

Linked Web APIs Dataset

Web APIs meet Linked Data

Editor(s): Name Surname, University, Country

Solicited review(s): Name Surname, University, Country

Open review(s): Name Surname, University, Country

Milan Dojchinovski * and Tomas Vitvar

Web Intelligence Research Group, Czech Technical University in Prague, Czech Republic

E-mail: {milan.dojchinovski,tomas.vitvar}@fit.cvut.cz

Abstract.

Web APIs enjoy significant increase in popularity and usage in the last decade. They **have become** the core technology for exposing functionalities and data. Nevertheless, due to the lack of semantic Web API descriptions their discovery, sharing, integration, and assessment of their quality and consumption is limited. In this paper, we present the Linked Web APIs dataset, an RDF dataset with semantic descriptions about Web APIs. It provides semantic descriptions for 11,339 Web APIs, 7,415 mashups and 7,717 developers profiles, which makes it the largest available dataset from the Web APIs domain. It captures the provenance, temporal, technical, functional, and non-functional aspects. We describe the Linked Web APIs Ontology, a minimal model which **build** on top of several well-known ontologies. The dataset has been interlinked and published according to the Linked Data principles. We describe several possible usage scenarios for the dataset and show its potential.

Keywords: Web APIs, Linked Data, Web services, Linked Web APIs, ontology

1. Introduction

The Web APIs have become the first-class citizen on the Web and the core functionality of any Web application. They primarily target and benefit developers and Web API providers **in terms of entry barrier**, data and functionality integration, and reusability of tools. Back in late 2008, ProgrammableWeb.com¹, the largest Web APIs and mashup directory, reported only 1,000 Web APIs, while 5,000 APIs in Feb 2014 and over 13,000 APIs in June 2015. The benefits of having these Web API descriptions provided as Linked Data are several. The Web API descriptions are contextualized, they can be referenced, re-used and combined. The Web APIs data is linked so API consumers can effectively discover new Web APIs. Least but not last,

on one side the developers can benefit from a sophisticated queries for discovery and selection of APIs of interest, on the other hand, the Web API providers can execute queries to get better insight and analysis of the Web APIs ecosystem.

To achieve these goals, we have developed the *Linked Web APIs* dataset. It provides information about Web APIs, mashups which **utilize** Web APIs in compositions, and mashups developers. The primary source for the dataset is ProgrammableWeb.com directory which acts as central repository for Web APIs descriptions. The dataset re-uses several well-known ontologies developed by the Semantic Web community. In order to conform to the Linked Data principles² we have also linked the dataset with four central LOD

*Corresponding author. E-mail: milan.dojchinovski@fit.cvut.cz, tel: +420 776 519502

¹<http://www.programmableweb.com>

²<http://www.w3.org/DesignIssues/LinkedData.html>

dataset: DBpedia³, Freebase⁴, LinkedGeoData⁵ and GeoNames⁶.

The reminder of this paper is structured as follows. We first, in Section 2 explain the source of information and how the data was collected. Section 3 describes the ontology developed for modelling relevant Web APIs information. Created Linked Web APIs dataset and its technical details are described in Section 4. The approach for interlinking the dataset with other LOD datasets is described in Section 5. Section 6 presents selected use cases to illustrate the potential of the dataset. Finally, Section 7 concludes the paper and provides future directions.

2. The Data Source

In our work, we have considered ProgrammableWeb as a primary source of information for creating the dataset. It adopts characteristics of a social Web platform where Web API providers can publish and share information about offered Web APIs and consequently increase the visibility of them. The API directory also allows developers to search and find appropriate APIs for their projects or see and learn from showcases of existing mashup applications.

Each web page describing Web API, mashup or a mashup developer was parsed and valuable information was extracted. An example of such Web page describing a Web API is the one describing the Twitter API⁷. For each Web API we extracted its title, short summary describing its functionalities, tags and categories assigned, technical information such as supported formats and protocols, as well as non-functional properties such as its homepage, usage limits, usage fees, security, etc. Similarly, for each mashup we extracted its title, short free-text description of its functionalities, assigned tags, and the homepage of the mashup. From each page describing a developer we extracted its username, homepage and short bio about the developer. Also, the city and country of residence, its given and family name and the gender were extracted, if these information were available as public information.

We also captured the *relationships* between the Web APIs, mashups and developers. In other words, for each mashup we extracted the list of Web APIs which were used by the mashup and also the information about the list of mashups created by each developer. The dataset also captures the *temporal aspects* - the creation time of the Web APIs, mashups and the time a user registered his profile.

3. The Ontology

The Linked Web APIs ontology⁸ is a **minimal model** that captures the most relevant information related to Web APIs and mashups. The ontology builds on top of existing and well established ontologies and appropriately extends them. The selection of appropriate ontologies for integration was driven by the following four crucial requirements:

- *Provenance:* It is important to keep information about **Who** (developers) created *What* (mashups) and *How* (using which APIs). What APIs a provider provides also needs to be captured.
- *Functional and Non-functional Properties:* What functionalities a Web API or mashup offers is more than important, as well as their usage limits and fees, supported security or authentication mechanisms.
- *Technical Properties:* Information about the supported protocols and formats and the Web APIs endpoint location is as important, as it allows a Web API consumer to search only for APIs with preferred technical capabilities.
- *Temporal Information:* When a mashup or Web API was created provides valuable information. For example, Web API analysts to analyze the recent trends in the API ecosystem, or Web API consumers to discover most recent Web APIs or mashups.

Figure 1 shows the overall Linked Web APIs ontology. The ontology contains three central classes: *iso:WebAPI* - to describe Web APIs, *iso:Mashup* - to describe mashup compositions which utilizes one or more Web APIs, and *iso:Agent* - to represents all kinds of entities involved in creation and/or consumption of Web APIs and mashups.

In order to capture the provenance information, the Linked Web APIs ontology integrates the PROV-

³<http://dbpedia.org/>

⁴<https://www.freebase.com/>

⁵<http://linkedgeodata.org/>

⁶<http://www.geonames.org/>

⁷<http://www.programmableweb.com/api/twitter>

⁸<http://linked-web-apis.fit.cvut.cz/ns/core/index.html>

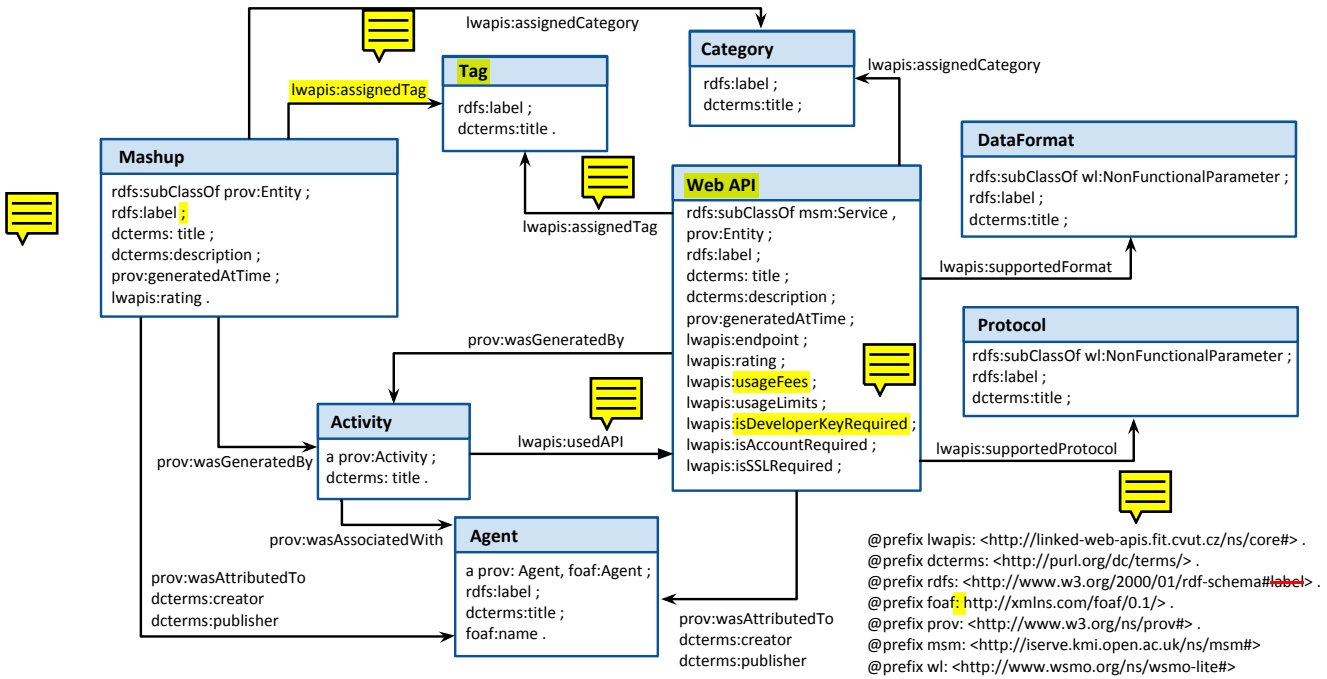


Fig. 1. The Linked Web APIs Ontology.

O ontology⁹ by incorporating its classes *prov:Entity*, *prov:Activity* and *prov:Agent*, and their related properties. The *prov:Entity* class serve as super-class of *lso:WebAPI* and *lso:Mashup* classes. Activities convey information about the process of consumption of Web APIs and generation of mashups by the agents. We introduce the *lso:usedAPI* property which refines the semantics of the *prov:used* property so it can be used to explicitly identify usage of a Web API in a mashup creation. The temporal information about the time of creation of a mashup or Web API is expressed using the *prov:generatedAtTime* property.

For the functional (tags and categories) and non-functional (formats and protocols) properties of the Web APIs and mashups we introduce new classes in our namespace. The ontology also integrates the *wl:NonFunctionalParameter* class from the WSMO-lite ontology¹⁰ [7], developed by the semantic Web services community, to explicitly identify non-functional properties. The Minimal Service Model (MSM)¹¹ ontology, initially defined for the hRESTS microformat [5] is also considered and the class *msm:Service* is integrated as super-class of the *lso:WebAPI*. This allows

to attach additional Web API information, such as operations, inputs and outputs, which is relevant for execution of Web APIs.

General metadata information such as Web API and mashup title, or their short textual description we describe using the Dublin Core vocabulary¹².

4. The Linked Web APIs Dataset

The Linked Web APIs is the largest and the first Linked Data dataset with Web API descriptions. It provides descriptions for 11,339 Web APIs, 7,415 mashups and 7,717 mashup creators and it contains over 550K RDF triples. For all the resources we mint URIs in our own namespace (<http://linked-web-apis.fit.cvut.cz/resource/{name}>). The name part from the URIs is a normalized form of the label of the resource, which is lowercased and each space is replaced with underscore sign. Further, since two different resources can have same name (e.g., the label XML can occur as a tag and also as a format) to each minted URI we attach its type as suffix to the URI. For example, *_api* for Web API URIs or *_tag* for tags. An example of a URI minted for

⁹<http://www.w3.org/TR/prov-o/>

¹⁰<http://www.wsmo.org/ns/wsmo-lite/>

¹¹<http://iserve.kmi.open.ac.uk/ns/msm#>

¹²<http://purl.org/dc/terms/>

the Google Maps API is `http://linked-web-apis.fit.cvut.cz/resource/google-maps_api`. Similar approach is employed by Wikipedia to distinguish between pages which have same title. For example, `/wiki/Food_(band)` for a page describing the musical band “Food” and `/wiki/Food_(film)` for a page describing the movie with the same name.

All URIs are dereferenceable and served according to the Linked Data principles in RDF/XML and Turtle format. The dataset is also available through a Virtuoso SPARQL endpoint and also as a dump. The landing page for the dataset is `http://linked-web-apis.fit.cvut.cz/` and it provides information about the latest news, releases and changes. Technical details about the dataset are listed in Table 1.

Table 1
Details of the Linked Web APIs dataset.

Name	Linked Web APIs dataset
URL	<code>http://linked-web-apis.fit.cvut.cz/</code>
Endpoint	<code>http://linked-web-apis.fit.cvut.cz/sparql</code>
Ontology	<code>http://linked-web-apis.fit.cvut.cz/ns/core#</code>
Version	0.1
Ver. Date	05.08.2015
Datahub	Linked Web APIs dataset

5. Dataset Linking

In order to assure maximal reusability and integrability, we linked the dataset with four central LOD datasets. Two multi-domain datasets, DBpedia and Freebase, and two geographical datasets GeoNames and LinkedGeoData. From the information we linked the Web APIs supported data formats, supported protocols, developers’ city and country of residence. Since GeoNames and LinkedGeoData are geographical datasets, only users’ city and country of residence was linked to those datasets. DBpedia and Freebase are multi-domain datasets and therefore we linked all information to these datasets. The links to DBpedia, and respectively to Freebase, were generated following the most-frequent-sense based approach used as entity linking method in EntityClassifier.eu NER system [2]. The linking to LinkedGeoData was governed by the intuition that the names of the cities and countries in our dataset have same names in the LinkedGeoData dataset. The approach was supported by a SPARQL query which retrieves resources with a given label. Following this linking methodology we generated 1,440 links out of which 722 are DBpedia links, 299 Free-

base links, 326 GeoNames links and 93 LinkedGeoData links. Table 2 provides more information about the linking.

Table 2
Number of linked resource per type and dataset.

	DBpedia	Freebase	LGD	GeoNames
Formats	283	208	/	/
Protocols	123	91	/	/
Cities	263	/	47	276
Countries	53	/	46	50
Total	722	299	93	326

6. Use Cases

The availability of a dataset with Web APIs descriptions in RDF can support various use cases, including, but not limited within, personalised Web API provisioning, API ecosystem analysis, and automated Web API descriptions processing contexts. In this section, we describe selected use cases and existing applications of the Linked Web APIs dataset.

6.1. Personalised Recommendations

The Linked Web APIs dataset contains links between the mashups and the developers, which is pertinent source of information for developing Web API recommendation methods. A simple scenario is when a user has already picked a Web API for his/her mashup and searches for other compatible Web APIs. Such scenario can be supported with the SPARQL query from Listing 1 which returns the top 5 most used Web APIs.

A developer can further customize the query to fit his/her needs, for example, to narrow down the results only to Web APIs which support particular data format (e.g., JSON) or APIs from a specific category (e.g., social, government, etc.).

```

1 SELECT ?api (COUNT(?api) as ?count)
2 WHERE {
3   ?mashup prov:wasGeneratedBy ?activity.
4   ?activity lso:usedAPI ls:google-maps_api .
5   ?activity lso:usedAPI ?api .
6   FILTER (!strends(str(?api), "google-maps_api"))
7 }
8 ORDER BY DESC(?count)
9 LIMIT 5

```

Listing 1: Top 5 most used Web APIs with Google Maps API.

In the context of personalised recommendations, the dataset has been recently employed in several works around **personalised recommendation of Web APIs** [3] and Linked Data resources [4]. Both works focus on developing graph based algorithms on top of the Linked Web APIs dataset and utilizing the provisioning information (who developed what), functional properties (tags and categories) and temporal information (when a mashup or Web API was developed).

6.2. Temporal Analysis

The dataset also captures the temporal aspect, i.e., the time when a mashup or a Web API was developed. Such information can help Web APIs providers to get better insights about the recent developments and study the consumption of a Web API, or the whole Web API ecosystem over time. The benefits from having temporal information can be illustrated with the SPARQL query from Listing 2.

```
1 SELECT COUNT(?mashup) as ?count
2 WHERE {
3   ?mashup prov:wasGeneratedBy ?activity.
4   ?activity lso:usedAPI ls:google-maps-api .
5   ?mashup prov:generatedAtTime ?date .
6   FILTER (?date >= "2013-01-01"^^xsd:dateTime
7     && ?date < "2014-01-01"^^xsd:dateTime)
8 }
```

Listing 2: Number of mashups **utilizing the Google Maps API in 2013.**

The SPARQL query in the listing gives information about the total number of mashups which utilized the Google Maps API in 2013. Figure 2 visualizes the results from such analysis for three popular APIs and their utilization over time.

The Web API provider might be interested in what kind of mashups their API was used. An answer to

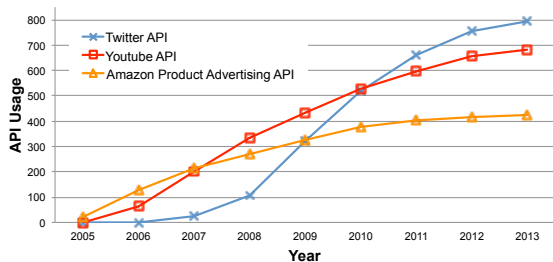


Fig. 2. Web API utilization over time.

such question can be answered with the SPARQL query in Listing 3.

```
1 SELECT ?category (COUNT(?category) as ?count)
2 WHERE {
3   ?mashup prov:wasGeneratedBy ?activity.
4   ?activity lso:usedAPI ls:google-maps-api .
5   ?mashup prov:generatedAtTime ?date .
6   ?mashup lso:assignedTag ?category .
7   FILTER (?date >= "2013-01-01"^^xsd:dateTime
8     && ?date < "2014-01-01"^^xsd:dateTime)
9 } ORDER BY DESC(?count)
```

Listing 3: The number of mashup categories the Google Maps API was used in 2013.

Further, a Web API analyst might be interested in the latest trends in the API ecosystem. Questions such as “*What protocols and formats are the most supported by the APIs?*” or “*Which domains provided most APIs in 2013?*” are likely to occur. Using the SPARQL query in Listing 4 we can get the top 5 most popular protocols in year 2013, which is also illustrated in Figure 3 for the two most used protocols REST and SOAP for a period of ten years.

```
1 SELECT ?protocol (COUNT(?api) as ?count)
2 WHERE {
3   ?api rdf:type lso:WebAPI .
4   ?api prov:generatedAtTime ?date .
5   ?api lso:supportedProtocol ?protocol .
6   FILTER (?date >= "2013-01-01"^^xsd:dateTime
7     && ?date < "2014-01-01"^^xsd:dateTime)
8 } ORDER BY DESC(?count)
9 LIMIT 5
```

Listing 4: The most popular API protocols in 2013.

An answer to the question “*Which domains provided most APIs in the 2013?*” can be answered with the SPARQL query in Listing 5. The results show that

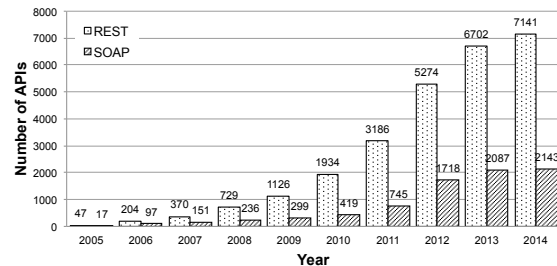


Fig. 3. Popularity of REST and SOAP protocols over time.

the most popular API category is “tools”, followed by the “science”, “internet”, “enterprise” and “financial” categories. It is interesting the fact that the “financial” and “enterprise” categories are among the top five most popular API categories, which indicates that APIs are already understood as relevant technology also by other domains than the internet and social networks domain.

```

1 SELECT ?category (COUNT(?api) as ?count)
2 WHERE {
3   ?api rdf:type lso:WebAPI .
4   ?api prov:generatedAtTime ?date .
5   ?api lso:assignedCategory ?category .
6   FILTER (?date > "2012-01-01"^^xsd:dateTime
7     && ?date < "2013-01-01"^^xsd:dateTime)
8 } ORDER BY DESC(?count)
9 LIMIT 5

```

Listing 5: The most popular API protocols in 2013.

A more in-depth analysis using the Linked Web APIs dataset has been conducted in [6]. In particular, the dataset has been used as a reference dataset for time-aware link prediction in RDF graphs.

7. Conclusion and Future Work

A growing number of available Web APIs requires new mechanisms to support the process of sharing, discovery, integration and re-use of Web APIs at large scale. In this paper, we have presented the Linked Web APIs dataset, the first and largest Linked Data dataset providing Web APIs descriptions. The dataset supports i) *API consumers*-in the process of discovery, selection and use of Web APIs, ii) *API providers*-in increasing the visibility and tracking the popularity of their Web APIs, and iii) *API analysts*-in analyzing the API ecosystem. The dataset will also help to raise the awareness about the importance of providing semantic Web API descriptions and publishing them as Linked Data. The dataset has been validated in several recent works in the context of personalized recommendations and link analysis. Also, on a set of usage scenarios we have shown the potential of the dataset.

In our future work, we want to enrich the dataset with Web API descriptions from other data sources. We also plan to integrate ontologies such as the SPARQL Service Description¹³ ontology and the

DataID¹⁴ dataset description model [1] which will in turn allow description of SPARQL processing services and corresponding Linked Data datasets. We also plan to enrich the dataset with user profiles from traditional social networks. We want to interlink the tags and categories information with relevant datasets from the LOD cloud such as the Wikidata¹⁵, Wiktionary¹⁶ and Dbnary¹⁷. Last but not least we want to explore other applications using the dataset and assess its potential.

Acknowledgement. We thank ProgrammableWeb.com for supporting this research.

References

- [1] M. Brümmer, et al. DataID: Towards Semantically Rich Metadata for Complex Datasets. In *Proceedings of the 10th International Conference on Semantic Systems, SEM '14*, pp. 84–91. ACM, New York, NY, USA, 2014.
- [2] M. Dojchinovski and T. Kliegr. Entityclassifier.eu: Real-time Classification of Entities in Text with Wikipedia. In H. Blockeel, K. Kersting, S. Nijssen, and F. Železný, (eds.) *Machine Learning and Knowledge Discovery in Databases*, vol. 8190 of *Lecture Notes in Computer Science*, pp. 654–658. Springer Berlin Heidelberg, 2013.
- [3] M. Dojchinovski, J. Kuchar, T. Vitvar, and M. Zaremba. Personalised Graph-Based Selection of Web APIs. In P. Cudré-Mauroux, et al., (eds.) *The Semantic Web ISWC 2012*, vol. 7649 of *Lecture Notes in Computer Science*, pp. 34–48. Springer Berlin Heidelberg, 2012.
- [4] M. Dojchinovski and T. Vitvar. Personalised Access to Linked Data. In K. Janowicz, S. Schlobach, P. Lambrix, and E. Hyvönen, (eds.) *Knowledge Engineering and Knowledge Management*, vol. 8876 of *Lecture Notes in Computer Science*, pp. 121–136. Springer International Publishing, 2014.
- [5] J. Kopecký, K. Gomadam, and T. Vitvar. hRESTS: An HTML Microformat for Describing RESTful Web Services. In *Web Intelligence and Intelligent Agent Technology, 2008. WI-IAT '08. IEEE/WIC/ACM International Conference on*, vol. 1, pp. 619–625. Dec 2008.
- [6] J. Kuchar, M. Dojchinovski, and T. Vitvar. Time-aware Link Prediction in RDF Graphs. In *Knowledge Engineering and Knowledge Management*, vol. 8876 of *Lecture Notes in Computer Science*. Springer International Publishing, 2014.
- [7] T. Vitvar, J. Kopecký, J. Viskova, and D. Fensel. WSMO-Lite Annotations for Web Services. In S. Bechhofer, M. Hauswirth, J. Hoffmann, and M. Koubarakis, (eds.) *The Semantic Web: Research and Applications*, vol. 5021 of *Lecture Notes in Computer Science*, pp. 674–689. Springer Berlin Heidelberg, 2008.

¹⁴<http://wiki.dbpedia.org/projects/dbpedia-dataid-unit>

¹⁵<https://www.wikidata.org>

¹⁶<https://www.wiktionary.org/>

¹⁷<http://kaiko.getalp.org/about-dbnary/>

¹³<http://www.w3.org/TR/sparql11-service-description/>