

Describing a network of live datasets with the SDS vocabulary

Arthur Vercruysse

Sitt Min Oo

Wout Slabbinck

Pieter Colpaert

Abstract—Datasets can be transformed by activities and often originates from a different live dataset. Multiple published datasets can be the result of the same initial dataset, we want to give query processors transparency in how datasets are linked. In this paper, we introduce the Smart Data Specification for Semantically Describing Streams (SDS) to annotate datasets with provenance information, describing the consumed stream and the applied transformations on that stream. We demo the execution of a pipeline that transforms an LDES and publishes the data in a different structure as described in the SDS description. The SDS vocabulary and application profile bring together DCAT-AP, LDES and P-Plan. In future work, we will create a source selection strategy for federated query processors that take into account this provenance information when selecting a dataset and interface to query the dataset.

I. INTRODUCTION

Data portals publish datasets and datasets can be derived from different datasets. This leads to a problem: when datasets are not accompanied by provenance information, the data portal has difficulty knowing whether or not this data is already published. Also, query agents are clueless as to what interfaces provide the same underlying dataset and as a fallback, they query the same underlying dataset multiple times.

In an application that notifies the user about changes in routes due to construction sites, multiple interfaces can be used. A time-based interface makes it easy to track the latest changes. On the other hand, a geospatial-based interface makes it easy to calculate whether or not a construction site will be encountered on a route. A SPARQL endpoint can be the solution if low availability and operation cost are of no concern[1]. Another possibility is to publish multiple Linked Data Event Streams (LDES)[2]. Each stream can be fragmented in a different way to accommodate the needs of the users.

Using different LDESes to publish different fragmentations of the same data exposes a plethora of interfaces which makes decent provenance a necessity. The provenance should cover the two steps a query agent takes to query an interface[3], [4]:

1. *Dataset discovery*: based on whether a dataset is going to contain statements that contribute to solving the query
2. *Interface discovery*: selecting the right interface that publishes the dataset

In this paper, we define a dataset as the accumulation of all updates that were published on a stream. When dataset A is derived from dataset B, there will be a stream A with updates that are a function of the updates of stream B. We introduce the Smart Data Specification ...

We set up a new interface for the Belgian street name registry and demonstrate that a query agent can find the optimal interface to execute particular queries.

II. RELATED WORK

DCTerms, DCAT and VoID Exposing metadata about datasets is long established. Dublin Core Terms (DCTerms) can be used to provide basic information about resources, providing terms like *title*, *author*, *description* and *subject*[5]. Data Catalog Vocabulary (DCAT) is designed to facilitate interoperability between data catalogs published on the web[6]. DCAT also provides terms like *license*, which makes it possible to define a new license for an interface. The Vocabulary of Interlinked Datasets (VoID) focuses on explicitly linking datasets together on some predicate and defining subset datasets[7].

LDES Linked Data Event Streams is a way of exposing immutable objects with HTTP resources. These resources can be divided into fragments that are linked together. Fragmentations are used to spread the items over different HTTP resources. Each HTTP resource can, for example, hold all items that start with a particular letter. A view description describes the meaning of the fragments and their links[2].

VoCaLS: Vocabulary of interoperable streams & On a Web of Data Streams (Dell Aglio): extends the ideas of DCAT with more information about streaming data[8]. The work defines a stream slightly differently than in this paper. VoCaLS focuses on streams that generate high throughput updates, this requires processors to use a windowing mechanism. In this paper, a stream is seen more broadly as a growing collection of objects, updates or otherwise.

P-Plan and PROV-O: The Ontology for Provenance and Plans (P-Plan) is an extension of the PROV-O ontology [9] created to represent the plans that guided the execution of scientific processes. P-Plan describes how the plans are composed and their correspondence to provenance records that describe the execution itself [10].

III. THE SMART DATA SPECIFICATION FOR SEMANTICALLY DESCRIBING STREAMS (SDS)

A stream in the context of SDS is a *physical* live channel that carries updates or items. A dataset can be derived from a stream as the collection of all updates or items. A *physical* channel can be any medium like a Kafka stream, WebSocket stream or even a file where updates are appended. A stream can carry any data: CSV rows, mutable or immutable linked data objects, video stream bytes, etc.

A stream can be derived from a transformation applied to items on a different stream. This transformation should be described with **p-plan** in the SDS description that is part of the resulting stream. The stream and the transformation correspond with **p-plan:Entity** and **p-plan:Activity** respectively. This is shown as the pink part of Figure 1. The transformation can **prov:used** a different stream. With the power of the **p-plan**, query agents can understand how datasets are linked and what interface fits a specific query the best.

The SDS description can be expanded with metadata about the dataset collected from the stream with the **sds:dataset** predicate. This way parts of the datasets' metadata can be changed after a transformation. This is represented as the green part in Figure 1.

Linking specific items to the correct stream is done with **sds:Record**. An **sds:Record** points to the data (**sds:payload**) and the corresponding stream (**sds:stream**). These small objects make it possible for multiple streams to use the same channel. Each transformation can thus push **sds:Record**'s and leave the original stream intact. A stream of immutable objects can still be transformed, this transformation can, for example, calculate a hash or add a fragment id to the **sds:Record** object. The yellow part of Figure 1 gives a visual overview of **sds:Record**.

IV. DEMO

Data published with Linked Data Event Streams can be partitioned or fragmented in a multitude of ways. This helps query agents resolve their queries with as few web requests as possible. A default fragmentation constitutes a timestamp fragmentation, this allows clients to replicate and synchronize the dataset efficiently. A substring fragmentation, on the other hand, makes autocompletion more efficient[11].

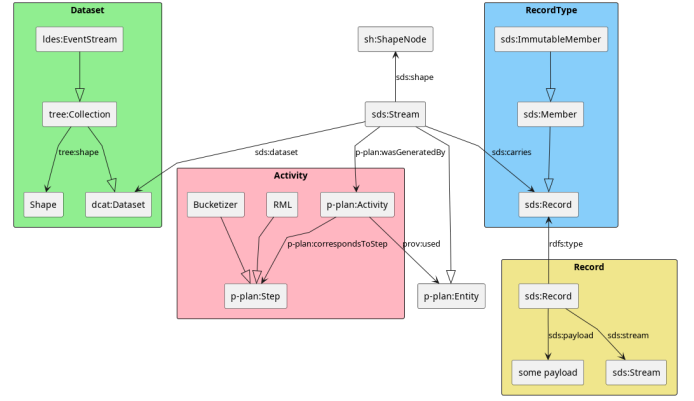


Figure 1: SDS Ontology

In this demo, we set up a pipeline starting from an existing LDES that exposes the registry of street names with a timestamp fragmentation. The pipeline calculates a substring fragmentation based on the name of the street and exposes a new LDES with the corresponding SDS Description.

When asking a query agent “What are the 10 latest updated street names?” starting from the newly created LDES, the query agent can derive from the SDS description that the current LDES is not suitable for this query. This query would require the query agent to request the entire LDES tree and manually find the 10 latest updates, whereas following the links from the SDS description back to the original LDES, this query would only require one or two HTTP requests.

Note that the original LDES does not expose an SDS description, so this has to be bootstrapped in the pipeline.

To execute this pipeline we use a proof of concept pipeline runner called Nautirust[12]. This makes it easy to start the three required processes with the correct arguments. The three required steps are: read the original LDES with an LDES client, add buckets to the SDS Records and ingest the new SDS records in an LDES server.

V. CONCLUSION

With the introduction of the SDS ontology it is possible to add a description to a stream and the resulting dataset, that provides provenance. The provenance links together stream and transformations applied to those stream. The SDS ontology aligns well with well established ontologies like DCAT and P-Plan to maximize interoperability.

The SDS description makes it possible for query agents to automatically select the right dataset and interface based on a given query.

Federated query processors, that utilize source selection based on this provenance information when selecting a

dataset and interface to query the dataset, is still future work.

VI. ACKNOWLEDGMENTS

Funded by the Flemish government's recovery fund VSDS project: the "Vlaamse Smart Data Space".

REFERENCES

- [1] C. Buil-Aranda, A. Hogan, J. Umbrich, and P.-Y. Vandenbussche, "SPARQL web-querying infrastructure: Ready for action?" in *The semantic web – iswc 2013*, 2013, pp. 277–293.
- [2] D. Van Lancker *et al.*, "Publishing base registries as linked data event streams," in *Web engineering*, 2021, pp. 28–36.
- [3] M. Ben Ellefi *et al.*, "RDF dataset profiling – a survey of features, methods, vocabularies and applications," *Semantic Web*, vol. 9, pp. 677–705, 2018.
- [4] F. Michel, C. Faron-Zucker, O. Corby, and F. Gandon, "Enabling automatic discovery and querying of web apis at web scale using linked data standards," in *Companion proceedings of the 2019 world wide web conference*, 2019, pp. 883–892.
- [5] T. Baker, "Libraries, languages of description, and linked data: A dublin core perspective," *Library Hi Tech*, vol. 30, no. 1, pp. 116–133, Jan. 2012.
- [6] A. G. Beltran, S. Cox, D. Browning, A. Perego, R. Albertoni, and P. Winstanley, "Data catalog vocabulary (DCAT) - version 2," W3C, W3C Recommendation, Feb. 2020.
- [7] K. Alexander, R. Cyganiak, M. Hausenblas, and J. Zhao, "Describing linked datasets," in *LDOW*, 2009.
- [8] R. Tommasini *et al.*, "VoCaLS: Vocabulary and catalog of linked streams," in *The semantic web – iswc 2018*, 2018, pp. 256–272.
- [9] T. Lebo, S. Sahoo, and D. McGuinness, "PROV-o: The PROV ontology," W3C, W3C Recommendation, Apr. 2013.
- [10] D. Garijo and Y. Gil, "The P-Plan ontology," Mar. 2014.
- [11] B. Van de Vyvere *et al.*, "Publishing cultural heritage collections of ghent with linked data event streams," in *Metadata and semantic research*, 2022, pp. 357–369.
- [12] A. Vercruysse and S. M. Oo, "Nautitrust connector architecture orchestrator," 2022. [Online]. Available: <https://github.com/ajuvercr/nautitrust>. [Accessed: 24-Aug-2022].