

Describing a network of live datasets with the SDS vocabulary

Arthur Vercruyse¹, Sitt Min Oo¹ and Pieter Colpaert¹

¹IDLab, Department of Electronics and Information Systems, Ghent University – imec, Technologiepark-Zwijnaarde 122, 9052 Ghent, Belgium

Abstract

Data publishers can provide multiple interfaces per dataset. Each interface has its own merits and drawbacks, SPARQL endpoints are expensive to host and clients find it difficult to work with static data dumps. Furthermore, query agents can only select the most fitting interface and dataset if provenance information is provided. In this paper, we introduce the Smart Data Specification for Semantically Describing Streams (SDS) to annotate dataset interfaces with provenance information, describing the consumed stream and the applied transformations on that stream. We focus on Linked Data Event Streams that can publish the same dataset with different fragmentations and demonstrate a pipeline that transforms a LDES and publishes the data with a different fragmentation as described in the accompanying provenance information. The SDS vocabulary is built upon the DCAT-AP, LDES and P-Plan vocabularies. 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.

Keywords

LDES, Linked Data, Provenance, Dataset selection

1. Introduction

The world is ever-changing, so data portals publish live datasets derived from streams of data (events, deltas, etc.). Streams are often related to other streams by transforming their data. This leads to problems when datasets and streams are not accompanied by provenance information. Query agents often use this information to query each dataset only once, while no or inadequate provenance information forces the query agents to query the same underlying data multiple times. Once a required dataset is found the query agent can go further and discover the most fitting interface. Different interfaces determine how data is accessed.

For example in a route planning 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, whereas a geospatial-based interface makes it easy to calculate whether or not a construction site will be encountered on a particular route. A SPARQL endpoint can fulfil both interfaces, but SPARQL carries high operation costs and may leave the users with


Managing the Evolution and Preservation of the Data Web (MEPDaW 2022)

✉ arthur.vercruyse@ugent.be (A. Vercruyse); x.sittminoo@ugent.be (S. M. Oo); pieter.colpaert@ugent.be (P. Colpaert)

🌐 <https://pietercolpaert.be> (P. Colpaert)

🆔 0000-0003-1586-5122 (A. Vercruyse); 0000-0001-9157-7507 (S. M. Oo); 0000-0001-6917-2167 (P. Colpaert)

© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

low availability [1]. Constantly updating a data dump is also an option, but leaves the client with high operation and bandwidth costs [2]. Linked Data Event Streams (LDES) alleviate these high costs by fragmenting the data in different web resources and linking these resources with semantic links in multiple dimensions[2], focussing on easy replication and synchronization. Each stream can be fragmented in a different way to accommodate to the needs of the users.

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.

2. 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 an evergrowing set of immutable objects. These objects can be divided into fragments as HTTP resources that are linked together with the TREE specification. Fragmentations are used to provide semantic meaning to links between HTTP resources. Each HTTP resource can, for example, hold all items that start with a particular letter. A view description is used to define the meaning of the fragments and their links[2].

VoCaLS: Vocabulary for Cataloging and Linking Streams and streaming services on the web 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 guide the execution of scientific processes. P-Plan describes how plans are composed. This information can be used to keep track of provenance tree[10].

3. The Smart Data Specification for Semantically Describing Streams (SDS)

A stream in the context of the Smart Data Specification for Semantically Describing Streams (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 such as a Kafka stream, a WebSocket stream or even a file where updates are appended. A stream can carry any kind of data: CSV rows, mutable or immutable linked data objects, video stream chunks, etc.

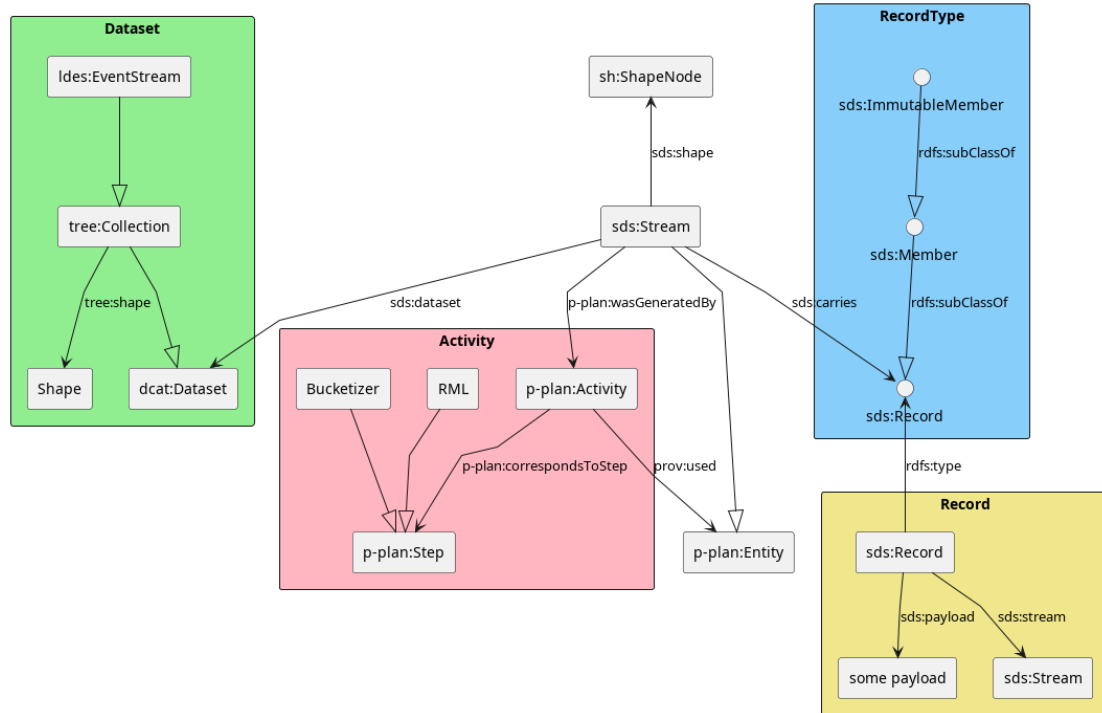


Figure 1: Visual representation of the SDS Ontology. The colored boxes have no semantic meaning and are only used for clarity.

A stream can be derived from a transformation applied to items of a different stream. This transformation is described with p-plan and resides in the SDS description 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 links to the previous stream with the prov:used predicate. 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 also links to metadata about the resulting dataset with the sds:dataset predicate. This makes modifying the datasets' metadata possible after a transformation. This is represented as the green part in Figure 1.

A sds:Stream should describe the type of objects it is carrying. First, the stream can use the sds:carries predicate that points to the type of the sds:Records it is carrying(see below). Next, the stream can describe the shacl shape with sds:shape that these members adhere to, which can only be used if the stream carries linked data objects.

Linking specific items to the correct stream is done with sds:Record objects. 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 objects and leave the original stream intact. A stream of immutable objects can still be transformed, for example, to 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.

4. 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 as fast as possible whilst ingesting as little data 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].

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 and substring fragmentation. The published Views can then be associated with their respective viewDescription as described in Listing 1.

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 a few HTTP requests. One HTTP request gets the SDS description and another request gets the latest updates due to the timestamp-based fragmentation.

5. Conclusion

The SDS ontology makes it possible to add a description to a stream and the resulting dataset, this adds provenance information. The provenance links streams together and transformations applied to the streams. The SDS ontology aligns well with long-established ontologies like DCAT and P-Plan to maximize interoperability.

With the SDS description, a query agent can now automatically select the right dataset and interface based on a given query. This can make a LDES more like a *querylike* interface because the best fragmentation for a particular problem can be found.

SDS descriptions also make it possible to track applied transformations down to the source. This enables the user that normally only extracts data from an interface to request changes at the data source. These changes will then be propagated through the streams and result in the intended change over all interfaces.

Federated query processors, that utilize source selection based on this provenance information when selecting a dataset and interface to query the dataset, are still future work.

6. Acknowledgments

Funded by the Flemish government’s recovery fund VSDS project: the “Vlaamse Smart Data Space”.

```

@prefix ex:      <http://example.org/ns#>.
@prefix rdfs:   <http://www.w3.org/2000/01/rdf-schema#> .
@prefix ldes:   <https://w3id.org/ldes#>.
@prefix p-plan: <http://purl.org/net/p-plan#> .
@prefix prov:   <http://www.w3.org/ns/prov#> .
@prefix sds:    <https://w3id.org/sds#> .

ex:MyLDES a ldes:EventStream;
  dcat:title "An example LDES".

# One fragmentation is based on an existing LDES
ex:BasicFragmentation a tree:ViewDescription ;
  dcat:endpointURL </basic> ; # A rootnote from which you can access all members
  dcat:servesDataset ex:MyLDES ; # the LDES
  # This viewDescription was created by importing this stream
  ldes:managedBy sds:LDESStream.

# The other fragmentation bucketizes the members before publishing
ex:SubstringFragmentation a tree:ViewDescription ;
  dcat:endpointURL </substring> ; # A rootnote from which you can access all members
  dcat:servesDataset ex:MyLDES ; # the LDES
  # This viewDescription was created by importing this stream
  ldes:managedBy sds:BucketizedStream.

ex:ImportLDES a p-plan:Activity;
  rdfs:comment "Reads csv file and converts to rdf members";
  prov:used <https://smartdata.dev-vlaanderen.be/base/gemeente>.

ex:LDESStream a sds:Stream;
  p-plan:wasGeneratedBy ex:ImportLDES;
  sds:carries sds:Member.

ex:BucketizeStrategy a ldes:BucketizeStrategy;
  ldes:bucketType ldes:SubstringFragmentation; # zegt aan de client dat ik een timestampfragmentation a
  tree:path rdfs:label;
  ldes:pageSize 50.

ex:BucketizeStream a p-plan:Activiy;
  rdfs:comment "Execute a substring bucketization on the incoming stream";
  prov:used ex:LDESStream, ex:BucketizeStrategy.

ex:BucketizedStream a sds:Stream;
  p-plan:wasGeneratedBy ex:BucketizeStream;
  sds:carries sds:Member.

```

Listing 1: RDF sample creating two LDES Views from different Streams

References

- [1] Buil-Aranda, C., Hogan, A., Umbrich, J., Vandenbussche, P.-Y.: SPARQL web-querying infrastructure: Ready for action? In: Alani, H., Kagal, L., Fokoue, A., Groth, P., Biemann, C., Parreira, J.X., Aroyo, L., Noy, N., Welty, C., and Janowicz, K. (eds.) *The semantic web – iswc 2013*. pp. 277–293. Springer Berlin Heidelberg, Berlin, Heidelberg (2013).
- [2] Verborgh, R., Vander Sande, M., Colpaert, P., Coppens, S., Mannens, E., Van de Walle, R.: Web-scale querying through Linked Data Fragments. In: Bizer, C., Heath, T., Auer, S., and Berners-Lee, T. (eds.) *Proceedings of the 7th workshop on linked data on the web* (2014).
- [3] Van Lancker, D., Colpaert, P., Delva, H., Van de Vyvere, B., Meléndez, J.R., Dedecker, R., Michiels, P., Buyle, R., De Craene, A., Verborgh, R.: Publishing base registries as linked data event streams. In: Brambilla, M., Chbeir, R., Frasincar, F., and Manolescu, I. (eds.) *Web engineering*. pp. 28–36. Springer International Publishing, Cham (2021).
- [4] Ben Ellefi, M., Bellahsene, Z., Breslin, J.G., Demidova, E., Dietze, S., Szymański, J., Todorov, K.: RDF dataset profiling – a survey of features, methods, vocabularies and applications. *Semantic Web*. 9, 677–705 (2018).
- [5] Michel, F., Faron-Zucker, C., Corby, O., Gandon, F.: Enabling automatic discovery and querying of web apis at web scale using linked data standards. *Companion proceedings of the 2019 world wide web conference*. pp. 883–892. Association for Computing Machinery, New York, NY, USA (2019).
- [6] Baker, T.: Libraries, languages of description, and linked data: A dublin core perspective. *Library Hi Tech*. 30, 116–133 (2012).
- [7] Beltran, A.G., Cox, S., Browning, D., Perego, A., Albertoni, R., Winstanley, P.: *Data catalog vocabulary (DCAT) - version 2*. W3C (2020).
- [8] Alexander, K., Cyganiak, R., Hausenblas, M., Zhao, J.: *Describing linked datasets*. LDOW (2009).
- [9] Tommasini, R., Sedira, Y.A., Dell’Aglia, D., Balduini, M., Ali, M.I., Le Phuoc, D., Della Valle, E., Calbimonte, J.-P.: VoCaLS: Vocabulary and catalog of linked streams. In: Vrandečić, D., Bontcheva, K., Suárez-Figueroa, M.C., Presutti, V., Celino, I., Sabou, M., Kaffee, L.-A., and Simperl, E. (eds.) *The semantic web – iswc 2018*. pp. 256–272. Springer International Publishing, Cham (2018).
- [10] Lebo, T., Sahoo, S., McGuinness, D.: *PROV-o: The PROV ontology*. W3C (2013).
- [11] Garijo, D., Gil, Y.: *The P-Plan ontology*. (2014).
- [12] Van de Vyvere, B., D’Huynslager, O.V., Ataulil, A., Segers, M., Van Campe, L., Vandekeybus, N., Teugels, S., Saenko, A., Pauwels, P.-J., Colpaert, P.: Publishing cultural heritage collections of ghent with linked data event streams. In: Garoufallou, E., Ovalle-Perandones, M.-A., and Vlachidis, A. (eds.) *Metadata and semantic research*. pp. 357–369. Springer International Publishing, Cham (2022).