling historical versioning and managing frequent changes. This is evident for Knowledge Graphs (KGs), materialized from such datasets, where synchronization happens through frequent regeneration. However, this is time-consuming, loses history, and wastes computing resources through redundant processing. We present a KG generation approach that efficiently handles evolving data sources with different change signaling strategies. We investigate change signaling strategies of real-world datasets, propose corresponding change detection algorithms, and introduce a declarative approach based on the RDF Mapping Language (RML) and Function Ontology to materialize changes for evolving KGs. Detected changes can be automatically published as a Linked Data Event Stream (LDES), using the Activity Streams 2.0 vocabulary to describe changes and communicate them over the Web. We implement our approach in the RMLMapper as Incremental RML and evaluate it both functionally, and quantitatively using a modified version of the GTFS Madrid Benchmark and several real-world data sources. Our approach reduces storage and computing requirements for generating and storing multiple KG versions (up to 315.83x less storage, 4.59x less CPU time, and 1.51x less memory) and reduces KG construction time up to 4.41x. Performance gains are more pronounced for larger datasets, while our approach's overhead partially offsets benefits for smaller ones. Overall, our approach lowers the cost of publishing and maintaining KGs and, via LDES, supports timely, Web-native dissemination of changes. We plan to optimize our change detection algorithms and use windowing to support streaming data.

## Keywords

resource description framework mapping language, linked data event stream, incremental knowledge graph construction

Received: March 18, 2024; accepted: October 22, 2025

## 1 Introduction

Most real-world datasets are not static. Whether it be the inclusion of new data points, the rectification of errors, or the evolution of data collection methodologies, datasets are, by their very nature, subject to continuous change. The change frequency may vary, with some datasets undergoing rapid transformations while others experiencing a more gradual evolution. Regardless of the particular change frequency, data consumers find themselves impacted by these alterations. Consumers must take this continuous change flow into account when using the data (Valencio et al., 2013) and incorporate these changes to stay synchronized with the current state of the original dataset. This leads them to face challenges such as storing multiple versions to keep historic records, having enough capacity to handle any dataset's size when processing

<sup>1</sup>IDLab, Department of Electronics and Information Systems, Ghent University – imec, Belgium

### Corresponding Author:

Dylan Van Assche, IDLab, Department of Electronics and Information Systems, Ghent University – imec, Belgium.  
Email: [dylan.vanassche@ugent.be](mailto:dylan.vanassche@ugent.be)



changes, and keeping up with dataset’s change frequency. Moreover, data source changes impact data integration processes, such as Knowledge Graph (KG) generation, which are forced to fully repeat their computations *every time one of the original data sources changes* (Denny et al., 2017; Rojas et al., 2021). This approach could become unsustainable when dealing with large volumes of data. It could also affect applications depending on (near) real-time information: A generated KG supporting the application might be already outdated if the source data sources change faster than the (re-)generation process. In general, this may result into: (i) Wasted computing resources and time due to unnecessary reprocessing of all data, especially when large parts remain unchanged, (ii) loss of historic records given that commonly only the latest version remains available, and (iii) outdated KGs when the data integration process is not able to keep up with the change frequency of the original data sources.

Change Data Capture (CDC) techniques (Gupta & Giri, 2018; Hao et al., 2023) try to address this problem by capturing and characterizing data changes using, for example, database transaction logs, database triggers, comparing snapshots of data dumps, etc. Denny et al. (2017). However, current approaches for integrating heterogeneous data sources into KGs (Daga et al., 2021; Das et al., 2012; Dimou et al., 2014; Lefrançois et al., 2017; Vu et al., 2019) cannot use existing CDC techniques, as they mainly focus on a specific type of data source (e.g., relational database logs). In this paper, we introduce a *declarative and incremental KG construction* approach, capable of handling heterogeneous and changing data sources based upon a CDC-based technique. Our approach can cope with different data change signaling strategies, allowing to detect and incorporate these changes in an incremental and continuous fashion into an existing KG. We enumerate and discuss different *data change signaling strategies* (both explicit and implicit) observed in real-world datasets, that are used to communicate changes to data consumers. We also propose a set of algorithms to handle these change signaling strategies for and during KG construction. Our approach only regenerates the parts of a KG that were changed in the original data source(s) (created and updated), while deleted data entities are detected, labelled as such and made explicit.

We opt for a declarative approach to incrementally generate KGs given its engine-agnostic and extensible nature to handle different types of data sources (Van Assche et al., 2022a). We implement our approach using resource description framework mapping language (RML) as declarative mapping language to generate resource description framework (RDF), and function ontology (FnO) for describing data transformations and state management functions, which together constitute our approach called Incremental RML (*IncRML*). However, our approach is not dependent on these particular technologies and could be implemented with any other declarative mapping language and data transformation vocabulary. Furthermore, we allow to optionally and automatically publish these changes on the Web in the form of a Linked Data Event Stream (LDES) (Van Lancker et al., 2021). LDES allows publishing data changes as an append-only event log that consumers can read to synchronize their KG with the latest data available. We show how it is possible to produce an LDES from the detected data changes, materialized with our IncRML approach, which we choose to describe semantically with the W3C Activity Streams 2.0 vocabulary (Snell & Prodromou, 2017).

We extend our previous work (Van Assche et al., 2022b) on aligning LDES and RML, where we introduced a preliminary approach for continuously generating a KG by detecting and materializing changes in its data sources. Our previous work focussed only on handling data creations and updates. In this paper, we extend that work to also detect and communicate deletions. We also include extensive functional and performance evaluations using (i) different types of real-world datasets, (ii) an extended version of the GTFS Madrid Benchmark (Chaves-Fraga et al., 2020) that allows to control the type and amount of changes in a data source, and (iii) a set of test cases that guarantee the support for different change signaling strategies. To keep this paper self-contained, we include the discussions of our previous work along with its extensions. Concretely, our contributions are the following:

1. An overview of different history and change signaling strategies observed in real-world data sources.
2. A set of CDC-based algorithms to detect changes in any of these change signaling categories.
3. The integration of these algorithms into a KG construction pipeline using RML, FnO, and (optionally) LDES.
4. An extensive evaluation and benchmark on the impact and performance of incrementally generating and updating KGs.

Overall, our approach is more efficient in terms of execution time and computing resources for generating the RDF quads of a KG, in exchange for a small overhead when the KG is constructed the first time. We observe that by materializing only detected changes when generating KGs, our approach uses 3.24–315.83x less storage to store multiple versions of a KG while also consuming less computing resources (0.85–4.59x less CPU time, 0.72–1.51x less memory consumption) depending on the dataset. The RDF generation time is reduced by 0.97–4.41x depending on the data size. To guarantee a fair comparison with a traditional full KG re-materialization approach, we measured the total time that it takes to go from existing data to an updated triplestore using SPARQL `UPDATE` queries, so that a consistent and up-to-date KG is achieved for both cases. For our approach, an additional step is required to produce the corresponding SPARQL query that

either creates, updates or deletes data entities. Despite this additional step, we observe that our approach is 11.07–57.66x faster to ingest all versions of a dataset, and on average 4.5–28.5x faster to ingest individual dataset updates. However it was only possible to measure this improvement on datasets whose updates were not larger than 20K quads, since the triplestore would fail due to reaching internal query length limits, to execute the SPARQL UPDATE queries for both the traditional and our approach.

IncRML is usable for any kind of data source, as supported by the underlying declarative mapping language. While we use RML as mapping language in our implementation, the logical definitions and algorithms of IncRML may be implemented with other mapping languages since they rely on widely supported features such as IRI templates and data transformations (Van Assche et al., 2022a). IncRML also allows for semantically describing data changes using any ontology of choice (e.g., W3C Activity Streams 2.0) and then publishing such changes with a structured approach such as LDES or any other data publishing strategy. Thanks to our work, data source changes can be integrated faster and with less resources into live and replicated KGs while allowing to keep access to the historical records in the form of an LDES. Although we observed that further optimizations are needed for triplestores to effectively support larger data updates through standard SPARQL UPDATE queries.

The remainder of this paper is structured as follows: Section 2 discusses related works, Section 3 introduces the main technological concepts used in this work, namely RML, FnO, and LDES. Section 4 presents different identified change signaling strategies and describes the rationale of our approach. Section 5 shows how we implement our approach. In Section 6, we present our evaluation design. Section 7 discusses our results, and Section 8 concludes.

## 2 Related Work

In this section, we discuss related work on (i) mapping rules for declarative KG generation describing how a KG can be generated, (ii) CDC approaches for detecting changes in data sources, (iii) versioning strategies for KGs to store the history of data sources and its impact on storage for producers and consumers, (iv) versioned generation of KGs, and (v) incremental mapping rules execution for optimizing execution time and resource usage.

### 2.1 Mapping Rules for Declarative KG Generation

Declarative mapping rules for KG generation is an active research domain since the introduction of the R2RML W3C Recommendation (Das et al., 2012) for transforming relational databases into an RDF KG (Cyganiak et al., 2014). R2RML was extended as the RML (Dimou et al., 2014; Iglesias-Molina et al., 2023) to support heterogeneous data sources (e.g., JSON, XML, CSV) while keeping backwards compatibility with R2RML. Recently, a survey (Van Assche et al., 2022a) was performed of existing approaches and systems for declarative KG generation from heterogeneous data. RML is widely used for declarative KG generation and was extended to support exporting RDF to various targets such as files and SPARQL endpoints (Van Assche et al., 2021), RDF Collections and Containers (Debruyne et al., 2017; Michel et al., 2017), and access to Web APIs (Chortaras & Stamou, 2018; Van Assche et al., 2021).

Besides RML, other declarative mapping languages were proposed to transform heterogeneous data sources into RDF such as xR2RML (Michel et al., 2015, 2017), SPARQL-Generate (Lefrançois et al., 2017), SPARQL-Anything (Daga et al., 2021), ShExML (García-González et al., 2020), D-REPR (Vu et al., 2019), and OTTR (Skjaveland et al., 2018). xR2RML also extends R2RML with support for heterogeneous data sources and RDFS Collections and Containers (Brickley & Guha, 2014). SPARQL-Generate and SPARQL-Anything do not extend R2RML, but SPARQL instead to transform heterogeneous data sources into a KG. Therefore, they can reuse existing SPARQL engines and syntax. Similar to SPARQL-based approaches, ShExML uses Shape Expressions (ShEx) (Prud'hommeaux et al., 2014) as syntax, while D-REPR defines its own syntax.

Currently, the W3C Community Group on KG Construction<sup>1</sup> is working on standardizing RML as a W3C Recommendation (Iglesias-Molina et al., 2023) and is supported by multiple implementations such as the RMLMapper (Dimou et al., 2014), Morph-KGC (Arenas-Guerrero et al., 2022), or SDM-RDFizer (Iglesias et al., 2020). Therefore, we decided to implement our approach with RML as the declarative mapping language in this work, but any mapping language may be used.

Data transformation support is an important requirement when generating KGs from heterogeneous data sources (Van Assche et al., 2022a), using, for example, the FnO (De Meester et al., 2020), SPARQL Functions (Harris & Seaborne, 2013), or FunUL (Junior et al., 2017). FnO is a popular vocabulary for describing functions to perform data transformations, and it is integrated with RML through FNML<sup>2</sup>. This way, RML+FnO mapping rules can perform both the generation of a KG and data transformations without requiring ad hoc or use case specific scripts. We use FnO in this paper to integrate the implementation of our change detection algorithms with RML to incrementally generate a KG from heterogeneous data sources.

## 2.2 Change Data Capture

CDC (Gupta & Giri, 2018; Hao et al., 2023) refers to a technique, primarily used in databases, to identify and capture data changes so that those changes can be tracked, recorded, and propagated to other systems or applications. CDC's primary purpose is identifying and capturing creations, updates, and deletions of data in a dataset to enable (near-)real-time synchronization of data across different systems. Most approaches focus on data sources which provide some sort of change signaling mechanism such as transaction logs, snapshots, triggers, or timestamps. Moreover, data sources may offer different granularity regarding change signaling, for example, change signaling on parts of the data source or the whole data source (Umbrich et al., 2010).

*Log-based approaches* (Hao et al., 2023; Ma & Yang, 2015) use database transaction logs to determine which changes were performed to the underlying data. As log systems are implementation-specific, these approaches are specific to each database. *Snapshot-based approaches* (Denny et al., 2017) compare the current version of a data source with previous versions to extract changes, requiring sufficient resources to store and compare these versions. *Trigger-based approaches* (Hu et al., 2019; MadeSukarsa et al., 2012; Valencio et al., 2013) hook into a data source to execute a trigger on each change, requiring support for triggers from the data source (e.g., stored procedures in relational databases). *Timestamp-based approaches* (Goyal & Dyreson, 2019) analyze a last-modified timestamp of a data source to detect and extract changes, requiring data sources to provide timestamp-annotated data.

Although these approaches clearly have their merit, many data sources do not signal their changes nor support triggers when a data record is changed (e.g., data streams, files, or Web APIs). Therefore, existing CDC approaches are insufficient to cover heterogeneous data sources which do not signal their changes to consumers. In this work, we combine and extend *Timestamp-based* and *Snapshot-based* approaches for detecting implicit and explicit changes.

## 2.3 Versioning of KGs

Several approaches for versioning of (RDF) KGs have been proposed (Papakonstantinou et al., 2016). Three main RDF archive storage strategies can be identified (Fernández et al., 2015): (i) Change-Based, (ii) Timestamp-Based, and (iii) Independent Copies. *Change-Based* only stores the changes; *Timestamp-Based* uses timestamps to define when a specific version is valid; and *Independent Copies* stores a copy of the data source each time it is updated.

*Change-Based* approaches include: R&Wbase (Vander Sande et al., 2013), based on Git<sup>3</sup>; a Version Control Based RDF storage approach using patches, similar to Frommhold et al. (2016); Cassidy and Ballantine (2007), SemVersion Völkel et al. (2005), R43ples Graube et al. (2014), and Im et al. (2012). *Timestamp-Based* approaches for accessing different data source versions include: Memento (HTTP) (de Sompel et al., 2013) and x-RDF-3X (Neumann & Weikum, 2010) (SPARQL). OSTRICH (Taelman et al., 2018) and TailR (Meinhardt et al., 2015) are both *hybrid* approaches, combining all three strategies for efficient query operations. All approaches implementing the 3 strategies put the burden of resolving versions on the data producer.

LDES (Van Lancker et al., 2021) uses an *Independent Copies* approach on an entity level (aka. member in LDES terminology) for versioning. An LDES entity/member may be defined as a named node and its properties, as defined by its Concise-Bounded Description<sup>4</sup>, or as named graph and its contained triples. However, Independent Copies suffers from scalability problems as storage is not infinite (Fernández et al., 2015), LDES addresses this by dropping the oldest versions of entities in the event stream, according to a specific and configurable *retention policy*. Since LDES consumers are aware of the retention policy, they can decide to store the LDES members themselves if they need to for their particular use case. For example, if an LDES producer's retention policy is 7 days and a consumer requires at least the last 30 days of data, the consumer must store a copy. If a consumer does not require history information longer than 7 days, it can solely rely on the LDES producer's data. LDES allows consumers to synchronize their local copy of the member collection, similar to a Copy and Log approach (Salzberg & Tsotras, 1999). This way, versioning is resolved on the consumer side and several versioning strategies can be applied independent of the publisher. In this work, we allow our approach to use LDES as a publishing strategy for communicating KG changes on the Web.

## 2.4 Versioned Generation of KGs

Ontologies (Change Detection Ontology (Randles & O'Sullivan, 2022, 2023)) and benchmarks (EvoGen Benchmark Suite (Meimaris & Papastefanatos, 2016), BEAR benchmark (Fernández et al., 2015)) were proposed for versioned KGs. However, they are tied to a specific ontology or focus on querying the versions while we focus on the generation in this work. The EvoGen Benchmark Suite (Meimaris & Papastefanatos, 2016) allows generating synthetic versioned RDF data for benchmarking purposes. BEAR (Fernández et al., 2015) proposed a benchmark for Semantic Web archiving systems to evaluate full materialization of different versioned KGs, only materializing the changes, or annotating triples when they

were created, updated, and deleted. However, both focus on materialized RDF for benchmarking query systems, while in this work, we focus on the generation of different KG versions from non-RDF heterogeneous data. The Change Detection Ontology (Randles & O’Sullivan, 2022, 2023) allows describing changes inside the original data sources as a changelog and is used in the MQ framework (Randles et al., 2022) to generate new KGs if changes are detected, using R2RML mappings from CSV, XML, and relational databases. However, this solution is tightly coupled with a specific ontology, while we aimed for a more generalized approach, that allowed to use any ontology to describe the detected changes.

## 2.5 Incremental Mapping Rules Execution

Execution planning of mapping rules has seen uptake in the KG community (Van Assche et al., 2022a) as seen in tools such as Morph-KGC (Arenas-Guerrero et al., 2022) and SDM-RDFizer (Iglesias et al., 2020). Both systems plan their execution and remove duplicates before they execute RML mapping rules to reduce the number of data records and improve performance. However, they consider that executing these mapping rules only happens once: If the datasets change, all the mapping rules must be fully executed again. Thus, incremental KG generation is not considered. Besides execution planning with Morph-KGC and SDM-RDFizer, existing work on incremental KG generation does not consider heterogeneous data as it focus on relational databases such as using log files to determine changes (Konstantinou et al., 2014), triggers to update virtualized views (Vidal et al., 2013), or indexing triples (Pu et al., 2014). There is no approach which supports heterogeneous data – besides relational databases – for example, JSON, XML, CSV. In this work, we present a novel approach to also detect changes in heterogeneous data and making them available to consumers, even if the data sources do not signal their changes. While our experiments focus on incremental KG generation from single data sources, multiple and heterogenous data source cases are also supported through the use of RML.

## 3 Background

In this section, we provide an introduction to the (i) RML, (ii) FnO, and (iii) LDES. These technologies came forward from Section 2 as prevalent methods to describe how a KG could be constructed from input datasets (RML), performing data transformations and generic functions during KG construction (FnO), and semantically describing and continuously publishing the KG changes (LDES) to allow for KG replication and synchronization.

### 3.1 RDF Mapping Language

RML (Dimou et al., 2014; Iglesias-Molina et al., 2023)<sup>5</sup> is an extension of W3C Recommendation R2RML (Das et al., 2012) to support heterogeneous data sources besides relational databases. RML mapping rules consist of Triples Maps (Listing 1: Lines 1–17) which define how the terms (subject, predicate, and object) of an RDF triple are generated. A named graph can also be specified using a Graph Map for generating RDF quads. Each Triples Map has one Logical Source (Listing 1: Lines 2–4), one Subject Map, and zero or more Predicate Object Maps. The Subject Map (Listing 1: Lines 6–9) defines how the subject IRIs are generated from the data source as defined by the Logical Source. This Subject Map also includes a Graph Map to specify the named graph of the RDF quad (Listing 1: Line 8). Predicate Object Maps (Listing 1: Lines 11–16) consist of Predicate Maps (Listing 1: Line 12) to specify the quad’s predicate and (Referencing) Object Maps (Listing 1: Lines 13–15) for the quad’s object. The Subject Map, Predicate Map, Object Map, and Graph Map are all Term Maps, generating an RDF term (an IRI, blank node, or literal). A Term Map may always generate the same RDF term with `rr:constant`, a referenced value from the data source with `rml:reference`, or construct a value based on a template with `rr:template`. If a Subject Map, Predicate Map, or Object Map uses an `rr:constant` for its RDF term generation, a shortcut can be used, for example, `rr:predicate`.

### 3.2 Function Ontology

The FnO (De Meester et al., 2020) semantically describes and declares implementation-independent functions and their relations to related concepts such as input parameters, outputs, mappings to concrete implementations, and executions. The alignment between RML and FnO (via FNML (De Meester et al., 2020)) specifies how a data transformation must be performed by specifying the function to execute (Listing 2: line 5) and its values (Listing 2: Lines 8–11). FnO is standalone which allows describing data transformations in a declarative way with or without RML. FnO is integrated in RML through an `fnml:FunctionMap` (Listing 2: Lines 1–14) which is an RML Term Map. Therefore, an FnO function can be used as Subject Map, Predicate Map, Object Map, or other RML Term Maps.

```

<#RMLMapping> a rr:TriplesMap;
  rml:logicalSource [ a rml:LogicalSource;
    rml:source "/path/to/data.csv";
  ];

  rr:subjectMap [ a rr:SubjectMap;
    rr:template "http://example.org/{ID}";
    rr:graphMap [ rr:constant ex:MyGraph ];
  ];

  rr:predicateObjectMap [ a rr:PredicateObjectMap;
    rr:predicate foaf:name;
    rr:objectMap [ a rr:ObjectMap;
      rml:reference "name";
    ];
  ];
.

```

**Listing 1.** RML uses Triples Maps with a Logical Source, a Subject Map, Graph Map, and zero or more Predicate Object Maps to specify how RDF quads must be generated from the referenced data source.

```

<#FunctionMap> a fnml:FunctionMap
  fnml:functionValue [
    rr:predicateObjectMap [ a rr:PredicateObjectMap;
      rr:predicate fno:executes;
      rr:object grel:toUppercase;
    ];
    rr:predicateObjectMap [ a rr:PredicateObjectMap;
      rr:predicate grel:inputString;
      rr:objectMap [ a rr:ObjectMap;
        rml:reference "name";
      ];
    ];
  ];
.

```

**Listing 2.** FnO defines a data transformation by specifying the function and the function's values. FnO is aligned with RML through a `fnml:FunctionMap` which is an RML Term Map. The function `toUppercase` is executed on all referenced name data values of the data source.

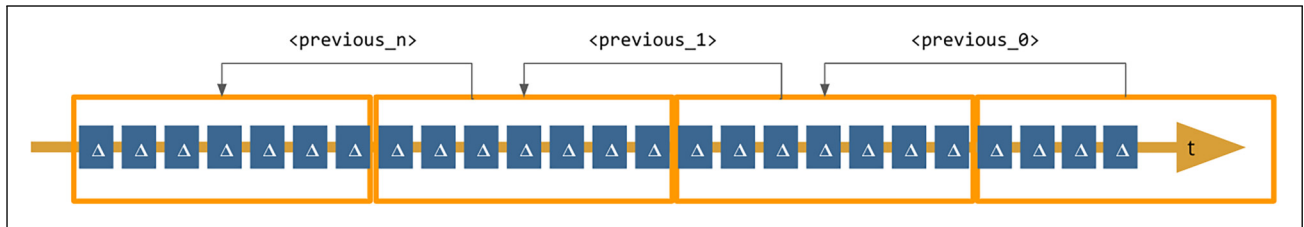
### 3.3 Linked Data Event Streams

LDES is an RDF data publishing approach fostered by the EU Semantic Interoperability Community<sup>6</sup> and officially adopted as a standard specification by the Flemish government through its Flemish Smart Data Space project<sup>7</sup>. LDES defines datasets in terms of a collection of immutable objects (a.k.a. members) such as versioned entities or observations (Listing 3), where every member must have its own unique IRI (Van Lancker et al., 2021). The main goal of a LDES is to enable efficient replication and synchronization of datasets over the Web. LDES allows data consumers to traverse the collection by relying on the TREE specification (Colpaert, 2022)<sup>8</sup> to semantically describe hypermedia relations among subsets or fragments of the data (Figure 1). These hypermedia relations can be defined in multiple ways, for example, by publishing fragments organized by time or by version, configuring tree-like data structures that can be efficiently traversed by clients. TREE also allows further describing the content of each member in an LDES collection through a SHACL shape (`tree:shape`, Listing 3: Line 4). Consumers may use the related SHACL shape to validate but also to understand the type and properties of the LDES members (`tree:member`, Listing 3: Line 5), for example, for discoverability and source selection purposes. Additional metadata on how a collection of LDES members is structured, is available via properties such as `ldes:timestampPath` (Listing 3: Line 2) or `ldes:versionOfPath` (Listing 3: Line 3), which describe the property paths that provide timestamp and versioning information in every member, similarly to the SHACL `sh:path` property. This way, LDES remains domain model agnostic and members are not limited to a specific ontology for describing timestamp-based versions or referring to other members' versions. Consumers may poll or subscribe to an LDES to obtain the latest (versions of) members published, which then can be further processed and materialized into a consumer's local environment to remain synchronized with the original data source.

## 4 Approach

In this section, we describe our incremental KG construction approach, which consists of the following high-level steps:

1. Detect changes in any heterogeneous dataset;
2. construct RDF data from these changes;
3. explicitly publish the changes with explicit semantics.



**Figure 1.** Linked data event stream (LDES)’s structure allows to traverse the collection by clients with semantic descriptions of the collection’s relations.

```

<#LDESDataCollection> a ldes:EventStream;
ldes:timestampPath sosa:resultTime;
ldes:versionOfPath dcterms:isVersionOf;
tree:shape <http://example.org/shacl/shape/>;
tree:member <LDESMember>;
.

<#LDESMember> a sosa:Observation;
sosa:resultTime "2023-01-01T00:00:00Z"^^xsd:dateTime;
sosa:versionOfPath <http://example.org/sensor/result/>;
sosa:hasSimpleResult "5"^^xsd:integer;
.

```

**Listing 3.** A data collection as an LDES with one member. The LDES member provides a sensor value at a given time and version.

Section 4.1 introduces the overall approach. How to detect changes in any heterogeneous dataset requires us to first understand how a particular dataset may signal changes (step 1a), hence, Section 4.2 discusses the various change signaling strategies along 3 dimensions: History, change communication, and change types. Then, Section 4.3 describes the algorithms to detect both implicit and explicit changes in data collections (step 1b). Where Sections 4.1–4.3 focus on detecting and describing changes in existing datasets, Section 4.4 contextualizes this work as part of an entire processing pipeline, from original data to a published event stream of RDF data, integrated in a mature KG construction and publication pipeline (steps 2 and 3). The overall approach is shown in Figure 3, top (green) row.

#### 4.1 Incremental RML

IncRML refers to the implementation of our proposed approach for incrementally generating a KG from heterogeneous data sources. It consists mainly of 2 high-level steps: (i) Detect changes in data sources, independent of whether they are signaled explicitly or implicitly, and (ii) enriching the original target ontology mapping to include additional metadata that makes explicit which quads are created, updated, and deleted. In general, our approach maximally relies on existing standards and specifications, both for the generation and the publishing of KGs. We aim at reducing execution time and resource consumption during the RDF generation process, as only the changes between consecutive data source versions are materialized. Through a CDC-based approach, we detect and extract changes during the KG generation process and (re-)generate the RDF triples/quads of all the entities (or members in LDES terminology) affected by such change. Each materialized member may include additional metadata specifying, for example, the type of change, the time of change, etc., as specified by LDES. Our approach does not limit the description of changes in data members to a specific ontology, thus it may be applied to any data collection modelled by an ontology with the semantics and expressivity to guarantee unique identification and describe member changes, or extensions thereof. This allows consumers to keep their local version of a KG in sync with the original producer across the Web, by interpreting the change semantics present in the materialized members and performing the corresponding create, update, and delete operations instead of fully re-fetching and re-ingesting a complete version of the KG every time there is an update.

Currently, RDF triplestores lack support for integrating LDES data directly, which poses the need for an additional intermediate step to interpret and execute the corresponding SPARQL UPDATE operation for each LDES member (e.g., INSERT DATA, DELETE WHERE, etc.), in order to keep the triplestore up to date. In this work we focus on studying the impact of our approach on the generation aspect of a CDC-based KG publishing system that adheres to the LDES specification, and provide a proof-of-concept implementation of an interpreter library that translates the change semantics of LDES members into the corresponding SPARQL UPDATE queries that ingest data updates into a target triplestore.

**Table 1.** Change Signaling Strategies According to Their History Availability and Change Communication.

History Availability	Change Communication	Change Types
Latest state	Explicit	Create Update
	Implicit	Delete
Latest changes	Explicit	Create Update
	Implicit	Delete
Full history	Explicit	Create Update
	Explicit	Delete
	Implicit	Delete

Change communication can differ per type of change, even in the same dataset.

## 4.2 Change Signaling Strategies

We identified a set of change signaling strategies differing along three dimensions: (i) availability of historical records, (ii) how changes are communicated to consumers, and (iii) type of change; by analyzing real-world datasets from various domains for example, bike-sharing data, public transport timetables, geographical data, traffic data, and meteorological data (Section 6.1.3).

*History.* We identified 3 different types of history availability:

- *latest state:* Latest state refers to datasets that publish only the latest version of all its members on every change.
- *latest changes:* Latest changes refers to datasets that publish only changed members (aka delta updates). Thus, consumers must have access to an initial complete version of the data upfront and reconcile updates over it.
- *full history:* Full history refers to datasets that are published including both historical and current versions of its members. The number of available versions is defined by the data publisher.

*Change communication.* We identified 2 change communication strategies:

- *explicit:* Dataset changes are *explicitly* communicated if metadata is also provided to point that a change has occurred (e.g., via uniquely identified members using timestamps, hashes, or logs).
- *implicit:* Dataset changes are *implicitly* communicated if members are changed without providing any kind of metadata indicating that a change happened, for example changes such as property updates, member deletions, or member creations, all happening silently across new versions of the data collection.

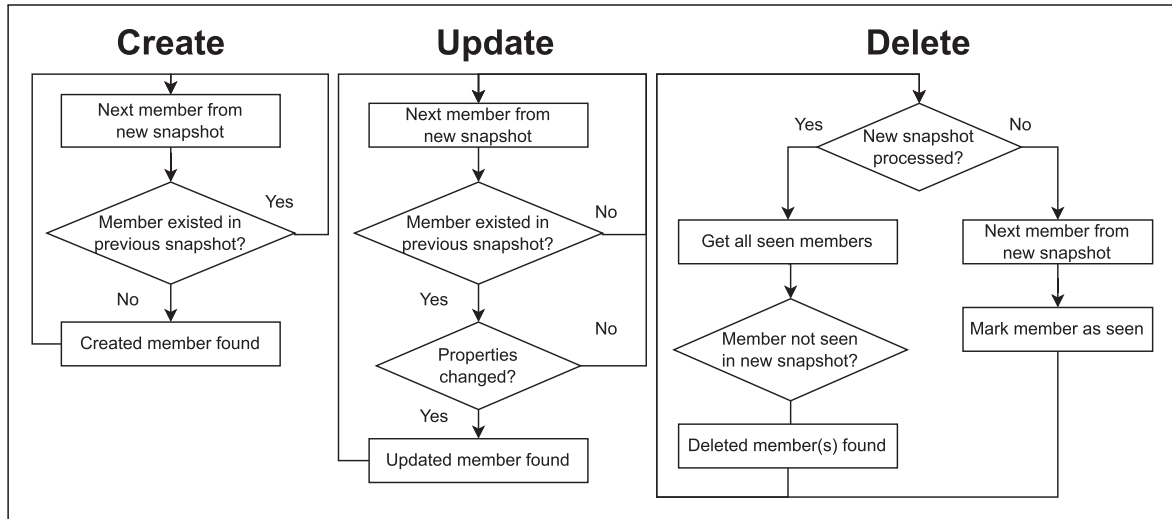
*Change types.* We identified 3 change types:

- *create:* A member is added to the dataset and it didn't exist before.
- *update:* An existing member of the dataset is modified. New properties are added to the member or existing properties are modified.
- *delete:* An existing member of the dataset is deleted.

Moreover, change communication may differ depending on the type of change. For example: Created and updated members may have a unique identifier (explicit change), while deleted members are simply removed in newer versions of a data collection (implicit change). Table 1 shows a summary of the different change signaling strategy combinations with respect to history availability, change communication, and change types identified on the set of real-world datasets analyzed in this work.

## 4.3 Implicit and Explicit Change Detection

Datasets communicate their changes mainly in 2 ways, regardless of historical records availability: (i) Explicitly through uniquely identified members, logs, etc., and (ii) implicitly by *silently* changing dataset members. The latter imposes the need for consumers to detect these changes themselves. Our approach (Figure 2) combines Timestamp-based and Snapshot-based CDC approaches (Denny et al., 2017) to handle both explicit and implicit changes. In our approach, explicit changes are detected by relying on uniquely identified dataset members (e.g., via subject IRIs that depend on last



**Figure 2.** Implicit updates must be detected by a change data capture (CDC) algorithm. Creates and updates are detected by checking if the member already exists in a previous snapshot. Members in a dataset are identified through their generated IRI. If the IRI does not exist in a previous snapshot, we conclude that the member was created. If the IRI does exist and the properties of the member have changed, we conclude that the member was updated. Deletions are detected by checking if the member is no longer present compared to the previous snapshot of the dataset.

modified timestamps). We detect implicit changes by comparing consecutive snapshots of dataset members, checking their subject IRIs, and (a subset of) their corresponding properties for changes.

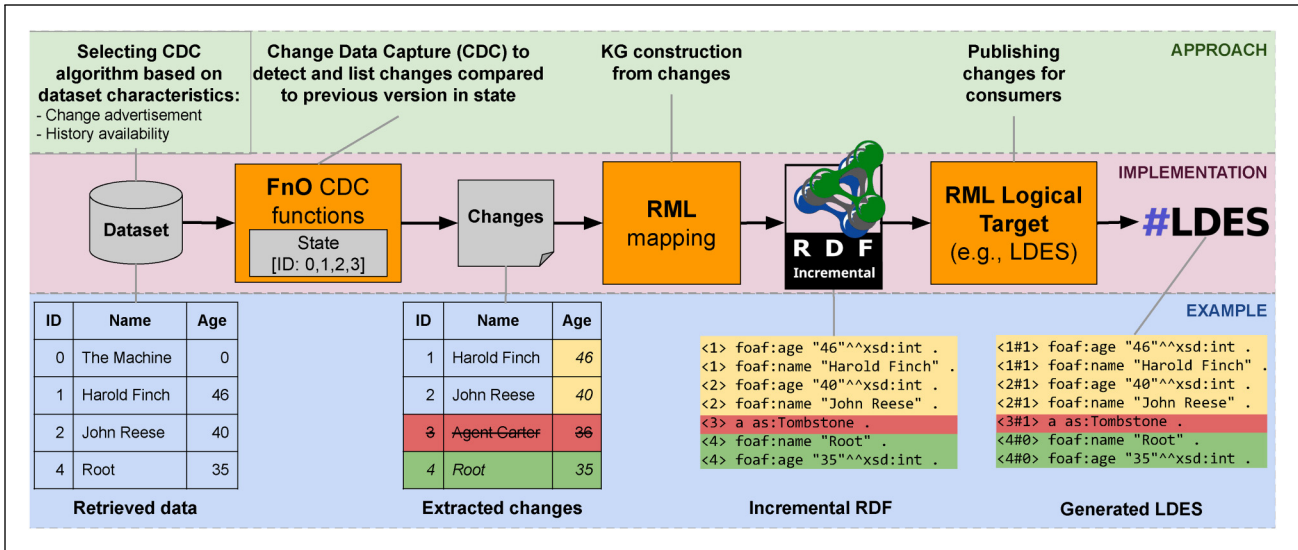
It is important to note that implicit deletion communication is not possible when combined with full history or latest changes. *Full history* datasets must communicate deletions explicitly otherwise a data consumer cannot determine that a member was deleted implicitly, since it will encounter it as part of the historical records also present in the dataset. Even when a data consumer can tell apart historical records from new data, if a member stops occurring in a full history dataset, it could be ambiguously interpreted either as being deleted or as a non-updated member. Assuming that the member is to be deleted could result in both false-positive and false-negative deletion detection. *Latest changes* datasets are similar to *full history* regarding deletions. If such dataset do not explicitly communicate that a member is deleted, consumers cannot determine with full certainty if a member was deleted or simply remains unchanged.

Explicit changes can be directly detected and do not require additional processing effort during the KG generation process. Implicit changes, however, require a stateful processing approach, to keep track of members' state across KG generation executions. Our approach introduces 3 algorithms to handle implicit changes, which scale in direct proportion to the number of members in the dataset. Each detection algorithm corresponds with a particular type of change: (i) *create*, detects when a new member is added to the dataset, (ii) *update*, detects when an existing member is modified in the dataset, and (iii) *delete*, detects when an existing member is deleted from the dataset. Next, we describe the logical flow of each algorithm in the case of implicit change communication.

*Implicit Create.* (Figure 2, left) For every dataset member being processed in the current version, the algorithm checks if the member was already present in a previous version of the dataset, based on its subject IRI. If not, a *created* member is found.

*Implicit Update.* (Figure 2, middle) Similar to create, the algorithm checks if the member was already present through its subject IRI and its properties. The algorithm does not necessarily check all properties of a member. It can simply check a predetermined subset which were labelled as *watched properties*. These properties can be used to, for example, compute a hash that allows to determine if a change has occurred. If the member's subject IRI was present and its watched properties' values were changed, a *updated* member is found.

*Implicit Delete.* (Figure 2, right) At KG generation time, the algorithm marks every generated member as seen. Once all members have been processed in the current version of the dataset, it identifies the members which were not seen in the current execution, with respect to the previous KG generation execution (if any). These are marked as the *deleted* members in the current version of the dataset.



**Figure 3.** Our approach IncRML (top green row) for which we use RML+FnO (middle pink row). We use FnO described functions to perform CDC based on the characteristics of the dataset and RML for constructing RDF from the detected changes. Changed RDF quads may be published as an LDES via an LDES Event Stream Logical Target. The pipeline is continuously executed to extract changes from new versions of the datasets. Our approach can be used by any RML engine with support for FnO (orange squares). Example data (bottom blue row) shows how data creations (green), data updates (yellow), and data deletions (red) are detected through CDC FnO functions. It is assumed that the previous state contains info on data rows with IDs 0, 1, 2, and 3. The extracted changes are then incrementally transformed into RDF and published as LDES members. IncRML: incremental resource description framework mapping language; RML: resource description framework mapping language; FnO: function ontology; CDC: change data capture; LDES: linked data event stream.

#### 4.4 Incremental KG Construction and Publishing Pipeline

Our approach effectively brings together (i) RML for declaratively defining generation rules for a KG; (ii) FnO for defining change detection functions across data source versions (Section 4.3); and (iii) LDES as an optional publishing strategy to publish semantically described changes as an event stream for consumers on the Web. Figure 3 presents a schematic view of a data processing pipeline using our approach to incrementally construct and publish a KG. If changes are detected, the corresponding RML mapping(s) is/are executed to generate and publish only the changed members to be later integrated into the KG. In practice, a member is generated by a RML Triples Map<sup>9</sup> (`rr:TriplesMap`) with one Subject Map (`rr:SubjectMap`) and zero or more Predicate Object Maps (`rr:PredicateObjectMap`).

Our CDC FnO functions monitor the Subject Maps and their correspondent Predicate Object Maps at execution time, to determine if there were any changes with respect to the previous execution. If so, the correspondent member is materialized and (optionally) published as a typed event in an LDES. As an example, we show how we can semantically annotate changed members using the W3C Activity Streams 2.0 vocabulary: `as:Create`<sup>10</sup> for created members, `as:Update` for updated members, `as>Delete` for deleted members. Created and updated members are fully re-materialized while for deleted members only a tombstone (a simplified version of the member) which indicates that the member was deleted without any other properties is generated. Annotating these changes with the type of change, allows consumers to interpret and incorporate these changes over their local copies of the KG. If published as an LDES, data consumers can subscribe to it to receive the changes and remain in sync over the Web. Our approach generates a unique IRI for version of a member as defined by the LDES specification<sup>11</sup>. Note that this requirement is solely for publishing a valid LDES and it is not required to perform the change detection process discussed in Section 4.3. Ultimately, both LDES and non-LDES data consumers can interpret the change semantics present in the materialized members produced by our incremental generation approach and use, for example, SPARQL UPDATE queries to incorporate the changed members into their KG.

## 5 Implementation

In this section, we describe how we implement our proposed approach (Section 4) in the RMLMapper<sup>12</sup> based on a set of CDC FnO functions, RML mapping rules to construct a KG, and LDES for publishing an incrementally constructed KG as an event stream. Section 5.1 explains how we implement our CDC FnO functions, Section 5.2 discusses how we

**Table 2.** Function Ontology (FnO) Functions for Explicitly and Implicitly Communicated Changes in Members.

FnO function	Purpose
<code>idlab-fn:explicitCreate</code>	Detect explicitly created members by checking if the member IRI existed already.
<code>idlab-fn:explicitUpdate</code>	Detect explicitly updated members by checking if the member IRI existed already.
<code>idlab-fn:explicitDelete</code>	Detect explicitly deleted members by checking if the member IRI existed already.
<code>idlab-fn:implicitCreate</code>	Detect implicitly created members by checking if the member IRI existed already.
<code>idlab-fn:implicitUpdate</code>	Detect implicitly updated members by checking if the member IRI already existed and its watched properties have changed.
<code>idlab-fn:implicitDelete</code>	Detect implicitly deleted members by marking each member IRI as seen and after processing all the members, returning the set of member IRIs which were not seen compared to the previous version.

Each type of change (creation, update, and deletion) has its own dedicated function.

integrated RML with the CDC FnO functions, Section 5.3 explains how we implemented an LDES Logical Target for directly generating LDES annotated data, and Section 5.4 demonstrates how our implementation is applied on an example dataset.

### 5.1 CDC FnO Functions

We implement our CDC algorithms from Section 4.3<sup>13</sup> as FnO described functions. A dedicated FnO function was implemented for each change type and change communication strategy (Table 2). This way, we semantically describe which type of detection is applied on the dataset, even if the underlying detection algorithms operate similarly. Listing 4 shows the FnO description of a CDC FnO function for detecting implicit updates in a dataset. It detects dataset updates based on the subject IRI of a member (`idlab-fn:iri`) and a set of watched properties (`idlab-fn:watchedProperty`), coalesced into as a structured string. By using a structured string, our FnO functions remain independent of the source data structure and allows us to use RML constructs such as `rml:template` for capturing the values of potentially complex watched properties through the used reference formulation (e.g., XPath or JSONPath) in the RML rules.

The function persists the state of the properties per member in the location specified by `idlab-fn:state`. Storing of state is implementation dependent, thus it can be stored as a plain file but also in a database for example. The state is a lookup table of the member’s properties with the member subject IRI as key. Through the state, we can keep track of any changes to members of a data source and detect added and deleted members by monitoring the presence of the member’s IRI.

Each function requires an `rr:template` to construct the IRI of the member being checked for changes. For explicitly communicated changes, each function only needs to check if a member IRI was already generated in the past: Since IRIs are guaranteed to be unique when changes are explicit, there is no need to materialize previously seen members. The implicit creation function `idlab-fn:implicitCreate` is identical to its explicit counterpart, (`idlab-fn:explicitCreate`), given that implicit member creation also occurs when new IRIs are detected. However, the `idlab-fn:implicitUpdate` and `idlab-fn:implicitDelete` functions are different. `idlab-fn:implicitUpdate` requires an additional list of properties (watched properties) for each member, to be compared with the previous version of the dataset and determine if anything changed. `idlab-fn:implicitDelete` keeps a state (e.g., a hashmap) of all seen member IRIs. This state is used at the end of every processing round to detect missing – thus, deleted – members from the dataset. It returns a list with all the deleted member IRIs.

### 5.2 KG Generation with RML and CDC FnO Functions

We integrate our FnO functions with RML through FNML in a RML Triples Map (Listing 4). For each type of change we have an independent RML Triples Map with a conditional Subject Map referencing one of these FnO functions. When a function returns the IRI that it received as input, the Triples Map is executed completely, thus materializing the member and its properties. If no IRI is returned by the FnO function, the Triples Map is not executed, which means that the member did not change in the dataset. A separate Triples Map can be used (Listing 5) to generate additional metadata in the form of an event log to describe which type of changes took place in the dataset. A possible ontology to describe these change types is the W3C Recommendation Activity Streams 2.0 (Snell & Prodromou, 2017), but other ontologies can be used in the mappings as well, that is, Change Detection Ontology (Randles & O’Sullivan, 2023; Randles et al., 2022), or PROV-O<sup>14</sup>.

```

<#FunctionMap> a fnml:FunctionMap
  fnml:functionValue [
    # CDC FnO function to use.
    rr:predicateObjectMap [
      rr:predicate fno:executes ;
      rr:objectMap [ rr:constant idlab-fn:implicitUpdate; ];
    ];
    # IRI template of a member
    rr:predicateObjectMap [
      rr:predicate idlab-fn:iri ;
      rr:objectMap [ rr:template "https://example.org/{id}"; ];
    ];
    # Properties to watch, can be one or multiple, might differ from IRI template.
    rr:predicateObjectMap [
      rr:predicate idlab-fn:watchedProperty ;
      rr:objectMap [ rr:template "prop1={prop1}&prop2={prop2}"; ];
    ];
    # Directory path to store the function's state
    rr:predicateObjectMap [
      rr:predicate idlab-fn:state ;
      rr:objectMap [ rr:constant "/path/to/state"; rr:dataType xsd:string; ];
    ];
  ];
.

```

**Listing 4.** RML usage of a CDC FnO function for detecting implicit updates in a dataset.

### 5.3 LDES Logical Target

We use LDES to publish a stream of dataset member changes events. LDES allows publishing only changed members as a stream which can be consumed by third-parties to replicate and synchronize with a dataset. We incrementally generate the detected changes (Section 5.2) and publish them as an LDES through an *Event Stream Target* in the RML mapping rules. An Event Stream Target is a subclass of an RML Logical Target with extra properties such as `rmlt:ldesBaseIRI`, `rmlt:ldes`, and `rmlt:generateImmutableIRI`. The Event Stream Target allows to indicate if unique and immutable IRIs are to be generated for each materialized member, as required by the LDES specification<sup>15</sup>. This way, consumers can uniquely identify individual changes in a dataset over time. We use an Event Stream Target for all member change types and one for generating an event log.

Listing 5 provides an example of our alignment between LDES and RML with CDC FnO functions. In this example, we use `as:Create`, `as:Update`, and `as>Delete` of W3C Activity Streams 2.0 (Snell & Prodromou, 2017) to semantically annotate the type of change for LDES consumers in the event log. Each change type has a dedicated RML Triples Map with an FnO function (Listing 5: Lines 65–82, 84–106, 108–121) which outputs the generated members to a LDES Event Stream Target (Listing 5: Lines 13–23). The W3C Activity Streams 2.0 event log is outputted to another LDES Event Stream Target (Listing 5: Lines 1–11).

### 5.4 Example Demonstrating Our Approach

We demonstrate our approach through an example dataset (Figure 3) which consists of an initial version (Table 3(a)) and an implicitly updated version (Table 3(b)). Figure 3 also visualizes this example in the pipeline. The following changes (Table 3(c)) are extracted between the initial data (Table 3(a)) and the newer version (Table 3(b)) through CDC:

- **ID 0:** Excluded from the changes as it is unchanged.
- **IDs 1,2:** Members are updated because the 'age' property changed, marked as 'update'.
- **ID 3:** A tombstone is generated since it is removed, marked as 'delete'.
- **ID 4:** New member, marked as 'create'

The resulting RDF quads of the initial (Listing 6) and updated (Listing 7) datasets consist of 3 named graphs<sup>16</sup>: `Create`, `Update`, and `Delete`, indicating the change type of each member. These named graphs are described in our examples (Listing 8) using W3C Activity Streams 2.0 (Snell & Prodromou, 2017)<sup>17</sup>. On each execution of our approach, the resulting RDF quads (Listing 7) are generated depending on the type of change compared to the previous execution.

## 6 Evaluation

In this section, we describe our methodology to evaluate our approach on several datasets using both synthetic and real-world data. First, we discuss our datasets and how they can be classified according to the different change signaling strategies identified in this work (Section 6.1). Then, we introduce our evaluation setup (Section 6.2).

```
# Logical Target for outputting W3C ActivityStreams 2.0 event log as an LDES
<LDESLogicalTargetAS> a rml:EventStreamTarget ;
  rml:target [ a void:Dataset ; void:dataDump <file:///eventlog.nq ; ] ;
  rml:serialization formats:N-Quads ;
  rml:ldes [ a ldes:EventStream ;
    ldes:timestampPath dct:created ; ldes:versionOfPath dct:isVersionOf ;
    tree:shape <https://example.org/shape/> ;
  ] ;
  rml:ldesBaseIRI <https://example.org/ldes/eventlog/> ;
  rml:ldesGenerateImmutableIRI "true"^^xsd:boolean .

# Logical Target for outputting data collection member changes as an LDES
<LDESLogicalTargetMember> a rml:EventStreamTarget ;
  rml:target [ a void:Dataset ; void:dataDump <file:///members.nq ; ] ;
  rml:serialization formats:N-Quads ;
  rml:ldes [ a ldes:EventStream ;
    ldes:timestampPath dct:created ; ldes:versionOfPath dct:isVersionOf ;
    tree:shape <https://example.org/shape/> ;
  ] ;
  rml:ldesBaseIRI <https://example.org/ldes/members/> ;
  rml:ldesGenerateImmutableIRI "true"^^xsd:boolean .

# Input CSV file as datasource
<DataSource> a rml:LogicalSource ;
  rml:source "data.csv" ;
  rml:referenceFormulation ql:CSV .

# Dedicated named graph for each change type
# W3C ActivityStreams 2.0 eventlog generation of created members
<TriplesMapASCreate> a rr:TriplesMap ;
  rml:logicalSource <DataSource> ;
  rr:subjectMap [
    rr:constant "http://blue-bike.be/event/create" ; rr:class as:Create ;
    rml:logicalTarget <LDESLogicalTargetAS> ;
  ] .

# W3C ActivityStreams 2.0 eventlog generation of updated members
<TriplesMapASUpdate> a rr:TriplesMap ;
```

**Listing 5.** RML mapping for generating an LDES from a CSV file as data source (data.csv). Each change type has a separate Triples Map with an FnO function as Subject Map. These functions return an IRI when changes are detected, thus triggering the full execution of the Triples Map. The generated RDF triples are written to the LDES Event Stream Target with the necessary LDES properties specified in the LDES Logical Target. W3C ActivityStreams 2.0 is used to indicate the type of change through an LDES-based event log. (Continued)

**Table 3.** Example Dataset with the Initial Dataset (Left), a Newer Version of the Dataset (Middle), and Extracted Changes by the Change Data Capture Functions Between the Initial Dataset and the Newer Version.

ID	Name	Age
<i>(a) Initial dataset produces 4 members with change type ‘create’.</i>		
0	The Machine	0
1	Harold Finch	44
2	John Reese	38
3	Agent Carter	36
<i>(b) Changed dataset has 2 updated, 1 unchanged, 1 created, and 1 deleted member(s).</i>		
0	The Machine	0
1	Harold Finch	46
2	John Reese	40
4	Root	35
<i>(c) Extracted changes by the Change Data Capture functions.</i>		
1	Harold Finch	<u>46</u>
2	John Reese	<u>40</u>
<u>3</u>	<u>Agent Carter</u>	<u>36</u>
<u>4</u>	<u>Root</u>	<u>35</u>

## 6.1 Methodology

We apply our approach on (i) a set of artificial test cases (Section 6.1.1) to verify if our approach is able to handle all change signaling strategies listed in Section 4.2; (ii) an extended version of the GTFS Madrid Benchmark (Section 6.1.2) to measure the scalability, performance and resource consumption of our approach for incrementally generating KGs; and (iii) 5 different real-world datasets (Section 6.1.3) to investigate its performance and resource impact on different types of datasets (Section 4.2),

```

rml:logicalSource <DataSource> ;
rr:subjectMap [
  rr:constant "http://blue-bike.be/event/update" ; rr:class as:Update ;
  rml:logicalTarget <LDESLogicalTargetAS> ;
] .

# W3C ActivityStreams 2.0 eventlog generation of deleted members
<TriplesMapASDelete> a rr:TriplesMap ;
rml:logicalSource <DataSource> ;
rr:subjectMap [
  rr:constant "http://blue-bike.be/event/delete" ; rr:class as:Delete ;
  rml:logicalTarget <LDESLogicalTargetAS> ;
] .

# Data collection member
<PersonName> a rr:PredicateObjectMap ;
rr:predicate schema:name ;
rr:objectMap [ rml:reference "name" ; rr:datatype xsd:string ; ] .

# Dedicated Triples Map per change type
# Detection of explicit member creations with FnO function,
# if the member IRI is not found in the state, a new created member is generated.
<TriplesMapObjectCreate> a rr:TriplesMap ;
rml:logicalSource <DataSource> ;
rr:subjectMap [
  fnml:functionValue [
    rr:predicateObjectMap [ rr:predicate fno:executes ; rr:object idlab-fn:explicitCreate ; ] ;
    rr:predicateObjectMap [
      rr:predicate idlab-fn:iri ;
      rr:objectMap [ rr:template "https://example.org/member/{id}" ]
    ] ;
  ] ;
  rr:graph <http://example.org/event/create> ; rr:class foaf:Person ;
  rml:logicalTarget <LDESLogicalTargetMember> ;
] ;
rr:predicateObjectMap <PersonName> .

# Detection of implicit member updates with FnO function
# Looks up the property 'name' of a member with the IRI of the member,
# if changed, an updated member is generated.
<TriplesMapObjectUpdate> a rr:TriplesMap ;
rml:logicalSource <DataSource> ;
rr:subjectMap [
  fnml:functionValue [
    rr:predicateObjectMap [ rr:predicate fno:executes ; rr:object idlab-fn:implicitUpdate ; ] ;
    rr:predicateObjectMap [
      rr:predicate idlab-fn:iri ;
      rr:objectMap [ rr:template "https://example.org/member/{id}" ]
    ] ;
    # Watch property 'name' of member for changes
    rr:predicateObjectMap [
      rr:predicate idlab-fn:watchedProperty ;
      rr:objectMap [ rr:template "name={name}" ]
    ] ;
  ] ;
  rr:graph <http://blue-bike.be/event/update> ; rr:class foaf:Person ;
  rml:logicalTarget <LDESLogicalTargetMember> ;
] ;
rr:predicateObjectMap <PersonName> .

# Detection of implicit member deletions with FnO function by IRI
# If member IRI is removed in the new version, a member as tombstone is generated.
<TriplesMapObjectDelete> a rr:TriplesMap ;
rml:logicalSource <DataSource> ;
rr:subjectMap [
  fnml:functionValue [
    rr:predicateObjectMap [ rr:predicate fno:executes ; rr:object idlab-fn:implicitDelete ; ] ;
    rr:predicateObjectMap [
      rr:predicate idlab-fn:iri ; rr:objectMap [ rr:template "https://example.org/member/{id}" ] ] ;
  ] ;
  rr:graph <http://blue-bike.be/event/delete> ; rr:class foaf:Person ;
  rml:logicalTarget <LDESLogicalTargetMember> ;
] .

```

Listing 5. Continued.

```

:Created {
  <http://ex.org/Mbr0#0> a foaf:Person .
  <http://ex.org/Mbr0#0> foaf:name "The Machine" .
  <http://ex.org/Mbr0#0> foaf:age "0"^^xsd:int .

  <http://ex.org/Mbr1#0> a foaf:Person .
  <http://ex.org/Mbr1#0> foaf:name "Harold Finch" .
  <http://ex.org/Mbr1#0> foaf:age "44"^^xsd:int .

  <http://ex.org/Mbr2#0> a foaf:Person .
  <http://ex.org/Mbr2#0> foaf:name "John Reese" .
  <http://ex.org/Mbr2#0> foaf:age "38"^^xsd:int .

  <http://ex.org/Mbr3#0> a foaf:Person .
  <http://ex.org/Mbr3#0> foaf:name "Agent Carter" .
  <http://ex.org/Mbr3#0> foaf:age "36"^^xsd:int .
}

```

**Listing 6.** The materialized KG in TriG of the initial dataset. All dataset members are materialized because this is the first time the dataset is processed. Thus, they are part of the Create named graph. The subject IRIs are versioned as required by LDES to publish versioned members.

```

:Created {
  <http://ex.org/Mbr4#0> a foaf:Person .
  <http://ex.org/Mbr4#0> foaf:name "Root" .
  <http://ex.org/Mbr4#0> foaf:age "35"^^xsd:int .
}

:Updated {
  <http://ex.org/Mbr1#1> a foaf:Person .
  <http://ex.org/Mbr1#1> foaf:name "Harold Finch" .
  <http://ex.org/Mbr1#1> foaf:age "46"^^xsd:int .

  <http://ex.org/Mbr2#1> a foaf:Person .
  <http://ex.org/Mbr2#1> foaf:name "John Reese" .
  <http://ex.org/Mbr2#1> foaf:age "40"^^xsd:int .
}

:Deleted {
  <http://ex.org/Mbr3#1> a foaf:Person .
}

```

**Listing 7.** The materialized KG in TriG of the changed dataset. A member is deleted (ID 3), a member’s age property is updated (ID 1 and 2), and a new member is created (ID 4). Unchanged members (ID 0) are not materialized. The subject IRIs are versioned as required by LDES to publish versioned members.

```

# Named graph for created members of data collection
:Created a as:Create ;
as:actor <http://ex.org/data-collection> .
# Named graph for updated members of data collection
:Updated a as:Update ;
as:actor <http://ex.org/data-collection> .
# Named graph for deleted members of data collection
:Deleted a as>Delete ;
as:actor <http://ex.org/data-collection> .

```

**Listing 8.** The different named graphs are described using the W3C Activity Streams 2.0 ontology as as:Create, as:Update, and as>Delete.

**6.1.1 Functionality.** Through a set of test cases (Table 4), we evaluate 3 change signaling strategy dimensions (history availability, change communication, and change type (Section 4.2)) to determine if our approach: (i) Can handle all dimensions, and (ii) is feasible to implement. We combined the 3 types of history availability (latest state, latest changes, and full history) with the 2 types of change signaling (explicit and implicit), and 3 change types (create, update, delete) resulting into 18 test-cases. We also included a test case where no change is applied to verify if an unchanged dataset is handled correctly by our approach, which yields additionally 6 more test-cases (24 in total). Since 2 combinations are not possible (Section 4: Implicit deletion signaling for full history and latest changes), we removed these from the test cases, bringing the number of test cases to 22.

We validate that all combinations of these dimensions are possible and use these test cases to verify that our implementation covers all possible scenarios. The test cases are publicly available on GitHub<sup>18</sup>.

**6.1.2 GTFS Madrid Benchmark Extension.** We extended the GTFS Madrid Benchmark (Chaves-Fraga et al., 2020) data generator<sup>19</sup> to allow generating different versions of the benchmark data by supporting creations, updates, and deletions based on the GTFS data model. We implemented creations, updates, and deletions in the GTFS Madrid Benchmark similar to how real-world GTFS datasets are changed at the Belgian public transport agencies De Lijn and NMBS. This way, we keep the characteristics of GTFS Madrid Benchmark which aims at using real-world data from the metro in Madrid. Creations are applied by adding GTFS routes and their associated GTFS trips, stops, stoptimes, and service entities. For example: 25% creations will provide 25 % additional new routes with respect to the original amount of routes. Updates are performed by modifying the GTFS services. For example: 50% updates will modify 50% of the GTFS service entries in the GTFS calendar. Deletes are applied by removing GTFS routes and their associated trips and services. Example: 10% deletes will remove 10% of the routes in the original data, together with the associated data. We use a random number generator to randomly select GTFS routes and services.

We extend the GTFS Madrid Benchmark with 4 additional configuration parameters:

- *seed*: The random seed value used for configuring the random number generator.
- *additions*: The percentage defining how many creations must be added to the generated data.
- *modifications*: The percentage defining how many updates must be performed on the generated data.
- *deletions*: The percentage defining how many deletions must be applied on the generated data.

**Table 4.** 22 Test Cases for Detecting Changes in Datasets.

# Test Cases	History Availability	Signaling	Change Type	Purpose
8	Latest state	Explicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Tombstone deleted entity
		Implicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Tombstone deleted entity
7	Latest changes	Explicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Tombstone deleted entity
		Implicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Tombstone deleted entity
7	Full history	Explicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Tombstone deleted entity
		Implicit	No change	Empty output
			Create	Added new entity
			Update	Existing entity changed
			Delete	Existing entity changed

Our approach is feasible to implement and covers detecting all possible change types and change signaling strategies.

We aim on analyzing the impact of our approach on datasets of varying size and change characteristics, by measuring execution time and resource usage (storage, CPU time, and RAM usage) to detect and materialize changes into a KG (CHANGE) or materialize the complete KG (ALL). In particular, we analyze the overhead of our approach for detecting changes and the reduction in execution time and resource usage achieved by only materializing the changes instead of the complete KG. The reduction of our approach might be affected by the amount of changes, type of changes, and the dataset size. Therefore, we use data size scales (1, 10, and 100) with a fixed change percentage of 50%, equally divided among the different change types (16.67% creations, 16.67% updates, and 16.67% deletions), to obtain reference measurements of increasing dataset size scales. This way, we avoid that other dimensions for example, change percentage and change types impacts our analysis of results from the data size dimension. We divide the changes equally among change types and use scale 100 to analyze the change percentage dimension (0%, 25%, 50%, 75%, and 100%). We also experiment with the type of change by using a fixed change percentage of 50% and scale 100 for either creations, updates, or deletions to analyze the impact of each change type on execution time and resources. For each experiment we use 10 new versions of the GTFS dataset which we apply as updates over an initial base version. We select 50% as fixed change percentage to avoid outliers from 0% or 100% and scale 100 for analyzing the change type dimension.

**6.1.3 Real-World Datasets.** We apply our approach on 5 types of real-world datasets: (i) BlueBike & JCDecaux bike-sharing data; (ii) timetables in GTFS format from the Belgian public transport agencies NMBS and De Lijn; (iii) OpenStreetMap (OSM) geographical data; (iv) dynamic message boards from the Flemish traffic controller centre Vlaams Verkeerscentrum (VVC); and (v) meteorological sensor data from the Belgian meteorological institute Koninklijk Meteorologisch Instituut van België (KMI). These datasets stand as representative examples for all the identified change signaling strategies (Table 5) regarding change communication, change types, and history availability. For each dataset, we collected released versions during 24 hours, with the exception of De Lijn and NMBS, since they only provide a new version per day. In this case, we collect new versions during a week. We use a timeframe instead of the number of versions to have different frequencies when new versions of the datasets are published, for example, VVC is changed more frequently (every 2–3 secs) compared to KMI (every 10 mins). Next we describe in detail each of the analyzed datasets, which are also summarized in Table 5.

**BlueBike & JCDecaux.** These datasets provide information related to bike-sharing services including data about stations and currently available bikes. They are publicly accessible via HTTP APIs<sup>20</sup>. Both datasets communicate creations explicitly through unique identifiers, and updates implicitly by modifying the number of available bikes and a last updated

**Table 5.** Real-world Datasets Addressed in our Evaluation.

Dataset	Number of Collected Versions	Collection Frequency	Change Communication			History Availability
			Create	Update	Delete	
BlueBike	1440	1 min	Explicit	Implicit	Implicit	Latest state
JCDecaux	1440	1 min	Explicit	Implicit	Implicit	Latest state
NMBS	7	1 day	Explicit	Implicit	Implicit	Latest state
De Lijn	7	1 day	Explicit	Implicit	Implicit	Latest state
OSM	1440	1 min	Explicit	Explicit	Explicit	Latest changes
VVC	28760	3 secs	Explicit	NA	Explicit	Latest state
KMI	144	10 mins	Explicit	NA	NA	Full history

Each dimension of dataset types is covered: Change signaling, change type and history availability. Some datasets only provide certain change types, inapplicable change types are marked as NA. OSM: OpenStreetMap; VVC: Vlaams Verkeerscentrum; KMI: Koninklijk Meteorologisch Instituut van België; NA: not applicable.

timestamp property. However, for some stations, they do not provide this timestamp. Thus, we use implicit change detection by *watching* the available bikes at each station. Deletes are implicitly communicated by removing stations from the dataset. Regarding history availability, these datasets follow the *latest state* strategy, overwriting the whole dataset with the current state of each member on every new version. We use a GBFS-based vocabulary<sup>21</sup> to transform this bike-sharing data into RDF and retrieve the data every minute to check for changes.

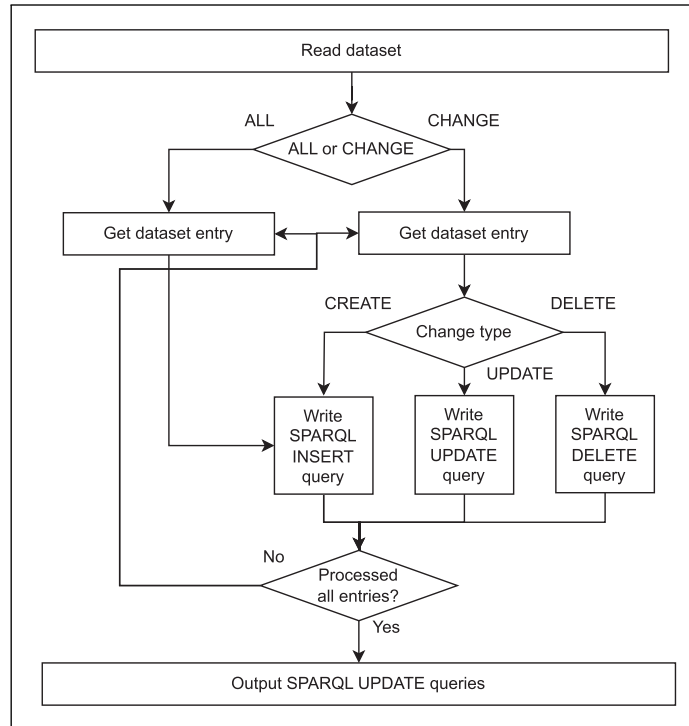
**NMBS & De Lijn.** These datasets contain public transport timetables in GTFS format, which are updated daily. They are publicly accessible as data dumps<sup>22</sup>. In GTFS timetables, creations are explicitly communicated by adding new entries with a new ID. Updates are implicit, silently changing properties of members. They do not provide unique identifiers nor timestamps for identifying changes. Deletes are also implicitly communicated by silently removing members. GTFS also follows the *latest state* approach for history availability. We use the Linked GTFS ontology<sup>23</sup> to transform the GTFS timetables into RDF, based on the RML mappings of the GTFS-Madrid-Bench (Chaves-Fraga et al., 2020).

**OpenStreetMap.** This is a crowdsourced geographical dataset which provides a feed<sup>24</sup> of changed data on a minutely basis. These updates publish OSM’s latest changes with explicit change communication for all types of change: Creations, updates, and deletions are all identified with a unique ID. OSM data is mapped using the Routable Tiles specification (Colpaert et al., 2019), which provides an ontology for OSM’s Nodes and Ways. OSM follows the *latest changes* approach regarding history availability in their replication feeds (although historical full latest-state data dumps are also available). Since OSM already explicitly provides unique IDs for each change, no additional processing is needed by our approach. Thus, semantically publishing the dataset changes can be achieved with regular RML mappings without applying the CDC functions on the dataset. We include this dataset in this work to show the wide range of dataset types regarding change signaling, history availability, and change types, supported by our approach.

**Vlaams Verkeerscentrum.** This dataset contains live information from traffic boards that include data about traffic jams, road works, accidents, and other incidents on the road across Flanders. These dynamic message boards are changed every 2–3 seconds and published as Open Data. We include this dataset due to its high update frequency, compared to the other datasets: The RDF generation process must keep up with this high frequency to avoid high latency. We use the DATEX II v3 VMS ontology<sup>25</sup> for the content of the dynamic message board messages and collect the data every 3 seconds. This dataset follows the *latest state* approach for history availability and contains a modified timestamp for each published message, which allows for explicit change communication. Messages are not modified in this dataset, but republished as new message instances instead. Therefore, there are no updates, only creations and deletions. Deletions are marked in the data as an out-of-service board. If a board is marked as out-of-service, we consider it deleted from the dataset. If the board is put in service again, it is considered as a creation.

**Koninklijk Meteorologisch Instituut van België.** This dataset is published by the Belgian meteorological institute which allows access to the measurement history of their weather sensors. The dataset contains the *full history* of sensor measurements labelled with unique IDs. However, it does not explicitly communicate deletions. Since detecting implicit deletions in a full history is not directly possible (Section 4), and each measurement has a unique ID, only creations are possible in this dataset. We use the W3C Semantic Sensor Network ontology<sup>26</sup> to transform the sensors’ measurements into RDF and collect the data every 10 minutes.

**6.1.4 Ingestion.** In practice, the outcome of a KG generation process is typically a set of RDF triples/quads that are subsequently ingested into an RDF triplestore, so that they can be queried and used in applications. In a traditional full re-materialization KG generation scenario (ALL), the ingestion process usually entails a complete deletion of the KG



**Figure 4.** `incrm12sparql` transforms the generated resource description framework (RDF) into SPARQL UPDATE queries for ingestion by a triplestore. ALL datasets always generate a SPARQL INSERT query for all data entries. CHANGE datasets generate a SPARQL INSERT/UPDATE/DELETE based on the type of change and only for the changed dataset entries.

in the triplestore (if any), followed by an insertion of all the newly generated RDF quads. However, the materialized members generated by our approach after the CDC process (CHANGE), do not constitute on their own a complete and standalone version of the KG (except for the initial generation), and require an additional interpretation step to determine the proper operations to be performed over the triplestore that will bring the KG to its latest state. We implemented a proof-of-concept library<sup>27</sup> that performs this interpretation step (Figure 4) and produces the corresponding SPARQL UPDATE queries for each materialized member, based on their change semantics (i.e., type of change). By using a triplestore as the same system to make the standalone KG accessible for both the ALL and CHANGE strategy, we make sure we are not monitoring side-effects when using reconciliation systems that are different for the ALL and CHANGE strategy.

To guarantee a fair comparison of our approach (CHANGE) with a traditional approach (ALL), we initially observe the number of RDF quads generated by each approach. Our premise is that less RDF quads translate into faster ingestion and thus more efficiently updated KGs. To verify this, we also measured the time that takes to ingest via SPARQL UPDATE (i) the fully materialized KGs (ALL); and (ii) to produce and execute the SPARQL UPDATE queries corresponding to the materialized members (CHANGE), for every real world dataset considered in this work (Section 6.1.3). We perform this measurement over a Virtuoso triplestore v7.2.14<sup>28</sup>.

## 6.2 Evaluation Setup

We evaluate our approach on the described datasets and a modified version of the GTFS Madrid Benchmark by measuring the following metrics: (i) Execution time to detect changes and materialize the complete or changed parts of a dataset into RDF; (ii) CPU usage to determine how CPU intensive is our generation approach; (iii) peak RAM memory usage to investigate the memory overhead of our approach; and (iv) storage usage of the generated KG on each new version to verify the impact of our approach for reducing storage of historical data. We execute this evaluation with and without our approach using the RMLMapper v6.3.0<sup>29</sup> to investigate how our approach impacts the measured metrics during RDF KG generation.

For each experiment, we materialize a KG from a given dataset and its corresponding updated versions as RDF quads. We manually verify if the output is complete and correct for each experiment. We compare 2 KG generation-and-version strategies:

- *ALL* is the traditional strategy to generate versioned RDF KGs, where a set of RDF quads are initially produced and completely re-materialized when a new version of the data source(s) is/are available.
- *CHANGE* corresponds to our approach to incrementally generated RDF KGs, where we initially materialize a complete version of the KG, but only materialize into RDF the actual changed members upon updates of the data source(s), not the complete KG. We also opt to use the LDES Logical Target, to add additional metadata that semantically describes how the KG changes to help consumers reconcile the changes into their own KG, for example, via SPARQL UPDATE queries over a triplestore.

In this work, we focus on measuring the generation step of both strategies to investigate the feasibility and impact of our approach over materializing the complete KG on each dataset version. Since our focus is the generation, we only perform a comparative experiment to have an indication of how our approach affects ingestion into triplestores using SPARQL UPDATE queries, given it is the standardized way to update KGs in any triplestore. Moreover, most triplestores offer custom approaches to update quads besides SPARQL UPDATE. Each of these ingestion approaches have their own benefits and drawbacks depending on the triplestore and vendor. Thus, we focus on comparing the amount of generated quads for each dataset version with and without our approach since we want to investigate how our approach affects KG generation itself.

The experiments are executed on an Ubuntu 22.04.1 LTS machine (Linux 5.15.0-83-generic, x86\_64) with an Intel(R) Xeon(R) CPU E5-2650 v2 @ 2.60Ghz, 48GB RAM memory and 2GB swap memory. Java JVM heap space is set to 90% of the available RAM memory. All experiments are executed 5 times, from which we report the median measurements of the execution. All resources and instructions to reproduce the experiments are available on Zenodo<sup>30</sup>.

## 7 Results

In this section, we present the results obtained during the execution of our evaluation with and without our approach. We first perform a functional evaluation through a set of test cases (Section 7.1). Then, we evaluate the performance in terms of execution time and computing resources of our approach on an extended version of the GTFS Madrid Benchmark (Section 7.2), and real-world datasets (Section 7.3). We then reconcile the results of both the ALL and CHANGE scenario (i.e., make the latest version of a KG accessible after a set of changes) and report the results of our comparative experiment after triplestore ingestion (Section 7.4). In Section 7.5, we discuss the impact of our approach on the measured metrics.

### 7.1 Functionality

We integrated our set of 22 test cases (Table 4 of Section 6.1.1) as unit tests in the RMLMapper v6.3.0<sup>31</sup> to validate if all possible change types and change signaling strategies are covered by our approach and are feasible to implement. Through these test cases, we confirm the feasibility of our approach and its coverage of all possible dimensions regarding change signaling and history availability.

### 7.2 GTFS Madrid Benchmark

In this subsection, we discuss the results obtained for the 2 generation-and-version strategies (Section 6.2) ALL and CHANGE per scenario: (i) Scaling data size, (ii) scaling amount of changes, and (iii) different change types. We compare the results of the initial version of the KG (Table 6), and the updates applied after the initial version (Table 7). We execute each experiment by starting with the base dataset to generate the initial KG and apply the corresponding dataset updates on top of the base dataset. Each experiment is executed 5 times from which the median value is taken for each metric of each dataset version. We report the average of these median metrics, for example, execution time, CPU time, and memory usage across multiple dataset updates in Tables 6 and 7. Storage usage is reported as the sum of KG sizes generated from the base dataset and all its updates (Table 8).

The results show that our approach **significantly reduces the storage requirements** for storing multiple versions of a KG and **reduces execution time and CPU time, without impacting memory usage** (Table 7). For the **initial KG generation** (Table 6), **the overhead of our approach causes an increase in execution time (2.22x in average), storage (1.57x in average), and CPU time (1.87x in average), while memory usage is mostly unaffected (1.15x in average)**. This is the result of constructing the initial KG but also initializing the state to keep track of dataset members, as the state is completely empty during the initial KG construction. Only the tracking of dataset members by our approach has a slight increase (+2.30% compared to average) in memory for CHANGE. Our approach remains unaffected by the amount of changes or type of change, being solely impacted by the dataset size.

**Table 6.** Initial Execution Results for All GTFS Madrid Benchmark for Strategies ALL (No Change Detection) and CHANGE (with Change Detection).

Scenario	Execution Time (s)			CPU Time (s)			Peak Memory (GB)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
GTFS <sub>SCALE</sub> 1	12.29	16.18	0.76	26.80	36.31	0.74	4.25	4.44	0.96
GTFS <sub>SCALE</sub> 10	82.15	124.84	0.66	129.45	206.05	0.63	7.57	14.00	0.54
GTFS <sub>SCALE</sub> 100	922.43	1943.38	0.47	2260.62	4626.85	0.49	46.30	51.78	0.89
GTFS <sub>CHANGE</sub> 0%	919.25	2736.11	0.34	2259.27	5002.98	0.45	45.29	52.30	0.87
GTFS <sub>CHANGE</sub> 25%	914.98	1575.10	0.58	2303.53	4521.50	0.51	46.38	51.64	0.90
GTFS <sub>CHANGE</sub> 50%	922.43	1943.38	0.47	2260.62	4626.85	0.49	46.30	51.78	0.89
GTFS <sub>CHANGE</sub> 75%	924.12	2357.18	0.39	2224.46	4603.18	0.48	46.77	52.09	0.90
GTFS <sub>CHANGE</sub> 100%	912.37	1420.53	0.64	2336.65	4431.36	0.53	46.32	51.70	0.90
GTFS <sub>TYPE</sub> CREATE	990.30	1397.49	0.71	2448.92	4414.65	0.55	48.07	51.73	0.93
GTFS <sub>TYPE</sub> UPDATE	924.16	1910.59	0.48	2239.66	4507.61	0.50	47.13	51.64	0.91
GTFS <sub>TYPE</sub> DELETE	929.74	1671.17	0.56	2302.00	4702.38	0.49	46.58	52.30	0.89

We scale one dimension for each experiment to analyze its impact. GTFS<sub>SCALE</sub> and GTFS<sub>TYPE</sub> use a fixed change percentage (50%). GTFS<sub>CHANGE</sub> and GTFS<sub>TYPE</sub> use a fixed data size (scale 100) which results into the same storage usage for the initial generation. GTFS<sub>SCALE</sub> and GTFS<sub>CHANGE</sub> spread the changes equally among the change types. lower is better.

**Table 7.** Execution Results for All GTFS Madrid Benchmark for Strategies ALL (No Change Detection) and CHANGE (with Change Detection).

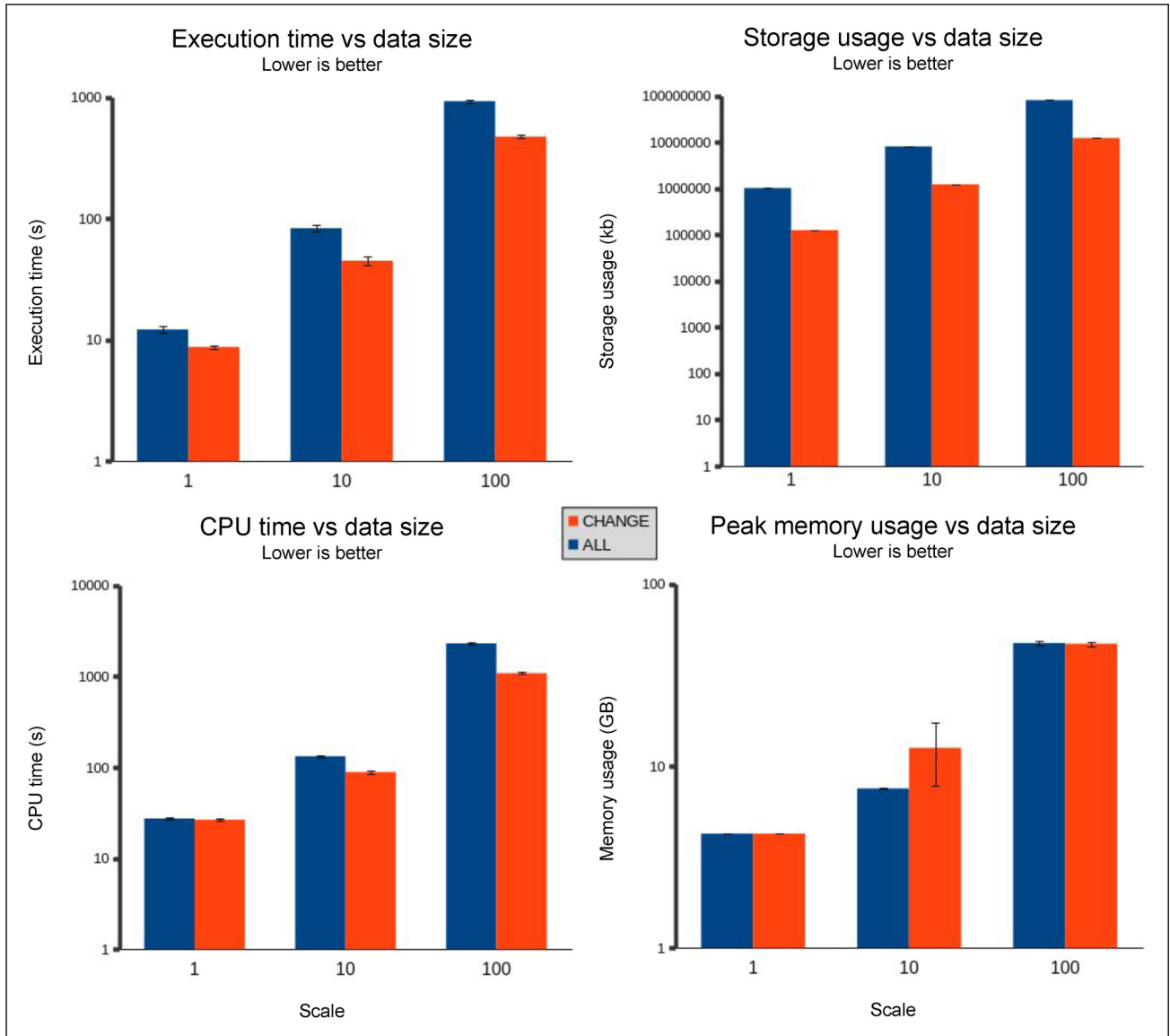
Scenario	Execution Time (s)			CPU Time (s)			Peak Memory (GB)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
GTFS <sub>SCALE</sub> 1	12.20	8.74	1.40	27.40	26.58	1.03	4.25	4.25	1.00
GTFS <sub>SCALE</sub> 10	83.41	44.93	1.87	132.63	88.73	1.49	7.53	12.58	0.60
GTFS <sub>SCALE</sub> 100	928.93	473.55	1.96	2306.50	1083.96	2.13	47.41	47.01	1.00
GTFS <sub>CHANGE</sub> 0%	927.82	470.73	1.97	2292.39	1072.47	2.14	47.19	47.33	1.00
GTFS <sub>CHANGE</sub> 25%	912.50	481.62	1.89	2303.02	1090.74	2.11	46.86	47.29	0.99
GTFS <sub>CHANGE</sub> 50%	928.93	473.55	1.96	2306.50	1083.96	2.13	47.41	47.01	1.00
GTFS <sub>CHANGE</sub> 75%	929.53	466.87	1.99	2292.81	1056.88	2.17	47.21	47.19	1.00
GTFS <sub>CHANGE</sub> 100%	914.68	460.30	1.98	2301.77	1063.46	2.16	47.05	47.35	0.99
GTFS <sub>TYPE</sub> CREATE	1021.04	454.06	2.25	2386.46	1063.90	2.24	48.55	46.61	1.04
GTFS <sub>TYPE</sub> UPDATE	953.32	479.87	1.99	2298.29	967.63	2.38	46.77	47.32	0.99
GTFS <sub>TYPE</sub> DELETE	922.04	464.78	1.98	2266.11	1056.56	2.14	46.69	47.33	0.99

Only results of dataset updates are included, initial execution is not included. We scale one dimension for each experiment to analyze its impact. GTFS<sub>SCALE</sub> and GTFS<sub>TYPE</sub> use a fixed change percentage (50%). GTFS<sub>CHANGE</sub> and GTFS<sub>TYPE</sub> use a fixed data size (scale 100) which results into the same storage usage for the initial generation. GTFS<sub>SCALE</sub> and GTFS<sub>CHANGE</sub> spread the changes equally among the change types. lower is better.

**Table 8.** Storage Usage for Initial and All Updates per Strategy of the GTFS Madrid Benchmark.

Scenario	Initial Storage Usage (kb)			Total Storage Usage (kb)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
GTFS <sub>SCALE</sub> 1	92,537.13	116,662.92	0.79	1023,192.74	125,962.43	8.12
GTFS <sub>SCALE</sub> 10	732,959.66	1168,374.14	0.63	8100,043.47	1234,254.86	6.56
GTFS <sub>SCALE</sub> 100	7347,846.47	11,701,971.43	0.62	81,178,874.08	12,323,348.17	6.59
GTFS <sub>CHANGE</sub> 0%	7347,846.47	11,701,971.43	0.62	80,826,311.20	11,703,851.89	6.91
GTFS <sub>CHANGE</sub> 25%	7347,846.47	11,701,971.43	0.62	81,005,737.97	12,019,442.49	6.74
GTFS <sub>CHANGE</sub> 50%	7347,846.47	11,701,971.43	0.62	81,178,874.08	12,323,348.17	6.59
GTFS <sub>CHANGE</sub> 75%	7347,846.47	11,701,971.43	0.62	81,355,130.62	12,631,882.64	6.44
GTFS <sub>CHANGE</sub> 100%	7347,846.47	11,701,971.43	0.62	81,530,274.97	12,934,632.82	6.30
GTFS <sub>TYPE</sub> CREATE	7347,846.47	11,701,971.43	0.62	82,031,098.47	13,530,545.51	6.06
GTFS <sub>TYPE</sub> UPDATE	7347,846.47	11,701,971.43	0.62	80,826,311.20	11,712,313.42	6.90
GTFS <sub>TYPE</sub> DELETE	7347,846.47	11,701,971.43	0.62	80,689,654.29	11,727,476.53	6.88

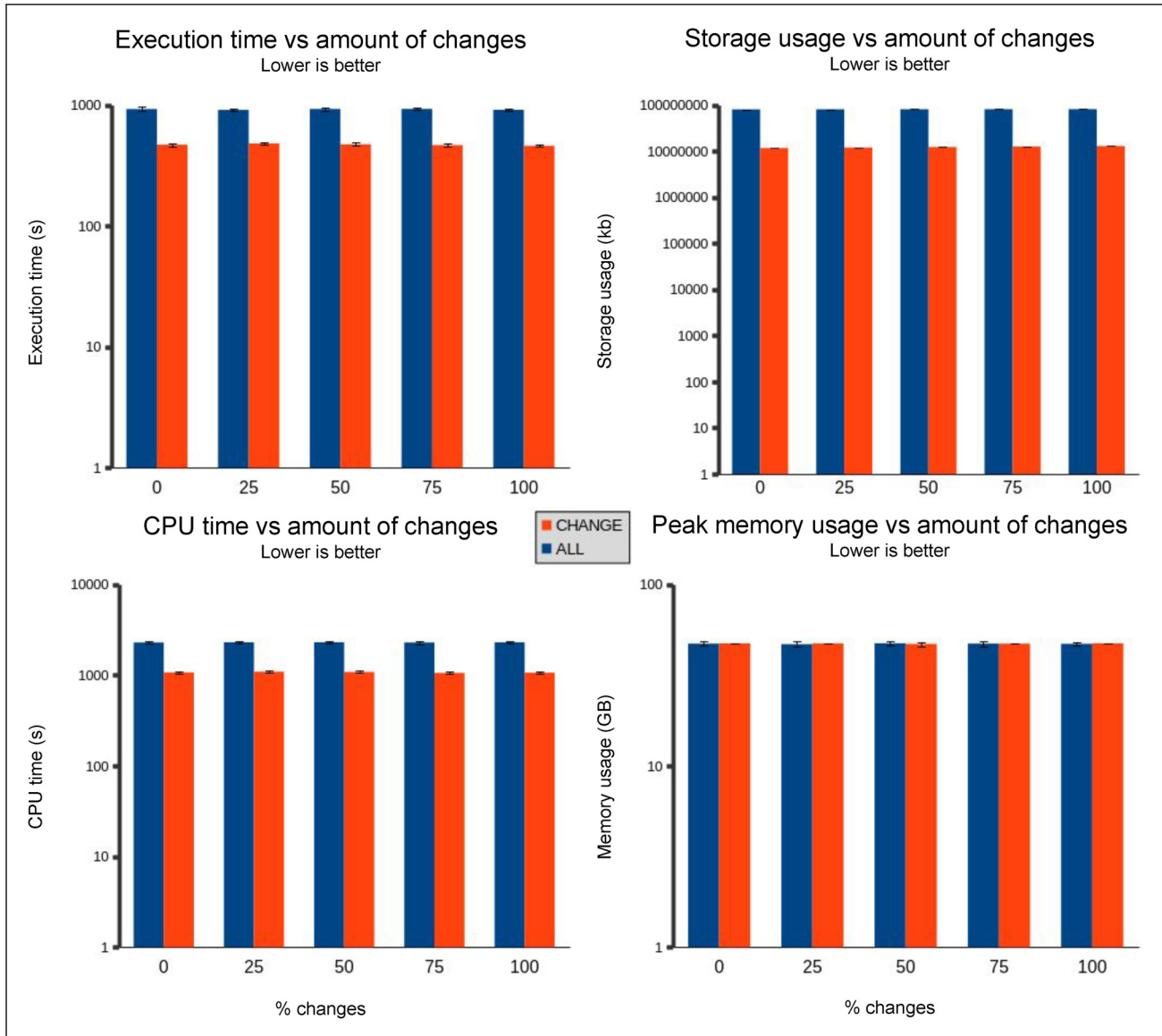
We scale one dimension for each experiment to analyze its impact. GTFS<sub>SCALE</sub> and GTFS<sub>TYPE</sub> use a fixed change percentage (50%). GTFS<sub>CHANGE</sub> and GTFS<sub>TYPE</sub> use a fixed data size (scale 100) which results into the same storage usage for the initial generation. GTFS<sub>SCALE</sub> and GTFS<sub>CHANGE</sub> spread the changes equally among the change types. Total storage usage is the sum of the base dataset and all applied updates upon it. lower is better.



**Figure 5.** GTFS Madrid Benchmark results for different data scales with a fixed amount of changes (50%). Only results of dataset updates are included, initial execution is not included.

*Scaling data size.* Our approach **reduces the execution time and resource usage (storage, CPU, and memory) of the different GTFS Madrid Benchmark scales (1, 10, 100)** by only materializing the actual changed members of the KG (Figure 5). However, **the initial generation of the KG from the base dataset has a longer execution time and higher resource consumption since all dataset members must be checked to initialize the internal state for tracking the members.** Moreover, the storage usage increases because of the metadata we generate to indicate that all members were created. When solely considering the dataset changes, the overhead of detecting changes and LDES event stream metadata is mostly noticeable with  $\text{GTFS}_{\text{SCALE}} 1$  (execution time reduced by 30%, no CPU time reduction). For larger GTFS scales the impact of the overhead is less noticeable, since they achieve an overall reduction in execution time (50% faster) and resource usage (CPU time is reduced up to 53%, and up to 8.10x less storage usage). Only for  $\text{GTFS}_{\text{SCALE}} 10$ , the memory usage increases since the Java VM does not perform garbage collection, given that the Java heap space still has enough free space. This is not the case for  $\text{GTFS}_{\text{SCALE}} 100$  where memory usage is similar again compared to no change detection.

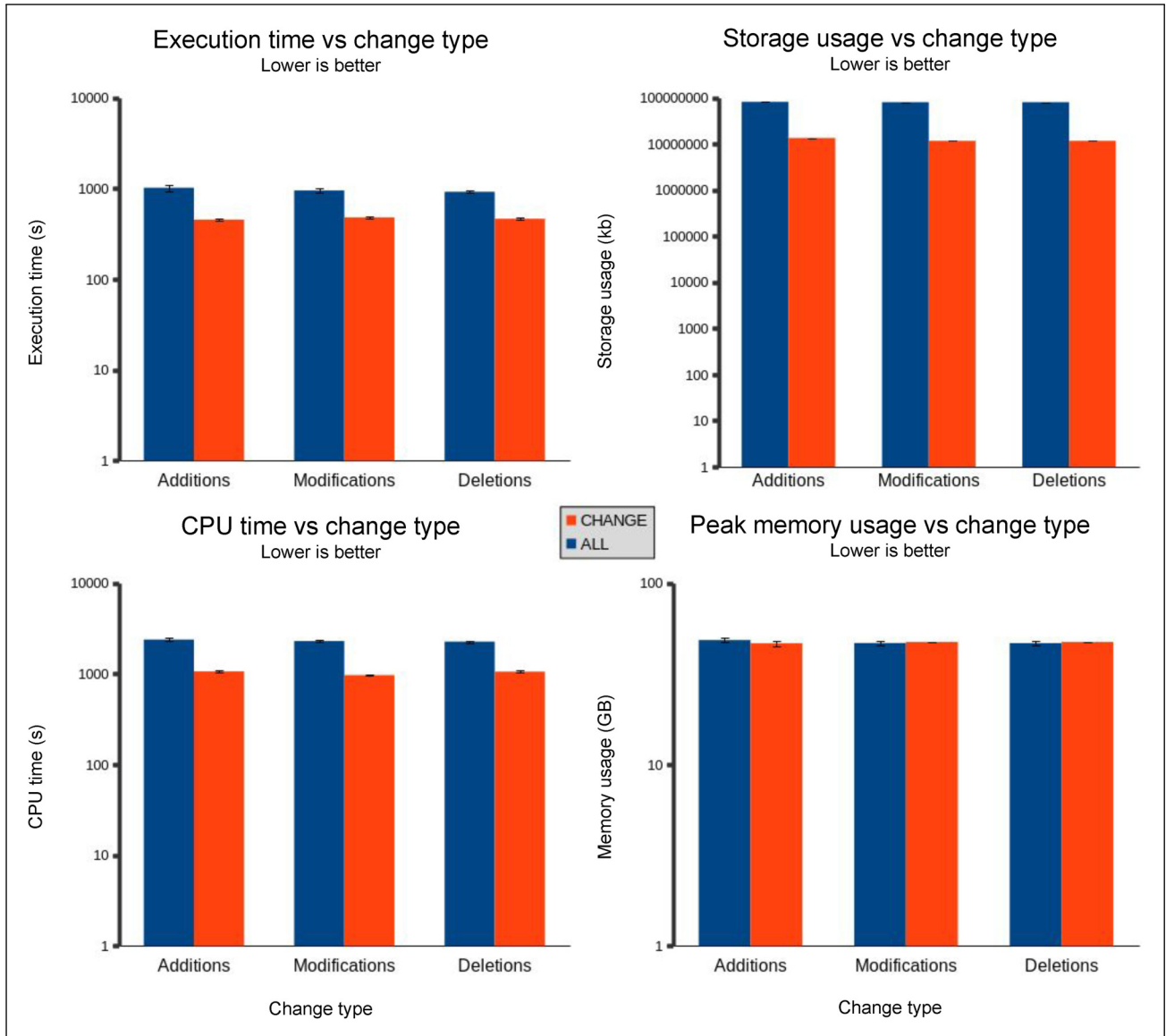
*Scaling amount of changes.* We observe an **increase in storage usage when the amount of changes increases**, while it has no impact on the resource consumption (Figure 6). This is due to each dataset member being evaluated regardless of the amount of changes. Similar to scaling the data size, the initial KG generation from the base dataset introduces overhead



**Figure 6.** GTFS Madrid Benchmark results for different amount of changes with a fixed data size (scale 100). Only results of dataset updates are included, initial execution is not included.

to initialize the state for tracking changes. When only considering the dataset updates, if no changes are found between 2 versions (0%), only the overhead of generating LDES event stream metadata is affecting storage usage, resulting in the lowest storage usage (11.7GB in total). Scaling up the amount of changes, increases the storage usage (12.93GB in total).

*Type of changes.* **The type of change in GTFS Madrid Benchmark mostly affects storage usage** because our algorithms evaluate each dataset member, for each change type (Figure 7). Therefore, CPU and memory usage is unaffected by the type of change. The same overhead of initially generating the KG applies as mentioned in the previous sections. Creations cause a higher storage usage because they result in a higher number of materialized RDF quads. In the case of GTFS-Madrid-Benchmark, properties of new GTFS routes and associated data such as GTFS trips, shapes, etc., must be generated as new triples. Deletions only keep a tombstone of the member, for example, a deleted GTFS route or trip. Updates also trigger the generation of new RDF quads, for example, the GTFS trip's service dates are modified in the GTFS Madrid Benchmark, which causes fewer changes among the GTFS datasets compared to creates or deletes. GTFS trip's service dates do not affect other information about a GTFS trip such as its route, shapes, or stops. Thus, updates have a lower storage usage for GTFS Madrid Benchmark since fewer changes happened in the GTFS dataset.



**Figure 7.** GTFS Madrid Benchmark results for different change types with a fixed data size (scale 100) and amount of changes (50%). Only results of dataset updates are included, initial execution is not included.

### 7.3 Real-World Datasets

In this subsection, we discuss the results obtained from applying our approach over the set of real-world datasets described in Section 6.1.3, with respect to storage usage, execution time, CPU time, and memory usage. For each dataset, we measured these metrics following strategies **ALL** (no change detection) and **CHANGE** (using our change detection approach). Table 9 reports the initial execution results, while Table 10 shows the impact of our approach on multiple dataset updates on the measured metrics. We execute each experiment by starting with the base dataset to generate the initial KG and apply the corresponding dataset updates on top of it. Each experiment is executed 5 times from which the median value is taken for each metric of each dataset update. We report the average of these median metrics for example, execution time, CPU time, and memory usage across multiple dataset updates in Tables 9 and 10. Storage usage is reported as the sum of KG sizes generated from the base dataset and all its updates (Table 11). OSM signals explicitly all changes, thus there is no  $OSM_{ALL}$ .

We observe that our approach **reduces execution time and CPU time without impacting memory usage** (Table 10), while also **reducing the storage needed to store multiple versions of a KG**. For the initial KG generation (Table 9), the

**Table 9.** Initial Execution Results for All Real-world Use-cases for Strategies ALL (No Change Detection) and CHANGE (with Change Detection).

Scenario	Execution Time (s)			CPU Time (s)			Peak Memory (GB)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
BlueBike	2.59	2.80	0.93	7.87	9.12	0.86	2.39	2.42	0.99
JCDecaux	3.73	4.65	0.80	16.82	21.63	0.77	2.13	2.33	0.91
NMBS	151.02	177.11	0.85	232.11	292.61	0.79	17.41	23.58	0.74
De Lijn	1156.67	1542.23	0.75	3101.31	6142.43	0.50	49.98	50.09	1.00
KMI	534.93	637.28	0.84	850.56	1147.12	0.74	36.11	38.86	0.93
VVC	19.94	17.94	1.11	44.67	43.09	1.04	7.61	6.46	1.18
OSM	NA	3.10	NA	NA	11.22	NA	NA	2.67	NA

Lower is better. OSM: OpenStreetMap; VVC: Vlaams Verkeerscentrum; KMI: Koninklijk Meteorologisch Instituut van België.

**Table 10.** Execution Results for All Real-life Datasets for Strategies ALL (No Change Detection) and CHANGE (with Change Detection).

Scenario	Execution Time (s)			CPU Time (s)			Peak Memory (GB)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
BlueBike	2.56	2.63	0.97	7.68	8.21	0.94	2.39	2.40	1.00
JCDecaux	3.79	3.78	1.00	16.74	19.77	0.85	2.13	2.16	0.99
NMBS	146.98	86.17	1.71	227.29	144.21	1.58	20.08	27.88	0.72
De Lijn	1072.17	385.53	2.78	2523.12	1268.39	1.99	48.88	48.61	1.00
KMI	527.95	119.63	4.41	845.98	184.25	4.59	34.03	31.14	1.09
VVC	17.47	6.36	2.75	41.10	22.73	1.94	7.35	4.87	1.51
OSM	NA	3.56	NA	NA	16.64	NA	NA	2.92	NA

Only results of dataset updates are included, initial execution is not included. lower is better. OSM: OpenStreetMap; VVC: Vlaams Verkeerscentrum; KMI: Koninklijk Meteorologisch Instituut van België.

**Table 11.** Storage Usage for Initial and All Updates per Strategy of the Real-World Datasets.

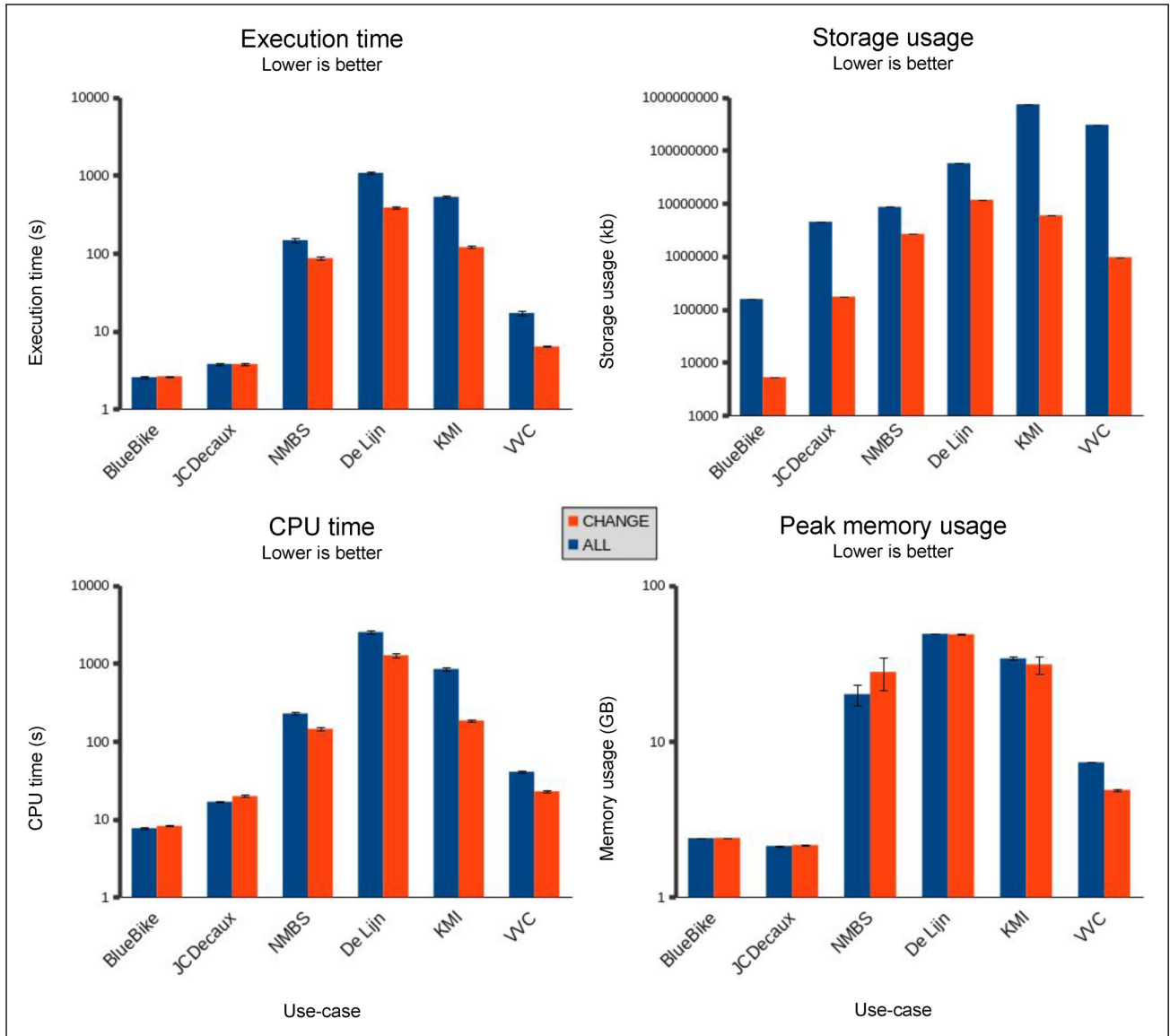
Scenario	Initial Storage Usage (kb)			Total Storage Usage (kb)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
BlueBike	108.60	164.00	0.66	156 381.58	5 241.14	29.84
JCDecaux	3120.69	3665.46	0.85	4493,804.23	173,649.86	25.88
NMBS	1242,672.71	1319,257.81	0.94	8564,018.79	2646,384.42	3.24
De Lijn	8705,837.26	10,050,766.41	0.87	57,132,528.26	11,638,171.53	4.91
KMI	5,130,955.63	5925,769.79	0.87	739,036,278.96	5928,882.90	124.66
VVC	10,581.32	10,957.26	0.97	304,588,651.83	964,395.21	315.83
OSM	NA	541.25	NA	NA	12,745,567.48	NA

Total storage usage is sum of the base dataset and the applied updates upon it. lower is better. OSM: OpenStreetMap; VVC: Vlaams Verkeerscentrum; KMI: Koninklijk Meteorologisch Instituut van België.

overhead of change detection causes a higher execution time, storage, and CPU time, while memory usage is unaffected, similar to the results observed and discussed in Section 7.2.

**Storage usage.** Our approach reduces storage usage by a factor ranging from 3.24 to 315.83 depending on the dataset (Figure 8). We observe a large reduction in storage for *full history* datasets such as KMI: The history itself does not change across new version releases and is considerably larger than the size of the members created/updated in new versions of the dataset. The same applies to datasets that have few changed members per version (e.g., VVC). In contrast, datasets with a large number of changes between versions, for example, NMBS, have a lower reduction in storage usage. Initial KG generation results in a higher storage usage for all datasets, as all members are generated along with the additional metadata indicating their creation.

**CPU time.** Similar to execution time, our approach reduces CPU time depending on the dataset size (Figure 8), with a factor up to 4.59. On larger datasets (e.g., KMI, De Lijn, NMBS, or VVC) more members avoid unnecessary



**Figure 8.** Our approach reduces the necessary resources to generate different versions of a knowledge graph (KG). Storage is reduced 3.24–315.83 times depending on dataset type. Execution time is reduced between 0.97–4.41 times depending on the dataset type. CPU time usage is reduced 0.85–4.59 times. Memory usage is mostly unaffected.

materialization, thus obtaining a higher reduction in CPU time. However, smaller datasets (e.g., BlueBike, JCDecaux) have a higher CPU time due to the overhead of continuously applying our change detection approach (7%–18% CPU time increase). The initial KG generation causes a higher CPU consumption for all datasets, as the state for each dataset is initialized to track each dataset member.

**Memory usage.** Results show that **our approach does not have a significant impact on the memory usage** (Figure 8) during the execution of KG generation processes. We attribute this to the compensation effect in our approach: Avoiding materializing unchanged members compensates for the memory overhead of detecting changes. NMBS is an exception in this case because the Java VM does not execute garbage collection while processing this dataset, since the heap space does not reach its limits yet. Bigger datasets, for example, KMI and De Lijn reach the heap space limit, triggering garbage collection. Therefore, memory usage grows and varies more for NMBS compared to the other datasets. The VVC dataset has a higher reduction in memory because processing XML data is costly regarding memory: When only the changes are processed, only a fraction of the XML is processed and kept in memory, causing a lower memory consumption. The initial

KG generation has no impact on memory consumption, similar to processing dataset updates because the same amount of memory is needed to evaluate each data member of the base dataset and its corresponding updates.

## 7.4 Ingestion

Comparing the size of each update between ALL and CHANGE indicates that processing CHANGE updates requires less resources and time because they are smaller. To fully compare our approach with the ALL strategy (i.e., having access to the latest standalone version of a KG), we measure the total processing time of the ingestion step by ingesting the real-world datasets, via SPARQL UPDATE queries, into a Virtuoso triplestore.

Table 12 presents the total time of processing all collected versions for each real-world dataset (as presented in Table 5), for each step of the pipeline: (i) KG generation time (as presented in Section 7.3), (ii) interpretation of the generated changes as SPARQL updates (using `incrml2sparql`, Section 6.1.4), and (iii) ingestion time in a triplestore, for both the ALL and CHANGE strategy. We present these metrics for the successfully completed real-world cases and the total execution time of the pipeline.

Unfortunately, most of these SPARQL queries are immediately rejected if the number of the changes increases, due to reaching Virtuoso's internal query length limits (10Mb). We overcame this problem by splitting up queries into multiple smaller ones, using the following pragmatic method: We started with 100 triples as upper limit and divided the queries in 2 parts each time a query failed, thus minimizing the number of queries needed. Depending on how the data is structured in RDF, Virtuoso was able to handle bigger queries or not, for example: Large string literals have a higher memory impact on Virtuoso. Therefore, only BlueBike, JCDecaux, and VVC were successfully ingested into Virtuoso. We further tried to optimize the queries by splitting them into smaller queries. This helped already for the JCDecaux dataset to ingest it, but not for bigger datasets with larger change sets such as NMBS, De Lijn, KMI, or GTFS Madrid Benchmark.

The results show a reduction on the overall ingestion time for the datasets that were successfully ingested, which includes also the time needed to interpret the change semantics and generate the corresponding SPARQL UPDATE queries for the CHANGE strategy. **We observe that the overall execution time (i.e., the time required to generate, interpret, and ingest all versions of a dataset) is reduced 1.05–2.23 times, depending on the dataset.** This results provide already a very promising indication of the usefulness of our approach and highlight a broader need for more effective implementations of SPARQL UPDATE query execution in triplestores.

## 7.5 Lessons Learned

In this subsection, we discuss the lessons learned (Table 13) from analyzing different real and synthetic datasets with varying change signaling strategies, change types, and history availability, for incremental KG generation.

**Functionality.** Our approach handles all possible dimensions and is feasible to implement. We implemented our approach by extending the RMLMapper v6.3.0 with **100% coverage for all the test cases** (Table 4). Therefore, we confirmed that our approach can detect all change types, both explicit and implicit, with all identified history signaling types and for all change types.

**GTFS Madrid Benchmark.** Our **approach reduces the storage (~6x), execution time (~2x), CPU time (~2x) and memory consumption is mostly unaffected** for any of the benchmark scales (data size, amount of changes, and change types). However, the initial KG generation has a higher execution time and consumes more resources because all members are created and tracked by our approach, which causes this initial overhead.

Scaling the data size (1, 10, 100) has the most impact for  $GTFS_{SCALE\ 1}$  because the dataset is rather small resulting in a more visible impact of the overhead from tracking dataset members. CPU time usage remains the same for all scales, while execution time has a larger increase for  $GTFS_{SCALE\ 1}$  (–30%) compared to  $GTFS_{SCALE\ 10}$  and  $GTFS_{SCALE\ 100}$ . Only for  $GTFS_{SCALE\ 10}$  the peak memory usage increases by 50%, due to the Java VM not performing garbage collection, as the Java heap space still has enough free memory.  $GTFS_{SCALE\ 100}$  is larger which causes the Java VM to trigger garbage collection, thus lowering memory usage. Storage usage is lower for any data size scale (6.56–8.12 times).

Scaling the amount of changes (0%, 25%, 50%, 75%, 100%) increases storage usage when more changes are present (11.70GB–12.93GB) since more new versions of dataset members are generated. CPU time usage and memory usage are unaffected because the number of dataset members remains unchanged, each dataset member is assessed to determine if it was changed or not.

Change types (create, update, delete) mostly affects the storage usage, for created members the storage usage (13.5GB) is higher compared to deleted (11.72GB) or updated members (11.71GB) because the new members that are added to the dataset need to be materialized. Deletions reduce the storage usage, but not significantly because deleting a GTFS Route only affects a subset of the dataset. Moreover, a tombstone comprising a few triples is still materialized to indicate the deletion of the member for historical purposes.

**Table 12.** Results for Generating, Interpreting, and Ingesting of the Different Real-World Datasets into a Triplestore (virtuoso).

Dataset	Generation Time (s)			Interpretation Time (s)			Ingestion Time (s)			Total Execution Time (s)		
	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C	ALL	CHANGE	Ratio A/C
BlueBike	3794.13	3691.84	1.02	601.74	607.65	0.99	123.24	111.13	11.07	4519.11	4310.72	1.05
JCDecaux	5450.90	5461.16	1.00	994.05	653.36	1.52	5746.05	189.84	30.27	12191.00	6304.46	1.93
VVC	503328.91	242419.45	2.08	22967.47	11721.36	1.96	43322.92	751.29	57.66	569619.30	254892.10	2.23
OSM	NA	5907.71	NA	NA	1852.16	NA	NA	6235.21	NA	NA	13995.08	NA

Failed datasets are caused by rejections of the SPARQL UPDATE queries by the triplestore when the initial graph or updates are too large. Larger datasets such as KMI fails to ingest, same with GTFS-based datasets, that is, NMBS and de lijn. GTFS Madridbenchmark is already larger than the failing GTFS-based datasets, thus it fails as well. Note that the number of versions is the same for ALL and CHANGE, the reported number of versions is lower for CHANGE since unchanged versions are skipped. lower is better: OSM: OpenStreetMap; VVC: Vlaams Verkeerscentrum; KMI: Koninklijk Meteorologisch Instituut van België.

**Table 13.** Analysis Goals on Change Signaling Strategies, Change Types, and History Availability from Datasets Using Incremental Knowledge Graph (KG) Generation.

Experiment	Goal
Functionality	Feasibility and coverage of our approach
GTFS Madrid Benchmark	Execution time and resource consumption analysis in multiple dimensions for example, data size, change percentage per version, and change types.
Real-world datasets	Real-world execution time and resource consumption analysis, including end-to-end comparison.

*Real-world datasets.* Depending on the data size, **our approach's RDF generation process runs faster (up to a factor of 4.41x)**. Besides faster RDF generation, the amount of resources is also **reduced in terms of storage** (factor 3.24–315.83) and **CPU time** (up to a factor of 4.59) depending on the dataset size. **No significant impact is observed for memory usage** because the additional memory required by our approach for tracking changes is compensated by avoiding materializing the unchanged members. Detecting changes has an overhead, but it is mitigated when only a certain part of the KG must be regenerated, which is usually the case when processing new versions of datasets in practice. The type of change does not affect our approach in execution time, CPU and memory usage because every dataset member is checked for changes in any case. However, storage usage is affected since additional created members need to be stored, deleted members release occupied storage, and updated members could increase or decrease storage usage based on the change. However, in a complete pipeline, **interpreting and ingesting the changes into a triplestore reduces the efficiency of our approach, but we are still faster by a factor of 1.05–2.23** compared to ingested complete datasets. Smaller datasets are more affected since they have more overhead.

**For smaller datasets** (e.g., BlueBike or JCDecaux), **the overhead of our approach increases CPU time (+6.96% to +18.12%) and execution time (−0.20% to +2.76%) while storage usage is significantly reduced (−96.14% to −96.65%)**. For the initial KG generation from the base dataset, the execution time is longer and more resources are consumed due to the initialization of the internal state used for tracking all dataset members. Once the dataset members are processed, the execution time is almost the same, but the CPU time is still higher. This is the overhead of tracking each member which impacts smaller datasets more visibly compared to larger ones. Storage usage is still heavily reduced. Despite the slight overhead on execution and CPU time introduced by our approach, the important reduction on storage usage may be decisive motivating factor for data publishers, which now could also store and offer historical records efficiently. Moreover, the semantic annotation of the different type of changes happening in a dataset, could also enable different and independent application behaviours both for the publishers and for remote data consumers if published via LDES.

**For larger datasets** (e.g., VVC, NMBS, KMI, or De Lijn), **our approach has a much larger impact on reducing CPU time usage (−33.57% to −78.22%), storage usage (−69.10% to −99.68%), and execution time (−41.37% to −77.34%)** with respect to its overhead, and compared to (re-)materializing the complete dataset. Note that memory consumption has an increase for NMBS (+38.8%) because the Java VM did not perform garbage collection as still enough heap memory is free for this dataset. Other datasets achieve a reduction of −0.55% to −31.31%.

In general, this results show a clear advantage of our approach, in terms of required computational resources, with respect to the traditional way of generating RDF KGs. Even considering the additional step required to fully ingest and integrate the materialized member changes into a KG (Table 12), we observe a clear reduction on the total time execution (SPARQL generation + ingestion). The beneficial impact of our approach increases with the size of the original datasets and provided that the number of changes remain relatively low compared to the total size, which is usually the case in most practical scenarios. This positions our approach as a scalable solution with promising potential to be used in production.

*Use cases benefiting from incremental KG generation.* Our approach is the most beneficial for use cases where data changes frequently or requires historical data. For example: Real-time information for public transport, IoT sensor data for smart cities, or weather history for analyzing the impact of climate change (Section 7). Real-time data needs to be integrated quickly to provide up-to-date information and the time to regenerate a KG for each data change is longer and consumes more computing resources when dealing with large datasets, compared to incrementally generating KGs. If historical information is important, it is beneficial to use our approach as only changes throughout history are incorporated into the KG which saves storage. This way, consumers can access and store the complete history for analytical purposes. However, if the data is rather small, for example, BlueBike or JCDecaux bicycle data, our approach's overhead to detect changes does not reduce the execution time or computing resources, but still provides an important reduction on storage requirements while providing access to historical information.

## 8 Conclusion

In this paper, we investigated how to detect and materialize only dataset updates towards establishing an incremental publishing approach for KGs. To achieve this, we designed an approach that combines established KG generation technologies and a novel KG publishing approach (IncRML), which we implemented by extending the RMLMapper, and evaluated on 5 types of heterogeneous datasets. We observed that in general, our IncRML achieves a reduction in execution time for RDF generation (up to 4.41x), CPU time (up to 4.59x), memory usage (up to 1.51x), and storage (up to 315.83x). In terms of ingesting and fully updating a KG hosted in a triplestore, we also observed faster ingestion in IncRML both overall for all updates (up to 57.66x) and on average for individual updates (up to 28.5x), although we were only able to measure this for smaller datasets in one triplestore (Virtuoso) due to internal query size limitations.

Through this work, we establish a trade-off for generating and publishing KGs, where following our approach can lead to significant time and computing resource savings to generate the raw RDF quads of a KG, at the cost of introducing an additional step for change reconciliation. On the other hand, a traditional KG generation is capable of producing an already updated and integrated KG, at the cost of additional computing resources and processing time. Nevertheless, our experiments indicate that despite the additional processing step required by IncRML, the overall processing time to update KGs is still lower (by a factor of 1.05–2.23) than with a fully re-materialization approach.

We evaluated our approach on a heterogeneous set of real-world datasets, including weather sensor data, public transport timetables, bike sharing data, live road traffic information and crowdsourced geospatial data. We show that our *IncRML* implementation is able to cover the different change communication strategies used by these set of real-world datasets (explicitly by the data source, or implicitly by silently changing the data), change types (creations, updates, and deletions), and history availability in the dataset (latest state, latest changes, or full history datasets).

Thanks to our approach, we provide the means to generate semantically annotated data and metadata from both implicitly and explicitly changed datasets. Therefore, we help to bring more transparency/provenance, in particular to implicitly signalled datasets. We integrate LDES as a Web native alternative to publish a semantically and structurally described stream of events that could enable data consumers to replicate and continuously synchronize with a KG, whether their changes are explicitly or implicitly signalled in the original data sources. Furthermore, our approach facilitates and reduces the cost of storing and publishing historic data, which is commonly an important requirement for for example, data analysis and machine learning applications. Existing LDES clients<sup>32</sup> can already take advantage of our incrementally generated KGs given that only changed members need to be processed, which is usually lower compared to the total size of a dataset. In general, IncRML show great potential to be further developed into production ready solutions that could lower the costs of creating and consuming KGs, with the goal of increasing adoption of Semantic Web technologies.

Further research includes expanding the pipeline to investigate the impact on *end-to-end performance* with different triplestores and increasing amount of consumers. On this direction, triplestores could be extended to allow direct native ingestion of incremental KGs (e.g., as an LDES) instead of performing full re-ingestion, thus benefiting from the usual lower amount of data that needs to be processed. On the other hand, we also highlight the need for more effective SPARQL UPDATE query processing in triplestores, in order to remain standard compliant and not having to rely only on custom and vendor-dependent data update techniques.

Investigating *optimizations for the proposed change detection algorithms* is also possible given that currently the implementation of our CDC algorithms as FnO functions uses in-memory lookup tables to detect changes. Avoid keeping the lookup tables in memory (e.g., using key-value stores) could reduce the memory footprint of IncRML even further. Also, exploring *windowing techniques for streaming data*, could allow to handle deletions on unbounded data streams. Our approach can perform change detection on streaming data, but requires a window definition for detecting deletions. Specifying a window is not supported yet in any declarative mapping language, but is considered by the W3C Community Group on KG Construction<sup>33</sup>. Another potential path for future work is *querying of incremental KGs*. Performing SPARQL queries on incrementally generated KGs requires further investigation to optimize query execution as we only tackled incremental generation. The study of approaches to *handle schema-level (ontology and mappings) changes* efficiently, is an important aspect to be consider (Conde-Herreros et al., 2024) since also the ontology and mappings can evolve besides the data itself. Lastly, performing a survey on dataset change signalling could be performed to validate our change signalling and communication strategies.


### Funding


The described research activities were supported by the Special Research Fund (Ghent University grant BOF20/DOC/132) and SolidLab Vlaanderen (Flemish Government, EWI and RRF project VV023/10).

## Declaration of Conflicting Interests


The authors declared no potential conflicts of interest with respect to the research, authorship and/or publication of this article.

## ORCID iDs

Dylan Van Assche  <https://orcid.org/0000-0002-7195-9935>

Julián Andrés Rojas  <https://orcid.org/0000-0002-6645-1264>

Ben De Meester  <https://orcid.org/0000-0003-0248-0987>

Pieter Colpaert  <https://orcid.org/0000-0001-6917-2167>

## Notes

1. <https://www.w3.org/community/kg-construct/>
2. We use the most implemented version of FnO with RML (<https://fno.io/rml>, De Meester et al., 2023), but our approach can be used with the latest version as well (<http://w3id.org/rml/fnml>).
3. <https://git-scm.com/>
4. <https://www.w3.org/submissions/CBD/>
5. We use the most implemented version of RML (<https://rml.io/spec>), but our approach can be used with the latest version as well (<http://w3id.org/rml/core>).
6. <https://joinup.ec.europa.eu/collection/semic-support-centre/linked-data-event-streams-ldes>
7. <https://www.vlaanderen.be/vlaamse-smart-data-space-portaal>
8. <https://w3id.org/tree/specification/>
9. The `rr` prefix expands to <http://www.w3.org/ns/r2rml#>
10. The `as` prefix expands to <https://www.w3.org/ns/activitystreams#>
11. LDES specification: <https://w3id.org/ldes/specification>
12. RMLMapper repository: <https://github.com/RMLio/rmlmapper-java>
13. FnO functions for CDC: <https://github.com/FnOio/idlab-functions-java/blob/main/src/main/java/be/ugent/knowns/idlabFunctions/IDLabFunctions.java>
14. Activities of the PROV-O ontology: <https://www.w3.org/TR/prov-o/>
15. LDES specification: <https://w3id.org/ldes/specification#introduction>
16. We split up in 3 named graphs by change type for showcasing purposes. However, some use cases are only interested in a specific change type. For example within public transport: History of cancelled routes requires only the Delete named graph while providing travelling information requires all named graphs to indicate added, updated, and deleted routes.
17. <https://www.w3.org/ns/activitystreams>
18. Repository: <https://github.com/RMLio/rml-ldes-testcases>, DOI: <https://doi.org/10.5281/zenodo.10171394>
19. Repository: <https://github.com/oeg-upm/gtfs-bench>, DOI: <https://doi.org/10.5281/zenodo.10256865>
20. BlueBike: <https://api.blue-bike.be/pub/location>; JCDecaux: <https://developer.jcdecaux.com/#/home>
21. <https://github.com/jiaoxlong/gbfs-json-schema/tree/gbfs-ld/GBFS-LD/v2.3>
22. NMBS: <https://gtfs.irail.be/nmbs/gtfs/latest.zip>, De Lijn: [https://gtfs.irail.be/de-lijn/de\\_lijn-gtfs.zip](https://gtfs.irail.be/de-lijn/de_lijn-gtfs.zip)
23. <http://vocab.gtfs.org/terms#>
24. <https://planet.openstreetmap.org/replication/>
25. <https://datex2.eu/vocab/3/Vms/>
26. <https://www.w3.org/TR/vocab-ssn/>
27. <https://github.com/julianrojas87/incrml2sparql.git>
28. OpenLink's Virtuoso triplestore: <https://virtuoso.openlinksw.com/>
29. Our approach is engine-agnostic, thus can be implemented in any KG construction engine that support a declarative mapping language and data transformations.
30. DOI: <https://doi.org/10.5281/zenodo.10171156>
31. Repository: <https://github.com/RMLio/rmlmapper-java>, DOI: <https://doi.org/10.5281/zenodo.10142511>
32. <https://github.com/rdf-connect/ldes-client>
33. Windowing operation for streaming data sources: <https://github.com/kg-construct/rml-core/issues/85>

## References

- Arenas-Guerrero, J., Chaves-Fraga, D., Toledo, J., Pérez, M. S., & Corcho, O. (2022). Morph-KGC: Scalable knowledge graph materialization with mapping partitions. *Semantic Web*, 15(1), 1–20. <https://doi.org/10.3233/sw-223135>
- Brickley, D., & Guha, R. V. (2014). RDF Schema 1.1. recommendation, world wide web consortium (W3C). <http://www.w3.org/TR/rdf-schema/>.

- Cassidy, S., & Ballantine, J. (2007). Version control for rdf triple stores. In *international conference on software and data technologies*. <https://api.semanticscholar.org/CorpusID:12177206>.
- Chaves-Fraga, D., Priyatna, F., Cimmino, A., Toledo, J., Ruckhaus, E., & Corcho, O. (2020). Gtfs-madrid-bench: A benchmark for virtual knowledge graph access in the transport domain. *Journal of Web Semantics*, 65, 100596. <https://doi.org/10.1016/j.websem.2020.100596>
- Chortaras, A., & Stamou, G. (2018). mapping diverse data to RDF in practice. In *The Semantic Web – ISWC 2018: 17<sup>th</sup> international semantic web conference, Monterey, CA, USA, October 8–12, 2018, Proceedings, Part I*, (pp. 441–457). <https://doi.org/10.1007/978-3-030-00671-6>.
- Colpaert, P. (2022). Building materializable querying interfaces with the TREE hypermedia specification. In D. Graux, F. Orlandi, E. Niazmand, G. Ydler & M. Vidal (Eds.), *Proceedings of the 8th Workshop on managing the evolution and preservation of the data Web (MEPDaW) co-located with the 21st International Semantic Web Conference (ISWC 2022), Virtual event, October 23rd, 2022, CEUR Workshop Proceedings*, (Vol. 3339, pp. 8–18). CEUR-WS.org. <https://ceur-ws.org/Vol-3339/paper2.pdf>.
- Colpaert, P., Abelshausen, B., Melendez, J. A. R., Delva, H., & Verborgh, R. (2019). Republishing OpenStreetMap’s roads as linked routable tiles. In P. Hitzler, S. Kirrane, O. Hartig, V. de Boer, M.-E. Vidal, M. Maleshkova, S. Schlobach, K. Hammar, N. Lasierra, S. Stadtmüller, K. Hose, R. Verborgh (Eds.), *Semantic web: ESWC 2019 Satellite events*, (Vol. 11762, pp. 13–17). Springer. ISBN 9783030323264. [http://doi.org/10.1007/978-3-030-32327-1\\_3](http://doi.org/10.1007/978-3-030-32327-1_3).
- Conde-Herrerros, D., Stork, L., Pernisch, R., Poveda-Villalón, M., Corcho, O., & Chaves-Fraga, D. (2024). Propagating ontology changes to declarative mappings in construction of knowledge graphs. In *Proceedings of the 5th international workshop on knowledge graph construction co-located with 21th extended semantic web conference (ESWC 2024), CEUR Workshop Proceedings*, (Vol. 3718, pp. 1–16). CEUR-WS.org. <https://ceur-ws.org/Vol-3718/paper1.pdf>.
- Cygniak, R., Wood, D., & Lanthaler, M. (2014). RDF 1.1 Concepts and Abstract Syntax. Recommendation, world wide web consortium (W3C). <http://www.w3.org/TR/rdf11-concepts/>.
- Daga, E., Asprino, L., Mulholland, P., & Gangemi, A. (2021). Facade-X: An opinionated approach to SPARQL anything. In *Further with knowledge graphs – proceedings of the 17<sup>th</sup> international conference on semantic systems, 6–9 September 2021, Amsterdam, The Netherlands*. (pp. 58–73). <https://doi.org/10.3233/SSW210035>.
- Das, S., Sundara, S., & Cygniak, R. (2012). R2RML: RDB to RDF Mapping Language. Working group recommendation, world wide web consortium (W3C). <http://www.w3.org/TR/r2rml/>.
- Debruyne, C., McKenna, L., & O’Sullivan, D. (2017). Extending r2rml with support for rdf collections and containers to generate mads-rdf datasets. In *Research and advanced technology for digital libraries: 21<sup>st</sup> international conference on theory and practice of digital libraries, TPDL 2017, Thessaloniki, Greece, September 18-21, 2017, Proceedings*. (pp. 531–536). [https://doi.org/10.1007/978-3-319-67008-9\\_42](https://doi.org/10.1007/978-3-319-67008-9_42)
- De Meester, B., Jozashoori, S., Maria, P., Chaves-Fraga, D., & Dimou, A. (2023). RML-FNML. Technical report, Knowledge graph construction community group. <https://kg-construct.github.io/rml-fnml/spec/docs/>.
- De Meester, B., Seymoens, T., Dimou, A., & Verborgh, R. (2020). Implementation-independent function reuse. *Future Generation Computer Systems*, 110, 946–959. <https://doi.org/10.1016/j.future.2019.10.006>
- Denny, A. I. P. M., Saptawijaya, A., & Aminah, S. (2017). Implementation of change data capture in etl process for data warehouse using hdfs and apache spark. In *2017 International workshop on big data and information security (IWBIS)*. (pp. 49–55). <https://doi.org/10.1109/IWBIS.2017.8275102>
- de Sompel, H. V., Nelson, M., & Sanderson, R. (2013). HTTP Framework for Time-based access to resource states – Memento. RFC 7089. <https://doi.org/10.17487/RFC7089>
- Dimou, A., Vander Sande, M., Colpaert, P., Verborgh, R., Mannens, E., & Van de Walle, R. (2014). RML: A Generic language for integrated rdf mappings of heterogeneous data. In *Proceedings of the 7<sup>th</sup> workshop on linked data on the web*.
- Fernández, J. D., Polleres, A., & Umbrich, J. (2015). Towards efficient archiving of dynamic linked open data.
- Frommhold, M., Piris, R. N., Arndt, N., Tramp, S., Petersen, N., & Martin, M. (2016). Towards versioning of arbitrary rdf data. In *Proceedings of the 12th international conference on semantic systems* <https://api.semanticscholar.org/CorpusID:14113981>.
- García-González, H., Boneva, I., Staworko, S., Labra-Gayo, J. E., & Lovelle, J. M. C. (2020). ShExML: Improving the usability of heterogeneous data mapping languages for first-time users. *PeerJ Computer Science*, e318.
- Goyal, A., & Dyreson, C. (2019). Temporal json. In *2019 IEEE 5th International conference on collaboration and internet computing (CIC)*. (pp. 135–144). <https://doi.org/10.1109/CIC48465.2019.00025>
- Graube, M., Hensel, S., & Urbas, L. (2014). R43ples: Revisions for triples - an approach for version control in the semantic web. In *LDQ@SEMANTiCS*. <https://api.semanticscholar.org/CorpusID:14184753>.
- Gupta, S., & Giri, V. (2018). Capture streaming data with change-data-capture. In *practical enterprise data lake insights: Handle data-driven challenges in an enterprise big data lake* (pp. 87–123). Apress. [https://doi.org/10.1007/978-1-4842-3522-5\\_3](https://doi.org/10.1007/978-1-4842-3522-5_3)
- Hao, L., Jiang, T., Lin, Y., & Lu, Y. (2023). Methods for solving the change data capture problem. In *Advances in natural computation, fuzzy systems and knowledge discovery*. (pp. 781–788).

- Harris, S., & Seaborne, A. (2013). SPARQL 1.1 Query Language. Recommendation, World Wide Web Consortium (W3C). <https://www.w3.org/TR/sparql11-query/>.
- Hu, Q., Gan, Z., & Zhang, B. (2019). Design and implementation of oracle database incremental data capture based on trigger and identification table. *Journal of Physics: Conference Series*, (), 0.
- Iglesias, E., Jozashoori, S., Chaves-Fraga, D., Collarana, D., & Vidal, M. E. (2020). SDM-RDFizer: An RML Interpreter for the Efficient Creation of RDF Knowledge Graphs. In *Proceedings of the 29<sup>th</sup> ACM international conference on information & knowledge management*. <https://doi.org/10.1145/3340531.3412881>
- 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. (2023). The RML Ontology: A community-driven modular redesign after a decade of experience in mapping heterogeneous data to rdf. In T. R. Payne, V. Presutti, G. Qi, M. Poveda-Villalón, G. Stoilos, L. Hollink, Z. Kaoudi, G. Cheng & J. Li (Eds.), *The Semantic Web – ISWC 2023* (pp. 152–175). Cham: Springer Nature Switzerland. ISBN 978-3-031-47243-5.
- Im, D. H., Lee, S. W., & Kim, H. J. (2012). A version management framework for rdf triple stores. *International Journal of Software Engineering and Knowledge Engineering*, 22(01), 85–106. <https://doi.org/10.1142/S0218194012500040>
- Junior, A. C., Debruyne, C., Brennan, R., & O’Sullivan, D. (2017). An evaluation of uplift mapping languages. *International Journal of Web Information Systems*, (), 0–0.
- Konstantinou, N., Kouis, D., & Mitrou, N. (2014). Incremental export of relational database contents into rdf graphs. In *Proceedings of the 4th international conference on web intelligence, mining and semantics (WIMS14)*, WIMS ’14. Association for Computing Machinery. <https://doi.org/10.1145/2611040.2611082>
- Lefrançois, M., Zimmermann, A., & Bakerally, N. (2017). A SPARQL extension for generating RDF from heterogeneous formats. In *The Semantic Web 14<sup>th</sup> international conference, ESWC 2017, Portorož, Slovenia, May 28 – June 1, 2017, Proceedings*. (pp. 35–50). [https://doi.org/10.1007/978-3-319-58068-5\\_3](https://doi.org/10.1007/978-3-319-58068-5_3)
- Ma, K., & Yang, B. (2015). Log-based change data capture from schema-free document stores using mapreduce. In *2015 International conference on cloud technologies and applications (CloudTech)*. (pp. 1–6). <https://doi.org/10.1109/CloudTech.2015.7336969>
- MadeSukarsa, I., Wisswani, N., Putra, I., & Linawati, L. (2012). Change data capture on oltp staging area for nearly real time data warehouse base on database trigger. *International Journal of Computer Applications*, 52(11), 32–37. <https://doi.org/10.5120/8248-1762>
- Meimaris, M., & Papastefanatos, G. (2016). The EvoGen Benchmark Suite for Evolving RDF Data. In *MEPDAW/LDQ@ESWC*. <https://api.semanticscholar.org/CorpusID:12789745>.
- Meinhardt, P., Knuth, M., & Sack, H. (2015). Tailr: A platform for preserving history on the web of data. In *Proceedings of the 11<sup>th</sup> international conference on semantic systems*. (pp. 57–64). <https://doi.org/10.1145/2814864.2814875>
- Michel, F., Djimenou, L., Faron-Zucker, C., & Montagnat, J. (2015). Translation of heterogeneous databases into RDF, and application to the construction of a SKOS taxonomical reference. In *International conference on web information systems and technologies*. (pp. 275–296). [https://doi.org/10.1007/978-3-319-30996-5\\_14](https://doi.org/10.1007/978-3-319-30996-5_14)
- Michel, F., Djimenou, L., Faron-Zucker, C., & Montagnat, J. (2017). xR2RML: Relational and non-relational databases to RDF Mapping Language. Rapport de recherche, Laboratoire d’Informatique, Signaux et Systèmes de Sophia-Antipolis (I3S). <https://hal.archives-ouvertes.fr/hal-01066663/document/>.
- Neumann, T., & Weikum, G. (2010). X-rdf-3x: Fast querying, high update rates, and consistency for rdf databases. *Proceeding of the VLDB Endow*, 3(1–2), 256–263. <https://doi.org/10.14778/1920841.1920877>
- Papakonstantinou, V., Flouris, G., Fundulaki, I., Stefanidis, K., & Roussakis, G. (2016). Versioning for linked data: Archiving systems and benchmarks. *BLINK@ ISWC, 1700*, 46–61.
- Prud’hommeaux, E., Labra Gayo, J. E., & Solbrig, H. (2014). Shape expressions: an RDF validation and transformation language. In H. Sack, A. Filipowska, J. Lehmann & S. Hellmann (Eds.), *Proceedings of the 10<sup>th</sup> international conference on semantic systems* (pp. 32–40). ACM, New York, NY, United States: Association for Computing Machinery. <https://doi.org/10.1145/2660517.2660523>
- Pu, X., Wang, J., Song, Z., Luo, P., & Wang, M. (2014). Efficient incremental update and querying in aweto rdf storage system. *Data & Knowledge Engineering*, 89, 55–75. <https://doi.org/10.1016/j.datak.2013.11.003>
- Randles, A., & O’Sullivan, D. (2022). Modelling & analyzing changes within ld source data. In *MEPDAW@ISWC*. <https://api.semanticscholar.org/CorpusID:257081241>.
- Randles, A., & O’Sullivan, D. (2023). Preserving the alignment of ld with source data. In *KGCW@ESWC*. <https://api.semanticscholar.org/CorpusID:259266370>.
- Randles, A., O’Sullivan, D., Keeney, J., & Fallon, L. (2022). Applying a mapping quality framework in cloud native monitoring. In *International conference on semantic systems*. <https://api.semanticscholar.org/CorpusID:252919576>.
- Rojas, J. A., Aguado, M., Vasilopoulou, P., Velitchkov, I., Assche, D. V., Colpaert, P., & Verborgh, R. (2021). Leveraging semantic technologies for digital interoperability in the European Railway domain <https://julianrojas.org/papers/iswc2021-in-use/>.
- Salzberg, B., & Tsotras, V. J. (1999). Comparison of access methods for time-evolving data. *ACM Computing Surveys*, (), 0–0.
- Skjæveland, M. G., Lupp, D. P., Karlsen, L. H., & Forssell, H. (2018). Practical ontology pattern instantiation, discovery, and maintenance with reasonable ontology templates. In D. Vrandečić, K. Bontcheva, M. C. Suárez-Figueroa, V. Presutti, I. Celino, M. Sabou,

- L. A. Kaffee & E. Simperl (Eds.), *The Semantic Web – ISWC 2018* (pp. 477–494). Cham: Springer International Publishing. ISBN 978-3-030-00671-6.
- Snell, J. M., & Prodomou, (2017). Activity Streams 2.0. Recommendation, World Wide Web Consortium (W3C). <http://www.w3.org/TR/activitystreams-core/>.
- Taelman, R., Vander Sande, M., Van Herwegen, J., Mannens, E., & Verborgh, R. (2018). Triple storage for random-access versioned querying of RDF archives. *Journal of Web Semantics*, 54(C), 4–28. <https://doi.org/10.1016/j.websem.2018.08.001>
- Umbrich, J., Villazón-Terrazas, B., & Hausenblas, M. (2010). Dataset dynamics compendium: A comparative study. In *COLD*. <https://api.semanticscholar.org/CorpusID:15551988>.
- Valencio, C. R., Marioto, M. H., Donega Zafalon, G. F., Machado, J. M., & Momente, J. C. (2013). Real time delta extraction based on triggers to support data warehousing. In *2013 International conference on parallel and distributed computing, applications and technologies*. (pp. 293–297). <https://doi.org/10.1109/PDCAT.2013.52>
- Van Assche, D., Delva, T., Haesendonck, G., Heyvaert, P., De Meester, B., & Dimou, A. (2022a). Declarative RDF graph generation from heterogeneous (semi-)structured data: A systematic literature review. *Journal of Web Semantics* <https://doi.org/10.1016/j.websem.2022.100753>
- Van Assche, D., Haesendonck, G., De Mulder, G., Delva, T., Heyvaert, P., De Meester, B., & Dimou, A. (2021). Leveraging Web of Things W3C Recommendations for Knowledge Graphs Generation. In *Web Engineering, 21<sup>st</sup> international conference, ICWE 2021, Biarritz, France, May 18–21, 2021*, (pp. 337–352). [https://doi.org/10.1007/978-3-030-74296-6\\_26](https://doi.org/10.1007/978-3-030-74296-6_26)
- Van Assche, D., Oo, S. M., Rojas, J. A., & Colpaert, P. (2022b). Continuous generation of versioned collections’ members with RML and LDES. In *Proceedings of the 3<sup>rd</sup> international workshop on knowledge graph construction (KGCW 2022) co-located with 19<sup>th</sup> Extended Semantic Web Conference (ESWC 2022)*.
- Vander Sande, M., Colpaert, P., Verborgh, R., Coppens, S., Mannens, E., & Van de Walle, R. (2013). R&Wbase: git for triples. In *Proceedings of the 6<sup>th</sup> workshop on linked data on the web*.
- Van Lancker, D., Colpaert, P., Delva, H., Van de Vyvere, B., Rojas Meléndez, J., Dedecker, R., Michiels, P., Buyle, R., De Craene, A., & Verborgh, R. (2021). Publishing base registries as linked data event streams. In *Web engineering, 21<sup>st</sup> International Conference, ICWE 2021, Biarritz, France, May 18–21, 2021*, (pp. 28–36). [https://doi.org/10.1007/978-3-030-74296-6\\_3](https://doi.org/10.1007/978-3-030-74296-6_3)
- Vidal, V. M. P., Casanova, M. A., & Cardoso, D. S. (2013). Incremental maintenance of rdf views of relational data. In *On the Move to meaningful internet systems: OTM 2013 conferences*. (pp. 572–587).
- Völkel, M., Winkler, W., Sure, Y., Kruk, S. R., & Synak, M. (2005). Semversion: A versioning system for rdf and ontologies. <https://api.semanticscholar.org/CorpusID:14892100>.
- Vu, B., Pujara, J., & Knoblock, C. A. (2019). D-repr: A language for describing and mapping diversely-structured data sources to rdf. In *Proceedings of the 10<sup>th</sup> international conference on knowledge capture*. (pp. 189–196). <https://doi.org/10.1145/3360901.3364449>