

RDF-based Deployment Pipelining for Efficient Dataset Release Management

Claus Stadler¹, Lisa Wenige¹, Sebastian Tramp², Kurt Junghanns¹, and Michael Martin¹

¹ Institute for Applied Informatics (InfAI), Goedelerring 9, Leipzig, Germany, D-04109

E-mail: (stadler, wenige, junghanns, martin)@infai.org

² eccenca GmbH, Hainstrasse 8, Leipzig, Germany, D-04109

E-mail: info@eccenca.com

Abstract. Open Data portals often struggle to provide release features (i.e., stable versioning, up-to-date download links, rich metadata descriptions) for their datasets. By this means, wide adoption of publicly available data collections is hindered, since consuming applications cannot access fresh data sources or might break due to data quality issues. While there exists a variety of tools to efficiently control release processes in software development, the management of dataset releases is not as clear. This paper proposes a deployment pipeline for efficient dataset releases that is based on automated enrichment of DCAT/DataID metadata and is a first step towards efficient deployment pipelining for Open Data publishing.

Keywords: Deployment, Open Data, DCAT, DataID, Data Quality

1 Introduction

With the advent of the Open Data movement, a multitude of datasets have been made available on public repositories.^{3,4} Therefore, researchers have developed methodologies to manage the data publishing process efficiently [9, 13]. There also exist software tools that (semi-)automatically assist data publishers during data conversion [2, 15], quality assurance [5, 16], linking [11, 14], metadata enrichment and data provision [8, 12, 13]. However, the integration of these separate publishing phases has not been entirely addressed by academia. Instead, the combination of the required tasks is often handled by individuals who need to manually execute hand-crafted transformations to make up for the missing technical links between the existing data processing units. This hinders frequent releases of fresh datasets as well as the reuse of data publishing pipelines across different application domains. Since quality assurance, timeliness and discoverability are among the most important preconditions for stakeholders and applications to actually consume the data [9], these issues should have priority when

³ <http://lodstats.aksw.org/>

⁴ <https://datahub.io>

publishing data collections. Hence, efficient software tools need to be in place to speed up publishing workflows and minimize manual adaptations thereby reducing error-proneness. In this context, the provision of automatic tools for bridging the phases of data conversion/preprocessing and the final dataset release is particularly relevant. Here, methods from the field of software engineering, such as *versioning, stable download links, automated testing and metadata enrichment* can serve as guiding best practices. However, the process of data publication also exhibits some specifics that need to be considered when developing deployment tools. In this paper, we present a methodology and a prototypical implementation of a deployment pipeline for automated release management of data collections based on RDF. We use a strict definition of *dataset* as the foundation for data management: A dataset is an instance of a data model. Hence, any procedures that yield (syntactic) representations of that instance, such as a CSV-to-RDF mapping and its materialization, are considered as different distribution forms of the *same* dataset.

2 Related Work

The research on efficient and integrated data release management is still in its infancy [8]. For instance, although the DataGraft application by Roman et al. enables data conversion and publication [13], the tool does not provide features to handle different dataset versions. On the other hand, while Rojas Meléndez et al. formulate the requirement of versioning for their Linked Data publication platform, the approach does not consider DevOps best practices [12]. In contrast, in software development, strategies of continuous integration (CI) are widespread. During CI processing, code repositories are frequently checked for updates and troubleshooting operations are triggered in case errors are found. Thus, it is ensured that released software fulfills qualitative requirements and that other applications can make use of the most recent updates [7]. However, current CI methods can not be transferred to data publishing 1:1. Proposals for adaptations have been made by Cirulli et al., who suggest to maintain the code for data conversion on a Jenkins server [4]. Meissner et al. present an integration pipeline that gets triggered upon new repository commits [10]. But these approaches are missing a clear versioning concept on the data level as well as methods to automatically add metadata to dataset releases. The issue of enhancing metadata is addressed by Frey et al., who present a novel strategy of publishing large datasets with their FlexiFusion approach [6]. However, a holistic method for deployment pipelining that considers all relevant aspects of the dataset release process has not yet been presented.

3 Use Case - The LIMBO Project

The German Federal Ministry of Transport and Digital Infrastructure is providing the mCLOUD Open Data portal, where currently more than 1070 datasets

on roads, rails, air traffic and waterways can be found. Additionally, the platform contains collections of space, climate and weather data. The datasets are published by the ministry and associated publicly funded agencies. By providing an Open Data portal, the ministry intends to harness research and development projects working on novel navigational services, smart travel and route planning as well as applications for highly precise weather forecasting.⁵ However, the mCLOUD data is provided in heterogeneous formats, not semantically described and poorly integrated with other data collections. Hence, the research project *Linked Data Services for Mobility (LIMBO)* has been started to convert selected datasets from the mCLOUD portal to RDF and link them to other datasets on the Linked Open Data cloud.⁶ As the original data is subject to changes, conversion processes are re-run frequently thus making efficient deployment and release pipelines necessary.

4 Dataset Deployment Pipeline

Repository management services, such as Github, Gitlab or Bitbucket provide numerous useful features, which can be configured for dataset publication. With the successful *Git branching model*⁷ there also exists a best practice for organizing release workflows. In a nutshell, only releases should be committed to the *master* branch. Currently, the LIMBO datasets including the deployment build tools can be accessed on the public gitlab instance.⁸ New (versions of) datasets are published by running Makefile build commands in a Git repository containing the raw RDF files. The deployment pipeline can be invoked in a broadly predefined, but configurable order of execution. The build process triggers *data transformation*, *RDFUnit-based testing*⁹ and *metadata enrichment* tasks (Fig. 1). During the test state, the data is checked for both schema alignment and general data quality. The corresponding error reports are automatically generated by RDFUnit. Thus, third-party consumers can assess the quality of datasets. During *metadata enrichment*, the deployment process assembles a local DCAT model of the dataset project and allows modifications to be specified with SPARQL update statements utilizing the tool *Sparql-Integrate*.¹⁰ This tool is a thin command-line wrapper for Jena's ARQ SPARQL engine¹¹ which registers several extensions, such as for passing environment variables from shell scripts to SPARQL queries. Upon release, first a *release branch* is created. Local file references in the assembled DCAT model are converted to public download URLs, yielding a DCAT record ready for publication in a dataset catalog with Maven-style group, artifact and version identifiers as proposed by [3] and [6]

⁵ <https://www.mcloud.de/>

⁶ <https://www.limbo-project.org/>

⁷ <https://nvie.com/posts/a-successful-git-branching-model/>

⁸ <https://gitlab.com/limbo-project/dataset-includes>

⁹ <https://github.com/AKSW/RDFUnit>

¹⁰ <https://github.com/SmartDataAnalytics/SparqlIntegrate>

¹¹ <http://jena.apache.org/>

for the DataID model. One of the most important benefits of our approach is that links to prior releases of any datasets contained in the Git repository can automatically be made part of a release record by inspecting whether the master branch already contains a prior one. If so, *owl:priorVersion* links can be created between all dataset identifiers with matching attributes, such as group and artifact id. Upon finalization, the release branch is merged into master and subsequently deleted, and a Git tag is created. In typical cases, this process only requires invoking the make goals *release-start* and *release-finish*. By this means, applications that rely on stable data can determine current releases from the metadata catalog.¹²

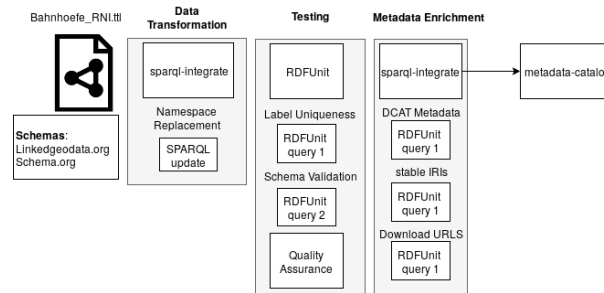


Fig. 1: Data and metadata transformations using Sparql-Integrate

5 Conclusion

In this paper, we have presented a Git-based methodology for publishing validated stable releases of datasets and corresponding DCAT records. Stable releases are beneficial for applications in many ways, such as when it comes to predictability and reproducibility of results, testing, and caching of derived information for performance-reasons. Up to date, we published a couple of datasets of the mCLOUD with our approach. As this approach does not depend on hosting new custom services, maintenance overhead is minimized. The adoption of Maven’s and DataID’s artifact identification scheme promises future interoperability with other existing infrastructures for data asset management, such as the Quit diff tool [1], which provides functions to determine changes between dataset versions.

Acknowledgement

This work has been supported by the Federal Ministry of Transport and Digital Infrastructure (BMVI) for the LIMBO project under the grant numbers 19F2029A and 19F2029G.

¹² <https://gitlab.com/limbo-project/metadata-catalog>

References

1. Arndt, N. and Radtke, N. (2016). Quit diff: Calculating the delta between rdf datasets under version control. In Proceedings of the 12th International Conference on Semantic Systems (pp. 185-188). ACM.
2. Auer, S., Dietzold, S., Lehmann, J., Hellmann, S. and Aumueller, D. (2009). Triplify: light-weight linked data publication from relational databases. In Proceedings of the 18th international conference on World wide web (pp. 621-630). ACM.
3. Brümmer, M., Baron, C., Ermilov, I., Freudenberg, M., Kontokostas, D. and Hellmann, S. (2014). DataID: Towards semantically rich metadata for complex datasets. In Proceedings of the 10th International Conference on Semantic Systems (pp. 84-91). ACM.
4. Cirulli, S. (2015). Continuous integration for XML and RDF Data. XML LONDON, 52-60.
5. Dimou, A., Kontokostas, D., Freudenberg, M., Verborgh, R., Lehmann, J., Manens, E., Hellmann, S. and Van de Walle, R. (2015). Assessing and Refining Mappings to RDF to Improve Dataset Quality. In International Semantic Web Conference (pp. 133-149). Springer, Cham.
6. Frey, J., Hofer, M., Hellmann, S. and Obraczka, D. DBpedia FlexiFusion Best of Wikipedia } Wikidata } Your Data.
7. Fowler, M. and Foemmel, M. (2006). Continuous integration. Thought-Works. <https://www.martinfowler.com/articles/continuousIntegration.html>.
8. Klímek, J., Skoda, P. and Necaský, M. (2016). Requirements on Linked Data Consumption Platform. In LDOW at WWW.
9. Kucera, J., Chlapek, D., Klímek, J. and Necaský, M. (2015). Methodologies and Best Practices for Open Data Publication. DATESO.
10. Meissner, R. and Junghanns, K. (2016). Using devOps principles to continuously monitor RDF data quality. In Proceedings of the 12th International Conference on Semantic Systems (pp. 189-192). ACM.
11. Ngomo, A. C. N. and Auer, S. (2011). LIMES a time-efficient approach for large-scale link discovery on the web of data. In Twenty-Second International Joint Conference on Artificial Intelligence.
12. Rojas Meléndez, J. A., Van de Vyvere, B., Gevaert, A., Taelman, R., Colpaert, P. and Verborgh, R. (2018). A Preliminary Open Data Publishing Strategy for Live Data in Flanders. In WWW2018, the International World Wide Web Conference (pp. 1847-1853). ACM Press.
13. Roman, D., Dimitrov, M., Nikolov, N., Putlier, A., Sukhobok, D., Elvester, B. and Petkov, Y. (2016). Datagraft: Simplifying open data publishing. In European Semantic Web Conference (pp. 101-106). Springer, Cham.
14. Schultz, A., Matteini, A., Isele, R., Bizer, C. and Becker, C. (2011). LDIF-linked data integration framework. In Proceedings of the Second International Conference on Consuming Linked Data-Volume 782 (pp. 125-130). CEUR-WS.org.
15. Unbehauen, J., Stadler, C. and Auer, S. (2012). Accessing relational data on the web with sparqlmap. In Joint International Semantic Technology Conference (pp. 65-80). Springer, Berlin, Heidelberg.
16. Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J. and Auer, S. (2016). Quality assessment for linked data: A survey. Semantic Web, 7(1), 63-93.