

LIMES - A Framework for Link Discovery on the Semantic Web

Axel-Cyrille Ngonga Ngomo, Mohamed Ahmed Sherif

Paderborn University, Data Science Group, Pohlweg 51, D-33098 Paderborn, Germany

E-mail: {firstname.lastname}@upb.de

Kleanthi Georgala, Mofeed Mohamed Hassan, Kevin Dreßler, Klaus Lyko, Daniel Obraczka, Tommaso Soru

University of Leipzig, Institute of Computer Science, AKSW Group, Augustusplatz 10, D-04009 Leipzig, Germany

E-mail: {lastname}@informatik.uni-leipzig.de

Abstract

The Linked Data paradigm builds upon the backbone of distributed knowledge bases connected by typed links. The mere volume of current knowledge bases as well as their sheer number pose two major challenges when aiming to support the computation of links across and within them. The first is that tools for link discovery have to be time-efficient when they compute links. Secondly, these tools have to produce links of high quality to serve the applications built upon Linked Data well. The current version of the LIMES framework is the product of seven years of research on these two challenges. The framework combines diverse algorithms for link discovery within a generic and extensible architecture. In this system paper, we give an overview of the 1.0 open-source release of the framework. In particular, we focus on an overview of the architecture of the framework, an intuition of its inner workings and a brief overview of the approaches it contains. Some of the applications within which the framework was used complete the paper. Our framework is open-source and available under a dual license at <http://github.com/aksw/limes-dev> together with a user manual and a developer manual.

Keywords: Link Discovery, Record Linking, System Description, Machine Learning, Semantic Web, OWL, RDF

1. Introduction

Establishing links between knowledge bases is one of the key steps of the Linked Data publication process.¹ A plethora of approaches has thus been devised to support this process [1]. In this paper, we present the LIMES framework, which was designed to accommodate a large number of link discovery approaches within a single extensible architecture. LIMES was designed as a declarative framework (i.e., a framework that processes link specifications, see Section 2) to address two main challenges.

1. *Time-efficiency*: The mere size of existing knowledge bases (e.g., 30+ billion triples in LinkedGeoData [2], 20+ billion triples in LinkedTCGA [3]) makes efficient solutions indispensable to the use of link discovery frameworks in real application scenarios. LIMES addresses this challenge by providing time-efficient approaches based on the characteristics of metric spaces [4, 5], orthodromic spaces [6] and on filter-based paradigms [7].
2. *Accuracy*: Efficient solutions are of little help if the results they generate are inaccurate. Hence, LIMES also accommodates (especially machine-learning) solutions that allow the generation of links between knowledge bases with a high accuracy. These solutions abide by paradigms

such as batch and active learning [8, 9, 10], unsupervised learning [10] and even positive-only learning [11].

The main goal of this paper is to give an overview of the LIMES framework and some of the applications within which it was used. We begin by presenting the link discovery problem and how we address this problem within a declarative setting (see Section 2). Then, we present our solution to the link discovery problem in the form of LIMES and its architecture. The subsequent sections present the different families of algorithms implemented within the framework. We begin by a short overview of the algorithms to ensure the efficiency of the framework (see Section 5). Thereafter, we give an overview of algorithms that address the accuracy problem (see Section 6). We round up the core of the paper with some of the applications within which LIMES was used, including benchmarking and the publication of 5-star linked datasets (Section 7). A brief overview of related work (Section 8), followed by an overview of the evaluation results of algorithms included in LIMES (Section 9) and a conclusion (Section 10) complete the paper.

2. The Link Discovery Problem

The formal specification of Link Discovery (LD) adopted herein is akin to that proposed in [12]. Given two (not necessarily distinct) sets S resp. T of source resp. target resources as well as a relation R , the goal of LD is to find the set $M = \{(s, t) \in$

¹<http://www.w3.org/DesignIssues/LinkedData.html>

$S \times T : R(s, t)$ of pairs $(s, t) \in S \times T$ such that $R(s, t)$. In most cases, computing M is a non-trivial task. Hence, a large number of frameworks (e.g., SILK [13], LIMES [12] and KnoFuss [14]) aim to approximate M by computing the *mapping* $M' = \{(s, t) \in S \times T : \sigma(s, t) \geq \theta\}$, where σ is a similarity function and θ is a similarity threshold. For example, one can configure these frameworks to compare the dates of birth, family names and given names of persons across census records to determine whether they are duplicates. We call the equation which specifies M' a *link specification* (short LS; also called linkage rule in the literature, see e.g., [13]). Note that the LD problem can be expressed equivalently using distances instead of similarities. As follows: Given two sets S and T of instances, a (complex) distance measure δ and a distance threshold $\theta \in [0, \infty[$, determine $M' = \{(s, t) \in S \times T : \delta(s, t) \leq \tau\}$.² Consequently, distance and similarities are used interchangeably within this paper.

Under this so-called *declarative paradigm*, two entities s and t are then considered to be linked by R if $\sigma(s, t) \geq \theta$. Naïve algorithms require $O(|S||T|) \in O(n^2)$ computations to output M' . Given the large size of existing knowledge bases, time-efficient approaches able to reduce this runtime are hence a central component of LIMES. In addition, note that the choice of appropriate σ and θ is central to ensure that M is approximated correctly by M' . Approaches that allow the computation of accurate σ and θ are thus fundamental for the LIMES framework [1].

3. Formal Overview

Several approaches can be chosen when aiming to define the syntax and semantics of LSs in detail [9, 13, 14]. In LIMES, we chose a grammar with a syntax and semantics based on set semantics. This grammar assumes that LSs consist of two types of atomic components:

- *similarity measures* m , which allow the comparison of property values or portions of the concise bound description of 2 resources and
- *operators* op , which can be used to combine these similarities into more complex specifications.

Without loss of generality, we define an *atomic similarity measure* a as a function $a : S \times T \rightarrow [0, 1]$. An example of an atomic similarity measure is the edit similarity dubbed `edit`³. Every atomic measure is a measure. A complex measure m combines measures m_1 and m_2 using measure operators such as `min` and `max`. We use *mappings* $M \subseteq S \times T$ to store the results of the application of a similarity measure to $S \times T$ or subsets thereof. We denote the set of all mappings as \mathcal{M} .

²Note that a distance function δ can always be transformed into a normed similarity function σ by setting $\sigma(x, y) = (1 + \delta(x, y))^{-1}$. Hence, the distance threshold τ can be transformed into a similarity threshold θ by means of the equation $\theta = (1 + \tau)^{-1}$.

³We define the edit similarity of two strings s and t as $(1 + lev(s, t))^{-1}$, where lev stands for the Levenshtein distance.

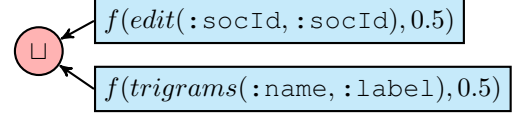


Figure 1: Example of a complex LS. The filter nodes are rectangles while the operator nodes are circles. `:socId` stands for social security number.

LS	$[[LS]]_M$
$f(m, \theta)$	$\{(s, t) (s, t) \in M \wedge m(s, t) \geq \theta\}$
$L_1 \sqcap L_2$	$\{(s, t) (s, t) \in [[L_1]]_M \wedge (s, t) \in [[L_2]]_M\}$
$L_1 \sqcup L_2$	$\{(s, t) (s, t) \in [[L_1]]_M \vee (s, t) \in [[L_2]]_M\}$
$L_1 \setminus L_2$	$\{(s, t) (s, t) \in [[L_1]]_M \wedge (s, t) \notin [[L_2]]_M\}$

Table 1: Link Specification Syntax and Semantics

We define a *filter* as a function $f(m, \theta)$. We call a specification *atomic* when it consists of exactly one filtering function. A complex specification can be obtained by combining two specifications L_1 and L_2 through an *operator* that allows the merging of the results of L_1 and L_2 . Here, we use the operators \sqcap , \sqcup and \setminus as they are complete w.r.t. the Boolean algebra and frequently used to define LS. An example of a complex LS is given in Figure 1. We denote the set of all LS as \mathcal{L} .

We define the semantics $[[L]]_M$ of a LS L w.r.t. a mapping M as given in Table 1. Those semantics are similar to those used in languages like SPARQL, i.e., they are defined extensionally through the mappings they generate. The mapping $[[L]]$ of a LS L with respect to $S \times T$ contains the links that will be generated by L .

4. Architecture of the Framework

4.1. Layers

As shown in Figure 2, the LIMES framework consists of 6 main layers, each of which can be extended to accommodate new or improved functionality. The input to the framework is a *configuration file*, which describes how the sets S and T are to be retrieved from two knowledge bases K_1 and K_2 (e.g., remote SPARQL endpoints). Moreover, the configuration file describes how links are to be computed. To this end, the user can choose to either provide a LS explicitly or a configuration for a particular machine learning algorithm. The result of the framework is a (set of) mapping(s). In the following, we give an overview of the inner workings of each of the layers.

4.1.1. Controller Layer

The processing of the configuration file by the different layers is orchestrated by the *controller layer*. The controller instantiates the different implementations of input modules (e.g., reading data from files or from SPARQL endpoints), the data modules (e.g., file cache, in-memory), the execution modules (e.g., planners, number of processors for parallel processing) and the output modules (e.g., the serialization format and the output location) according to the configuration file or using default values. Once the layers have been instantiated, the configuration is forwarded to the input layer.

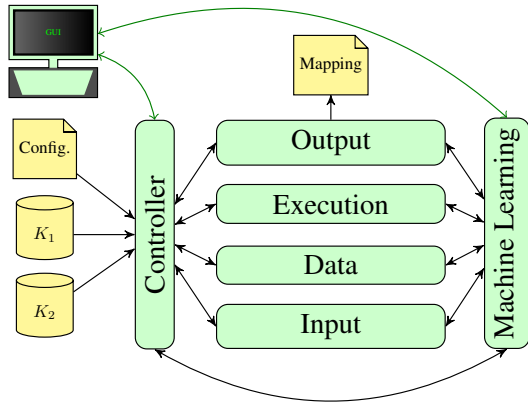


Figure 2: General architecture of LIMES.

4.1.2. Input Layer

The input layer reads the configuration and extracts all the information necessary to execute the specification or the machine learning approach chosen by the user. This information includes (1) the location of the input knowledge bases K_1 and K_2 (e.g., SPARQL endpoints or files), (2) the specification of the sets S and T , (3) the measures and thresholds or the machine learning approach to use. The current version of LIMES supports RDF configuration files based on the LIMES Configuration Ontology (LCO)⁴ (see Figure 3) and XML configuration files based on the LIMES Specification Language (LSL) [5]. If the configuration file is valid (w.r.t. the LCO or LSL), the input layer then calls its *query module*. This module uses the configuration for S and T to retrieve instances and properties from the input knowledge bases that adhere to the restrictions specified in the configuration file. All data retrieved is then forwarded to the data layer via the controller.

4.1.3. Data Layer

The *data layer* stores the input data gathered by the input layer by using memory, file or hybrid storage techniques. Currently, LIMES relies on a *hybrid cache* as default storage implementation. This implementation generates a hash for any set of resources specified in a configuration file and serializes this set into a file. Whenever faced with a new specification, the cache first determines whether the data required for the computation is already available locally by using the hash aforementioned. If the data is available, it is retrieved from the file into memory. In other cases, the hybrid cache retrieves the data as specified, generates a hash and caches it on the hard drive. By these means, the data layer addresses the practical problem of data availability, especially from remote data sources. Once all data is available, the controller layer chooses (based on the configuration file) between execution of complex measures specified by the user of running a machine learning algorithm.

⁴<https://github.com/AKSW/LIMES-dev/blob/master/limes-core/resources/lco.owl>

4.1.4. Execution Layer

If the user specifies the LS to execute manually, the LIMES controller layer calls the *execution layer*. The execution layer then re-writes, plans and executes the LS specified by the user. The *rewriter* aims to simplify complex LS by removing potential redundancy so as to eventually speed up the overall execution. To this end, it exploits the subset relations and the Boolean algebra which underpin the syntax and semantics of complex specifications. The *planner module* maps a LS to an equivalent execution plan, which is a sequence of atomic execution operations from which a mapping results. These operations include `EXECUTE` (running a particular atomic similarity measure, e.g., the trigram similarity) and operations on mappings such as `UNION` (union of two mappings) and `DIFF` (difference of two mappings). For an atomic LS, i.e., a LS such that σ is an atomic similarity measure (e.g., Levenshtein similarity, qgrams similarity, etc.), the planner module generates a *simple plan* with one `EXECUTE` instruction. For a complex LS, the planner module determines the order in which atomic LS should be executed as well as the sequence in which intermediary results should be processed. For example, the planner module may decide to first run some atomic LS L_1 and then filter the results using another atomic LS L_2 instead of running both L_1 and L_2 independently and merge the results. The *execution engine* is then responsible for executing the plan of an input LS and returns the set of potential links as a mapping.

4.1.5. Machine Learning layer

In case the user opts for using a machine learning approach, the LIMES controller calls the *machine learning layer*. This layer instantiates the machine learning algorithm selected by the user and executes it in tandem with the execution layer. To this end, the machine learning approaches begin by retrieving the training data provided by the user if required. The approaches then generate LSs, which are run in the execution layer. The resulting mappings are evaluated using quality measures (e.g., F-measure, pseudo-F-measure [10]) and returned to the user once a termination condition has been fulfilled. Currently, LIMES implements the EAGLE [9], COALA [15], EUCLID [10] and WOMBAT [11] algorithms. A brief overview of these approaches is given in Section 5. More details on the algorithms can be found in the corresponding papers.

4.1.6. Output Layer

The *output layer* serializes the output of LIMES. Currently, it supports the serialization into files in any RDF serialization also chosen by the user. In addition, the framework support CSV and XML as output formats.

4.2. Graphical User Interface

In addition to supporting configurations as input files, LIMES provides a graphical user interface (GUI) to assist the end user during the LD process [16]. The framework supports the manual creation of LS as shown in Figure 4 for users who already know which LS they would like to execute to achieve a certain linking task. However, most users do not know which LS suits their linking task best and therefore need help throughout this

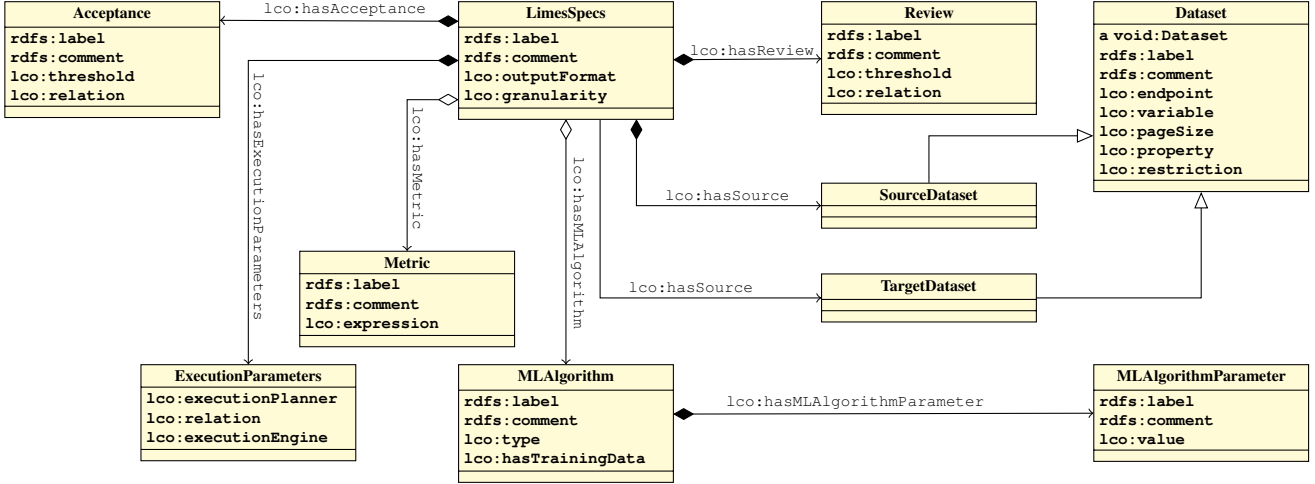


Figure 3: LCO ontology

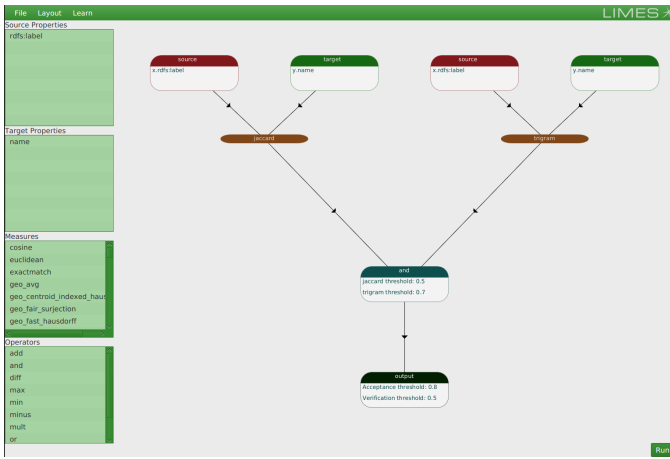


Figure 4: Link Specification Interface

Property	Value	Property	Value
rdfs:label	isopropamide	name	isopropamide
Source URI			
Target URI			
value			
Match?			
http://dbpedia.org/resource/Butoconazole	http://dbpedia.org/resource/Butoconazole	3.7442161...	<input checked="" type="checkbox"/>
http://dbpedia.org/resource/Flecainide	http://dbpedia.org/resource/Flecainide	3.7466641...	<input type="checkbox"/>
http://dbpedia.org/resource/Tamsulosin	http://dbpedia.org/resource/Tamsulosin	3.74544E7	<input type="checkbox"/>
http://dbpedia.org/resource/Onchidal	http://dbpedia.org/resource/Onchidal	3.7478884...	<input type="checkbox"/>
http://dbpedia.org/resource/Isopropamide	http://dbpedia.org/resource/Isopropamide	3.7491129...	<input checked="" type="checkbox"/>
http://dbpedia.org/resource/Anagestone	http://dbpedia.org/resource/Anagestone	3.7515625...	<input type="checkbox"/>
http://dbpedia.org/resource/Suritozole	http://dbpedia.org/resource/Suritozole	3.7503376...	<input type="checkbox"/>
http://dbpedia.org/resource/Diphenamil_m...	http://dbpedia.org/resource/Diphenamil_m...	3.7527876...	<input type="checkbox"/>
http://dbpedia.org/resource/Chromium(III)	http://dbpedia.org/resource/Chromium(III)	3.7552384...	<input type="checkbox"/>

Figure 5: User feedback interface for active learning

process. Hence, the GUI provides wizards which ease the LS creation process and allow configuration of the machine learning algorithms. After the selection of an algorithm, the user can modify the default configurations if necessary and run it. If a batch or unsupervised learning approach is selected, the execution is carried out once and the results are presented to the user. In the case of active learning, the LIMES GUI presents the most highly informative link candidates to the user, who must now label these candidates as either a match or non-match (see Figure 5). This procedure is repeated until the user is satisfied with the result. In all cases, the GUI allows the exportation of the configurations generated by machine learning approaches for future use.

5. Scalability Algorithms

As mentioned in the introduction, the scalability of LIMES is based on a number of time-efficient approaches for atomic LS that the framework supports. In the following, we summarize the approaches supported by LIMES version 1.0.0. By no means do we aim to give all technical details necessary to understand how

these approaches work in detail. Interested readers are referred to the publications mentioned in each subsection.

5.1. LIMES

The LIMES algorithm [4] is the first scalable approach developed explicitly for LD. The basic intuition behind the approach is based on regarding the LD problem as a bounded distance problem $\delta(s, t) \leq \tau$. If δ abides by the triangle inequality, the $\delta(s, t) \geq \delta(s, e) - \delta(e, t)$ for any e . Based on this insight, LIMES implements a space tiling approach that assigns each t to an exemplar e such that $\delta(e, t)$ is known for many t . Hence, for all s , LIMES can first compute a lower bound for the distance between s and t . If this lower bound is larger than τ , then the exact computation of $\delta(s, t)$ is not carried out, as (s, t) is guaranteed not to belong in the output of this atomic measure.

5.2. \mathcal{HR}^3

The \mathcal{HR}^3 algorithm [5] builds upon some of the insights of the LIMES algorithms within spaces with *Minkowski distances*. The basic insight behind the approach is that the equation $\delta(s, t) \leq \tau$ describes a hypersphere of radius τ around s . While

computing spheres in multidimensional spaces can be time-consuming, hypercubes can be computed rapidly. Hence, the approach relies on the approximation hyperspheres with hypercubes. In contrast to previous approaches such as HYPPO [17], the approach combines hypercubes with an indexing approach and can hence discard a larger number of comparisons. One of the most important theoretical insights behind this approach is that it is the first approach guaranteed to being able to achieve the smallest possible number of comparisons when given sufficient memory. It is hence the first *reduction-ratio-optimal* approach for link discovery.

5.3. ORCHID

ORCHID [6] is a link discovery algorithm for geo-spatial data. Like \mathcal{HR}^3 , ORCHID is reduction-ratio-optimal. In contrast to \mathcal{HR}^3 , it operates in orthodromic spaces, in which Minkowski distances do not hold. To achieve the goal of reduction-ratio optimality, ORCHID uses a space discretization approach which relies on tiling the surface of the planet based on latitude and longitude information. The approach then only compares polygons $t \in T$ which lie within a certain range of $s \in S$. Like in \mathcal{HR}^3 , the range is described via a set of hypercubes (i.e., squares in 2 dimensions). However, the shape that is to be approximated is a disk projected onto the surface of the sphere. The projection is accounted for by compensating for the curvature of the surface of the planet using an increased number of squares. See [18] for a survey of point-set distance measures implemented by LIMES and optimized using ORCHID.

5.4. AEGLE

AEGLE [19] is a time-efficient approach for computing temporal relations. Our approach supports all possible temporal relations between event data that can be modelled as intervals according to Allen’s Interval Algebra [20]. The key idea behind the algorithm is to reduce the 13 Allen relations to 8 simpler relations, which compare either the beginning or the end of an event with the beginning or the end of another event. Given that time is ordered, AEGLE reduces the problem of computing these simple relations to the problem of matching entities across sorted lists, which can be carried out within a time complexity of $O(n \log n)$. The approach is guaranteed to compute complete results for any temporal relation.

5.5. RADON

RADON [21] addresses the efficient computation of topological relations on geo-spatial datasets, which belong to the largest sources of Linked Data. The main innovation of the approach is a novel sparse index for geo-spatial resources based on *minimum bounding boxes* (MBB). Based on this index, it is able to discard unnecessary computations for DE-9IM relations.

Even small datasets can contain very large geometries that span over a large number of hypercubes and would lead to a large spatial index when used as source. Therefore, RADON applies a *swapping strategy* as its first step, where it swaps source and

target datasets and computes the reverse⁵ relation r' instead of r in case it finds that the *estimated total hypervolume* of the target is less than the one of the source dataset. Then in its second step, RADON utilizes a space tiling approach to insert all source and target geometries into an index I , which maps resources to sets of hypercubes. Finally, RADON implements the last speedup strategy using a MBB-based filtering technique. Like AEGLE, RADON is guaranteed to achieve a result completeness of 100% as it is able to compute all topological relations of the DE-9IM model.

5.6. ROCKER

Keys for graphs are the projection of the concept of primary keys from relational databases. A key is thus a set of properties such that all instances of a given class are distinguishable from each other by their property values. LIMES features ROCKER [22], a refinement operator-based approach for key discovery. The ROCKER algorithm first streams the input RDF graph to build a hash index including instance URIs and a concatenation of object values for each property. The indexes are then stored in a local database to enable an efficient computation of the discriminability score. Starting from the empty set, a refinement graph is visited, adding a property to the set at each refinement step. As keys abide by several monotonicities, additional optimizations allow scalability to large datasets. ROCKER has shown state-of-the-art results in terms of discovered keys, efficiency, and memory consumption [22].

5.7. REEDED

LD frameworks rely on similarity measures such as the Levenshtein (or edit) distance to discover similar instances. The edit distance calculates the number of operations (i.e., addition, deletion or replacement of a character) needed to transform a string to another. However, the assumption that all operations must have the same weight does not always apply. For instance, the replacement of a *z* with an *s* often leads to the same word written in American and British English, respectively. This operation should necessarily weigh less than, e.g., an addition of the character *I* to the end of *King George I*. While manifold algorithms have been developed for the efficient discovery of similar strings using edit distance [23, 24], REEDED was developed to support also weighted edit distances. REEDED joins two sets of strings applying three filters to avoid computing the similarity values for all the pairs. The (1) length-aware filter discards the pairs of strings having a too different length, the (2) character-aware filter selects only the pairs which do not differ by more than a given number of characters and the (3) verification filter finally verifies the weighted edit distance among the pairs against a threshold [7].

5.8. Jaro-Winkler

The Jaro-Winkler string similarity measure, which was originally designed for the deduplication of person names, is a

⁵Formally, the reverse relation r' of a relation r is defined as $r'(y, x) \Leftrightarrow r(x, y)$.

common choice to generate links between knowledge bases based on labels. We therefore developed a runtime-optimized Jaro-Winkler algorithm for LD [25]. The approach consists of two steps: indexing and tree-pruning. The indexing step itself also consists of two phases: (1) strings of both source and target datasets get sorted into buckets based on upper and lower bounds pertaining to their lengths and (2) in each bucket, strings of the target dataset get indexed by adding them to a tree akin to a *trie* [26]. Tree pruning is then applied to cut off subtrees for which the similarity threshold can not be achieved due to character mismatches. This significantly reduces the number of necessary Jaro-Winkler similarity computations.

5.9. HELIOS

In contrast to the algorithms above, HELIOS [27] addresses the scalability of link discovery by improving the planning of LSs. To this end, HELIOS implements a rewriter and a planner. The rewriter consists of a set of algebraic rules, which can transform portions of LS into formally equivalent expressions, which can be potentially executed faster. For example, expressions such as $AND(\sigma(s, t) \geq \theta_1, \sigma(s, t) \geq \theta_2)$ are transformed into $\sigma(s, t) \geq \max(\theta_1, \theta_2)$. In this particular simple example, the number of EXECUTE calls can hence be reduced from 2 to 1. The planner transforms a rewritten LS into a plan by aiming to find a sequence of execution steps that minimizes the total runtime of the LS. To this end, the planner relies on runtime approximations derived from linear regressions on large static dictionaries.

6. Accuracy Algorithms

In addition to pushing for efficiency, the LIMES framework also aims to provide approaches that compute links accurately. To this end, the framework includes a family of machine learning approaches designed specifically for the purpose of link discovery. The core algorithms can be adapted for unsupervised, active and batch learning. In the following, we present the main ideas behind these approaches while refraining from providing complex technical details, which can be retrieved from the corresponding publications.

6.1. RAVEN

RAVEN [8] is the first approach designed to address the active learning of LS. The basic intuition behind the approach is that LS can be regarded as classifiers. Links that abide by the LS then belong to the class +1 while all other links belong to -1. Based on this intuition, the approach tackles the problem of learning LSs by first using the solution of the hospital-resident problem to detect potentially matching classes and properties in K_1 and K_2 . The algorithm then assumes that it knows the type of LS that is to be learned (e.g., conjunctive specifications, in which all specifications operators are conjunctions). Based on this assumption and the available property matching, it learns the thresholds associated with each property pair iteratively by using an extension of the perceptron learning paradigm. RAVEN first applies perceptron learning to all training data that is made

available by the user. Then, most informative unlabeled pairs (by virtue of their proximity to the decision boundary of the perceptron classifier) are sent as queries to the user. The labeled answers are finally added to the training data, which closes the loop. The iteration is carried on until the user terminates the process or a termination condition such as a maximal number of questions is reached.

6.2. EAGLE

EAGLE [9] implements a machine learning approach for LS of arbitrary complexity based on genetic programming. To this end, the approach models LS as trees, each of which stands for the genome of an individual in a population. EAGLE begins by generating a population of random LS, i.e., of individuals with random genomes. Based on training data or on an objective function, the algorithm determines the fitness of the individuals in the population. Operators such as mutation and crossover are then used to generate new members of the population. The fittest individuals are finally selected as members of the next population. When used as an active learning approach, EAGLE relies on a committee-based algorithm to determine most informative positive and negative examples. In addition to the classical entropy-based selection, the COALA [15] extension of EAGLE also allows the correlation between resources to be taken into consideration during the selection process for most information examples. All active learning versions of EAGLE ensure that each iteration of the approach only demands a small number of labels from the user. EAGLE is also designed to support the unsupervised learning approach of LS. In this case, the approach aims to find individuals that maximize so-called *pseudo F-Measures* [10, 28].

6.3. ACIDS

The ACIDS approach [29] targets instance matching by joining active learning with classification using linear Support Vector Machines (SVM). Given two instances s and t , the similarities among their datatype values are collected in a feature vector of size N . A pair (s, t) is represented as a point in the similarity space $[0, 1]^N$. At each iteration, pairs are (1) selected based on their proximity to the SVM classifier, (2) labeled by the user, and (3) added to the SVM training set. Mini-batches of pairs are labeled as positive if the instances are to be linked with an R , or negative otherwise. The SVM model is then built after each step, where the classifier is a hyperplane dividing the similarity space into two subspaces. String similarities are computed using a weighted Levenshtein distance. Using an update rule inspired by perceptron learning, edit operation weights are learned by maximizing the training F-score.

6.4. EUCLID

EUCLID [10] is a deterministic learning algorithm loosely based on RAVEN. EUCLID reuses the property matching techniques implemented in RAVEN and the idea of having known classifier shapes. In contrast to previous algorithms, EUCLID learns combinations of atomic LS by using a hierarchical search

approach. While EUCLID was designed for unsupervised learning, it can also be used for supervised learning based on labeled training data[10].

6.5. WOMBAT

One of the most crucial tasks when dealing with evolving datasets lies in updating the links from these data sets to other data sets. While supervised approaches have been devised to achieve this goal, they assume the provision of both positive and negative examples for links [30]. However, the links available on the Data Web only provide positive examples for relations and no negative ones, as the *open-world assumption* underlying the Web of Data suggests that the non-existence of a link between two resources cannot be understood as stating that these two resources are not related. In LIMES, we addressed this drawback by proposing WOMBAT [11], the first approach for learning links based on positive examples only. WOMBAT is inspired by the concept of generalisation in quasi-ordered spaces. Given a set of positive examples, it aims to find a classifier that covers a large number of positive examples (i.e., achieves a high recall on the positive examples) while still achieving a high precision. The simple version of WOMBAT relies on a two-step approach, which learns atomic LS and subsequently finds ways to combine them to achieve high F-measures. The complete version of the algorithm relies on a full-fledged upward refinement operator, which is guaranteed to find the best specification possible but scales less well than the simple version.

7. Use Cases

LIMES was already used in a large number of use cases. In the following, we present some of the datasets and other applications, in which techniques developed for LIMES or the LIMES framework itself, were used successfully.

7.1. Semantic Quran

The SEMANTICQURAN dataset [31] is a multilingual RDF representation of translations of the *Quran*. The dataset was created by integrating data from two different semi-structured sources and aligning them to an ontology designed to represent multilingual data from sources with a hierarchical structure. The resulting RDF data encompasses 43 different languages which belong to the most under-represented languages in the Linked Data Cloud, including Arabic, Amharic and Amazigh. Using LIMES, we were able to generate links from SEMANTICQURAN to 3 versions of the RDF representation of *Wiktionary* and to *DBpedia*. The basic intuition behind our linking of SEMANTICQURAN was to link words in a given language in SEMANTICQURAN to words in the same language with exactly the same label. We provide 7617 links to the English version of *DBpedia*, which in turn is linked to non-English versions of *DBpedia*. In addition, we generated 7809 links to the English, 9856 to the French and 1453 to the German *Wiktionary*.

7.2. Linked TCGA

A second dataset that was linked with LIMES is LINKEDTCGA [3], an RDF representation of The Cancer Genome Atlas. The data source describes cancer patient data gathered from 9000 patients for 30 different cancer diseases. The dataset aims to support bio-medical experts who work on characterizing the changes that occur in genes due to cancer. LINKEDTCGA is one of the largest knowledge bases on the Linked Data Web and accounts for more than 20 billion triples. The linking task for this dataset consisted of connecting it with the biomedical portion of the Linked Open Data Cloud, in particular to Bio2RDF.⁶ LIMES was used to compute the more than 16 million links which connect LINKEDTCGA with chromosomes from OMIM and HGNC. The more than 500 thousand links that connect LINKEDTCGA with Homologene and similar knowledge bases were also computed using LIMES.

7.3. Knowledge Base Repair

The COLIBRI [28] approach attempts to repair instance knowledge in n knowledge bases. COLIBRI discovers links for transitive relations (e.g., `owl:sameAs`) between instances in knowledge bases while correcting errors in the same knowledge bases. In contrast to most of the existing approaches, COLIBRI takes an n -set⁷ of set of resources K_1, \dots, K_n with $n \geq 2$ as input. Thereafter, COLIBRI relies on the LIMES's unsupervised deterministic approach EUCLID (see Section 6.4) to link each pair (K_i, K_j) of sets of resources (with $i \neq j$). The resource mappings resulting from the link discovery are then forwarded to a voting approach, which is able to detect poor mappings. This information is subsequently used to find sources of errors in the mappings, such as erroneous or missing information in the instances. The errors in the mappings are finally corrected and the approach starts the next linking iteration.

7.4. Question Answering

LIMES was also used to create the knowledge base behind the question answering engine DEQA [32]. This engine was designed to be a comprehensive framework for deep web question answering. To this end, the engine provides functionality for the extraction of structured knowledge from (1) unstructured data sources (e.g., using FOX [33] and AGDISTIS [34]) and (2) the deep Web using the XPath language.⁸ The results of the extraction are linked with other knowledge sources on the Web using LIMES. In the example presented in [32], geo-spatial links such as `nearBy` were generated using LIMES' approach \mathcal{HR}^3 . The approach was deployed on a dataset combining data on flats in Oxford and geo-spatial information from LinkedGeoData and was able to support the answering of complex questions such as "Give me a flat with 2 bedrooms near to a school"

⁶<http://bio2rdf.org>

⁷An n -set is a set of magnitude n .

⁸<http://www.xpath.org/>

7.5. Benchmarking

While the creation of datasets for the Linked Data Web is an obvious use of the framework, LIMES and its algorithms can be used for several other purposes. A non-obvious application is the generation of benchmarks based on query logs [35, 36]. The DBpedia SPARQL benchmark [35] relies on DBpedia query logs to detect clusters of similar queries, which are used as templates for generating queries during the execution of the benchmark. Computing the similarity between the millions of queries in the DBpedia query log proved to be an impractical endeavor when implemented in a naïve fashion. With LIMES, the runtime of the similarity computations could be reduced drastically by expressing the problem of finding similar queries as a near-duplicate detection problem.

The application of LIMES in FEASIBLE [36] was also in the clustering of queries. The approach was designed to enable users to select the number of queries that their benchmark should generate. The exemplar-based approach from LIMES [4] was used as a foundation for detecting the exact number of clusters are required by the user while optimizing the homogeneity of the said clusters. The resulting benchmark was shown to outperform the state of the art in how well it encompasses the idiosyncrasies of the input query log.

7.6. Dataset Enrichment

Over the recent years, a few frameworks for RDF data enrichment such as LDIF⁹ and DEER¹⁰ have been developed. The frameworks provide enrichment methods such as entity recognition [33], link discovery [12] and schema enrichment [37]. However, devising appropriate configurations for these tools can prove a difficult endeavour, as the tools require the right sequence of enrichment functions to be chosen and these functions to be configured adequately. In DEER [38], we address this problem by presenting a supervised machine learning approach for the automatic detection of enrichment pipelines based on a refinement operator and self-configuration algorithms for enrichment functions. The output of DEER is an enrichment pipeline that can be used on whole datasets to generate enriched versions.

The DEER's *linking* enrichment function is used to connect further datasets using link predicates such as `owl:sameAs`. With the proper configurations, the *linking* enrichment function can be used to generate arbitrary predicate types. For instance, using the `gn:nearby` for linking geographical resources. The idea of the self configuration of linking enrichment function is to (1) first perform an automatic property matching, then (2) based on LIMES' WOMBAT algorithm [10] (see Section 6.5), perform supervised link discovery.

8. Related Work

There has been a significant amount of work pertaining to the two key challenges of LD: efficiency and accuracy. Here,

we focus on frameworks for link discovery. Dedicated approaches and algorithms can be found in the related work sections of [4, 5, 12, 9, 10, 15, 19, 25, 21]. As described in [1], most LD frameworks address the efficiency challenge by aiming to discover unnecessary computations of the similarity or distance measures efficiently. The accuracy challenge is commonly addressed using machine learning techniques.

One of the first LD frameworks is SILK [39], which uses a multi-dimensional blocking technique (*MultiBlock* [13]) to optimize the linking runtime through a rough index pre-matching. To parallelize the linking process, SILK relies on *MapReduce*. Like LIMES, SILK configuration supports input files in XML and is also able to retrieve RDF data by querying SPARQL endpoints. Both frameworks also allow user-specified link types between resources as well as `owl:sameAs` links. SILK incorporates element-level matchers on selected properties using string, numeric, temporal and geo-spatial similarity measures. SILK also supports multiple matchers, as it allows the comparison of different properties between resources that are combined together using match rules. SILK also implements supervised and active learning methods for identifying LS for linking.

KNOFUSS [14] incorporates blocking approaches derived from databases. Like LIMES, it supports unsupervised machine learning techniques based on genetic programming. However, in contrast to both aforementioned tools, KNOFUSS supports no other link types apart from `owl:sameAs` and implements only string similarity measures between matching entities properties. Additionally, its indexing technique for time complexity minimization is not guaranteed to achieve result completeness.

ZHISHI.LINKS [40] is another LD framework that achieves efficiency by using an index-based technique that comes at the cost of effectiveness. Similar to KNOFUSS, it only supports `owl:sameAs` links but implements geo-spatial relations between resources. In comparison with all previous three tools, it supports both property-based and semantic-based matching, using knowledge obtained by an ontology.

SERIMI [41] is a LD tool that retrieves entities for linking by querying SPARQL endpoints. In contrast to the previous frameworks, it only supports single properties to be used for linking resources. However, it incorporates an adaptive technique that weighs differently the properties of a knowledge base and chooses the most representative property for linking. Furthermore, SLINT+ [42] is a LD tool that is very similar to SERIMI but supports comparisons between multiple properties.

Finally, there are a set of frameworks (RIMOM [43], AGREEMENTMAKER [44], LOGMAP [45], CODI [46]) that initially supported ontology matching and then evolved to support LD between resources. RIMOM is based on Bayesian decision matching in order to link ontologies and transforms the problem of linking into a decision problem. Even though RIMOM utilizes dictionaries for ontology matching, it does not support them in entity linking. Similar to RIMOM, AGREEMENTMAKER also uses dictionaries in ontology level. AGREEMENTMAKER incorporates a variety of matching methods based on different properties considered for comparison and the different granularity of the components. LOGMAP is an efficient tool for ontology matching that scales up to tens (even hundreds) of thousands

⁹<http://ldif.wbsg.de/>

¹⁰<http://aksw.org/Projects/DEER.html>

of different classes included in an ontology. It propagates the information obtained from ontology matching to link resources, by using logical reasoning to exclude links between resources that do not abide by the restrictions obtained from ontology matching. CODI is a probabilistic-logical framework for ontology matching based on the syntax and semantics of Markov logic. It incorporates typical matching approaches that are joined together to increase the quality of ontology alignment.

The basic difference between the previous LD framework and LINES is that LINES provides both theoretical and practical guarantees for efficient and scalable LD. LINES is guaranteed to lead to exactly the same matching as a brute force approach while at the same time reducing significantly the number of comparisons. The approaches incorporated in LINES facilitate different approximation techniques to compute estimates of the similarity between instances. These estimates are then used to filter out a large number of those instance pairs that do not suffice the mapping conditions. By these means, LINES can reduce the number of comparisons needed during the mapping process by several orders of magnitude. LINES supports the first planning technique for link discovery HELIOS, that minimizes the overall execution of a link specification, without any loss of completeness. As shown in our previous works (see [27, 19, 47]), LINES is one of the most scalable link discovery frameworks currently available.

9. Evaluation

The approaches included in LINES are the result of seven years of research, during which the state of the art evolved significantly. In Table 2, we give an overview of the performance improvement (w.r.t. runtime and/or accuracy) of a selection of algorithms currently implemented in LINES. The improvements mentioned refer to improvements w.r.t. the state of art at the time at which the papers were written. We refer the readers to the corresponding research paper for the evaluation settings and the experimental results for each of the algorithms mentioned in the table.

10. Conclusion and Future Work

In LINES, we offer an extensible framework for the discovery of links between knowledge bases. The framework has already been used in several applications and shown to be a reliable, near industry-grade framework for link discovery. The graphical user interface and the manuals for users¹¹ and developers¹², which accompany the tool, make it a viable framework for novices and experts that aims to create and deploy linked data sets. While a number of challenges have been addressed in the framework, the increasing amount of Linked Data available on the Web and the large number of applications that rely on it demand that we address ever more challenges over the years to come. The intelligent use of memory (disk, RAM, etc.)

Table 2: Evaluation of various algorithms used in LINES. The numbers provided were retrieved from the experiments and results of the corresponding papers. These numbers are not upper bounds but merely reflect improvements observed empirically on particular datasets. Orders of magnitude are abbreviated to o.o.m..

Algorithm	Improvement		Sec.	Ref.
	Runtime	F-Measure		
LINES	Up to 60-fold	–	5.1	[4]
\mathcal{HR}^3	Up to 7%	–	5.2	[5]
ORCHID	Up to 2 o.o.m.	–	5.3	[6]
AEGLE	Up to 4 o.o.m.	Correct (100%)	5.4	[19]
RADON	Up to 3 o.o.m.	Correct (100%)	5.5	[21]
ROCKER	Up to 1 o.o.m.	Correct (100%)	5.6	[22]
REDED	Up to 15-fold	Correct (100%)	5.7	[7]
Jaro-Winkler	Up to 5.5-fold	–	5.8	[25]
HELIOS	Up to 300%	–	5.9	[27]
RAVEN	Up to 33%	90% – 100%	6.1	[8]
EAGLE	Up to 14-fold	> 90%	6.2	[9]
ACIDS	–	+7%	6.3	[29]
COALA	–	+25%	6.2	[15]
EUCLID	Up to 33%	+11%	6.4	[10]
WOMBAT	–	+23%	6.5	[11]
RAVEN	Up to 33%	90% – 100%	6.1	[8]
EAGLE	Up to 14-fold	Superior to 90%	6.2	[9]
ACIDS	–	Up to 7%	6.3	[29]
COALA	–	Up to 25%	6.2	[15]
EUCLID	Up to 33%	Up to 11%	6.4	[10]
WOMBAT	–	Up to 23%	6.5	[11]

becomes a problem of central importance when faced with small amounts of RAM that cannot fit the sets S and T [48] or when faced with streaming or complex data (e.g., 5D geospatial data). The appropriate use of hardware resources, which has already been shown to be of central importance for link discovery approaches [49, 50], remains an important challenge for the future. Big Data frameworks such as SPARK¹³ and FLINK¹⁴ promise to alleviate this problem when used correctly.

Acknowledgments

This work has been supported by H2020 projects SLIPO (GA no. 731581) and HOBBIT (GA no. 688227) as well as the DFG project LinkingLOD (project no. NG 105/3-2) and the BMWI Projects SAKE (project no. 01MD15006E) and GEISER (project no. 01MD16014).

References

- [1] M. Nentwig, M. Hartung, A. N. Ngomo, E. Rahm, A survey of current link discovery frameworks, *Semantic Web* 8 (3) (2017) 419–436. doi: 10.3233/SW-150210. URL <http://dx.doi.org/10.3233/SW-150210>
- [2] 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. doi:10.3233/SW-2011-0052. URL <http://dx.doi.org/10.3233/SW-2011-0052>

¹¹http://aksw.github.io/LINES-dev/user_manual/

¹²http://aksw.github.io/LINES-dev/developer_manual/

¹³<http://spark.apache.org/>

¹⁴<https://flink.apache.org/>

- [3] M. Saleem, S. S. Padmanabhuni, A. N. Ngomo, J. S. Almeida, S. Decker, H. F. Deus, Linked cancer genome atlas database, in: M. Sabou, E. Blomqvist, T. D. Noia, H. Sack, T. Pellegrini (Eds.), I-SEMANICS 2013 - 9th International Conference on Semantic Systems, ISEM '13, Graz, Austria, September 4-6, 2013, ACM, 2013, pp. 129–134. doi:10.1145/2506182.2506200.
URL <http://doi.acm.org/10.1145/2506182.2506200>
- [4] A. N. Ngomo, S. Auer, LINES - A time-efficient approach for large-scale link discovery on the web of data, in: T. Walsh (Ed.), IJCAI 2011, Proceedings of the 22nd International Joint Conference on Artificial Intelligence, Barcelona, Catalonia, Spain, July 16-22, 2011, IJCAI/AAAI, 2011, pp. 2312–2317. doi:10.5591/978-1-57735-516-8/IJCAI11-385.
URL <http://dx.doi.org/10.5591/978-1-57735-516-8/IJCAI11-385>
- [5] A. N. Ngomo, Link discovery with guaranteed reduction ratio in affine spaces with minkowski measures, in: The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference, Boston, MA, USA, November 11-15, 2012, Proceedings, Part I, 2012, pp. 378–393. doi:10.1007/978-3-642-35176-1_24.
URL http://dx.doi.org/10.1007/978-3-642-35176-1_24
- [6] A. N. Ngomo, ORCHID - reduction-ratio-optimal computation of geospatial distances for link discovery, in: Alani et al. [51], pp. 395–410. doi:10.1007/978-3-642-41335-3_25.
URL http://dx.doi.org/10.1007/978-3-642-41335-3_25
- [7] T. Soru, A. N. Ngomo, Rapid execution of weighted edit distances, in: Shvaiko et al. [52], pp. 1–12.
URL http://ceur-ws.org/Vol-1111/om2013_Tpaper1.pdf
- [8] A. N. Ngomo, J. Lehmann, S. Auer, K. Höffner, RAVEN - active learning of link specifications, in: Shvaiko et al. [53], pp. 25–36.
URL http://ceur-ws.org/Vol-814/om2011_Tpaper3.pdf
- [9] A. N. Ngomo, K. Lyko, EAGLE: efficient active learning of link specifications using genetic programming, in: The Semantic Web: Research and Applications - 9th Extended Semantic Web Conference, ESWC 2012, Heraklion, Crete, Greece, May 27-31, 2012. Proceedings, 2012, pp. 149–163. doi:10.1007/978-3-642-30284-8_17.
URL http://dx.doi.org/10.1007/978-3-642-30284-8_17
- [10] A. N. Ngomo, K. Lyko, Unsupervised learning of link specifications: deterministic vs. non-deterministic, in: Shvaiko et al. [52], pp. 25–36.
URL http://ceur-ws.org/Vol-1111/om2013_Tpaper3.pdf
- [11] M. Sherif, A.-C. Ngonga Ngomo, J. Lehmann, WOMBAT - A Generalization Approach for Automatic Link Discovery, in: 14th Extended Semantic Web Conference, Portorož, Slovenia, 28th May - 1st June 2017, Springer, 2017.
URL http://svn.aksw.org/papers/2017/ESWC_WOMBAT/public.pdf
- [12] A. N. Ngomo, On link discovery using a hybrid approach, J. Data Semantics 1 (4) (2012) 203–217. doi:10.1007/s13740-012-0012-y.
URL <http://dx.doi.org/10.1007/s13740-012-0012-y>
- [13] R. Isele, A. Jentzsch, C. Bizer, Efficient multidimensional blocking for link discovery without losing recall, in: Proceedings of the 14th International Workshop on the Web and Databases 2011, WebDB 2011, Athens, Greece, June 12, 2011, 2011.
URL <http://webdb2011.rutgers.edu/papers/Paper%2039/silk.pdf>
- [14] A. Nikolov, V. S. Uren, E. Motta, Knofuss: a comprehensive architecture for knowledge fusion, in: Proceedings of the 4th International Conference on Knowledge Capture (K-CAP 2007), October 28-31, 2007, Whistler, BC, Canada, 2007, pp. 185–186. doi:10.1145/1298406.1298446.
URL <http://doi.acm.org/10.1145/1298406.1298446>
- [15] A. N. Ngomo, K. Lyko, V. Christen, COALA - correlation-aware active learning of link specifications, in: The Semantic Web: Semantics and Big Data, 10th International Conference, ESWC 2013, Montpellier, France, May 26-30, 2013. Proceedings, 2013, pp. 442–456. doi:10.1007/978-3-642-38288-8_30.
URL http://dx.doi.org/10.1007/978-3-642-38288-8_30
- [16] A.-C. Ngonga Ngomo, D. Obraczka, K. Georgala, Using machine learning for link discovery on the web of data, in: Demos at the European Conference on Artificial Intelligence, 2016.
- [17] A. N. Ngomo, A time-efficient hybrid approach to link discovery, in: Shvaiko et al. [53].
URL http://ceur-ws.org/Vol-814/om2011_Tpaper1.pdf
- [18] M. A. Sherif, A.-C. N. Ngomo, A Systematic Survey of Point Set Distance Measures for Link Discovery, Semantic Web Journal.
URL <http://www.semantic-web-journal.net/content/systematic-survey-point-set-distance-measures-link-discovery-1>
- [19] K. Georgala, M. A. Sherif, A. N. Ngomo, An efficient approach for the generation of allen relations, in: ECAI 2016 - 22nd European Conference on Artificial Intelligence, 29 August-2 September 2016, The Hague, The Netherlands - Including Prestigious Applications of Artificial Intelligence (PAIS 2016), 2016, pp. 948–956. doi:10.3233/978-1-61499-672-9-948.
URL <http://dx.doi.org/10.3233/978-1-61499-672-9-948>
- [20] J. F. Allen, Maintaining knowledge about temporal intervals, Commun. ACM 26 (11) (1983) 832–843. doi:10.1145/182.358434.
URL <http://doi.acm.org/10.1145/182.358434>
- [21] 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, February 4-9, 2017, San Francisco, California, USA., 2017, pp. 175–181.
URL <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14199>
- [22] T. Soru, E. Marx, A. N. Ngomo, ROCKER: A refinement operator for key discovery, in: Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015, 2015, pp. 1025–1033. doi:10.1145/2736277.2741642.
URL <http://doi.acm.org/10.1145/2736277.2741642>
- [23] C. Xiao, W. Wang, X. Lin, Ed-join: an efficient algorithm for similarity joins with edit distance constraints, Proceedings of the VLDB Endowment 1 (1) (2008) 933–944.
- [24] G. Li, D. Deng, J. Wang, J. Feng, Pass-join: A partition-based method for similarity joins, Proceedings of the VLDB Endowment 5 (3) (2011) 253–264.
- [25] K. Dreßler, A. N. Ngomo, On the efficient execution of bounded jaro-winkler distances, Semantic Web 8 (2) (2017) 185–196. doi:10.3233/SW-150209.
URL <http://dx.doi.org/10.3233/SW-150209>
- [26] E. Fredkin, Trie memory, Commun. ACM 3 (9) (1960) 490–499. doi:10.1145/367390.367400.
URL <http://doi.acm.org/10.1145/367390.367400>
- [27] A. N. Ngomo, HELIOS - execution optimization for link discovery, in: Mika et al. [54], pp. 17–32. doi:10.1007/978-3-319-11964-9_2.
URL http://dx.doi.org/10.1007/978-3-319-11964-9_2
- [28] A. N. Ngomo, M. A. Sherif, K. Lyko, Unsupervised link discovery through knowledge base repair, in: The Semantic Web: Trends and Challenges - 11th International Conference, ESWC 2014, Anissaras, Crete, Greece, May 25-29, 2014. Proceedings, 2014, pp. 380–394. doi:10.1007/978-3-319-07443-6_26.
URL http://dx.doi.org/10.1007/978-3-319-07443-6_26
- [29] T. Soru, A. N. Ngomo, Active learning of domain-specific distances for link discovery, in: Semantic Technology, Second Joint International Conference, JIST 2012, Nara, Japan, December 2-4, 2012. Proceedings, 2012, pp. 97–112. doi:10.1007/978-3-642-37996-3_7.
URL http://dx.doi.org/10.1007/978-3-642-37996-3_7
- [30] A. N. Ngomo, S. Auer, J. Lehmann, A. Zaveri, Introduction to linked data and its lifecycle on the web, in: Reasoning Web. Reasoning on the Web in the Big Data Era - 10th International Summer School 2014,

- Athens, Greece, September 8-13, 2014. Proceedings, 2014, pp. 1–99. doi:10.1007/978-3-319-10587-1_1. URL http://dx.doi.org/10.1007/978-3-319-10587-1_1
- [31] M. A. Sherif, A. N. Ngomo, Semantic quran, *Semantic Web* 6 (4) (2015) 339–345. doi:10.3233/SW-140137. URL <http://dx.doi.org/10.3233/SW-140137>
- [32] J. Lehmann, T. Furche, G. Grasso, A. N. Ngomo, C. Schallhart, A. J. Sellers, C. Unger, L. Bühmann, D. Gerber, K. Höffner, D. Liu, S. Auer, deqa: Deep web extraction for question answering, in: *The Semantic Web - ISWC 2012 - 11th International Semantic Web Conference*, Boston, MA, USA, November 11-15, 2012, Proceedings, Part II, 2012, pp. 131–147. doi:10.1007/978-3-642-35173-0_9. URL http://dx.doi.org/10.1007/978-3-642-35173-0_9
- [33] R. Speck, A. N. Ngomo, Ensemble learning for named entity recognition, in: Mika et al. [54], pp. 519–534. doi:10.1007/978-3-319-11964-9_33. URL http://dx.doi.org/10.1007/978-3-319-11964-9_33
- [34] R. Usbeck, A. N. Ngomo, M. Röder, D. Gerber, S. A. Coelho, S. Auer, A. Both, AGDISTIS - graph-based disambiguation of named entities using linked data, in: Mika et al. [54], pp. 457–471. doi:10.1007/978-3-319-11964-9_29. URL http://dx.doi.org/10.1007/978-3-319-11964-9_29
- [35] M. Morsey, J. Lehmann, S. Auer, A.-C. Ngonga Ngomo, DBpedia SPARQL Benchmark—Performance Assessment with Real Queries on Real Data, in: *ISWC 2011*, 2011. URL <http://jens-lehmann.org/files/2011/dbpsb.pdf>
- [36] M. Saleem, Q. Mehmood, A.-C. Ngonga Ngomo, Feasible: A feature-based sparql benchmark generation framework, in: *International Semantic Web Conference (ISWC)*, 2015. URL http://svn.aksw.org/papers/2015/ISWC_FEASIBLE/public.pdf
- [37] L. Bühmann, J. Lehmann, Pattern based knowledge base enrichment, in: Alani et al. [51], pp. 33–48. doi:10.1007/978-3-642-41335-3_3. URL http://dx.doi.org/10.1007/978-3-642-41335-3_3
- [38] M. A. Sherif, A. N. Ngomo, J. Lehmann, Automating RDF dataset transformation and enrichment, in: *The Semantic Web. Latest Advances and New Domains - 12th European Semantic Web Conference, ESWC 2015, Portoroz, Slovenia, May 31 - June 4, 2015. Proceedings*, 2015, pp. 371–387. doi:10.1007/978-3-319-18818-8_23. URL http://dx.doi.org/10.1007/978-3-319-18818-8_23
- [39] J. Volz, C. Bizer, M. Gaedke, G. Kobilarov, Discovering and maintaining links on the web of data, in: *The Semantic Web - ISWC 2009, 8th International Semantic Web Conference, ISWC 2009, Chantilly, VA, USA, October 25-29, 2009. Proceedings*, 2009, pp. 650–665. doi:10.1007/978-3-642-04930-9_41. URL http://dx.doi.org/10.1007/978-3-642-04930-9_41
- [40] X. Niu, S. Rong, Y. Zhang, H. Wang, Zhishi.links results for OAEI 2011, in: Shvaiko et al. [53]. URL http://ceur-ws.org/Vol-814/oaei11_paper16.pdf
- [41] S. Araújo, J. Hidders, D. Schwabe, A. P. de Vries, SERIMI - resource description similarity, RDF instance matching and interlinking, in: *Proceedings of the 6th International Workshop on Ontology Matching*, Bonn, Germany, October 24, 2011, 2011. URL http://ceur-ws.org/Vol-814/om2011_poster6.pdf
- [42] K. Nguyen, R. Ichise, B. Le, SLINT: a schema-independent linked data interlinking system, in: *Proceedings of the 7th International Workshop on Ontology Matching*, Boston, MA, USA, November 11, 2012, 2012. URL http://ceur-ws.org/Vol-946/om2012_Tpaper1.pdf
- [43] J. Tang, B. Liang, J. Li, K. Wang, Risk minimization based ontology mapping, in: *Content Computing, Advanced Workshop on Content Computing*, AWCC 2004, ZhenJiang, JiangSu, China, November 15-17, 2004, Proceedings, 2004, pp. 469–480. doi:10.1007/978-3-540-30483-8_58. URL http://dx.doi.org/10.1007/978-3-540-30483-8_58
- [44] I. F. Cruz, F. P. Antonelli, C. Stroe, Agreementmaker: Efficient matching for large real-world schemas and ontologies, *PVLDB* 2 (2) (2009) 1586–1589. URL <http://www.vldb.org/pvldb/2/vldb09-1003.pdf>
- [45] E. Jiménez-Ruiz, B. C. Grau, Logmap: Logic-based and scalable ontology matching, in: *The Semantic Web - ISWC 2011 - 10th International Semantic Web Conference*, Bonn, Germany, October 23-27, 2011, Proceedings, Part I, 2011, pp. 273–288. doi:10.1007/978-3-642-25073-6_18. URL http://dx.doi.org/10.1007/978-3-642-25073-6_18
- [46] M. Niepert, C. Meilicke, H. Stuckenschmidt, A probabilistic-logical framework for ontology matching, in: *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*, 2010. URL <http://www.aaai.org/ocs/index.php/AAAI/AAAI10/paper/view/1744>
- [47] L. Zhu, M. Ghasemi-Gol, P. A. Szekely, A. Galstyan, C. A. Knoblock, Unsupervised entity resolution on multi-type graphs, in: *The Semantic Web - ISWC 2016 - 15th International Semantic Web Conference*, Kobe, Japan, October 17-21, 2016, Proceedings, Part I, 2016, pp. 649–667. doi:10.1007/978-3-319-46523-4_39. URL http://dx.doi.org/10.1007/978-3-319-46523-4_39
- [48] A.-C. N. Ngomo, M. M. Hassan, The lazy traveling salesman - memory management for large-scale link discovery., in: H. Sack, E. Blomqvist, M. d’Aquino, C. Ghidini, S. P. Ponzetto, C. Lange (Eds.), *ESWC*, Vol. 9678 of *Lecture Notes in Computer Science*, Springer, 2016, pp. 423–438. URL <http://dblp.uni-trier.de/db/conf/eswc/eswc2016.html#NgomoH16>
- [49] A.-C. Ngonga Ngomo, L. Kolb, N. Heino, M. Hartung, S. Auer, E. Rahm, When to reach for the cloud: Using parallel hardware for link discovery., in: *Proceedings of ESCW*, 2013.
- [50] M. A. Sherif, A. N. Ngomo, An optimization approach for load balancing in parallel link discovery, in: *Proceedings of the 11th International Conference on Semantic Systems, SEMANTICS 2015, Vienna, Austria, September 15-17, 2015*, 2015, pp. 161–168. doi:10.1145/2814864.2814872. URL <http://doi.acm.org/10.1145/2814864.2814872>
- [51] H. Alani, L. Kagal, A. Fokoue, P. T. Groth, C. Biemann, J. X. Parreira, L. Aroyo, N. F. Noy, C. Welty, K. Janowicz (Eds.), *The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference*, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part I, Vol. 8218 of *Lecture Notes in Computer Science*, Springer, 2013. doi:10.1007/978-3-642-41335-3. URL <http://dx.doi.org/10.1007/978-3-642-41335-3>
- [52] P. Shvaiko, J. Euzenat, K. Srinivas, M. Mao, E. Jiménez-Ruiz (Eds.), *Proceedings of the 8th International Workshop on Ontology Matching co-located with the 12th International Semantic Web Conference (ISWC 2013)*, Sydney, Australia, October 21, 2013, Vol. 1111 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2013. URL <http://ceur-ws.org/Vol-1111>
- [53] P. Shvaiko, J. Euzenat, T. Heath, C. Quix, M. Mao, I. F. Cruz (Eds.), *Proceedings of the 6th International Workshop on Ontology Matching*, Bonn, Germany, October 24, 2011, Vol. 814 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2011. URL <http://ceur-ws.org/Vol-814>
- [54] P. Mika, T. Tudorache, A. Bernstein, C. Welty, C. A. Knoblock, D. Vrandečić, P. T. Groth, N. F. Noy, K. Janowicz, C. A. Goble (Eds.), *The Semantic Web - ISWC 2014 - 13th International Semantic Web Conference*, Riva del Garda, Italy, October 19-23, 2014, Proceedings, Part I, Vol. 8796 of *Lecture Notes in Computer Science*, Springer, 2014. doi:10.1007/978-3-319-11964-9. URL <http://dx.doi.org/10.1007/978-3-319-11964-9>