

Towards Efficient Processing of RDF Data Streams - Short Paper

Alejandro Llaves, Javier D. Fernández, and Oscar Corcho

Ontology Engineering Group, Universidad Politécnica de Madrid, Madrid, Spain
{allaves,jdfernandez,ocorcho}@fi.upm.es

Abstract. In the last years, there has been an increase in the amount of real-time data generated. Sensors attached to *things* are transforming how we interact with our environment. Extracting meaningful information from these streams of data is essential for some application areas and requires processing systems that scale to varying conditions in data sources, complex queries, and system failures. This paper describes on-going research on the development of a scalable RDF streaming engine.

Keywords: RDF Stream Processing, Adaptive Query Processing, RDF Compression

1 Introduction to RDF Stream Processing

The Resource Description Framework (RDF) format is the W3C standard for data interchange on the Web.¹ With the growth of the Semantic Web,² many organizations started to publish their data as Linked Data³, e.g. the UK Government⁴ or the European Environmental Agency⁵.

The origins of Linked Stream Data [18,14] are in the enrichment of the Sensor Web with spatial, temporal, and thematic metadata [19]. Standardization initiatives such as the W3C Semantic Sensor Network Incubator Group⁶ and its SSN ontology [6] have fostered the publication of sensor datasets as Linked Stream Data and its consequent integration with other datasets. A more recent initiative is the W3C RDF Stream Processing (RSP) Community Group, which works on defining a common model for producing, transmitting, and continuously querying data streams encoded in RDF.⁷ Within the RSP group, several RDF stream processing systems that appeared during the last years are being studied, such as CQELS Cloud [13], C-SPARQL [4], INSTANS [17], EP-SPARQL [1], and morph-streams [5], among others.

¹ <http://www.w3.org/RDF/>

² <http://www.w3.org/standards/semanticweb/>

³ <http://www.w3.org/DesignIssues/LinkedData.html>

⁴ <http://data.gov.uk/>

⁵ <http://semantic.eea.europa.eu/>

⁶ <http://www.w3.org/2005/Incubator/ssn/>

⁷ <http://www.w3.org/community/rsp/>

Extracting information from data streams is complex because of the heterogeneity of the data, the rate of data generation, the high volumes, and the often unclear data provenance. One use case would be real-time monitoring of public transportation in a city. Here, decisions on unexpected events, such as a car crash, should be taken on short time slots based on a set of spatio-temporal data streams coming from different providers. For instance, by diverting a bus line route. This means reasoning over data in the temporal order it is ingested by the processing system. Linked Stream Data may help to integrate datasets from different providers and to solve interoperability problems. Yet, remaining challenges require solutions from a stream processing approach. In the remainder of the paper we discuss some of these challenges (section 2), a general description of our approach towards efficient processing of RDF streams (section 3), and open questions to debate during the workshop (section 4).

2 Challenges on RDF Stream Processing

The development of scalable services for real-time stream processing involves support for high throughput, management of complex queries, low latency response, fault-tolerance, and statistics extraction, among others. Nowadays, cloud services offer solutions to some of these problems, e.g. by applying elastic load balancing in presence of input data bursts. We mainly focus on the **efficient processing of user queries over RDF streams (C.1)**, which requires parallelization at the query operator level.

Distributed computing refers to the processing of data in distributed systems. The **continuous transmission of data (C.2)** between sources and processing nodes, and among nodes may cause response delays that should be minimized when possible.

The **integration of historical and real-time data with background knowledge (C.3)** is challenging in Web-scale environments. Many RSP systems combine background knowledge with real-time processing, but historical data management is often overlooked [15]. Traditional systems tend to store data aggregates in order to save space, but with the capacity of current systems it is possible (and recommended) to store all data in raw format and define views on data batches [16]. The efficient management of historical data is essential to detect trends in data, extract statistics, or compare old data with current data to identify anomalies [15]. For instance, by aggregating metro users on the last hours and compare the numbers to the average of users during the last days.

3 Efficient Processing of Queries over RDF Streams

This section gives an high level overview of our approach, which is motivated by the challenges described above. We address three aspects of stream processing: adaptivity in query processing, data compression, and the architecture choice for integrating real-time and historical data.

3.1 Adaptive query processing for data streams

Heterogeneous data streams are generated from different sources, at different rates, and include multiple domains. Our purpose is to build a distributed stream processing engine capable of adapting to changing conditions while serving complex continuous queries. Some sources already generate Linked Data streams [3]. Otherwise, we provide a layer serving an ontology-based access to non RDF data stream sources. Adapters for various input formats, such as CSV or REST APIs, are used to convert heterogeneous streams to RDF.

To reach an efficient query processing over data streams (section 2, C.1) we will focus on query execution planning. Traditional databases include a query optimizer that designs an execution plan based on the registered query and data statistics. In a distributed stream processing environment, there are several aspects to contemplate: changing rates of the input data, failure of processing nodes, and distribution of workload, among others. Adaptive Query Processing (AQP) techniques [7] allow adjusting the query execution plan to varying conditions of the data input, the incoming queries, and the system. Additionally, it is used to correct query optimizer mistakes and cope with unknown statistics [2].

First, we will analyze different strategies to process query operators, such as JOIN or FILTER. There are various examples in the literature that use ordering approaches for a more scalable stream processing, e.g. in the implementation of JOIN operators [12] or in eviction strategies [11]. We will design Storm⁸ topologies to efficiently process a set of common operators based on parallelizable tasks. Storm topologies are formed by spouts and bolts. Spouts are stream sources, whereas bolts are stream processors. Trident is a high-level abstraction on top of Storm that allows for stateful stream processing.⁹ For instance, a windowed join can be implemented using stream snapshots, which in Trident are called *states*. Depending on criteria such as the selectivity of the join or the size of snapshots, the processing engine can decide on the more appropriate join operator in terms of efficiency. Then, we will define a list of queries for a specific use case and will extend the topologies to fit the queries. While new data is entering the system, a dedicated bolt will manage stream data statistics to reassign a different topology if data stream conditions vary. In Storm, the coordination between the master node, which controls the assignation of tasks to spouts and bolts, and the worker nodes is managed by Zookeeper.¹⁰ However, Zookeeper does not provide elastic load balancing off-the-shelf. The use of cloud services for this purpose, such as the one offered by Amazon EC2, will be addressed in the near future.

3.2 Compressing RDF streams

To date, universal compressors (e.g. gzip) and specific RDF compressors (e.g. HDT [9]) are commonly used to reduce RDF exchange costs and delays on the network. These approaches, though, consider a static view of RDF datasets,

⁸ <http://storm.incubator.apache.org/>

⁹ <https://storm.incubator.apache.org/documentation/Trident-tutorial.html>

¹⁰ <http://zookeeper.apache.org/>

disregarding the dynamic nature of RDF stream management. A recent work [10] points out the importance of efficient RDF stream compression and proposes an initial solution leveraging the inherent data redundancy of the RDF data streams. Based on this proposal, we are currently working on a compressed data structure specifically designed for the particularities of dynamic RDF streams [8]. In particular, we aim at providing a lightweight serialization of RDF streams which i) minimizes the data exchange among processing nodes (section 2, C.2) while ii) serving a small set of operators on the compressed data.

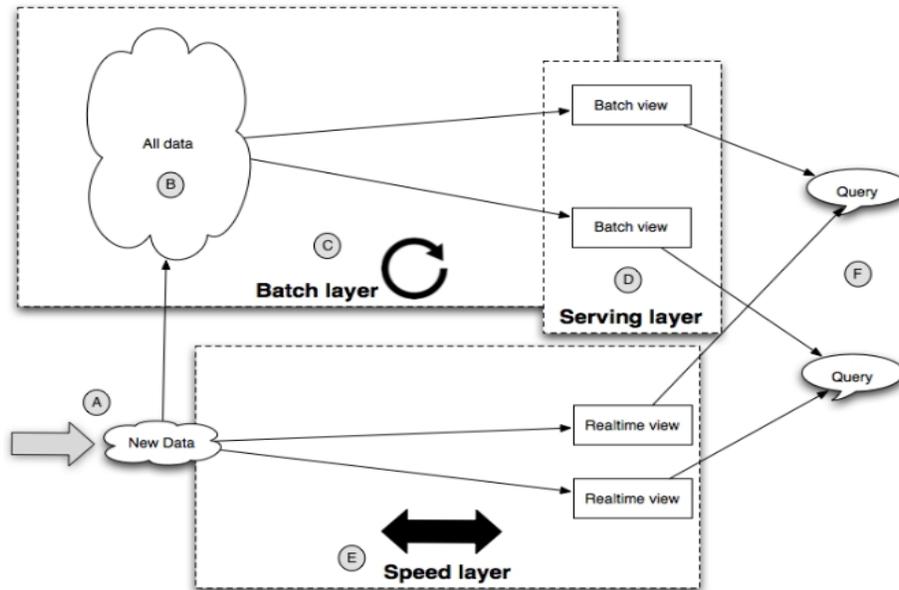


Fig. 1. Diagram of Lambda Architecture (from [16]). (A) Data input is consumed by the batch layer and the speed layer. The batch layer stores all the data raw (B) and continuously computes query functions to create batch views (C). The serving layer indexes batch views to access them quickly (D). The speed layer only processes recent data to produce real-time views (E). Batch and real-time views are integrated to respond queries (F).

3.3 Architecture overview

Our engine will address real-time processing on the Web of Things context following the Lambda principles [16]. Lambda is a 3-layer architecture designed to fit the requirements of Big Data: a batch layer (using Hadoop¹¹) stores all the incoming data in an immutable master dataset and pre-computes batch views on

¹¹ <http://hadoop.apache.org/>

historic data (section 2, C.3); a serving layer (NoSQL database) indexes views on the master dataset; and a speed layer manages the real-time processing issues and requests data views depending on incoming queries. The cloud platform to deploy our engine should provide computing services that scale on demand, as well as elastic load balancing.

4 Open Questions

Next steps on the short term will address the implementation of an scalable RDF stream processing system. Open questions that we would like to discuss at the workshop are:

- How does the order of tuple arrival affect the parallelization of join processing tasks?
- Are the spatial (or spatio-temporal) properties of a tuple a dimension to have into account for ordering? In this case, what influence does it have on reasoning tasks? And on parallelization tasks?
- How does the out-of-order tuples affect the processing of streams? In case of discarding them, how to communicate this decision in the results?

Acknowledgements

The research leading to these results has received funding from the European Union’s Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 257641, PlanetData network of excellence, and from Ministerio de Economía y Competitividad (Spain) under the project ”4V: Volumen, Velocidad, Variedad y Validez en la Gesti3n Innovadora de Datos” (TIN2013-46238-C4-2-R).

References

1. Anicic, D., Fodor, P., Rudolph, S., Stojanovic, N.: EP-SPARQL: A Unified Language for Event Processing and Stream Reasoning. In: Proceedings of the 20th International Conference on World Wide Web. pp. 635–644. WWW ’11, ACM, New York, NY, USA (2011)
2. Babu, S., Bizarro, P.: Adaptive Query Processing in the Looking Glass. In: Proceedings of the Second Biennial Conference on Innovative Data Systems Research (CIDR), Jan. 2005 (2005)
3. Balduini, M., Della Valle, E., Dell’Aglio, D., Tsytsarau, M., Palpanas, T., Confalonieri, C.: Social Listening of City Scale Events Using the Streaming Linked Data Framework. In: Alani, H., Kagal, L., Fokoue, A., Groth, P., Biemann, C., Parreira, J.X., Aroyo, L., Noy, N., Welty, C., Janowicz, K. (eds.) The Semantic Web ISWC 2013, pp. 1–16. Springer Berlin Heidelberg (2013)
4. Barbieri, D.F., Braga, D., Ceri, S., Valle, E.D., Grossniklaus, M.: C-sparql: a continuous query language for rdf data streams. *Int. J. Semantic Computing* 4(1), 3–25 (2010)

5. Calbimonte, J.P., Jeung, H., Corcho, O., Aberer, K.: Enabling Query Technologies for the Semantic Sensor Web. *International Journal on Semantic Web and Information Systems* 8(1), 43–63 (2012)
6. Compton, M., Barnaghi, P., Bermudez, L., Garcia-Castro, R., Corcho, O., Cox, S., Graybeal, J., Hauswirth, M., Henson, C., Herzog, A., Huang, V., Janowicz, K., Kelsey, W.D., Phuoc, D.L., Lefort, L., Leggieri, M., Neuhaus, H., Nikolov, A., Page, K., Passant, A., Sheth, A., Taylor, K.: The ssn ontology of the w3c semantic sensor network incubator group. *Web Semantics: Science, Services and Agents on the World Wide Web* 17(0) (2012)
7. Deshpande, A., Ives, Z., Raman, V.: Adaptive Query Processing. *Foundations and Trends in Databases* 1(1), 1–140 (Jan 2007)
8. Fernández, J.D., Llaves, A., Corcho, O.: Efficient RDF Interchange (ERI) Format for RDF Data Streams (accepted). In: *ISWC 2014* (2014)
9. Fernández, J.D., Martínez-Prieto, M.A., Gutiérrez, C., Polleres, A., Arias, M.: Binary RDF representation for publication and exchange (HDT). *Web Semantics: Science, Services and Agents on the World Wide Web* 19(0), 22–41 (2013)
10. Fernández García, N., Arias Fisteus, J., Sánchez Fernández, L., Fuentes-Lorenzo, D., Corcho, O.: RDSZ : An approach for lossless RDF stream compression. In: *In 11th European Semantic Web Conference (ESWC 2014)*. pp. 1–16. Crete, Greece (2014)
11. Gao, S., Scharrenbach, T., Bernstein, A.: The CLOCK Data-Aware Eviction Approach: Towards Processing Linked Data Streams with Limited Resources. In: *The 11th Extended Semantic Web Conference. Lecture Notes in Computer Science*, Springer (May 2014)
12. Golab, L., Özsu, M.T.: Processing Sliding Window Multi-joins in Continuous Queries over Data Streams. In: *Proceedings of the 29th International Conference on Very Large Data Bases - Volume 29*. pp. 500–511. VLDB '03, VLDB Endowment (2003)
13. Le-phuoc, D., Nguyen, H., Quoc, M., Van, C.L., Hauswirth, M.: Elastic and Scalable Processing of Linked Stream Data in the Cloud. In: Alani, H., Kagal, L., Fokoue, A., Groth, P., Biemann, C., Parreira, J.X., Aroyo, L., Noy, N., Welty, C., Janowicz, K. (eds.) *The Semantic Web ISWC 2013*, vol. 287305, pp. 280–297. Springer Berlin Heidelberg (2013)
14. Le-phuoc, D., Parreira, J.X., Hauswirth, M.: Linked Stream Data Processing. In: *Reasoning Web. Semantic Technologies for Advanced Query Answering*, pp. 245–289. Springer Berlin Heidelberg (2012)
15. Margara, A., Urbani, J., van Harmelen, F., Bal, H.: Streaming the web: Reasoning over dynamic data. *Web Semantics: Science, Services and Agents on the World Wide Web* 0(0) (0)
16. Marz, N., Warren, J.: *Big Data: principles and best practices of scalable realtime data systems*. Manning Publications (2013)
17. Rinne, M., Nuutila, E., Seppo, T.: INSTANS: High-Performance Event Processing with Standard RDF and SPARQL. In: Henson, C., Taylor, K., Corcho, O. (eds.) *Proceedings of the 5th International Workshop on Semantic Sensor Networks*. pp. 81–96. CEUR-WS, Boston, Massachusetts, USA (2012)
18. Sequeda, J.F., Corcho, O.: Linked Stream Data: A Position Paper. In: *Proceedings of the 2nd International Workshop on Semantic Sensor Networks, SSN 09*. CEUR-WS, Washington, USA (2009)
19. Sheth, A., Henson, C., Sahoo, S.S.: Semantic Sensor Web. *IEEE Internet Computing* 12(4), 78–83 (2008)