

Creating Knowledge Graphs for Geographic Data on the Web

ELENA DEMIDOVA, Data Science and Intelligent Systems Group (DSIS), University of Bonn, Germany

ALISHIBA DSOUZA, Data Science and Intelligent Systems Group (DSIS), University of Bonn, Germany

SIMON GOTTSCHALK, L3S Research Center, Leibniz University Hannover, Germany

NICOLAS TEMPELMEIER, L3S Research Center, Leibniz University Hannover, Germany

RAN YU, Data Science and Intelligent Systems Group (DSIS), University of Bonn, Germany

Abstract. Geographic data plays an essential role in various Web, Semantic Web and machine learning applications. OpenStreetMap and knowledge graphs are critical complementary sources of geographic data on the Web. However, data veracity, the lack of integration of geographic and semantic characteristics, and incomplete representations substantially limit the data utility. Verification, enrichment and semantic representation are essential for making geographic data accessible for the Semantic Web and machine learning. This article describes recent approaches we developed to tackle these challenges.

1 INTRODUCTION

Geographic data plays an essential role in a range of real-world applications on the Web, including machine learning models estimating travel time or charging demand for electric vehicles, recommending points of interest and predicting traffic accidents (e.g., [2], [9]). Such applications rely on rich representations of a variety of geographic entities including monuments, roads and charging stations.

Geographic Web Information Sources and Knowledge Graphs are major complementary Web sources providing information regarding geographic entities, their spatio-temporal context, characteristics, and relationships.

- *Geographic Web Information Sources* such as OpenStreetMap (OSM¹) provide characteristics of geographic entities and their relationships. Today, OSM is an essential source of free and open geographic Web information created by voluntary effort, containing over 6.8 billion entities from 188 countries².
- *Knowledge Graphs* such as Wikidata, DBpedia and EventKG [5] contain real-world entities (e.g., persons and places), events and their relationships in a graph-based format. Semantic interpretation of these facts is facilitated through ontologies.

These data representation paradigms have different focuses: On the one hand, knowledge graphs provide rich semantic information about real-world entities, facilitating querying, exploration, and reasoning. On the other hand, OSM provides rich geographic information, i.e., fine-grained coordinates of real-world locations, but does not possess a clear schema, thus lacking direct semantic interpretation.

Table 1 illustrates an example geographic entity (“Zugspitze”, a mountain in Germany) and its different representations in OSM and Wikidata³. While OSM provides the information in the form of heterogeneous key-value pairs, so-called

©Elena Demidova, Alishiba Dsouza, Simon Gottschalk, Nicolas Tempelmeier, Ran Yu, 2022. This is the author’s version of the work. It is posted here for your personal use. Not for redistribution. The definitive version was published in the ACM SIGWEB Newsletter, Issue Winter 2022 Article No.: 4 pp 1–8 <https://doi.org/10.1145/3522598.3522602>.

¹OpenStreetMap, OSM and the OpenStreetMap magnifying glass logo are trademarks of the OpenStreetMap Foundation, and are used with their permission. We are not endorsed by or affiliated with the OpenStreetMap Foundation.

²<https://osmstats.neis-one.org>

³wd and wtd are the prefixes of <http://www.wikidata.org/entity/> and <http://www.wikidata.org/prop/direct/>, respectively.

“tags”, e.g., `natural=peak`, entities in Wikidata are represented by Uniform Resource Identifier (URIs) and are connected to other entities via well-defined properties.

Table 1. Representations of the German mountain Zugspitze in OpenStreetMap and Wikidata. `wd:Q3375` identifies the Zugspitze in Wikidata.

Key	Value	Subject	Predicate	Object
name	Zugspitze	wd:Q3375	<code>rdfs:label</code> (<i>label</i>)	Zugspitze
natural	peak	wd:Q3375	<code>wdt:P279</code> (<i>instance of</i>)	mountain
summit:cross	yes	wd:Q3375	<code>wdt:P2044</code> (<i>elevation</i>)	2,962.06 metre
ele	2962	wd:Q3375	<code>wdt:P3137</code> (<i>parent peak</i>)	wd:Q15127 (<i>Finsteraarhorn</i>)

(a) OpenStreetMap representation

(b) Wikidata representation

Knowledge graphs that include data from geographic Web information sources have a significant potential to make geographic data more accessible for the Semantic Web and machine learning applications. However, several significant challenges have to be addressed, when it comes to creating such knowledge graphs. In particular, this includes:

- (C1) Data quality verification: Geographic Web Information Sources, including OSM, often come as volunteered geographic information (VGI), i.e., volunteers contribute to the data collection and modeling. While volunteered contributions increase worldwide coverage of geographic entities, the openness of these sources potentially leads to data quality issues, such as vandalism and misinformation. Therefore, data quality verification is an essential prerequisite for using VGI in knowledge graphs.
- (C2) Semantic enrichment on schema and data level: Alignments of geographic entities and their classes between geographic Web information sources and knowledge graphs are rare. Consequently, there is a need to establish entity and schema links across sources. The varying coverage and heterogeneous representations of geographic entities, attributes, and relationships make such enrichment particularly challenging.
- (C3) Creating meaningful representations of geographic entities for semantic and machine learning applications: The combined utilization of Geographic Web Information Sources and knowledge graphs for machine learning requires representations of geographic entities that reflect their semantic and spatial extent.

Recently, we developed several methods to tackle these challenges, making geographic entities available for downstream applications. An overview of our methods is shown in Fig. 1. We consider existing knowledge graphs and OSM as input data sources. To enable data quality verification (C1), we propose a supervised method for *vandalism detection* in OSM that can cope with crowdsourced geographic Web data. To tackle (C2), we establish links between geographic entities in OSM and a knowledge graph through *geographic entity linking*. Based on interlinked geographic entities, *schema alignment* aims at aligning their classes (e.g., “mountain” in Wikidata and “peak” in OSM). Finally, we make the resulting semantic geographic entities available in two forms: (i) as a WorldKG knowledge graph through *geographic knowledge graph creation*, and (ii) through *semantic geographic entity representation learning*, to represent geographic entities for the usage in machine learning applications (C3).

In the remainder of this article, we provide more details of these methods and the datasets built upon them. First, we present our methodology: Section 2 discusses the *OVID* approach for OSM vandalism detection, Section 3 presents *OSM2KG* for linking Wikidata and OSM entities, and Section 4 introduces the *NCA* approach for schema alignment between Wikidata, DBpedia and OSM. Then, we describe the creation of two resources: *WorldKG* presented in Section 5 is

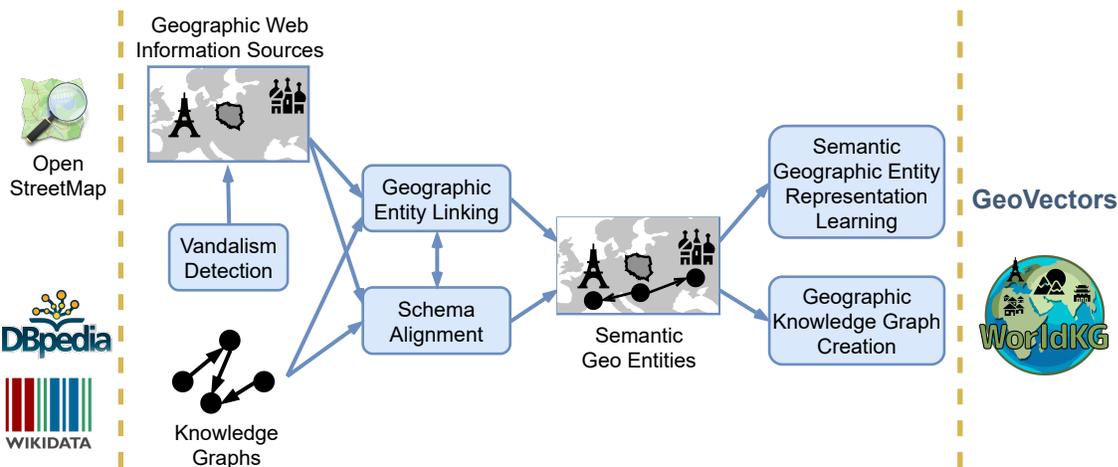


Fig. 1. Overview of the resources and methods for creating knowledge graphs for geographic data on the Web discussed in this article.

our geographic knowledge graph, and *GeoVectors* presented in Section 6 is our resource of geographic entity embeddings covering their spatial and semantic aspects. Finally, in Section 7, we provide a conclusion and discuss future work.

2 VANDALISM DETECTION – OVID

Ensuring a high quality of openly available geographic information is vital for adopting this information in real-world applications. In particular, protecting crowdsourced geographic data from vandalism is a challenging task. Whereas this task has been addressed in knowledge bases such as Wikipedia and Wikidata, existing vandalism detection approaches typically rely on textual features and do not consider the spatial dimension of the data.

In [11], we developed novel methods to automatically detect vandalism on OpenStreetMap. We proposed the OVID (OpenStreetMap Vandalism Detection) model. OVID combines user, edit, and changeset features with an attention-based neural architecture in a supervised classification model. Furthermore, we systematically extracted vandalism incidents from the OSM edit history from 2014 to 2019 and provided a dataset containing over 9,000 real-world vandalism examples. We made the dataset available under an open license⁴. Our evaluation results confirmed that the OVID model outperformed the existing models for vandalism detection.

3 GEOGRAPHIC ENTITY LINKING – OSM2KG

According to the five-star open data principle proposed by Tim Berners-Lee, the interlinking across open datasets is an essential data quality indicator. However, such entity links are rarely available in geographic Web information sources. For instance, as of January 2021, only 0.05% of all OSM entities provided a link to the Wikidata knowledge graph.

To systematically increase the number of links between OSM and knowledge graph entities, we developed the *OSM2KG* approach [10]. Given an OSM entity, *OSM2KG* aims to determine a geographic entity in a knowledge graph representing the same real-world entity. In contrast to a knowledge graph, OSM does not enforce a strict schema and represents its entities through key-value pairs (see Table 1a). This representation leads to sparse attributes and imposes significant challenges to feature extraction for traditional link discovery. *OSM2KG* is a supervised machine

⁴<https://github.com/NicolasTe/Ovid>

learning approach for linking OSM and knowledge graphs at the entity level. The core of OSM2KG is a novel key-value embedding for OSM entities. OSM2KG applies an unsupervised representation learning algorithm that captures OSM entity semantics. The resulting embedding represents OSM entities in a classification model for link discovery. Our evaluation conducted on OpenStreetMap and the Wikidata and DBpedia knowledge graphs demonstrates that OSM2KG reliably outperforms existing link discovery approaches. We make our code available under an open license⁵.

4 SCHEMA ALIGNMENT – NCA

Although OpenstreetMap contains numerous geographic entities, these entities are not directly accessible to semantic applications due to their key-value-based structure, as discussed in Section 3. Aligning schema elements such as classes and properties of knowledge graphs to OSM keys and key-value pairs (tags) can benefit semantic applications, providing direct access to geographic entities on a world scale. Existing approaches for schema alignment are typically based on structural similarity, or string similarity [7, 8]. Due to its flat schema and heterogeneous keys, such methods cannot be applied to OSM. To tackle these challenges, we introduced a two-step approach called NCA (Neural Class Alignment) for the tag-to-class alignment between OSM and knowledge graphs [3].

In the first step, NCA utilizes existing links between OSM and knowledge graphs as a supervision signal to classify OSM entities into the semantic classes of knowledge graphs. As a result of the classification, a shared latent space capturing the semantics of tag-to-class alignment is created. NCA systematically probes the classification model with specific tags in the second step to obtain the corresponding semantic classes. We made our code available under an open license⁶.

We evaluated the NCA approach on six country-specific datasets of different countries using manually annotated tag-to-class pairs. We observed that NCA outperforms state-of-the-art approaches. As the result of this alignment, we also observed an over 400% increase in the semantic class annotations for OSM data. NCA mapping of the schema elements is an essential step towards the semantic representation of OSM entities. This schema mapping enables us to enrich OSM entities with semantic class information using knowledge graphs. We utilize the tag-to-class matches obtained by NCA for the construction of the WorldKG knowledge graph, a new geographic knowledge graph built starting from OSM, presented in Section 5.

5 GEOGRAPHIC KNOWLEDGE GRAPH CREATION – WORLDCG

Knowledge graphs are vital sources of semantic information regarding real-world entities and their relations. However, the coverage of geographic entities in popular knowledge graphs is relatively poor. As discussed in Section 3, on the one hand, popular general-purpose knowledge graphs such as DBpedia and Wikidata only provide a limited number of geographic entities. On the other hand, specialized geographic knowledge graphs such as LinkedGeoData [1] and YAGO2geo [6] contain only a tiny subset of geographic classes. To bridge this gap, we develop a new geographic knowledge graph *WorldKG* [4], aiming to provide better coverage of semantic representations of geographic entities through a fusion of knowledge graphs and OpenStreetMap.

The WorldKG knowledge graph⁷ currently contains over 100 million geographic entities from 188 countries. This knowledge graph is based on a novel WorldKG ontology having over 1000 geographic classes created using the OSM schema. We convert the flat OSM schema into a hierarchical ontology structure using the schema alignment NCA

⁵<https://github.com/NicolasTe/osm2kg>

⁶<https://github.com/alishiba14/NCA-OSM-to-KGs>

⁷<https://www.worldkg.org/>

approach presented in Section 4. The results of this alignment can determine the correct classes of WorldKG entities in Wikidata and DBpedia ontology with over 99% accuracy, on average.

To enable direct semantic access to WorldKG, we provide a SPARQL endpoint and downloadable data files in standard RDF turtle format. We believe that the scale and accuracy of WorldKG can help various semantic data-driven applications. Examples include event-centric and geospatial question answering and geographic information retrieval.

6 SEMANTIC GEOGRAPHIC ENTITY REPRESENTATION LEARNING – GEOVECTORS

Usable representations of geographic entities are of utmost importance for various applications, including travel time estimation, location recommendation, and geographic information retrieval. Dependent on the specific application, such representations need to capture the spatial and semantic entity dimensions. The spatial dimension reflects the geographic extent and proximity to neighboring entities. The semantic dimension can include an entity type and further type-dependent attributes. In OSM, where geographic entities are represented using key-value pairs reflecting various entity types and their heterogeneous attributes, discrete vector-based entity representations are extremely sparse and high-dimensional. The computation of geographic entity representations that machine learning algorithms can effectively use is not trivial.

In our *GeoVectors* corpus [12], we provide a ready to use corpus containing embeddings of over 980 million geographic entities extracted from OpenStreetMap. *GeoVectors* corpus consists of two parts. The *GeoVectors-location* embeddings capture the spatial dimension of the geographic entities. To this extent, we adopted an established geographic representation learning algorithm and developed a framework that enables the algorithm training on a world scale. The *GeoVectors-tags* embeddings capture the semantic dimension of geographic entities. To obtain these embeddings, we adopted a pre-trained word embedding model to map OSM key-value pairs into a latent space and obtain the semantic representation of geographic entities. Our experiments conducted on real-world data demonstrate that the *GeoVectors* embeddings effectively capture the spatial and semantic dimensions of geographic entities.

We offer direct access to the *GeoVectors* corpus on our website⁸. On this website, we provide a detailed corpus description. Moreover, we provide a knowledge graph that integrates *GeoVectors* with popular knowledge graphs such as Wikidata and DBpedia, making our embeddings of geographic entities in these sources directly accessible. We make our knowledge graph accessible via a SPARQL endpoint. Finally, we provide a search function that enables name-based entity access to the *GeoVectors* knowledge graph.

7 CONCLUSION AND FUTURE WORK

In this article, we described three significant challenges towards the creation of knowledge graphs for geographic data on the Web, including data quality verification, semantic enrichment, and the creation of effective representations. We tackled these challenges with several new approaches, including OVID for vandalism detection, OSM2KG for geographic entity linking, and NCA for schema alignment. These approaches build the basis for creating the WorldKG knowledge graph for geographic data and the *GeoVectors* corpus containing distributional representations of geographic entities.

In future work, we aim to provide an increasingly complete and semantically rich representation of geographic data on the Web. We will further work on bringing these methods closer together to enable cross-fertilization. Furthermore, WorldKG and *GeoVectors* can benefit from interlinking and enrichment with additional datasets. We make the datasets openly available to enable their broad reuse.

⁸<https://geovectors.l3s.uni-hannover.de>

ACKNOWLEDGMENTS

This work was partially funded by DFG, German Research Foundation (“WorldKG”, 424985896), the Federal Ministry for Economic Affairs and Climate Action (BMWK), Germany (“d-E-mand”, 01ME19009B and “CampaNeo”, 01MD19007B), the Federal Ministry of Education and Research (BMBF), Germany (“Simple-ML”, 01IS18054), and the European Commission (EU H2020, “smashHit”, grant-ID 871477).

REFERENCES

- [1] Sören Auer, Jens Lehmann, and Sebastian Hellmann. 2009. LinkedGeoData: Adding a Spatial Dimension to the Web of Data. In *Proceedings of The 8th International Semantic Web Conference, ISWC 2009 (Lecture Notes in Computer Science, Vol. 5823)*. Springer, 731–746. https://doi.org/10.1007/978-3-642-04930-9_46
- [2] Rajjat Dadwal, Thorben Funke, and Elena Demidova. 2021. An Adaptive Clustering Approach for Accident Prediction. In *Proceedings of The 24th IEEE International Intelligent Transportation Systems Conference, ITSC 2021*. IEEE, 1405–1411. <https://doi.org/10.1109/ITSC48978.2021.9564564>
- [3] Alishiba Dsouza, Nicolas Tempelmeier, and Elena Demidova. 2021. Towards Neural Schema Alignment for OpenStreetMap and Knowledge Graphs. In *Proceeding of The 20th International Semantic Web Conference, ISWC 2021 (Lecture Notes in Computer Science, Vol. 12922)*. Springer, 56–73. https://doi.org/10.1007/978-3-030-88361-4_4
- [4] Alishiba Dsouza, Nicolas Tempelmeier, Ran Yu, Simon Gottschalk, and Elena Demidova. 2021. WorldKG: A World-Scale Geographic Knowledge Graph. In *Proceeding of The 30th ACM International Conference on Information and Knowledge Management, CIKM '21*. ACM, 4475–4484. <https://doi.org/10.1145/3459637.3482023>
- [5] Simon Gottschalk and Elena Demidova. 2019. EventKG - the Hub of Event Knowledge on the Web - and Biographical Timeline Generation. *Semantic Web* 10, 6 (2019), 1039–1070. <https://doi.org/10.3233/SW-190355>
- [6] Nikolaos Karalis, Georgios M. Mandilaras, and Manolis Koubarakis. 2019. Extending the YAGO2 Knowledge Graph with Precise Geospatial Knowledge. In *Proceedings of the 18th International Semantic Web Conference, ISWC 2019 (Lecture Notes in Computer Science, Vol. 11779)*. Springer, 181–197. https://doi.org/10.1007/978-3-030-30796-7_12
- [7] Jayant Madhavan, Philip A. Bernstein, and Erhard Rahm. 2001. Generic Schema Matching with Cupid. In *VLDB 2001, Proceedings of 27th International Conference on Very Large Data Bases, September 11-14, 2001, Roma, Italy*. Morgan Kaufmann, 49–58. <http://www.vldb.org/conf/2001/P049.pdf>
- [8] DuyHoa Ngo, Zohra Bellahsene, and Konstantin Todorov. 2013. Opening the Black Box of Ontology Matching. In *Proceedings of the ESWC 2013 (Lecture Notes in Computer Science, Vol. 7882)*. Springer, 16–30. https://doi.org/10.1007/978-3-642-38288-8_2
- [9] Ashutosh Sao, Nicolas Tempelmeier, and Elena Demidova. 2021. Deep Information Fusion for Electric Vehicle Charging Station Occupancy Forecasting. In *Proceedings of The 24th IEEE International Intelligent Transportation Systems Conference, ITSC 2021*. IEEE, 3328–3333. <https://doi.org/10.1109/ITSC48978.2021.9565097>
- [10] Nicolas Tempelmeier and Elena Demidova. 2021. Linking OpenStreetMap with Knowledge Graphs - Link Discovery for Schema-agnostic Volunteered Geographic Information. *Future Gener. Comput. Syst.* 116 (2021), 349–364. <https://doi.org/10.1016/j.future.2020.11.003>
- [11] Nicolas Tempelmeier and Elena Demidova. 2022. Attention-Based Vandalism Detection in OpenStreetMap. In *Proceedings of The Web Conference, WWW, 2022*. ACM. <https://doi.org/10.1145/3485447.3512224>
- [12] Nicolas Tempelmeier, Simon Gottschalk, and Elena Demidova. 2021. GeoVectors: A Linked Open Corpus of OpenStreetMap Embeddings on World Scale. In *Proceedings of The 30th ACM International Conference on Information and Knowledge Management, CIKM '21*. ACM, 4604–4612. <https://doi.org/10.1145/3459637.3482004>