

Linking OpenStreetMap with Knowledge Graphs - Link Discovery for Schema-Agnostic Volunteered Geographic Information

Nicolas Tempelmeier¹, Elena Demidova^{1,2}

¹*L3S Research Center, Leibniz Universität Hannover, Germany*

²*Data Science & Intelligent Systems Group (DSIS), University of Bonn, Germany*

Abstract

Representations of geographic entities captured in popular knowledge graphs such as Wikidata and DBpedia are often incomplete. OpenStreetMap (OSM) is a rich source of openly available, volunteered geographic information that has a high potential to complement these representations. However, identity links between the knowledge graph entities and OSM nodes are still rare. The problem of link discovery in these settings is particularly challenging due to the lack of a strict schema and heterogeneity of the user-defined node representations in OSM. In this article, we propose OSM2KG - a novel link discovery approach to predict identity links between OSM nodes and geographic entities in a knowledge graph. The core of the OSM2KG approach is a novel latent, compact representation of OSM nodes that captures semantic node similarity in an embedding. OSM2KG adopts this latent representation to train a supervised model for link prediction and utilises existing links between OSM and knowledge graphs for training. Our experiments conducted on several OSM datasets, as well as the Wikidata and DBpedia knowledge graphs, demonstrate that OSM2KG can reliably discover identity links. OSM2KG achieves an F1 score of 92.05% on Wikidata and of 94.17% on DBpedia on average, which corresponds to a 21.82 percentage points increase in F1 score on Wikidata compared to the best performing baselines.

1. Introduction

OpenStreetMap¹ (OSM) has recently evolved as the key source of openly accessible volunteered geographic information (VGI) for many parts of the world, building a backbone for a wide range of real-world applications on the Web and beyond [1]. Prominent examples of OSM applications include mobility and transportation services such as route planners [2], public transportation information sites² and Global Positioning System (GPS) tracking³, as

well as geographic information services⁴ and spatial data mining.

The OSM data is produced by a large number of contributors (approx. 5.6 million in August 2019⁵) and lacks a pre-defined ontology. The description of geographic entities in OSM (so-called “OSM nodes”) includes few mandatory properties such as an identifier and a location as well as a set of user-defined key-value pairs (so-called “tags”). As a result, the representations of OSM nodes are extremely heterogeneous. The tags provided for the individual OSM nodes vary highly [3].

Knowledge graphs (KGs), i.e., graph-based knowledge bases [4], including Wikidata [5], DBpedia [6], YAGO2 [7]

*©2020. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Email addresses: tempelmeier@l3s.de (Nicolas Tempelmeier¹), elena.demidova@cs.uni-bonn.de (Elena Demidova^{1,2})

¹OpenStreetMap is a trademark of the OpenStreetMap Foundation, and is used with their permission. We are not endorsed by or affiliated with the OpenStreetMap Foundation.

²<http://www.öpnv-karte.de>

³<https://gitlab.com/eneiluj/phonetrack-oc>

⁴<https://histosm.org/>

⁵Statistics from https://www.openstreetmap.org/stats/data_stats.html

and EventKG [8] are a rich source of semantic information for geographic entities, including for example cities and points of interest (POIs). This information, typically represented according to the RDF data model, has a high and so far, mostly unexploited potential for semantic enrichment of OSM nodes. An interlinking of OSM nodes and geographic entities in knowledge graphs can bring semantic, spatial, and contextual information to its full advantage and facilitate, e.g., spatial question answering [9] and semantic trip recommendation [10].

Interlinking of OSM and knowledge graphs has recently attracted interest in the Wikidata⁶ and OSM⁷ communities. Our analysis results, presented in Section 2, illustrate that the coverage of the existing interlinking between the OSM nodes and Wikidata entities varies significantly across entity types and geographic regions. For example, in a recent OSM snapshot of Germany (referred to as OSM-DE), cities are linked more often (73%) than less popular entities like mountains (5%). For another example, there are 42% more linked OSM nodes in the OSM snapshot of Germany than in that of Italy (OSM-IT). In practice, the interlinking of OSM nodes with semantic reference sources such as Wikidata or DBpedia is typically conducted manually by volunteers (and sometimes companies, see, e.g., [11]).

The problem of OSM link discovery is particularly challenging due to the heterogeneity of the OSM node representations. Other factors affecting the effectiveness of OSM node disambiguation in the context of link discovery include place name ambiguity and limited context [12]. Furthermore, geographic coordinates in the VGI sources such as OSM often represent the points of community consensus rather than being determined by objective criteria [13] and can thus vary significantly across sources. For example, an average geographic distance between the coordinates of the corresponding entities in Germany in the

OSM and Wikidata datasets is 2517 meters. This example illustrates that geographic coordinates alone are insufficient to effectively discover identity links between the corresponding entities in VGI sources.

Although research efforts such as the LinkedGeoData project [13] and Yago2Geo [14] have been conducted to lift selected parts of OSM data in the Semantic Web infrastructure to facilitate link discovery, these efforts typically rely on manually defined schema mappings. Maintenance of such mappings does not appear feasible or sustainable, given the large scale, and openness of the OSM schema. Therefore, link discovery approaches that can address the inherent heterogeneity of OSM datasets are required.

In this article, we propose the novel OSM2KG link discovery approach to establish identity links between the OSM nodes and equivalent geographic entities in a knowledge graph. OSM2KG addresses OSM’s heterogeneity problem through a novel latent representation of OSM nodes inspired by the word embedding architectures [15]. Whereas embeddings have recently gained popularity in several domains, their adoption to volunteered geographic information in OSM is mostly unexplored. In contrast to state-of-the-art approaches to link discovery in OSM (such as [14, 13]), OSM2KG does not require any schema mappings between OSM and the reference knowledge graph.

The core of the OSM2KG approach is a novel latent representation of OSM nodes that captures semantic node similarity in an embedding. OSM2KG learns this latent, compact node representation automatically from OSM tags. To the best of our knowledge OSM2KG is the first approach to address the heterogeneity of the OSM data by a novel embedding representation. This embedding representation is created in an unsupervised fashion and is task-independent. The embedding systematically exploits the co-occurrence patterns of the OSM’s key-value pairs to capture their semantic similarity. Building upon this embedding, along with spatial and semantic information in the target knowledge graph, OSM2KG builds a su-

⁶<https://www.wikidata.org/wiki/Wikidata:OpenStreetMap>

⁷https://wiki.openstreetmap.org/wiki/Proposed_features/Wikidata

perervised machine learning model to predict missing identity links. To train the proposed link prediction model, we exploit publicly available community-created links between OSM, Wikidata, and DBpedia as training data.

The key contribution of our work is the novel OSM2KG link discovery approach to infer missing identity links between OSM nodes and geographic entities in knowledge graphs, including:

- A novel unsupervised embedding approach to infer latent, compact representations that capture semantic similarity of heterogeneous OSM nodes.
- A supervised classification model to effectively predict identity links, trained using the proposed latent node representation, selected knowledge graph features, and existing links.
- We describe an algorithm for link discovery in the OSM datasets that uses the proposed supervised model and the latent representation to effectively identify missing links.

The results of the extensive experimental evaluation on three real-world OSM datasets for different geographic regions, along with the Wikidata and DBpedia knowledge graphs, confirm the effectiveness of the proposed OSM2KG link discovery approach. According to our evaluation results, OSM2KG can reliably predict links.

OSM2KG achieves an F1 score of 92.05% on Wikidata and of 94.17% on DBpedia on average, which corresponds to a 21.82 percentage points increase in F1 score on Wikidata compared to the best performing baselines.

The remainder of the article is organised as follows. In Section 2, we discuss the representation of geographic information in OSM and Wikidata and the existing inter-linking between these sources to motivate our approach. Then in Section 3, we formally introduce the link discovery problem addressed in this article. In Section 4, we present the proposed OSM2KG approach. Following that, we describe the evaluation setup in Section 5 and provide and

discuss our evaluation results in Section 6. Then in Section 7, we discuss related work. Finally, in Section 8, we provide a conclusion.

2. Motivation

Volunteered geographic information is a special case of user-generated content that represents information about geographic entities [16]. VGI is typically collected from non-expert users via interactive Web applications, with the OpenStreetMap project⁸ being one of the most prominent and successful examples. OSM is a rich source of spatial information available under an open license (Open Database License) and created collaboratively through an international community effort. Today OSM data has become available at an unprecedentedly large scale. While in 2006 OSM captured only 14.7 million GPS points, this number has increased to 7.4 billion by 2019. Similarly the number of users who contribute to OSM has grown from 852 in 2006 to 5.6 million in 2019⁹.

OSM includes information on *nodes* (i.e., points representing geographic entities such as touristic sights or mountain peaks), as well as *lines* (e.g. lists of points) and their topological *relations*. The description of nodes in OSM consists of few mandatory properties such as the node identifier and the location (provided as geographic coordinates) and an optional set of tags. *Tags* provide information about nodes in the form of key-value pairs. For instance, the tag “`place=city`” is used to express that a node represents a city. OSM does not provide a fixed taxonomy of keys or range restrictions for the values but encourages its users to follow a set of best practices¹⁰. For example, the node labels are often available under the “`name`” key, whereas the labels in different languages can be specified using the “`name:code=`” convention¹¹. The

⁸<https://www.openstreetmap.org>

⁹<https://blackadder.dev.openstreetmap.org/OSMStats/>.

¹⁰https://wiki.openstreetmap.org/wiki/Any_tags_you_like

¹¹https://wiki.openstreetmap.org/wiki/Multilingual_names

Table 1: Number of nodes, tags and distinct keys in the country-specific OSM snapshots (OSM-[country]) and their respective subsets linked to Wikidata (Wikidata-OSM-[country]).

	France			Germany			Italy		
	OSM-FR	Wikidata-OSM-FR	Ratio	OSM-DE	Wikidata-OSM-DE	Ratio	OSM-IT	Wikidata-OSM-IT	Ratio
No. Nodes	390,586,064	21,629	0.01%	289,725,624	24,312	< 0.01%	171,576,748	18,473	0.01%
No. Nodes with Name	1,229,869	20,507	1.67%	1,681,481	23,979	1.43%	557,189	18,420	3.31%
No. Tags	27,398,192	199,437	0.73%	37,485,549	212,727	0.56%	18,850,692	122,248	0.65%
No. Distinct Keys	6,009	1,212	20.17%	12,392	1,700	13.72%	4,349	892	20.51%

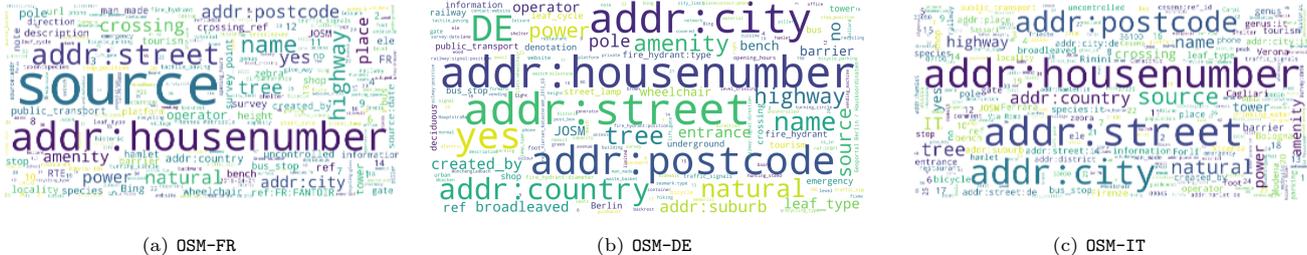


Figure 1: Tag clouds generated from the 1000 most frequent tags in each respective OSM dataset.

tags can also be used to specify identity links across datasets, e.g., to link OSM nodes to the equivalent entities in a knowledge graph.

For example, the link between the OSM node representing the city of Berlin and its Wikidata counterpart is established via the tag “wikidata=Q64” assigned to the OSM node. Here, “Q64”¹² denotes the identifier of the corresponding Wikidata entity. Recent studies indicate that the level of details provided for the individual OSM nodes is very heterogeneous [3]. Contextual information, e.g., regarding the historical development of the city population, is typically not available in OSM. Furthermore, the individual keys and tags do not possess any machine-readable semantics, which further restricts their use in applications.

Country-specific OSM snapshots are publicly available¹³. In the following, we refer to the country-specific snapshots as of September 2018 as the OSM-[country] dataset. E.g.m the snapshot for Germany is referred to as “OSM-DE”. The linked sets Wikidata-OSM-FR, Wikidata-OSM-DE, and Wikidata-OSM-IT are the subsets of the OSM-[country] datasets obtained by extracting all nodes that link to Wikidata entities from the respective OSM snapshot. Table 1 provides an overview of the number of

nodes, nodes with name, tags, and distinct key contained in the OSM-[country] datasets and the respective linked sets Wikidata-OSM-[country]. As we can observe, only a small fraction of nodes, tags, and distinct keys from the overall datasets appear in the linked sets. Furthermore, nearly all nodes contained in one of the linked sets exhibit a name tag. In addition, in Figure 1, we illustrate the most frequent keys of the OSM-FR, OSM-DE, and OSM-IT datasets in a tag cloud visualisation.

Figure 2 depicts the mean and the standard deviation of the number of tags contained in the OSM-DE dataset for the four most common entity types in Wikidata-OSM-DE, such as cities, train stations, castles, and mountains. Note that, unlike a knowledge graph, OSM does not define the node type information explicitly. To generate the statistics presented in this section, we used the existing links between the OSM nodes and the Wikidata entities to manually identify the tags in OSM indicative for the particular entity types in Wikidata and collected the OSM nodes annotated with these tags. We observe that the number of tags varies significantly with the entity type. Moreover the standard deviation is relatively high (between 35% and 63%) for all entity types. While for some entity types (e.g., mountains) the variation in the absolute number of tags

¹²<https://www.wikidata.org/wiki/Q64>

¹³OSM snapshots can be found at <http://download.geofabrik.de>.

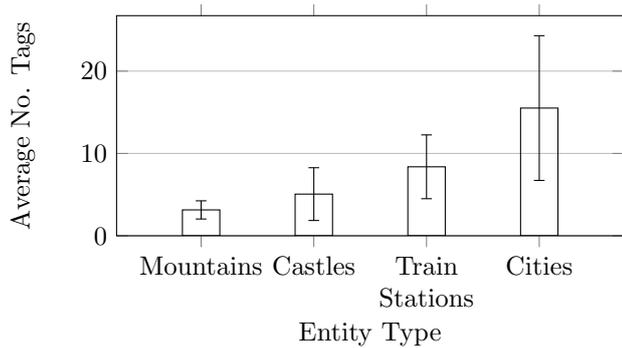


Figure 2: Average number of tags per entity type in Wikidata-OSM-DE. Error bars indicate the standard deviation.

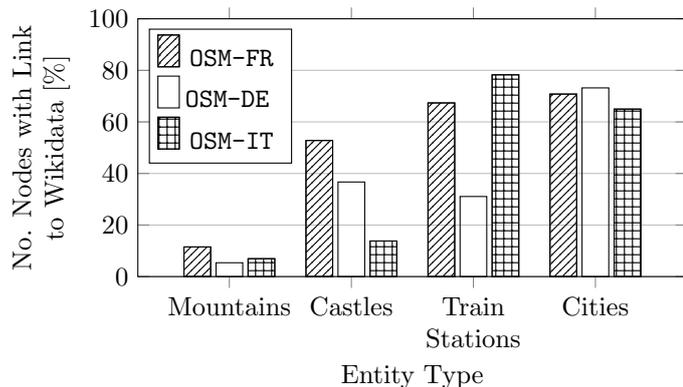


Figure 3: Percentage of frequent OSM node types with links to Wikidata entities within the OSM datasets for Germany (OSM-DE), France (OSM-FR), and Italy (OSM-IT) as of September 2018.

is rather small, other types (e.g., cities) exhibit more substantial variations, meaning that some of the cities possess more detailed annotations compared with the rest.

Knowledge graphs such as Wikidata [5], DBpedia [6], and YAGO [7] are a rich source of contextual information about geographic entities, with Wikidata currently being the largest openly available knowledge graph linked to OSM. In September 2018, Wikidata contained more than 6.4 million entities for which geographic coordinates are provided. Overall, the geographic information in OSM and contextual information regarding geographic entities in the existing knowledge graphs are highly complementary. As an immediate advantage of the existing effort to manually interlink OSM nodes and Wikidata entities, the names of the linked OSM nodes have become available in many languages [11].

The links between the OSM nodes and geographic en-

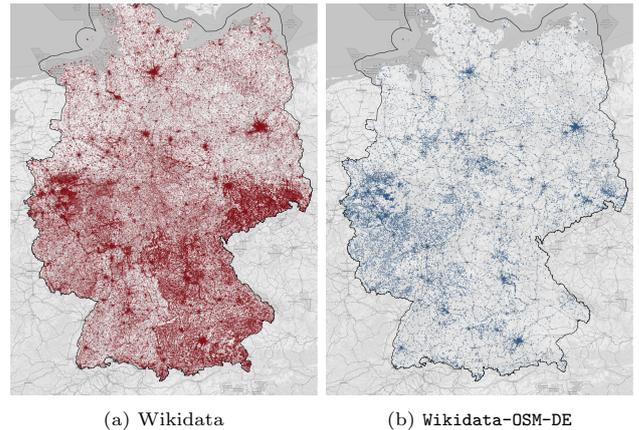


Figure 4: Wikidata geo-entities located within Germany and Wikidata geo-entities linked by OSM. Map image: ©OpenStreetMap contributors, ODbL.

tities in Wikidata are typically manually annotated by volunteers and community efforts and are still only rarely provided. Figure 3 illustrates the percentage of the four most frequent geographic entity types (i.e., cities, train stations, mountains, and castles) that link to Wikidata from the OSM datasets for Germany, France, and Italy. Here, entity types are obtained from Wikidata using existing links between the OSM nodes and Wikidata entities. As we can observe, the cities are linked most frequently, with a link coverage of approximately 70% for all datasets. The link coverage of the other entity types is significantly lower, with mountains having the smallest coverage across these four categories with approximately 5% in Germany. Figure 4 provides a visual comparison of the number of Wikidata entities located in Germany and the number of Wikidata entities to which links from OSM exist. While a significant fraction of links is still missing, existing links manually defined by volunteers reveal a high potential for being used as training data for supervised machine learning to increase link coverage automatically.

In summary, volunteered geographic information is a continually evolving large-scale source of heterogeneous spatial data, whereas knowledge graphs provide complementary, contextual information for geographic entities. The links between VGI and knowledge graphs are mainly manually specified and are still only rarely present in the

OSM datasets. The existing links represent a valuable source of training data for supervised machine learning methods to automatically increase the link coverage between OSM and knowledge graphs. This interlinking can provide a rich source of openly available semantic, spatial, and contextual information for geographic entities.

3. Problem Statement

In this work, we target the problem of identity link discovery between the nodes in a semi-structured geographic corpus such as OSM with equivalent entities in a knowledge graph.

Definition 1. Knowledge graph: Let E be a set of entities, R a set of labelled directed edges and L a set of literals. A *knowledge graph* $\mathcal{KG} = \langle E \cup L, R \rangle$ is a directed graph where entities in E represent real-world entities and the edges in $R \subseteq (E \times E) \cup (E \times L)$ represent entity relations or entity properties.

In this work, we focus on interlinking entities in a knowledge graph that possess geographic coordinates, i.e., longitude and latitude. We refer to such entities as *geo-entities*. Typical examples of geo-entities include cities, train stations, castles, and others.

Definition 2. Geo-entity: A *geo-entity* $e \in E$ is an entity for which a relation $r \in R$ exists that associates e with geographic coordinates, i.e., a longitude $lon \in L$ and a latitude $lat \in L$.

For instance, a geo-entity representing the city of Berlin may be represented as follows (the example illustrates an excerpt from the Wikidata representation of Berlin):

Entity	Property	Entity/Literal
Q64	<i>name</i>	<i>Berlin</i>
Q64	<i>instance of</i>	<i>Big City</i>
Q64	<i>coordinate location</i>	52°31'N, 13°23'E
Q64	<i>capital of</i>	<i>Germany</i>

We denote the subset of nodes representing geo-entities in the knowledge graph \mathcal{KG} as $E_{geo} \subseteq E$.

Definition 3. Geographic corpus: A *geographic corpus* \mathcal{C} is a set of nodes. A node $n \in \mathcal{C}$, $n = \langle i, l, T \rangle$ is represented as a triple containing an identifier i , a location l , and a set of tags T . Each tag $t \in T$ is represented as a key-value pair with the key k and a value v : $t = \langle k, v \rangle$.

For instance, the city of Berlin is represented as follows (the example illustrates an excerpt from the OSM representation):

<i>i</i>	240109189
<i>l</i>	52.5170365, 13.3888599
<i>name=</i>	<i>Berlin</i>
<i>place=</i>	<i>city</i>
<i>capital=</i>	<i>yes</i>

Let $sameAs(n, e) : \mathcal{C} \times E_{geo} \mapsto \{true, false\}$ be the predicate that holds iff $n \in \mathcal{C}$ and $e \in E_{geo}$ represent the same real-world entity. We assume that a node $n \in \mathcal{C}$ corresponds to at most one geo-entity in a knowledge graph \mathcal{KG} . Then the problem of link discovery between a knowledge graph \mathcal{KG} and a geographic corpus \mathcal{C} is defined as follows.

Definition 4. Link discovery: Given a node $n \in \mathcal{C}$ and the set of geo-entities $E_{geo} \subseteq E$ in the knowledge graph \mathcal{KG} , determine $e \in E_{geo}$ such that $sameAs(n, e)$ holds.

In the example above, given the OSM node representing the city of Berlin, we aim to identify the entity representing this city in E_{geo} .

4. OSM2KG Approach to Link Discovery

The intuition of the proposed OSM2KG approach is as follows:

1. Equivalent nodes and entities are located in geospatial proximity. Therefore, OSM2KG adopts geospa-

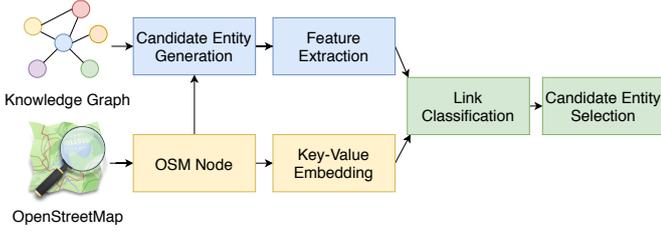


Figure 5: OSM2KG Link discovery pipeline overview.

tial blocking to identify candidate entities in large-scale datasets efficiently.

2. OSM nodes are schema-agnostic and heterogeneous. Therefore OSM2KG relies on an unsupervised model to infer latent, compact node representation that captures semantic similarity.
3. Equivalent nodes and entities can indicate common representation patterns. Therefore, OSM2KG adopts a supervised classification model for link prediction.

Figure 5 presents the OSM2KG link discovery pipeline. In the first blocking step, for each node $n \in \mathcal{C}$ in the geographic corpus \mathcal{C} , a set of candidates $E' \subseteq E_{geo}$ is generated from the set of geo-entities E_{geo} contained in the knowledge graph. In the next feature extraction step, representations of the node n and the relevant entities E' from the knowledge graph are extracted. A latent representation of the node $n \in \mathcal{C}$ is a *key-value embedding* that is learned in an unsupervised fashion. Representations of the knowledge graph entities in E' are generated using selected knowledge graph features. Furthermore, distance and similarity metrics for each candidate pair ($n \in \mathcal{C}$, $e \in E'$) are computed. Following that, each candidate pair is processed by a supervised machine learning model during the link classification step. The model predicts if the pair represents the same real-world entity and provides a confidence score for the link prediction. Finally, an identity link for the pair with the highest confidence among the positively classified candidate pairs for the node n is generated. In the following, we discuss these steps in more detail.

4.1. Candidate Entity Generation

Representations of a real-world geographic entity in different data sources may vary; this can be especially the case for the geographic coordinates in VGI, where the reference points represent typical points of community consensus rather than an objective metric [13]. The blocking step is based on the intuition that geographic coordinates of the same real-world entity representation in different sources are likely to be in a short geographic distance.

Given a node $n \in \mathcal{C}$ contained in a geographic corpus and a knowledge graph $\mathcal{KG} = \langle E \cup L, R \rangle$, with a set of geo-entities $E_{geo} \subseteq E$, in the blocking step we compute a set of candidate geo-entities $E' \subseteq E_{geo}$ from \mathcal{KG} , i.e., the geo-entities potentially representing the same real-world entity as n .

The set of candidates E' for a node n consists of all geographic entities $e \in E_{geo}$ that are in a short geographic distance to n . In particular, we consider all entities within the distance specified by the blocking threshold th_{block} :

$$E' = \{e \in E_{geo} \mid distance(n, e) \leq th_{block}\},$$

where $distance(n, e)$ is a function that computes the geographic distance between the node n and a geo-entity e . Here the geographic distance is measured as *geodisc distance* [17].

Note that E' can be computed efficiently by employing spatial index structures such as R-trees [18]. The value of the threshold th_{block} can be determined experimentally (see Section 6.5.2).

4.2. Key-Value Embedding for Geographic Corpus

In this work, we propose an unsupervised approach to infer novel latent representations of nodes in a geographic corpus. This representation aims at capturing the semantic similarity of the nodes by utilising typical co-occurrence patterns of OSM tags. Our approach is based on the intuition that semantic information, like for example entity

types, can be inferred using statistical distributions [19]. To realise this intuition in the context of a geographic corpus such as OSM, we propose a neural model inspired by the skip-gram model for word embeddings by Mikolov et al. [15]. This model creates latent node representations that capture the semantic similarity of the nodes by learning typical co-occurrences of the OSM tags.

In particular, we aim to obtain a latent representation of the node $n = \langle i, l, T \rangle, n \in \mathcal{C}$ that captures the semantic similarity of the nodes. To this extent, we propose a neural model that encodes the set of key-value pairs T describing the node in an embedding representation. Figure 6 depicts the architecture of the adopted model that consists of an input, a projection, and an output layer. The *input layer* encodes the identifier $n.i$ of each node $n = \langle i, l, T \rangle$. In particular, vector representations are obtained by applying one-hot-encoding¹⁴ of the identifiers, i.e., each identifier $n.i$ corresponds to one dimension of the input layer. The corresponding entry of the vector representation is set to 1, while other entries are set to 0. The *projection layer* computes the latent representation of the nodes. The number of neurons in this layer corresponds to the number of dimensions in the projection, i.e., the embedding size. The *output layer* maps the latent representation to the encoded keys and values using softmax [20]. The key-value pairs $\langle k, v \rangle \in n.T$ for each node n are encoded by applying one-hot-encoding to both keys and values separately. As the set of values might be highly diverse, we only consider the top-k most frequent values to be represented as an individual dimension. The non-frequent values are unlikely to be indicative for semantic similarity, whereas the information of the presence of a rare value can be discriminative. Thus, all non-frequent values are mapped to a single dimension.

The embedding aims to generate a similar representation for the nodes with similar properties, independent of

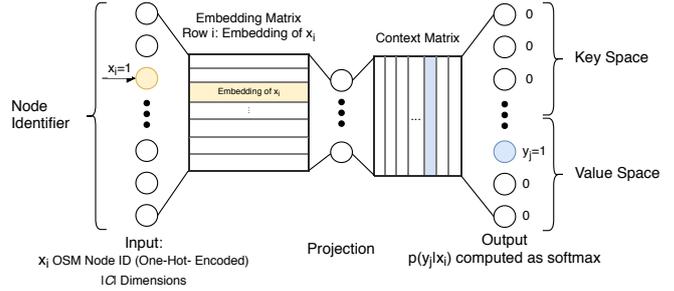


Figure 6: Architecture of the key-value embedding model. The input layer 1-hot encodes the node identifiers. The embedding matrix transforms the input to the latent representation in the projection layer. The output layer maps the latent representation to the encoded keys and values by applying the softmax function.

their location. Therefore, we do not include location information, such as geographic coordinates, in the embedding. Note that the value of name tags are typically not part of the embedding, as names typically have rare values.

The objective of the proposed model is to maximise the following log probability:

$$\sum_{n \in \mathcal{C}} \sum_{\langle k, v \rangle \in n.T} \log p(k|n.i) + \log p(v|n.i).$$

Here, the term $\log p(k|n.i) + \log p(v|n.i)$ expresses the node’s log probability with the identifier $n.i$ to be annotated with the key-value pair $\langle k, v \rangle$, i.e. $\langle k, v \rangle \in n.T$. The probabilities are calculated using softmax. The training of the network aims at minimising the key-value based loss function. This way, nodes that exhibit similar keys or values are assigned similar representations in the projection layer. Thus, we use the activation of the projection layer as a latent representation of each respective OSM node. This representation captures the latent semantics of the keys and values of the node. We refer to this feature as *KV-embedding*. We learn the *KV-embedding* for each OSM node. The training is conducted without any supervision. The resulting node representation is task-independent.

4.3. Feature Extraction from KG

This step aims at extracting features for the entities $e \in E'$, where E' denotes the set of candidate geo-entities in the knowledge graph for the target node $n \in \mathcal{C}$. We

¹⁴<https://www.kaggle.com/dansbecker/using-categorical-data-with-one-hot-encoding>

adopt the following features:

Entity Type: Entities and nodes that belong to the same category, for instance “city” or “train station”, are more likely to refer to the same real-world entity than the candidates of different types. In the knowledge graph, we make use of the *rdf:type*¹⁵ property as well as knowledge graph specific properties (e.g. *wikidata:instanceOf*) to determine the type of e . To encode the type, we create a vector of binary values in which each dimension corresponds to an entity type. For each type of e , the corresponding dimension is set to “1” while all other dimensions are set to “0”. Concerning the target node n , the node type is not expected to be explicitly provided in a geographic corpus. Nevertheless, we expect that the *KV-embedding* of the geographic corpora implicitly encodes type information, based on the intuition that types can be inferred using statistical property distributions [19].

Popularity: A similar level of entity popularity in the respective sources can provide an indication for matching. Popular entities are likely to be described with a higher number of relations and properties than less popular entities. To represent entity popularity, we employ the number of edges starting from e in \mathcal{KG} as a feature. More formally: $popularity(e) = |\{(e, x) \in R \mid x \in E \cup L\}|$. We expect that the *KV-embedding* implicitly encodes the node popularity information in the geographic corpora as popular nodes have a higher number of tags.

4.4. Similarity and Distance Metrics

This step aims at extracting features that directly reflect the similarity between an OSM node $n \in \mathcal{C}$ and a candidate geo-entity $e \in E'$. To this extent, we utilise name similarity and geographical distance.

Name Similarity: Intuitively, a geo-entity and an OSM node sharing the same name are likely to represent the same real-world object. Therefore, we encode the similarity between the value of the *name* tag of an OSM node

$n \in \mathcal{C}$ and the *rdfs:label*¹⁶ of a geo-entity $e \in E'$ as a feature. We compute the similarity using the Jaro-Winkler distance [21], also adopted by [13]. Jaro-Winkler distance assigns a value between [0,1], where 0 corresponds to no difference and 1 to the maximum dissimilarity. If a *name* tag or a *rdfs:label* is not available for a particular pair (n, e) , the value of this feature is set to 1.

Geo Distance: Based on the intuition that nodes and candidate entities that exhibit smaller geographic distance are more likely to refer to the same real-world entity, we employ geographic distance as a feature. To this extent, we utilise the logistic distance function proposed in [13]:

$$geo-distance(n, e) = 1/(1 + exp(-12d'(n, e) + 6)),$$

with $d' = 1 - d(n, e)/th_{block}$, where d denotes the so-called *geodisc distance* [17] between n and e and takes the spheroid form of the earth into account. th_{block} denotes the threshold that defines the maximum geographic distance at which the candidates are considered to be similar. To facilitate efficient computation, the th_{block} threshold is also utilised in the blocking step, described in Section 4.1. The intuition behind the logistic distance function is to allow for smaller differences of the geographic positions and to punish more significant differences. The Geo Distance feature directly encodes the geospatial similarity between the node n and the candidate geo-entity e .

4.5. Link Classification

We train a supervised machine learning model to predict whether the target node $n \in \mathcal{C}$ and a candidate geo-entity represent the same real-world entity. Each target node n and the set of candidates E' for this node are transformed into the feature space. Each node-candidate pair is interpreted as an instance for a supervised machine learning model by concatenating the respective feature vectors. For training, each pair is then labelled as correct or incor-

¹⁵rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns>

¹⁶rdfs: <http://www.w3.org/2000/01/rdf-schema>

rect, where labels are obtained from the existing links to the knowledge graph within the OSM corpus \mathcal{C} . Note that the number of pairs labelled as incorrect (i.e., negative examples) is typically higher than the number of correct pairs. To allow an efficient training of classification models, we limit the number of incorrect candidates for each node n to 10 candidates via random sampling. To address the imbalance of classes within the training data, we employ oversampling to level out the number of instances per class. In particular, we employ the state-of-the-art SMOTE algorithm [22]. The data is then normalised by removing the mean and scaling to unit variance. We use the normalised data as input to the classification model. We consider the following models: RANDOM FOREST, DECISION TREE, NAÏVE BAYES, and LOGISTIC REGRESSION. We discuss the model performance in Section 6.3. We optimise the hyperparameters using random search [23].

Finally, the candidate entity selection is based on the assumption that the knowledge graph contains at most one geo-entity equivalent to the target node. If at least one node within E' is classified as correct (with a confidence $> 50\%$), a link between node n and $e_{max} \in E'$ is created, where e_{max} denotes the entity with the highest confidence score of the model. If all entities are labelled as incorrect, no link for the node n is created.

4.6. Algorithm for Link Discovery

Finally, Algorithm 1 details the process of link discovery. The algorithm integrates the above described steps, namely *candidate entity generation* (line 1), *feature extraction* (lines 2-7), *link classification* (lines 9-12) and *candidate entity selection* (lines 12-17). Table 2 presents a description of the functions used in the algorithm.

4.7. Implementation

In this section, we provide implementation details of the OSM2KG components. We implemented our overall experimental framework and the proposed algorithm in

Algorithm 1 Link Discovery

Input: Node $n \in \mathcal{C}$
Knowledge graph \mathcal{KG}

Output: Entity $e_{link} \in \mathcal{KG}$ that should be linked to n
or **null** if no matching entity was found

```

1:  $E' \leftarrow \text{generateCandidates}(n, \mathcal{KG})$ 
2:  $\text{features} \leftarrow []$ 
3:  $\text{features}[n] \leftarrow \text{KV-embedding}(n)$ 
4: for all  $e \in E'$  do
5:    $\text{features}[e] \leftarrow \text{KG-features}(e, \mathcal{KG})$ 
6:    $\text{features}[e] \leftarrow \text{features}[e] \cup \text{similarity-features}(e, n)$ 
7: end for
8:  $\text{confidences} \leftarrow []$ 
9: for all  $e \in E'$  do
10:   $\text{confidences}[e] \leftarrow \text{link-classification}(\text{features}[n], \text{features}[e])$ 
11: end for
12:  $e_{link} \leftarrow \text{argmax}_{e \in E'}(\text{confidences}[e])$ 
13: if  $\text{classifiedAsCorrect}(e_{link})$  then
14:   return  $e_{link}$ 
15: else
16:   return null
17: end if

```

Table 2: Description of functions used in Algorithm 1.

Function Name	Returned Result	Section
generateCandidates	Candidate entities from \mathcal{KG} nearby n	4.1
KV-embedding	Latent representation of n	4.2
KG-features	Feature representation for e	4.3
similarity-features	Similarity features between e and n	4.4
link-classification	Confidence score for (n, e)	4.5
classifiedAsCorrect	True iff a link between (n, e) is classified to be correct	4.5

Java 8. We stored the evaluation results in a PostgreSQL¹⁷ database (version 9.6). In a pre-processing step, we extracted relevant data from OpenStreetMap using Python (version 3.6) and the osmium¹⁸ library (version 2.14). We extracted relevant knowledge graph entities from Wikidata with geographic coordinates using pyspark¹⁹ (version 2.2). The geographic data was stored in a PostgreSQL database (version 9.6) and indexed using the PostGIS²⁰ extension (version 2.3). The feature extraction is implemented in Java 8 within our experimental framework. We implemented the extraction of the KV-embedding in Python 3.6, using Tensorflow²¹ version 1.14.1. The machine learning algorithms were implemented in Python 3.7 using the scikit-learn²² (version 0.21) and the imbalanced-learn²³ (version 0.5) libraries. To facilitate the reproducibility, we make our code available under the open MIT license in a GitHub repository²⁴.

5. Evaluation Setup

In this section, we describe the datasets, metrics, baselines and OSM2KG configurations utilised in the evaluation.

5.1. Datasets and Metrics

We conduct our evaluation on three large-scale OSM datasets for France, Germany, and Italy as well as the Wikidata and DBpedia knowledge graphs.

Knowledge Graphs: In our experiments, we consider Wikidata snapshot from September 2018, as well as DBpedia in its German, French and Italian editions, snapshots from August 2019, as the target knowledge graphs.

¹⁷<https://www.postgresql.org/>

¹⁸<https://osmcode.org/libosmium/>

¹⁹<https://spark.apache.org/docs/latest/api/python/pyspark.html>

²⁰<https://postgis.net/>

²¹<https://www.tensorflow.org/>

²²<https://scikit-learn.org/stable/>

²³<https://imbalanced-learn.readthedocs.io/en/stable/api.html>

²⁴<https://github.com/NicolasTe/osm2kg>

Table 3: The number of geographic entities, distinct types and average statements per geo-entity in the considered knowledge graphs.

Knowledge Graph	No. Geo-Entities	No. Distinct Types	Average No. Edges/Entity
Wikidata	6,465,081	13,849	24.69
DBpedia-FR	317,500	185	18.33
DBpedia-DE	483,394	129	31.60
DBpedia-IT	111,544	11	31.13

Wikidata [5] is a publicly available collaborative knowledge graph. Wikidata is the central repository for structured information of the Wikimedia Foundation and the currently largest openly available knowledge graph. *DBpedia* [6] is a knowledge graph that extracts structured data from the information of various Wikimedia projects, e.g., the Wikipedia²⁵ encyclopedia. DBpedia is provided in language-specific editions. We refer to each language-specific edition of DBpedia as *DBpedia-[language]*. Table 3 presents the number of available geographic entities as well as the number of distinct types and the average number of edges per geo-entity in each knowledge graph. Note that we consider geo-entities in the knowledge graphs with valid geographic coordinates, i.e., coordinates that can be located on the globe.

OpenStreetMap: We consider OSM datasets extracted from the three largest country-specific OSM snapshots as of September 2018. In particular, we consider the snapshots of Germany, France, and Italy. We denote the country-specific snapshots as *OSM-[country]*. Furthermore, we extract all nodes that exhibit a link to a geo-entity contained in Wikidata or DBpedia. For DBpedia, we consider links to the DBpedia version of the language that corresponds to the country of the individual OSM snapshot, since the existing links in the country-specific snapshots target the respective language-specific edition of DBpedia in all cases for the considered datasets. We denote the considered link datasets as *[KG]-OSM-[language]*. For instance, *DBpedia-OSM-FR* denotes the dataset that interlinks the OSM snapshot of France with the French DBpe-

²⁵<https://www.wikipedia.org>

Table 4: The number of existing links between OpenStreetMap, Wikidata and DBpedia. *OSM-[country]* denote the country-specific snapshots of OSM as of September 2018. The existing links serve as ground truth for the experimental evaluation.

Knowledge Graph	OSM-FR	OSM-DE	OSM-IT
Wikidata	21,629	24,312	18,473
DBpedia-FR	12,122	-	-
DBpedia-DE	-	16,881	-
DBpedia-IT	-	-	2,353

dia.

Table 4 provides an overview of the number of existing links between OSM and the knowledge graphs. The existing links between the OSM datasets and knowledge graphs in these link datasets serve as ground truth for the experimental evaluation of all link discovery approaches considered in this work.

To assess the performance of link discovery approaches, we compute the following metrics:

Precision: The fraction of the correctly linked OSM nodes among all nodes assigned a link by the considered approach.

Recall: The fraction of the OSM nodes correctly linked by the approach among all nodes for which links exist in the ground truth.

F1 score: The harmonic mean of recall and precision. In this work, we consider the F1 score to be the most relevant metric since it reflects both recall and precision.

We apply the 10-fold cross-validation. We obtain the folds by random sampling the links from the respective link datasets. For each fold, we train the classification model on the respective training set. We report the macro average over the folds of each metric.

5.2. Baselines

We evaluate the link discovery performance of OSM2KG against the following unsupervised and supervised baselines:

BM25: This naive baseline leverages the standard BM25 text retrieval model [24] to predict links. We created an inverted index on English labels of all geo-entities

(i.e., for all $e \in E_{geo}$) in a pre-processing step to apply this model. Given the target node n , we query the index using the value of the name tag of n to retrieve geo-entities with similar labels. We query the index using either the English name tag of the node n (if available) or the name tag without the language qualifier. We create the link between n and the entity with the highest similarity score returned by the index. If the name tag is not available, we do not create any link.

SPOTLIGHT: This baseline employs the *DBpedia Spotlight* [25] model to determine the links. *DBpedia Spotlight* is a state-of-the-art model to perform entity linking, i.e., to link named entities mentioned in the text to the DBpedia knowledge graph. Given an OSM node n , we use the name tag of this node in the language native to the specific OSM dataset as an input to the DBpedia Spotlight model in the same language edition. The model returns a set of DBpedia entities out of which we choose the entity with the highest confidence score. To increase precision, we restrict the DBpedia Spotlight baseline to return only entities of type *dbo:Place*²⁶. DBpedia entities are resolved to the equivalent Wikidata entities using existing *wikidata:about* links.

GEO-DIST: This baseline predicts the links solely based on the geographic distance, measured as geodisc distance. For a target OSM node n , the link is created between n and $e_{min} \in E_{geo}$, where

$$e_{min} = \operatorname{argmin}_{e \in E_{geo}} (\operatorname{distance}(n, e)).$$

Here, $\operatorname{distance}(n, e)$ is a function that computes the geodisc distance between the OSM node n and the geo-entity e .

LGD: This baseline implements a state-of-the-art approach of interlinking OSM with a knowledge graph proposed in the *LinkedGeoData* project [13]. The LGD baseline utilises a combination of name similarity computed using the *Jaro-Winkler* string distance and geographic dis-

²⁶dbo: DBpedia Ontology

tance. It aims at computing links with high precision. For each OSM node n a link between n and $e \in E_{geo}$ is generated if the condition $\frac{2}{3}s(n, e) + \frac{1}{3}g(n, e, th_{block}) > th_{str}$ is fulfilled, where $th_{str} = 0.95$. Here, $s(n, e)$ denotes the Jaro-Winkler distance between the value of the name tag of n and the label of e . If the name tag is not available, an empty string is used to compute the distance. $g(n, e, th_{block})$ is a logistic geographic distance function specified in [13]. The parameter th_{block} denotes the maximum distance between a geo-entity and the node n . In our experiments, we use $th_{block} = 20000$ meter to allow for high recall.

LGD-SUPER: We introduce supervision into the LGD baseline by performing exhaustive grid search for $th_{block} \in \{1000, 1500, 2500, 5000, 10000, 20000\}$ meter and $th_{str} \in \{0.05 \cdot i \mid i \in \mathbb{N}, 1 \leq i \leq 20\}$. We evaluate each combination on the respective training set and pick the combination that results in the highest F1 score.

YAGO2GEO: This method was proposed in [14] to enrich the YAGO2 knowledge graph with geospatial information from external sources, including OpenStreetMap. Similar to LGD, this baseline relies on a combination of the Jaro-Winkler and geographic distance. In particular, a link between an OSM node n and $e \in E_{geo}$ is established if $s(n, e) < th_{str}$ and $distance(n, e) < th_{block}$ with $th_{str} = 0.82$, $th_{block} = 20000$ meter. $s(n, e)$ denotes the Jaro-Winkler distance between the value of the name tag of n and the label of e , and $distance(n, e)$ denotes the geographic distance between e and n .

YAGO2GEO-SUPER: We introduce supervision into the YAGO2GEO baseline by performing exhaustive grid search for $th_{block} \in \{1000, 1500, 2500, 5000, 10000, 20000\}$ meter and $th_{str} \in \{0.05 \cdot i \mid i \in \mathbb{N}, 1 \leq i \leq 20\}$. We evaluate each combination on the respective training set and pick the combination that results in the highest F1 score.

LIMES/Wombat: The Wombat algorithm, integrated within the LIMES framework [26], is a state-of-the-art approach for link discovery in knowledge graphs. The algo-

rithm learns rules, so-called link specifications, that rate the similarity of two entities. The rules conduct pairwise comparisons of properties, which are refined and combined within the learning process. As LIMES requires the data in the RDF format, we transformed the OSM nodes into RDF triples, in which the OSM id represents the subject, the key represents the predicate, and the value represents the object. We further added *geo:lat*²⁷ and *geo:long* properties representing geographic coordinates of the OSM nodes. LIMES requires all entities to contain all considered properties. Therefore we limit the properties to the geographic coordinates *geo:lat*, *geo:lon* as well as the name tag in OSM and the *rdfs:label*²⁸ in the knowledge graph. We use the default similarity metrics of LIMES, namely Jaccard, trigram, 4-grams, and cosine similarity and accept all links with a similarity score higher or equal to 0.7. Note that LIMES does not distinguish between data types when using machine learning algorithms. Therefore, it is not possible to simultaneously use string similarity and spatial similarity metrics (e.g. Euclidean distance).

5.3. OSM2KG Configurations

We evaluate our proposed OSM2KG approach in the following configuration: RANDOM FOREST as classification model (according to the results presented later in Section 6.3, RANDOM FOREST and DECISION TREE perform similarly on our datasets), dataset-specific embedding size of 3-5 dimensions (Section 6.5.1), and a blocking threshold of 20 km for DBpedia-OSM-IT and 2.5 km for all other datasets (Section 6.5.2).

Furthermore, we evaluate our proposed approach in the following variants:

OSM2KG: In this variant, we run OSM2KG as described in Section 4 using the features KV-embedding, Name Similarity, Geo Distance, Entity Type, and Popularity. To obtain latent representations of the OSM nodes, we train unsupervised embedding models as described in Section 4.2

²⁷geo: http://www.w3.org/2003/01/geo/wgs84_pos

²⁸rdfs: <http://www.w3.org/2000/01/rdf-schema>

on each of the OSM-FR, OSM-IT, OSM-DE datasets. During training, we consider the top-k most frequent values with $k=1000$ to be represented in the value space and compute 1000 epochs using a learning rate of $\alpha = 1.0$. We make the key-value embeddings of OpenStreetMap nodes created in our experiments publicly available²⁹. These key-value embeddings provide a task-independent compact representation of OSM nodes.

OSM2KG-TFIDF: To better understand the impact of the proposed embedding method on the link discovery performance, in this variant, we exchange the proposed KV-embedding with a simple TF-IDF representation of the keys and values (i.e., term frequency and inverse document frequency). To this extent, we computed the TF-IDF values of the top 1000 most frequent keys and values for each OSM dataset. In this representation, each of the keys and values is described by a single dimension, resulting in a 1000-dimension vector. All other features, such as Name Similarity, Geo Distance, Entity Type, and Popularity remain the same.

6. Evaluation

The main goal of the evaluation is to assess the link discovery performance of OSM2KG compared to the baselines. Moreover, we analyse the effectiveness of the classification model and the proposed features and perform parameter tuning.

6.1. Link Discovery Performance

Table 5 summarises the overall link discovery performance results of the BM25, SPOTLIGHT, GEO-DIST, LGD, LGD-SUPER, YAGO2GEO, YAGO2GEO-SUPER, and LIMES/WOMBAT baselines as well as our proposed approach in the OSM2KG and OSM2KG-TFIDF variants. Table 5a reports the results of the experiments conducted on the link datasets from Wikidata, while Table 5b

reports the result on the DBpedia datasets. We report the macro averages of the 10-fold cross-validation conducted on the corresponding link dataset concerning the precision, recall, and F1 score. In our experiments, we observed that the micro averages behave similarly.

Overall, we observe that in terms of F1 score, OSM2KG performs best on all Wikidata datasets, where it achieves an F1 score of 92.05% on average and outperforms the best performing LGD-SUPER baseline by 21.82 percentage points. Furthermore, we observe that OSM2KG achieves the best performance concerning the recall on all datasets. Moreover, OSM2KG maintains high precision, i.e., 94.62% on Wikidata and 97.94% on DBpedia, on average. Regarding the DBpedia datasets, we observe that OSM2KG outperforms the baselines on DBpedia-OSM-FR and DBpedia-OSM-IT, whereas the difference to the LGD-SUPER baseline is much smaller, compared to Wikidata. On DBpedia-OSM-DE, LGD-SUPER archives a slightly higher F1 score, compared to OSM2KG. This result indicates that, in contrast to Wikidata, the respective DBpedia and OSM datasets are well-aligned in terms of names and geographic coordinates, such that simple heuristics utilising name similarity and geographic distance can already yield good results in many cases. In contrast, the task of link discovery in Wikidata is more challenging. In these settings, the advantages of the OSM2KG approach become clearly visible.

The BM25 and SPOTLIGHT baselines adopt name similarity for matching, whereas SPOTLIGHT can also make use of the knowledge graph context, including entity types. As we can observe, BM25 shows relatively low performance in terms of both precision (on average 45.66% (Wikidata) and 53.94% (DBpedia)) and recall (on average 41.95% (Wikidata) and 62.61% (DBpedia)). The SPOTLIGHT baseline can improve on BM25 regarding precision and F1 score on Wikidata and DBpedia datasets. However, the absolute precision and F1 scores of SPOTLIGHT, with the maximum F1 score of 65.40% on Wikidata, are not competitive.

²⁹<http://13s.de/~tempelmeier/osm2kg/key-value-embeddings.zip>

Table 5: Macro averages for precision, recall and F1 score [%], best scores are bold. Statistically significant (according to paired t-tests with $p < 0.05$) F1 score results of OSM2KG compared to all baselines and OSM2KG-TFIDF are marked with *.

(a) Link prediction performance on the Wikidata datasets

Approach	Wikidata-OSM-FR			Wikidata-OSM-DE			Wikidata-OSM-IT			Average		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
BM25	45.22	42.59	43.86	47.28	41.60	44.26	44.49	41.67	43.04	45.66	41.95	43.72
SPOTLIGHT	65.17	32.26	43.15	69.79	51.03	58.95	54.79	26.89	36.08	63.25	36.73	46.06
GEO-DIST	74.46	74.46	74.46	62.16	62.16	62.16	72.80	72.80	72.80	69.81	69.81	69.81
LGD	100.00	44.09	61.20	100.00	47.46	64.37	100.00	43.59	60.71	100.00	45.05	62.09
LGD-SUPER	100.00	53.25	69.50	100.00	55.34	71.25	100.00	53.79	69.95	100.00	54.13	70.23
YAGO2GEO	63.66	44.98	52.71	64.48	48.61	55.43	58.40	47.36	52.30	62.18	46.98	53.48
YAGO2GEO-SUPER	78.49	47.38	59.09	73.49	48.96	58.76	72.25	48.73	58.20	74.74	48.36	58.69
LIMES/WOMBAT	74.03	17.50	28.31	78.54	17.01	27.97	65.28	17.22	27.25	72.62	17.25	27.84
OSM2KG-TFIDF	95.06	90.60	92.77	93.67	86.37	89.87	93.98	87.07	90.39	94.24	88.01	91.01
OSM2KG	95.51	91.90	93.67*	93.98	88.29	91.05*	94.39	88.68	91.45*	94.62	89.63	92.05

(b) Link prediction performance on the DBpedia datasets

Approach	DBpedia-OSM-FR			DBpedia-OSM-DE			DBpedia-OSM-IT			Average		
	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
BM25	70.04	69.32	69.68	47.28	76.84	75.58	44.49	41.67	43.04	53.94	62.61	62.77
SPOTLIGHT	72.40	49.42	58.74	79.08	62.31	69.70	85.38	56.17	67.76	78.95	55.97	65.40
GEO-DIST	85.94	85.94	85.94	66.49	66.49	66.49	86.17	86.17	86.17	79.53	79.53	79.53
LGD	100.00	61.81	76.40	100.00	60.72	75.56	100.00	64.94	78.74	100.00	62.49	76.90
LGD-SUPER	100.00	88.18	93.72	100.00	84.56	91.63	100.00	86.90	92.99	100.00	86.55	92.78
YAGO2GEO	77.52	70.40	73.78	87.41	75.84	81.22	94.74	78.47	85.84	86.56	74.90	80.28
YAGO2GEO-SUPER	84.74	82.47	83.59	93.62	80.14	86.36	97.46	81.28	88.64	91.94	81.30	86.19
LIMES/WOMBAT	82.34	60.33	69.64	79.00	68.00	73.09	97.38	70.89	82.05	86.24	66.41	74.93
OSM2KG-TFIDF	98.68	95.35	96.99	95.61	84.93	89.95	98.46	89.83	93.95	97.91	90.04	93.63
OSM2KG	99.06	96.25	97.63*	95.65	85.83	90.47	99.11	90.13	94.41	97.94	90.74	94.17

Table 6: Parameters learned by the LGD-SUPER and the YAGO2GEO-SUPER baselines

Data Set	LGD-SUPER		YAGO2GEO-SUPER	
	th_{block}	th_{str}	th_{block}	th_{str}
Wikidata-OSM-FR	1500	0.1	1000	0.70
Wikidata-OSM-DE	2000	0.1	2000	0.80
Wikidata-OSM-IT	1500	0.1	1000	0.70
DBpedia-OSM-FR	1000	0.1	1000	0.30
DBpedia-OSM-DE	5000	0.1	2000	0.75
DBpedia-OSM-IT	20000	0.3	1500	0.30

Overall, we conclude that name similarity, as adopted by these baselines, is not sufficient for effective link prediction.

The LGD and LGD-SUPER baselines that combine name similarity and geographic distance achieve the best precision of 100% on all datasets. However, LGD baselines suffer from lower recall. LGD-SUPER achieves on average 54.13% recall on Wikidata and 86.55% recall on DBpedia, overall resulting in lower F1 scores on average compared to OSM2KG. The YAGO2GEO baseline that uses similar features as LGD achieves higher recall scores

than LGD (46.98% on Wikidata, 74.90% on DBpedia on average) but cannot maintain the high precision of LGD (on average 62.18% on Wikidata, 86.56% on DBpedia). Overall, YAGO2GEO achieves lower F1 scores compared to OSM2KG.

Regarding the supervised baselines, Table 6 presents the parameters learned by LGD-SUPER and the YAGO2GEO-SUPER during the training process. We observe that YAGO2GEO-SUPER learns more restrictive parameters, whereas LGD-SUPER allows for less restrictive threshold values. This result indicates that the ranking function of LGD-SUPER that combines geographic distance and name similarity is more robust than the ranking function of YAGO2GEO-SUPER. YAGO2GEO-SUPER uses geographic distance exclusively for blocking and ranks the candidates based solely on the name similarity. We observe that both baselines achieve a reasonably good performance on the DBpedia datasets. On the contrary, both baselines can not

reach comparable performance on the Wikidata datasets and result in 70.23% F1 score for LGD-SUPER, and 58.69% F1 score for YAGO2GEO-SUPER, on average.

GEO-DIST, which solely relies on the geographic distance, achieves an F1 score of 69.81% on Wikidata, and 79.53% on DBpedia on average. Although a significant fraction of the OSM nodes can be correctly linked solely based on the geographic distance, still a significant fraction of nodes (on average 30.19% for Wikidata and 20.74% for DBpedia) can not be appropriately linked this way. We observe that the lower performance of GEO-DIST corresponds to densely populated areas (e.g., large cities), where we expect knowledge graphs to have a higher number of entities, making disambiguation based on geographic distance ineffective. OSM2KG overcomes this limitation and outperforms the GEO-DIST baseline by 22.24 percentage points (Wikidata) and 14.64 percentage points (DBpedia) on average concerning F1 score.

The LIMES/WOMBAT baseline that aims to learn rules for link discovery in a supervised fashion does not achieve competitive performance on any considered dataset and results in 27.84% F1 score for Wikidata and 74.93% F1 score for DBpedia on average. One of the main reasons for such low performance is that LIMES/WOMBAT requires all entities to contain all considered properties. As none of the OSM tags is mandatory, this baseline is de-facto limited to only frequently used properties, such as the name and the geo-coordinates. These properties alone are insufficient to extract the rules leading to competitive performance in the link discovery task on these datasets.

Comparing the performance of OSM2KG across the datasets, we observe that scores achieved on the Wikidata-OSM-FR and DBpedia-OSM-FR datasets (93.67%, and 97.63% F1 score) are higher than on the other language editions. This result can be explained through a more consistent annotation of the nodes within the OSM-FR dataset. For instance, in OSM-FR eight key-value combinations appeared more than 2000 times, whereas in OSM-DE and OSM-IT only

two to four combinations are that frequent.

Comparing the overall link discovery performance on the DBpedia and Wikidata datasets, we observe that higher F1 scores are achieved on DBpedia by all considered approaches. Furthermore, the LGD-SUPER and YAGO2GEO-SUPER baselines that utilise only geographic distance and name similarity heuristics can reach high performance on DBpedia (up to 92.78% F1 score on average). In contrast, their maximal performance on Wikidata is limited to 70.23% F1 score. This result indicates that, in general, geographic coordinates and entity names of OSM are better aligned with DBpedia than with Wikidata. This result also suggests that the link discovery task is more difficult on Wikidata. Our OSM2KG approach is particularly useful in these settings, where we achieve 21.82 percentage points increase in F1 score compared to the best performing LGD-SUPER baseline.

6.2. Comparison to OSM2KG-TFIDF

Comparing the performance of OSM2KG with the OSM2KG-TFIDF variant, we observe that the embedding of OSM2KG leads to better performance (1.04 percentage points of F1 score for Wikidata and 0.54 percentage points of F1 score for DBpedia on average).

We observe a statistically significant difference between the F1 scores of OSM2KG and OSM2KG-TFIDF on all Wikidata datasets and DBPEDIA-OSM-FR (paired t-tests with $p < 0.01$). Through a manual inspection of exemplary instances, we found that OSM2KG especially improves over OSM2KG-TFIDF on discovering links for nodes with name information and nodes corresponding to Wikidata types with a small number of instances. For example, a node corresponding to a private school³⁰ was wrongly assigned to a trade school³¹ instead of the entity³². In this example, the name of the OSM node and the geo-entity are identical. We believe that through the high number

³⁰<https://www.openstreetmap.org/node/2733503641>

³¹<https://www.wikidata.org/wiki/Q828825>

³²<https://www.wikidata.org/wiki/Q2344470>

of dimensions in the TF-IDF representation, the *name* dimension and the corresponding *name similarity* might lose importance, even though the name is typically a very effective feature in the context of link discovery. From the RANDOM FOREST models, we observe that the *name similarity* achieves a lower mean decrease impurity [27] in OSM2KG-TFIDF than in OSM2KG, indicating the lower contribution of the feature. Moreover, the *KV-embedding* poses a distributed representation of the OpenStreetMap tags. We believe that especially for Wikidata types with a small number of instances the distributed representation might be more robust, whereas in a TF-IDF representation single tags could introduce bias towards types with a higher number of instances. In the example above, the tag `toilets:wheelchair=yes` is likely to co-occur with both the private school and trade school types but might be biased towards the more populated type.

We do not observe statistically significant differences between OSM2KG and OSM2KG-TFIDF on the DBpedia-OSM-DE and DBpedia-OSM-IT datasets. On these datasets, baselines that exclusively make use of geographic distance and name similarity such as LGD-SUPER achieve the best or close-to-best F1 score. Therefore, the individual importance of the *KV-embedding* or the TF-IDF feature is not as high as for the other datasets.

Furthermore, the proposed *KV-embedding* provides a compact representation that consists of only 3-5 dimensions, whereas the corresponding TF-IDF representations consist of 1000 dimensions. Figure 7 contrasts the average memory consumption across the folds of the random forest models of OSM2KG and OSM2KG-TFIDF. We observe that the usage of the *KV-embedding* generally results in a lower memory footprint than the TF-IDF variant, which becomes particularly visible for larger datasets. The difference is largest on the Wikidata-OSM-FR dataset, where the *KV-embedding* (0.7 GB) requires only 5% of memory compared to the TF-IDF variant (14 GB). We observe the smallest difference on DBpedia-OSM-IT. This dataset has

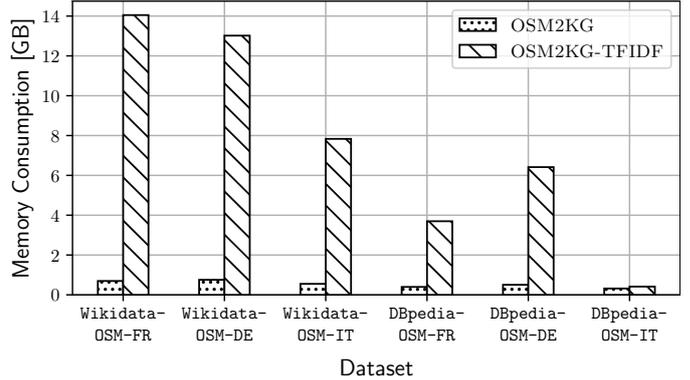


Figure 7: Average memory consumption across folds of the training of the RANDOM FOREST models used by OSM2KG and OSM2KG-TFIDF.

the smallest number of instances (2353), resulting in the small memory difference between the models (0.1 GB).

We conclude that KV-embedding is an effective, concise, and task-independent way to represent the OSM information. We believe that this representation makes OSM data more usable for models that may suffer from the curse of dimensionality or memory limitations.

6.3. Classification Model Performance

Table 7 presents the F1 scores achieved by OSM2KG with respect to each dataset while varying the classification model. In particular, we evaluate the performance of RANDOM FOREST, DECISION TREE, NAÏVE BAYES, and LOGISTIC REGRESSION. As we can observe, the performance of the classification models is consistent among the datasets. RANDOM FOREST and DECISION TREE achieve similar F1 scores and show the best performance, i.e., on average 92.05% (Wikipedia), 94.17% (DBpedia) F1 score using RANDOM FOREST, and 92.21% (Wikidata), 93.77% (DBpedia) using DECISION TREE. According to a paired t-test, the observed differences between the RANDOM FOREST and DECISION TREE are not statistically significant on our datasets. In contrast, the performance of NAÏVE BAYES and LOGISTIC REGRESSION is much lower, i.e., they achieve on average only 66.99% (Wikidata), 80.93% (DBpedia) F1 score using NAÏVE BAYES and 67.54% (Wikidata), 87.49% (DBpedia) using LOGISTIC REGRESSION.

Table 7: Comparison of OSM2KG F1 scores [%] with respect to the classification model, best scores are bold.

Classifier	Wikidata- OSM-FR	Wikidata- OSM-DE	Wikidata- OSM-IT	Wikidata- Average	DBpedia- OSM-FR	DBpedia- OSM-DE	DBpedia- OSM-IT	DBpedia- Average
RANDOM FOREST	93.67	91.05	91.45	92.05	97.63	90.47	94.41	94.17
DECISION TREE	94.45	91.17	91.01	92.21	97.12	89.62	94.56	93.77
NAÏVE BAYES	70.88	63.64	66.45	66.99	76.69	77.69	88.40	80.93
LOGISTIC REGRESSION	65.36	66.40	70.87	67.54	86.84	86.93	88.71	87.49

Table 8: Differences in OSM2KG F1 score [percentage points] when leaving out single features using RANDOM FOREST.

Left out Feature	Wikidata- OSM-FR	Wikidata- OSM-DE	Wikidata- OSM-IT	Wikidata- Average	DBpedia- OSM-FR	DBpedia- OSM-DE	DBpedia- OSM-IT	DBpedia- Average
KV-embedding	2.80	3.91	4.53	3.75	1.94	1.96	0	1.30
Geo Distance	15.28	14.72	11.98	13.99	2.81	2.19	8.67	4.56
Name	1.92	3.52	3.51	2.98	3.61	5.66	6.86	5.38
Entity Type	0.71	2.00	2.77	1.83	0.45	0.54	-0.08	0.30
Popularity	0.29	1.07	0.94	0.77	0.29	0.28	-0.02	0.18
Entity Type & Popularity	1.67	9.30	6.94	5.97	0.84	1.50	-0.08	0.75

We conclude that non-linear classification models such as RANDOM FOREST and DECISION TREE are better suited to the problem we address than the linear models. This result also suggests that the classification problem is not linearly separable. In our experiments in Section 6.1, we made use of RANDOM FOREST classification models.

6.4. Feature Evaluation

In this section, we assess the feature contributions of OSM2KG. To assess the contribution of the single features to link discovery, we conducted a leave-one-out feature evaluation. In particular, we removed each feature individually from the feature set and determined the difference in F1 score to quantify the feature importance.

Table 8 shows the differences in the F1 score of the OSM2KG model when a single feature is left out compared to the F1 score achieved when the entire feature set is used. Since no difference is negative, except for DBpedia-OSM-IT, we conclude that all features typically contribute to better classification performance. *Geo Distance* results in the most substantial difference of 13.99 percentage points on average for Wikidata. On DBpedia, *Geo Distance* results in the second-largest difference of 4.56 percentage points on average. The most considerable difference for DBpedia results from the *Name* feature, with 5.38 percentage points on average. For Wikidata,

the *Name* feature results in a variation of 2.98 percentage points on average. The importance of the *Name* feature on DBpedia indicates that the names of the OSM and DBpedia datasets are well-aligned. This result confirms our observations in Section 6.1, where we discussed the performance of the LGD-SUPER baseline that utilises both features.

The *KV-embedding* feature shows the second-largest difference on Wikidata (3.75 percentage points) and the third-largest difference on DBpedia (1.30 percentage points) on average. As expected, the contribution of this feature is higher for the more complex link discovery task in Wikidata, as opposed to DBpedia, where simple heuristics may suffice. As an extreme example, we do not observe any contribution of *KV-embedding* for DBpedia-OSM-IT. As discussed before, simple heuristics (e.g., geographic distance and name similarity) are sufficient to achieve relatively high performance on this dataset.

The *Entity Type* and *Popularity* show the smallest differences, where *Entity Type* has slightly larger differences than *Popularity*. For the Wikidata datasets, we observe that the individual contributions of the features are rather small, i.e. 1.83 percentage points (*Entity Type*) and 0.77 percentage points (*Popularity*) on average. When leaving both features out, we observe a difference of 5.97 percent-

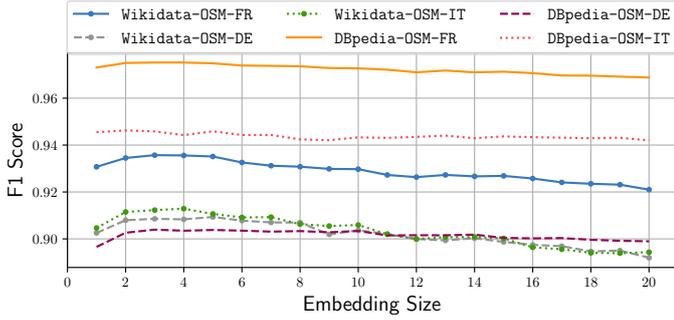


Figure 8: Influence of the embedding size on F1 score of the RANDOM FOREST classifier.

age points on average. We conclude that the information encoded in both features is partly redundant. Furthermore, this relatively large difference indicates feature importance. We conclude that for Wikidata datasets the information of the *Entity Type* is especially useful when combined with the *Popularity* feature. On the contrary, for the DBpedia datasets, we observe that the contribution of the *Popularity* feature is nearly identical to the joint contribution of *Entity Type* and *Popularity*. For DBpedia-OSM-IT we observe negative contributions for both features. Again, this indicates that geographic distance and name similarity are sufficient for link discovery in this dataset.

Although *Entity Type* and *Popularity* are correlated in many cases, they can provide complementary information for some instances. Intuitively, the joint information can help to disambiguate entities similar concerning one of the features, but dissimilar regarding the other. For example, two railway stations of different sizes are likely to be described with a different number of statements, whereas the type is identical. In such cases, in addition to the *Entity Type*, *Popularity* can help to disambiguate entities better.

6.5. Parameter Tuning

We evaluate the influence of the parameters such as embedding size and the blocking threshold value on the performance of OSM2KG.

6.5.1. Embedding Size

The embedding size corresponds to the number of dimensions (i.e. neurons) in the projection layer of the neural model presented in Section 4.2. Figure 8 shows F1 scores obtained with respect to the number of dimensions of the *KV-embedding* achieved by the RANDOM FOREST classifier on all datasets.

We observe similar trends for all datasets except for DBpedia-OSM-IT. Overall, we can observe a growth of the F1 score of the classifier with an increasing number of dimensions, between one and four dimensions for all datasets. We conclude that embeddings with an insufficient number of dimensions are not able to capture all relevant information. When the number of dimensions increases, more information can be encoded, which leads to better performance. As we can observe, the curve achieves its maximum at three dimensions for the Wikidata-OSM-FR, and DBpedia-OSM-FR datasets, at four dimensions for Wikidata-OSM-IT and at five dimensions for the Wikidata-OSM-DE and DBpedia-OSM-DE datasets. Further increase of the embedding size does not lead to an increase in performance. On the contrary, the performance can drop, indicating that no additional beneficial information is obtained by adding further dimensions.

For DBpedia-OSM-IT, we observe a near-constant performance around 94% F1 score of the classifier. As discussed in Section 6.4, here the contribution of the *KV-embedding* is not as high as for the other datasets. Thus the variation of the embedding size does not result in any significant performance changes for this dataset.

Overall, we conclude that 3-5 dimensions are most suited for the datasets that make effective use of the *KV-embedding* feature. Thus we adopted the following number of dimensions: Wikidata-OSM-FR: 3, Wikidata-OSM-DE:5, Wikidata-OSM-IT: 4, DBpedia-OSM-FR: 3, DBpedia-OSM-DE: 5, DBpedia-OSM-IT: 4.

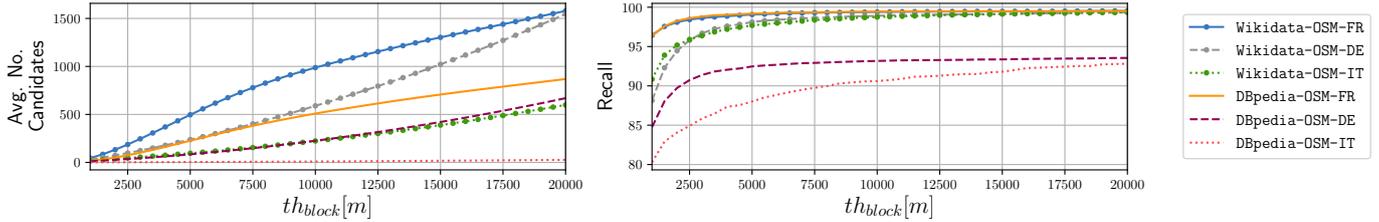


Figure 9: Influence of the threshold th_{block} on the average number of candidates and recall of the blocking step.

6.5.2. Blocking Threshold

The blocking threshold th_{block} represents the maximal geographic distance considered for candidate entity generation, as discussed in Section 4.1. For a single OSM node, all knowledge graph entities that are closer than th_{block} are considered as candidates. The value of th_{block} can be determined experimentally by evaluating the recall of the blocking step.

Figure 9 shows the influence of th_{block} on the average number of candidates and the recall of the blocking step. Considering the average number of candidates, we observe a linear-like rise (i.e., the slope of the curve is nearly constant) of the number of candidates concerning th_{block} for all datasets, whereas the datasets differ in slope. Due to the low geographic density of the DBpedia-OSM-IT dataset, the corresponding slope is especially low. Concerning recall, we observe that the curve starts with a steady incline, but quickly saturates with an increasing th_{block} . We conclude that in most cases, the correct candidate exhibits a geographic distance of about 2.5 km. Thus, in our experiments, we chose $th_{block} = 2.5$ km. This threshold value allows for more than 85% recall of correct candidates for the DBpedia datasets and 95% recall for the Wikidata datasets in the blocking step, while effectively limiting the number of candidates. For DBpedia-OSM-IT, we adopt a different th_{block} threshold of 20 km to increase recall on this dataset.

To make the impact of geospatial blocking comparable across the considered approaches, we assess the effect of the blocking step on the overall link discovery perfor-

mance. To this extent, we added an additional blocking step to the BM25 and GEO-DIST baselines and evaluate the models BM25, GEO-DIST, LGD, YAGEO2GEO and OSM2KG with the blocking thresholds $th_{block} \in \{1, 2.5, 5, 10, 20\}$ km. Figure 10 presents the F1 scores regarding the blocking threshold value th_{block} . As we can observe, the general link discovery performance is not very sensitive to the th_{block} value. However, if th_{block} value is chosen too low, e.g. 1 km, the link discovery performance can drop, as shown in Figure 10b. Overall, an optimal threshold value depends on the model as well as on the dataset. For example, LGD may benefit from a lower blocking threshold value, as shown in Figure 10e, whereas GEO-DIST works better with a higher threshold (Figure 10f). For OSM2KG we do not observe any significant impact for values of $th_{block} \geq 2.5$ km for most datasets. For the supervised variants of the baselines LGD and YAGO2GEO, LGD-SUPER and YAGO2GEO-SUPER, we observe that the appropriate threshold can be determined during the training process. The performance of the GEO-DIST baseline is degraded with the limitation of the additional blocking step, as this limitation does not contribute to precision, but potentially limits recall of this baseline. The BM25 baseline benefits from the blocking step but is still clearly outperformed by OSM2KG. In summary, as presented by Figure 10, we observe that OSM2KG outperforms all baselines for all values of the blocking threshold th_{block} on all considered datasets concerning F1 score.

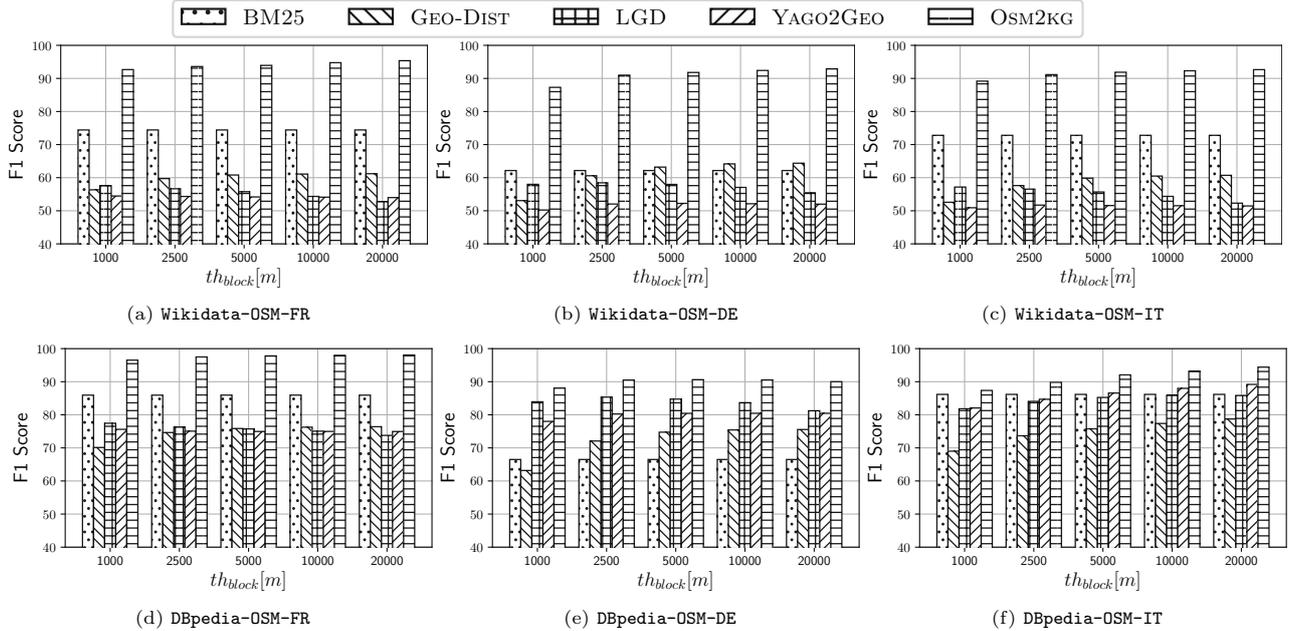


Figure 10: Link discovery performance concerning th_{block} value for OSM2KG and the baselines that can include a blocking step. X-axis presents the value of th_{block} in meter. Y-axis presents the F1 score.

Table 9: Distribution of error types on nodes for which no correct link could be found by OSM2KG.

Error Type	Wikidata-OSM-FR	Wikidata-OSM-DE	Wikidata-OSM-IT	Avg.
No candidate found	41%	54%	54%	49.67%
Wrong candidate selected	39%	37%	22%	32.67%
Duplicate entity in Wikidata	17%	4%	20%	13.67%
Wrong link in ground truth	3%	5%	4%	4.00%

6.6. Error Analysis

We conducted an error analysis through manual inspection of a random sample of 100 nodes for which OSM2KG identified no correct link for each of the Wikidata datasets. Table 9 presents the resulting error distribution. As we can observe, the most common reason for errors is a too restrictive candidate selection leading to an empty candidate set (in 49.67% of cases), followed by the selection of wrong candidates (in 32.67% of cases) and quality issues in Wikidata such as duplicate entities (in 13.67%) as well as wrong links in the ground truth data (in 4%). Note that the restrictive candidate selection is subject to the choice of the blocking threshold value. For this study, the threshold was chosen in such a way that 95% recall of the blocking step was achieved. In a small number of cases (3% on average), the candidate set is not empty, but the

correct candidate is not included in this set. This issue can be addressed by an adaptive increase of the threshold for the nodes without any candidates.

Furthermore, we observe that the selection of wrong candidates in most cases happens within the regions with a high geographic density of Wikidata entities, e.g., in cities where single houses can represent entities, resulting in a large candidate set. To further increase the precision of OSM2KG, a dedicated, supervised model for geographically dense regions can be trained. Such a model can follow a more restrictive policy, e.g., by requiring higher confidence to establish a link.

Finally, the detection of duplicate entities and wrong ground truth links indicates the potential to adopt OSM2KG for de-duplication of geo-entities in Wikidata to increase data quality. These observations provide a basis for an incremental tuning of OSM2KG in future work.

6.7. Discussion

Approaches that mainly rely on name similarity heuristics and do not leverage any geospatial features are not suitable for effective link prediction for the OSM nodes.

We can observe this by considering the relatively low performance of the BM25 and SPOTLIGHT baselines, where SPOTLIGHT achieved F1 scores of 46.06% (Wikidata) and 65.40% (DBpedia), on average. Geospatial features such as geographic distance are a reliable indicator to match OSM nodes with knowledge graph entities in our datasets. This observation is confirmed by the GEO-DIST baseline, which reached F1 scores of 69.81% (Wikidata) and 79.53% (DBpedia) by solely considering the geographic distance. However, in a significant fraction of cases, geospatial information alone is insufficient to disambiguate OSM nodes effectively. Heuristics using a combination of the name similarity and geospatial features, and in particular the supervised LGD-SUPER baseline, can achieve competitive performance on the DBpedia datasets. However, they are insufficient for link discovery in more complex datasets, such as Wikidata, where the entity names are not well-aligned with OSM.

The proposed OSM2KG approach combines the latent representation of OSM nodes that captures the semantic similarity of the nodes with geospatial information and is highly effective for link prediction. OSM2KG is of particular advantage for link discovery between OSM and Wikidata, where it significantly outperforms the baselines concerning the recall and F1 score. Overall, we observe that the proposed latent node representation as *key-value embedding* combined with geospatial distance is an effective way to facilitate link discovery in a schema-agnostic volunteered geographic dataset such as OSM. This representation, with only 3-5 dimensions, is compact and task-independent.

Limitations in link discovery can arise from the candidate generation step, where we consider the set of entities for which geographic coordinates are available in the knowledge graph. A promising direction for future research is to discover identity links between OSM nodes and geographic entities for which geographic coordinates are not available in the knowledge graph.

In this work, we focused the discussion and evaluation of OSM2KG on Wikidata and DBpedia as target knowledge graphs due to their openness, popularity, and availability of training data (i.e., the links between these knowledge graphs and OSM). Nevertheless, the proposed OSM2KG approach is applicable to other knowledge graphs provided a set of identity links between OSM and the target knowledge graph is available for training the OSM2KG classifier.

7. Related Work

Link Discovery is the task of identifying semantically equivalent resources in different data sources [28]. Nentwig et al. [28] provide a recent survey of link discovery frameworks, with prominent examples, including Silk [29] and LIMES [30].

In particular, the Wombat algorithm, integrated within the LIMES framework [26], is a state-of-the-art approach for link discovery in knowledge graphs. Link discovery approaches that operate on Linked Data typically expect datasets in Resource Description Framework (RDF) format having a schema defined by an underlying ontology and data exhibiting graph structure. This assumption does not apply to the OSM data represented as key-value pairs.

Besides the syntactic and structural differences, LIMES relies on several assumptions that severely limit its applicability to OSM datasets. First, LIMES assumes a one-to-one mapping between properties. In contrast, the required mappings between the Wikidata properties and the OSM keys are 1:n, as a Wikidata property can correspond to several OSM keys. For example, the “instanceOf” property in Wikidata corresponds to “place,” “natural,” “historic,” and many other keys in OSM. Second, LIMES requires all instances to contain all considered properties. Therefore LIMES is limited to utilise only frequently used properties, such as the name and the geo-coordinates. To this end, LIMES is not suited to utilise the information from

other infrequent properties for mapping. Finally, the current LIMES implementation does not adequately support a combination of different data types, such as strings and geo-coordinates. Given these differences, the application of LIMES to the OSM data is de-facto restricted to the name matching. We utilise Wombat/LIMES as a baseline for the evaluation. Our experimental results confirm that OSM2KG outperforms this baseline.

In the context of individual projects such as LinkedGeoData and Yago2Geo [13, 14], a partial transformation of OSM data to RDF was conducted using manually defined schema mappings for selected keys. In contrast, the proposed OSM2KG approach adopts an automatically generated latent representation of OSM data.

Entity linking (also referred to as entity disambiguation) is the task of linking mentions of real-world entities in unstructured sources (e.g., text documents) to equivalent entities in a knowledge base. A recent survey on entity linking approaches is provided in [31]. Entity linking approaches typically adopt Natural Language Processing (NLP) techniques and use the context of the entity mentions such as phrases or sentences. However, such a context is not available in OSM, where textual information is mainly limited to node labels (typically available as a specialised name tag). One of the most popular state-of-the-art models to automatically annotate mentions of DBpedia entities in natural language text is *DBpedia Spotlight* [25]. DBpedia Spotlight adopts NLP techniques to extract named entities (including locations) from text and uses a context-aware model to determine the corresponding DBpedia entities. This approach serves as a baseline in our experiments, whereas we use the name tag of an OSM node as its textual representation.

Linking geographic data: The most relevant projects in the context of our work are LinkedGeoData [13] and Yago2Geo [14]. LinkedGeoData is an effort to lift OSM data into semantic infrastructure. This goal is addressed through deriving a lightweight ontology from the OSM

tags and transforming OSM data to the RDF data model. LinkedGeoData interlinks OSM nodes represented as RDF with geo-entities in external knowledge sources such as DBpedia and GeoNames. Yago2Geo aims at extending the knowledge graph YAGO2 [7] with geographic knowledge from external data sources. To this extent, identity links between YAGO2 and OSM are computed. Both interlinking approaches rely on manually defined schema mappings and heuristics based on name similarity and geographic distance. The dependence of both approaches on manual schema mappings restricts the coverage of mapped entity types and can also negatively affect link maintenance. In contrast, the OSM2KG approach proposed in this article extracts latent representations of OSM nodes fully automatically. The LinkedGeoData and Yago2Geo interlinking approaches serve as a baseline in our experiments. Our experimental results confirm that OSM2KG outperforms both baselines. The applications of linked geographic data include, for example, the training of comprehensive ranking models [32] or the creation of linked data based gazetteers [33].

Geospatial link discovery [34, 35, 36, 37] refers to the problem of creating topological relations across geographic datasets. These links express the topographic relations between entities (e.g., intersects and overlaps). For example, [37] presented the problem of discovery of spatial and temporal links in RDF datasets. In Radon [36], efficient computation of topological relations between geospatial resources in the datasets published according to the Linked Data principles was presented. In contrast, in this work, we focus on link discovery for identity links.

Geographic representation learning: Recently, several approaches emerged that employ representation learning to encode geographic data. Typical data sources are point of interest and floating car data, where the proposed architectures include graph embeddings [38, 39, 40], metric embeddings [41], stacked autoencoders [42], generative models [43], and word2vec-like models [44, 45]. [46]

proposed neural embeddings for Geonames that explicitly takes the geospatial proximity into account. The proposed OSM2KG approach relies on an embedding architecture inspired by word2vec to automatically encode semantic similarity of the OSM nodes using key-value pairs. The embedding aims to generate a similar representation for the nodes with similar properties, independent of their location. Thus, we do not include location information in the embedding.

8. Conclusion

In this article, we proposed OSM2KG, a novel link discovery approach to predict identity links between OpenStreetMap nodes and geographic entities in knowledge graphs. OSM2KG combines latent representations of heterogeneous OSM nodes with a supervised classification model to predict identity links across large-scale, diverse datasets effectively. Our experiments conducted on three large-scale OSM datasets for Germany, France, and Italy and Wikidata and DBpedia knowledge graphs demonstrate that the proposed OSM2KG approach can reliably discover identity links.

OSM2KG achieves an F1 score of 92.05% on Wikidata and of 94.17% on DBpedia on average, which corresponds to a 21.82 percentage points increase in F1 score on Wikidata compared to the best performing baselines.

Whereas we conducted our evaluation on OSM, Wikidata and DBpedia, our approach can be applied to other VGI sources and knowledge graphs as long as a training set of identity links is available. In future work, we would like to develop novel applications that take advantage of integrated geographic and semantic information created by OSM2KG. Furthermore, we would like to explore the applicability of the proposed *KV-embedding* to further datasets and tasks.

Acknowledgement

This work is partially funded by the DFG, German Research Foundation (“WorldKG”, DE 2299/2-1, 424985896), the Federal Ministry of Education and Research (BMBF), Germany (“Simple-ML”, 01IS18054), (“Data4UrbanMobility”, 02K15A040), and the Federal Ministry for Economic Affairs and Energy (BMWi), Germany (“d-E-mand”, 01ME19009B).

References

- [1] J. Jokar Arsanjani, A. Zipf, P. Mooney, M. Helbich, An Introduction to OpenStreetMap in Geographic Information Science: Experiences, Research, and Applications, Lecture Notes in Geoinformation and Cartography, Springer International Publishing, 2015, pp. 1–15.
- [2] S. Huber, C. Rust, Calculate Travel Time and Distance with OpenStreetMap Data Using the Open Source Routing Machine (OSRM), The Stata Journal 16 (2) (2016).
- [3] G. Touya, A. Reimer, Inferring the scale of openstreetmap features, in: OpenStreetMap in GIScience - Experiences, Research, and Applications, Lecture Notes in Geoinformation and Cartography, Springer, 2015, pp. 81–99.
- [4] M. Färber, F. Bartscherer, C. Menne, A. Rettinger, Linked Data Quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO, Semantic Web 9 (1) (2018) 77–129.
- [5] D. Vrandečić, M. Krötzsch, Wikidata: a free collaborative knowledgebase, Commun. ACM 57 (10) (2014) 78–85.
- [6] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. van Kleef, S. Auer, C. Bizer, DBpedia - A large-scale, multilingual knowledge base extracted from wikipedia, Semantic Web 6 (2) (2015) 167–195.
- [7] J. Hoffart, F. M. Suchanek, K. Berberich, G. Weikum, YAGO2: A spatially and temporally enhanced knowledge base from wikipedia, Artif. Intell. 194 (2013) 28–61.
- [8] S. Gottschalk, E. Demidova, EventKG - the hub of event knowledge on the web - and biographical timeline generation, Semantic Web 10 (6) (2019) 1039–1070.
- [9] D. Punjani, K. Singh, A. Both, M. Koubarakis, I. Angelidis, K. Bereta, T. Beris, D. Bilidas, T. Ioannidis, N. Karalis, C. Lange, D. Pantazi, C. Papaloukas, G. Stamoulis, Template-based question answering over linked geospatial data, in: Proceedings of the 12th Workshop on Geographic Information Retrieval, GIR@SIGSPATIAL 2018, ACM, 2018, pp. 7:1–7:10.
- [10] D. Herzog, S. Sikander, W. Wörndl, Integrating route attractiveness attributes into tourist trip recommendations, in: Com-

- panion of The 2019 World Wide Web Conference, WWW 2019, ACM, 2019, pp. 96–101.
- [11] A. Ganesh, Scaling OpenStreetMap with Wikidata knowledge, Blogpost, 2016, <https://blog.mapbox.com/scaling-openstreetmap-with-wikidata-knowledge-675d4495815f>.
- [12] M. Gritta, M. T. Pilehvar, N. Limsopatham, N. Collier, What’s missing in geographical parsing?, *Lang. Resour. Evaluation* 52 (2) (2018) 603–623.
- [13] C. Stadler, J. Lehmann, K. Höffner, S. Auer, LinkedGeoData: A core for a web of spatial open data, *Semantic Web* 3 (4) (2012) 333–354.
- [14] N. Karalis, G. M. Mandilaras, M. Koubarakis, Extending the YAGO2 knowledge graph with precise geospatial knowledge, in: *Proceedings of the 18th International Semantic Web Conference, ISWC 2019*, Vol. 11779 of *Lecture Notes in Computer Science*, Springer, 2019, pp. 181–197.
- [15] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: *Proceedings of the 27th Annual Conference on Neural Information Processing Systems 2013*, 2013, pp. 3111–3119.
- [16] M. F. Goodchild, Citizens as sensors: the world of volunteered geography, *GeoJournal* 69 (4) (2007).
- [17] R. Kimmel, A. Amir, A. M. Bruckstein, Finding shortest paths on surfaces using level sets propagation, *IEEE Trans. Pattern Anal. Mach. Intell.* 17 (6) (1995) 635–640.
- [18] A. Guttman, R-trees: A dynamic index structure for spatial searching, in: *SIGMOD’84, Proceedings of Annual Meeting*, ACM Press, 1984, pp. 47–57.
- [19] H. Paulheim, C. Bizer, Improving the quality of linked data using statistical distributions, *Int. J. Semantic Web Inf. Syst.* 10 (2) (2014) 63–86.
- [20] I. J. Goodfellow, Y. Bengio, A. C. Courville, *Deep Learning*, Adaptive computation and machine learning, MIT Press, 2016.
- [21] W. E. Winkler, The state of record linkage and current research problems, Tech. rep., Statistical Research Division, U.S. Bureau of the Census (1999).
- [22] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE: synthetic minority over-sampling technique, *J. Artif. Intell. Res.* 16 (2002) 321–357.
- [23] J. Bergstra, Y. Bengio, Random search for hyper-parameter optimization, *Journal of Machine Learning Research* 13 (2012) 281–305.
- [24] C. D. Manning, P. Raghavan, H. Schütze, *Introduction to Information Retrieval*, Cambridge University Press, 2008.
- [25] J. Daiber, M. Jakob, C. Hokamp, P. N. Mendes, Improving efficiency and accuracy in multilingual entity extraction, in: *Proceedings of the 9th International Conference on Semantic Systems, I-SEMANTICS 2013*, ACM, 2013, pp. 121–124.
- [26] M. A. Sherif, A. Ngonga Ngomo, J. Lehmann, Wombat - A generalization approach for automatic link discovery, Vol. 10249 of *Lecture Notes in Computer Science*, Springer, 2017, pp. 103–119.
- [27] G. Louppe, L. Wehenkel, A. Sutera, P. Geurts, Understanding variable importances in forests of randomized trees, in: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1, NIPS’13*, 2013, p. 431–439.
- [28] M. Nentwig, M. Hartung, A. N. Ngomo, E. Rahm, A survey of current link discovery frameworks, *Semantic Web* 8 (3) (2017) 419–436.
- [29] J. Volz, C. Bizer, M. Gaedke, G. Kobilarov, Silk - A link discovery framework for the web of data, in: *Proceedings of the WWW2009 Workshop on Linked Data on the Web, LDOW 2009*, Vol. 538 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2009.
- [30] A. N. Ngomo, S. Auer, LIMS - A time-efficient approach for large-scale link discovery on the web of data, in: *Proceedings of the 22nd International Joint Conference on Artificial Intelligence, IJCAI 2011, IJCAI/AAAI*, 2011, pp. 2312–2317.
- [31] W. Shen, J. Wang, J. Han, Entity linking with a knowledge base: Issues, techniques, and solutions, *IEEE Trans. Knowl. Data Eng.* 27 (2) (2015) 443–460.
- [32] A. Dessi, M. Atzori, A machine-learning approach to ranking rdf properties, *Future Generation Computer Systems* 54 (2016) 366 – 377.
- [33] S. D. Cardoso, F. K. Amanqui, K. J. Serique, J. L. dos Santos, D. A. Moreira, Swi: A semantic web interactive gazetteer to support linked open data, *Future Generation Computer Systems* 54 (2016) 389 – 398.
- [34] T. Saveta, G. Flouris, I. Fundulaki, A. N. Ngomo, Benchmarking link discovery systems for geo-spatial data, in: *Joint Proceedings of BLINK2017: 2nd International Workshop on Benchmarking Linked Data and NLIWoD3: Natural Language Interfaces for the Web of Data co-located with (ISWC2017)*, Vol. 1932 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2017.
- [35] A. F. Ahmed, M. A. Sherif, A. N. Ngomo, On the effect of geometries simplification on geo-spatial link discovery, in: *Proceedings of the 14th International Conference on Semantic Systems, SEMANTICS 2018*, Vol. 137 of *Procedia Computer Science*, Elsevier, 2018, pp. 139–150.
- [36] M. A. Sherif, K. Dreßler, P. Smeros, A. N. Ngomo, Radon - rapid discovery of topological relations, in: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017, AAAI Press, 2017, pp. 175–181.
- [37] P. Smeros, M. Koubarakis, Discovering spatial and temporal links among RDF data, in: *Proceedings of the Workshop on*

Linked Data on the Web, LDOW 2016, Vol. 1593 of CEUR Workshop Proceedings, CEUR-WS.org, 2016.

- [38] M. Xie, H. Yin, H. Wang, F. Xu, W. Chen, S. Wang, Learning graph-based POI embedding for location-based recommendation, in: Proceedings of the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016, ACM, 2016, pp. 15–24.
- [39] H. Wang, Z. Li, Region representation learning via mobility flow, in: Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, ACM, 2017, pp. 237–246.
- [40] D. Yang, B. Qu, J. Yang, P. Cudré-Mauroux, Revisiting user mobility and social relationships in lbsns: A hypergraph embedding approach, in: Proceedings of The World Wide Web Conference, WWW 2019, ACM, 2019, pp. 2147–2157.
- [41] W. Liu, J. Wang, A. K. Sangaiah, J. Yin, Dynamic metric embedding model for point-of-interest prediction, *Future Generation Computer Systems* 83 (2018) 183 – 192.
- [42] C. Ma, Y. Zhang, Q. Wang, X. Liu, Point-of-interest recommendation: Exploiting self-attentive autoencoders with neighbor-aware influence, in: Proceedings of the 27th ACM International Conference on Information and Knowledge Management, CIKM 2018, ACM, 2018, pp. 697–706.
- [43] L. Huang, Y. Ma, Y. Liu, A. K. Sangaiah, Multi-modal bayesian embedding for point-of-interest recommendation on location-based cyber-physical-social networks, *Future Generation Computer Systems* 108 (2020) 1119 – 1128.
- [44] S. Feng, G. Cong, B. An, Y. M. Chee, Poi2vec: Geographical latent representation for predicting future visitors, in: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI Press, 2017, pp. 102–108.
- [45] C. Yang, D. H. Hoang, T. Mikolov, J. Han, Place deduplication with embeddings, in: The World Wide Web Conference, WWW '19, ACM, 2019, p. 3420–3426.
- [46] M. Kejriwal, P. A. Szekely, Neural embeddings for populated geonames locations, in: Proceedings of the 16th International Semantic Web Conference, ISWC 2017, Vol. 10588 of Lecture Notes in Computer Science, Springer, 2017, pp. 139–146.