

# A Generalization Approach for Automatic Link Discovery

Mohamed Ahmed Sherif and Axel-Cyrille Ngonga Ngomo and Jens Lehmann<sup>1</sup>

**Abstract.** A significant portion of the evolution of Linked Data datasets lies in updating the links to other datasets. An important challenge when aiming to update these links automatically under the open-world assumption is the fact that usually only positive examples for the links exist. We address this challenge by presenting and evaluating WOMBAT, a novel approach for the discovery of links between knowledge bases that relies exclusively on positive examples. WOMBAT is based on generalisation via an upward refinement operator to traverse the space of link specification. We study the theoretical characteristics of WOMBAT and evaluate it on 8 different benchmark datasets. Our evaluation suggests that WOMBAT outperforms state-of-the-art supervised approaches while relying on less information. Moreover, our evaluation suggests that WOMBAT’s pruning algorithm allows it to scale well even on large datasets.

## 1 Introduction

The Linked Open Data Cloud has grown from a mere 12 datasets at its beginning to a compendium of more than 9,000 public RDF<sup>2</sup> data sets.<sup>3</sup> In addition to the number of the datasets published growing steadily, we also witness the size of single datasets growing with each new edition. For example, *DBpedia* has grown from 103 million triples describing 1.95 million things (DBpedia 2.0) to 583 million triples describing 4.58 million things (DBpedia 2014) within 7 years. This growth engenders an increasing need for automatic support when maintaining evolving datasets. 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 [1]. However, the links available on the Data Web only provide positive examples for relations and no negative examples, 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 these two resources are not related. Consequently, state-of-the-art supervised learning approaches for link discovery can only be employed if the end users are willing to provide the algorithms with information that is generally not available on the Linked Open Data Cloud, i.e., with negative examples.

We address this drawback by proposing the first approach for learning links based on positive examples only. Our approach, dubbed WOMBAT, is inspired by the concept of generalisation in

quasi-ordered spaces. Given a set of positive examples we aim 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. We use Link Specifications (LS, see Section 2) as classifiers [1, 7, 14]. We are thus faced with the challenge of using various similarity metrics, acceptance thresholds and nested logical combinations of those when learning. The contributions of this paper are: 1) We provide the first approach for learning LS that is able to learn links from positive examples only. 2) Our approach is based on an upward refinement operator for which we analyse its theoretical characteristics. 3) We use the characteristics of our operator to devise a pruning approach and improve the scalability of WOMBAT. 4) We evaluate WOMBAT on 8 benchmark datasets and show that in addition to needing less training data, it also outperforms the state of the art in most cases.

## 2 Preliminaries

The aim of link discovery (LD) is to discover the set  $\{(s, t) \in S \times T : Rel(s, t)\}$  provided an input relation *Rel* and two sets *S* (source) and *T* (target) of RDF resources. To achieve this goal, declarative LD frameworks rely on LS, which describe the conditions under which  $Rel(s, t)$  can be assumed to hold for a pair  $(s, t) \in S \times T$ . Several grammars have been used for describing LS in previous works [15, 5, 19]. In general, these grammars assume that LS consist of two types of atomic components: *similarity measures* *m*, which allow comparing property values of input resources and *operators* *op*, which can be used to combine these similarities to more complex specifications. Without loss of generality, we define a similarity measure *m* as a function  $m : S \times T \rightarrow [0, 1]$ . An example of a similarity measure is the edit similarity dubbed *edit*<sup>4</sup> which allows computing the similarity of a pair  $(s, t) \in S \times T$  with respect to the properties  $p_s$  of *s* and  $p_t$  of *t*. We use *mappings*  $M \subseteq S \times T$  to store the results of the application of a similarity function to  $S \times T$  or subsets thereof. We denote the set of all mappings as  $\mathcal{M}$  and the set of all LS as  $\mathcal{L}$ . 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 merging the results of  $L_1$  and  $L_2$ . Here, we use the operators  $\sqcap$ ,  $\sqcup$  and  $\setminus$  as they are complete and frequently used to define LS. An example of a complex LS is given in Figure 1.

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

<sup>1</sup> Department of Computer Science, University of Leipzig, 04109 Leipzig, Germany, email: {sherif, ngonga, lehmann}@informatik.uni-leipzig.de

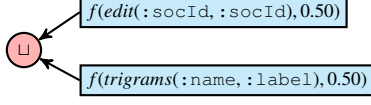
<sup>2</sup> See <https://www.w3.org/RDF/> for the specification of the RDF standard.

<sup>3</sup> <http://lodstats.aksw.org>

<sup>4</sup> We define the edit similarity of two strings *s* and *t* as  $(1 + lev(s, t))^{-1}$ , where *lev* stands for the Levenshtein distance.

to  $S \times T$  contains the links that will be generated by  $L$ . A LS  $L$  is *subsumed* by  $L'$ , denoted by  $L \sqsubseteq L'$ , if for all mappings  $M$ , we have  $[[L]]_M \subseteq [[L']]_M$ . Two LS are *equivalent*, denoted by  $L \equiv L'$  iff  $L \sqsubseteq L'$  and  $L' \sqsubseteq L$ . Subsumption ( $\sqsubseteq$ ) is a partial order over  $\mathcal{L}$ .

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



**Table 1.** Link Specification Syntax and Semantics

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

### 3 Constructing and Traversing Link Specifications

The goal of our learning approach is to learn a specification  $L$  that generalizes a mapping  $M \subseteq S \times T$  which contains a set of pairs  $(s, t)$  for which  $Rel(s, t)$  holds. Our approach consists of two main steps. First, we aim to derive initial atomic specifications  $A_i$  that achieve the same goal. In a second step, we combine these atomic specifications to the target complex specification  $L$  by using the operators  $\sqcap$ ,  $\sqcup$  and  $\setminus$ . In the following, we detail how we carry out these two steps.

#### 3.1 Learning Atomic Specifications

The goal here is to derive a set of initial atomic specifications  $\{A_1, \dots, A_n\}$  that achieves the highest possible F-measure given a mapping  $M \subseteq S \times T$  which contains all known pairs  $(s, t)$  for which  $Rel(s, t)$  holds. Given a set of similarity functions  $m_i$ , the set of properties  $P_s$  of  $S$  and the set of properties  $P_t$  of  $T$ , we begin by computing the subset of properties from  $S$  and  $T$  that achieve a coverage above a threshold  $\tau \in [0, 1]$ , where the coverage of a property  $p$  for a knowledge base  $K$  is defined as

$$coverage(p) = \frac{|\{s : (s, p, o) \in K\}|}{|\{s : \exists q : (s, q, o) \in K\}|}. \quad (1)$$

Now for all property pairs  $(p, q) \in P_s \times P_t$  with  $coverage(p) \geq \tau$  and  $coverage(q) \geq \tau$ , we compute the mappings  $M_{ij} = \{(s, t) \in S \times T : m_{ij}(s, t) \geq \theta_{ij}\}$ , where  $m_{ij}$  compares  $s$  and  $t$  w.r.t.  $p$  and  $q$  and  $M_{ij}$  is maximal w.r.t. the F-measure it achieves when compared to  $M$ . To this end, we apply an iterative search approach. Finally, we select  $M_{ij}$  as the atomic mapping for  $p$  and  $q$ . Thus, we return as many atomic mappings as property pairs with sufficient coverage. Note that this approach is not quintessential for WOMBAT and can thus be replaced with any approach of choice which returns a set of initial LS that is to be combined.

#### 3.2 Combining Atomic Specifications

After deriving atomic LS as described above, WOMBAT computes complex specifications by using an approach based on generalisation operators. The basic idea behind these operators is to perform an iterative search through a solution space based on a score function. Formally, we rely on the following definitions:

**Definition 1** ((Refinement) Operator). *In the quasi-ordered space  $(\mathcal{L}, \sqsubseteq)$ , we call a function from  $\mathcal{L}$  to  $2^{\mathcal{L}}$  an (LS) operator. A downward (upward) refinement operator  $\rho$  is an operator, such that for all  $L \in \mathcal{L}$  we have that  $L' \in \rho(L)$  implies  $L' \sqsubseteq L$  ( $L \sqsubseteq L'$ ).  $L'$  is called a specialisation (generalisation) of  $L$ .  $L' \in \rho(L)$  is usually denoted as  $L \rightsquigarrow_{\rho} L'$ .*

**Definition 2** (Refinement Chains). *A refinement chain of a refinement operator  $\rho$  of length  $n$  from  $L$  to  $L'$  is a finite sequence  $L_0, L_1, \dots, L_n$  of LS, such that  $L = L_0$ ,  $L' = L_n$  and  $\forall i \in \{1 \dots n\}$ ,  $L_i \in \rho(L_{i-1})$ . This refinement chain goes through  $L''$  iff there is an  $i$  ( $1 \leq i \leq n$ ) such that  $L'' = L_i$ . We say that  $L''$  can be reached from  $L$  by  $\rho$  if there exists a refinement chain from  $L$  to  $L''$ .  $\rho^*(L)$  denotes the set of all LS which can be reached from  $L$  by  $\rho$ .  $\rho^m(L)$  denotes the set of all LS which can be reached from  $L$  by a refinement chain of  $\rho$  of length  $m$ .*

**Definition 3** (Properties of refinement operators). *An operator  $\rho$  is called (1) **(locally) finite** iff  $\rho(L)$  is finite for all LS  $L \in \mathcal{L}$ ; (2) **redundant** iff there exists a refinement chain from  $L \in \mathcal{L}$  to  $L' \in \mathcal{L}$ , which does not go through (as defined above) some LS  $L'' \in \mathcal{L}$  and a refinement chain from  $L$  to  $L'$  which does go through  $L''$ ; (3) **proper** iff for all LS  $L \in \mathcal{L}$  and  $L' \in \mathcal{L}$ ,  $L' \in \rho(L)$  implies  $L \not\equiv L'$ . An LS upward refinement operator  $\rho$  is called **weakly complete** iff for all LS  $\perp \sqsubseteq L$  we can reach a LS  $L'$  with  $L' \equiv L$  from  $\perp$  (most specific LS) by  $\rho$ .*

We designed two different operators for combining atomic LS to complex specifications: The first operator takes an atomic LS and uses the three logical connectors to append further atomic LS. Assuming that  $(A_1, \dots, A_n)$  is the set of atomic LS found,  $\varphi$  can be defined as follows:

$$\varphi(L) = \begin{cases} \bigcup_{i=1}^n A_i & \text{if } L = \perp \\ (\bigcup_{i=1}^n L \sqcup A_i) \cup (\bigcup_{i=1}^n L \sqcap A_i) \cup (\bigcup_{i=1}^n L \setminus A_i) & \text{otherwise} \end{cases}$$

This naive operator is not a refinement operator (neither upward nor downward). Its main advantage lies in its simplicity allowing for a very efficient implementation. However, it cannot reach all specifications, e.g., a specification of the form  $(A_1 \sqcup A_2) \sqcap (A_3 \sqcup A_4)$  cannot be reached. Examples of chains generated by  $\varphi$  are as follows:

1.  $\perp \rightsquigarrow_{\varphi} A_1 \rightsquigarrow_{\varphi} A_1 \sqcup A_2 \rightsquigarrow_{\varphi} (A_1 \sqcup A_2) \setminus A_3$
2.  $\perp \rightsquigarrow_{\varphi} A_2 \rightsquigarrow_{\varphi} A_2 \sqcap A_3 \rightsquigarrow_{\varphi} (A_2 \sqcap A_3) \setminus A_4$

The second operator,  $\psi$ , uses a more sophisticated expansion strategy in order to allow learning arbitrarily nested LS and is shown in Figure 2. Less formally, the operator works as follows: It takes a LS as input and makes a case distinction on the type of LS. Depending on the type, it performs the following actions:

- The  $\perp$  LS is refined to the set of all combinations of  $\setminus$  operations. This set can be large and will only be built iteratively (as required by the algorithm) with at most approx.  $n^2$  refinements per iteration (see the next section for details).
- In LS of the form  $A_1 \setminus A_2$ ,  $\psi$  can drop the second part in order to generalise.
- If the LS is a conjunction or disjunction, the operator can perform a recursion on each element of the conjunction or disjunction.
- For LS of any type, a disjunction with an atomic LS can be added.

Below are two example refinement chains of  $\psi$ :

1.  $\perp \rightsquigarrow_{\psi} A_1 \setminus A_2 \rightsquigarrow_{\psi} A_1 \rightsquigarrow_{\psi} A_1 \sqcup A_2 \setminus A_3$

$$\psi(L) = \begin{cases} \{A_{i_1} \setminus A_{j_1} \sqcap \dots \sqcap A_{i_m} \setminus A_{j_m} \mid A_{i_k}, A_{j_k} \in P \\ \text{for all } 1 \leq k \leq m\} & \text{if } L = \perp \\ \{L \sqcup A_i \setminus A_j \mid A_i \in P, A_j \in P\} & \text{if } L = A \text{ (atomic)} \\ \{L_1\} \cup \{L \sqcup A_i \setminus A_j \mid A_i \in P, A_j \in P\} & \text{if } L = L_1 \setminus L_2 \\ \{L_1 \sqcap \dots \sqcap L_{i-1} \sqcap L' \sqcap L_{i+1} \sqcap \dots \sqcap L_n \mid L' \in \psi(L_i)\} \\ \cup \{L \sqcup A_i \setminus A_j \mid A_i \in P, A_j \in P\} & \text{if } L = L_1 \sqcap \dots \sqcap L_n (n \geq 2) \\ \{L_1 \sqcup \dots \sqcup L_{i-1} \sqcup L' \sqcup L_{i+1} \sqcup \dots \sqcup L_n \mid L' \in \psi(L_i)\} \\ \cup \{L \sqcup A_i \setminus A_j \mid A_i \in P, A_j \in P\} & \text{if } L = L_1 \sqcup \dots \sqcup L_n (n \geq 2) \end{cases}$$

**Figure 2.** Definition of the refinement operator  $\psi$ .

$$2. \perp \rightsquigarrow_\psi A_1 \setminus A_2 \sqcap A_3 \setminus A_4 \rightsquigarrow_\psi A_1 \sqcap A_3 \setminus A_4 \rightsquigarrow_\psi A_1 \sqcap A_3 \rightsquigarrow_\psi (A_1 \sqcap A_3) \sqcup (A_5 \setminus A_6)$$

$\psi$  is an upward refinement operator with the following properties.

**Proposition 1.**  $\psi$  is an upward refinement operator.

*Proof.* For an arbitrary LS  $L$ , we have to show for any element  $L' \in \psi(L)$  that  $L \sqsubseteq L'$  holds. The proof is straightforward by showing that  $L'$  cannot generate less links than  $L$  via case distinction and structural induction over LS:

- $L = \perp$ : Trivial.
- $L$  is atomic: Adding a disjunction cannot result in less links (this also holds for the cases below).
- $L$  is of the form  $L_1 \setminus L_2$ :  $L' = L_1$  cannot result in less links.
- $L$  is a conjunction / disjunction:  $L'$  cannot result in less links by structural induction.  $\square$

**Proposition 2.**  $\psi$  is weakly complete.

*Proof.* To show this, we have to show that an arbitrary LS  $L$  can be reached from the  $\perp$  LS. First, we convert everything to negation normal form by pushing  $\setminus$  inside, e.g. LS of the form  $L_1 \setminus (L_2 \sqcap L_3)$  are rewritten to  $(L_1 \setminus L_2) \sqcup (L_1 \setminus L_3)$  and LS of the form  $L_1 \setminus (L_2 \sqcup L_3)$  are rewritten to  $(L_1 \setminus L_2) \sqcap (L_1 \setminus L_3)$  exhaustively. We then further convert the LS to conjunction normal including an exhaustive application of the distribute law, i.e., conjunctions cannot be nested within disjunctions. The resulting LS is dubbed  $L'$  and equivalent to  $L$ . We show that  $L'$  can always be reached from  $\perp$  via induction over its structure:

- $L' = \perp$ : Trivial via the empty refinement chain.
- $L' = A$  (atomic): Reachable via  $\perp \rightsquigarrow_\psi A \setminus A' \rightsquigarrow_\psi A$ .
- $L' = A_1 \setminus A_2$  (atomic negation): Reachable directly via  $\perp \rightsquigarrow_\psi A_1 \setminus A_2$ .
- $L'$  is a conjunction with  $m$  elements:  $\perp \rightsquigarrow_\psi A_{i_1} \setminus A_{j_1} \sqcap \dots \sqcap A_{i_m} \setminus A_{j_m}$  where an element  $A_{i_k} \setminus A_{j_k}$  is chosen as follows: Let the  $k$ -th element of conjunction  $L'$  be  $L''$ .
  - If  $L''$  is an atomic specification  $A$ , then  $A_{i_k} = A$  ( $A_{j_k}$  can be arbitrarily).
  - If  $L''$  is an atomic negation  $A_1 \setminus A_2$ , then  $A_{i_k} = A$  and  $A_{j_k} = A_2$ .
  - If  $L''$  is a disjunction, the first element of this disjunction falls into one of the above two cases and  $A_{i_k}$  and  $A_{j_k}$  can be set as described there.

Each element of  $L''$  is then further refined to  $L'$  as follows:

- If  $L''$  is an atomic specification  $A$ :  $A \setminus A_{j_k}$  is refined to  $A$ .
- If  $L''$  is an atomic negation  $A_1 \setminus A_2$ : No further refinements are necessary.

- If  $L''$  is a disjunction. The first element of the disjunction is first treated according to the two cases above. Subsequent elements of the disjunction are either atomic LS or atomic negation and can be added straightforwardly as the operator allows adding disjunctive elements to any non- $\perp$  LS.

Please note that the case distinction is exhaustive as we assume  $L'$  is in conjunctive negation normal form, i.e., there are no disjunctions on the outer level, negation is always atomic, conjunctions are not nested within other conjunction and elements of disjunctions within conjunctions cannot be conjunctions.  $\square$

**Proposition 3.**  $\psi$  is finite, not proper and redundant.

*Proof. Finiteness:* There are only finitely many atomic LS. Hence, there are only finitely many atomic negations and, consequently, finitely many possible conjunctions of those. Consequently,  $\psi(\perp)$  is finite. The finiteness of  $\psi(L)$  with  $L \neq \perp$  is straightforward.

*Properness:* The refinement chain  $\perp \rightsquigarrow_\psi^* A_1 \sqcap A_2 \rightsquigarrow_\psi^* (A_1 \sqcup A_2) \sqcap A_2$  is a counterexample.

*Redundancy:* The two refinement chains  $A_1 \sqcap A_3 \rightsquigarrow_\psi^* (A_1 \sqcup A_2) \sqcap A_3 \rightsquigarrow_\psi^* (A_1 \sqcup A_2) \sqcap (A_3 \sqcup A_4)$  and  $A_1 \sqcap A_3 \rightsquigarrow_\psi^* A_1 \sqcap (A_3 \sqcap A_4) \rightsquigarrow_\psi^* (A_1 \sqcup A_2) \sqcap (A_3 \sqcup A_4)$  are a counterexample.  $\square$

Naturally, the restrictions of  $\psi$  (being redundant and not proper) raise the question whether there are LS refinement operators satisfying all theoretical properties:

**Proposition 4.** There exists a weakly complete, finite, proper and non-redundant refinement operator in  $\mathcal{L}$ .

*Proof.* Let  $C$  be the set of LS in  $\mathcal{L}$  in conjunctive negation normal form without any LS equivalent to  $\perp$ . We define the operator  $\alpha$  as  $\alpha(\perp) = C$  and  $\alpha(L) = \emptyset$  for all  $L \neq \perp$ .  $\alpha$  is obviously complete as any LS has an equivalent in conjunctive negation normal form. It is finite as  $S$  can be shown to be finite with an extended version of the argument in the finiteness proof of  $\psi$ .  $\alpha$  is trivially non-redundant and it is proper by definition.  $\square$

The existence of an operator which satisfies all considered theoretical criteria of a refinement operator is an artifact of only finitely many semantically inequivalent LS existing in  $\mathcal{L}$ . This set is however extremely large and not even small fractions of it can be evaluated in all but very simple cases. For example, the operator  $\alpha$  as  $\alpha(\perp) = C$  and  $\alpha(L) = \emptyset$  for all  $L \neq \perp$  is trivially non-redundant and it is proper by definition. Such an operator  $\alpha$  is obviously not useful as it does not help *structuring the search space*. Providing a useful way to structure the search space is the main reason for refinement operators being successful for learning in other complex languages as it allows to gradually converge towards useful solutions while being

able to prune other paths which cannot lead to promising solutions (explained in the next section). This is a reason why we sacrificed properness and redundancy for a better structure of the search space.

## 4 The WOMBAT Algorithm

---

### Algorithm 1: WOMBAT Learning Algorithm

---

**Input:** Sets of resources  $S$  and  $T$ ; examples  $E \subseteq S \times T$ ; property coverage threshold  $\tau$ ; set of similarity functions  $\mathbf{F}$

```

1  $\mathbf{A} \leftarrow \text{null}$  (the list of initial atomic metrics);
2  $i \leftarrow 1$ ;
3 foreach property  $p_s \in S$  do
4   if  $\text{coverage}(p_s) \geq \tau$  then
5     foreach property  $p_t \in T$  do
6       if  $\text{coverage}(p_t) \geq \tau$  then
7         Find atomic metric  $m(p_s, p_t)$  that leads to
          highest F-measure;
8         Optimize similarity threshold for  $m(p_s, p_t)$  to
          find best mapping  $A_i$ ;
9         Add  $A_i$  to  $\mathbf{A}$ ;
10         $i \leftarrow i + 1$ ;
11  $\Gamma \leftarrow \perp$  (initiate search tree  $\Gamma$  to the root node  $\perp$ );
12  $F_{\text{best}} \leftarrow 0, L_{\text{best}} \leftarrow \text{null}$ ;
13 while termination criterion not met do
14   Choose the node with highest scoring LS  $L$  in  $\Gamma$ ;
15   if  $L == \perp$  then
16     foreach  $A_i, A_j \in \mathbf{A}$ , where  $i \neq j$  do
17       Only add refinements of form  $A_i \setminus A_j$ ;
18   else
19     Apply operator to  $L$ ;
20     if  $L$  is a refinement of  $\perp$  then
21       foreach  $A_i, A_j \in \mathbf{A}$ , where  $i \neq j$  do
22         In addition to refinements, add conjunctions
          with specifications of the form  $A_i \setminus A_j$  as
          siblings;
23   foreach refinement  $L'$  do
24     if  $L'$  is not already in the search tree  $\Gamma$  then
25       Add  $L'$  to  $\Gamma$  as children of the node containing  $L$ ;
26   Update  $F_{\text{best}}$  and  $L_{\text{best}}$ ;
27   if  $F_{\text{best}}$  has increased then
28     foreach subtree  $t \in \Gamma$  do
29       if  $F_{\text{best}} > F_{\text{max}}(t)$  then
30         Delete  $t$ ;
31 Return  $L_{\text{best}}$ ;

```

---

We have now introduced all ingredients necessary for defining the WOMBAT algorithms. The first algorithm, which we refer to as *simple* version, uses the operator  $\varphi$ , whereas the second algorithm, which we refer to as *complete*, uses the refinement operator  $\psi$ . The complete algorithm has the following specific characteristics: First, while  $\psi$  is finite, it would generate a prohibitively large number of refinements when applied to the  $\perp$  concept. For that reason, those refinements will be computed stepwise as we will illustrate below. Second, as  $\psi$  is an upward refinement operator it allows to prune parts of the

search space, which we will also explain below. We only explain the implementation of the complex WOMBAT algorithm as the other is a simplification excluding those two characteristics.

Algorithm 1 shows the individual steps of WOMBAT complete. Our approach takes the source dataset  $S$ , the target dataset  $T$ , examples  $E \subseteq S \times T$  as well as the property coverage threshold and the set of considered similarity functions as input. In Line 3, the property matches are computed by optimizing the threshold for properties that have the minimum coverage (Line 7) as described in Section 3.1. The main loop starts in Line 13 and runs until a termination criterion is satisfied, e.g. 1) a fixed number of LS has been evaluated, 2) a certain time has elapsed, 3) the best F-score has not changed for a certain time or 4) a perfect solution has been found. Line 14 states that a heuristic-based search strategy is employed.

By default, we employ the F-score directly. The intuition behind the use of the F-score is derived from the concept of *least general generalization*: We aim to find a specification that (1) abides by the formal model of specifications (our language) while (2) that contains as many of the positive examples as possible while (3) containing the smallest possible number of unknown link, thus being least general. This intuition is well captured by the F-measure and works practically as shown by our results (see Section 5). Still, more complex heuristics introducing a bias towards specific types of LS (for example specifications which consist of a small number of atomic specifications) could be encoded here.

In Line 15, we make a case distinction: Since the number of refinements of  $\perp$  is extremely high and not feasible to compute in most cases, we perform a stepwise approach: In the first step, we only add simple LS of the form  $A_i \setminus A_j$  as refinements (Line 17). Later, in Line 22, we add more complex conjunctions if the simpler forms are promising. Apart from this special case, we apply the operator directly. Line 24 updates the search tree by adding the nodes obtained via refinement. Moreover, the algorithm contains a redundancy elimination procedure: We only add those nodes to the search tree which are not already contained in it.

The subsequent part starting from Line 26 defines our *pruning procedure*: Since  $\psi$  is an upward refinement operator, we know that the set of links generated by a child node is a superset of or equal to the set of links generated by its parent. Hence, while both precision and recall can improve in subsequent refinements, they cannot rise arbitrarily. Precision is bound as false positives cannot disappear during generalisation. Furthermore, the achievable recall  $r_{\text{max}}$  is that of the most general constructable LS, i.e.,  $\mathcal{A} = \bigcup A_i$ . This allows to compute an upper bound on the achievable F-score. In order to do so, we first build a set  $S'$  with those resources in  $S$  occurring in the input examples  $E$  as well as a set  $T'$  with those resources in  $T$  occurring in  $E$ . The purpose of those is to restrict the computation of F-score to the fragment  $S' \times T' \subseteq S \times T$  relevant for example set  $E$ . We can then compute an upper bound of precision of a LS  $L$  as follows:

$$p_{\text{max}}(L) = \frac{|E|}{|E| + |\{(s, t) \mid (s, t) \in [[L]], s \in S' \text{ or } t \in T' \setminus E\}|}$$

$F_{\text{max}}$  is then computed as the F-measure obtained with recall  $r_{\text{max}}$  and precision  $p_{\text{max}}$ , i.e.,  $F_{\text{max}} = \frac{2p_{\text{max}}r_{\text{max}}}{p_{\text{max}}+r_{\text{max}}}$ . It is an upper bound for the maximum achievable F-measure of any node reachable via refinements. We can disregard all nodes in the search tree which have a maximum achievable F-score that is lower than the best F-score already found. This is implemented in Line 28. The pruning is conservative in the sense that no solutions are lost. In the evaluation, we give statistics on the effect of pruning. WOMBAT ends by returning  $L_{\text{best}}$  as the best LS found, which is the specification with the highest F-score. In case of ties, we prefer shorter specifications over long

ones. Should the tie persist, then we prefer specifications that were found early.

**Proposition 5.** *WOMBAT is complete, i.e., it will eventually find the LS with the highest F-measure within  $\mathcal{L}$ .*

*Proof.* This is a consequence of the weak completeness of  $\psi$  and the fact that the algorithm will eventually generate all refinements of  $\psi$ . For the latter, we have to look at the refinement of  $\perp$  as a special case since otherwise a straightforward application of  $\psi$  is used. For the refinements of  $\perp$  it is easy to show via induction over the number of conjunctions in refinements that any element in  $\psi(\perp)$  can be reached via the algorithm. (The pruning is conservative and only prunes nodes never leading to better solutions.)  $\square$

## 5 Evaluation

We evaluated our approach using 8 benchmark datasets. Five of these benchmarks were real-world datasets while three were synthetic. The real-world interlinking tasks used were those in [9]. The synthetic datasets were from the OAEI 2010 benchmark<sup>5</sup>. All experiments were carried out on a 64-core 2.3 GHz PC running *OpenJDK* 64-Bit Server 1.7.0.75 on *Ubuntu* 14.04.2 LTS. Each experiment was assigned 20 GB RAM.

For testing WOMBAT against the benchmark datasets in both its simple and complete version, we used the `jaccard`, `trigrams`, `cosine` and `qgrams` similarity measures. We used two termination criteria: Either a LS with F-measure of 1 was found or a maximal depth of refinement (10 resp. 3 for the simple resp. complete version) was reached. This variation of the maximum refinement trees sizes between the simple and complete version was because WOMBAT complete adds a larger number of nodes to its refinement tree in each level. The coverage threshold  $\tau$  was set to 0.6. A more complete list of evaluation results are available at the project web site.<sup>6</sup> Altogether, we carried out 6 sets of experiments to evaluate WOMBAT.

In the *first set of experiments*, we compared the average F-Measure achieved by the simple and complete versions of WOMBAT to that of four other state-of-the-art LS learning algorithms within a 10-fold cross validation setting. The other four LS learning algorithms were EAGLE [15] as well as the *linear*, *conjunctive* and *disjunctive* versions of EUCLID [16]. EAGLE was configured to run 100 generations. The mutation and crossover rates were set to 0.6 as in [15]. To address the non-deterministic nature of EAGLE, we repeated the whole process of 10-fold cross validation 5 time and present the average results. EUCLID's grid size was set to 5 and 100 iterations were carried out as in [16]. The results of the evaluation are presented in Table 2. The simple version of WOMBAT was able to outperform the state-of-the-art approaches in 4 out of the 8 data sets and came in the second position in 2 datasets. WOMBAT complete was able to achieve the best F-score in 4 data sets and achieve the second best F-measure in 3 datasets. On average, both versions of WOMBAT were able to achieve an F-measure of 0.9, by which WOMBAT outperforms the three version of EUCLID by an average of 11%. While WOMBAT was able to achieve the same performance of EAGLE in average, WOMBAT is still to be preferred as (1) WOMBAT only requires positive examples and (2) EAGLE is indeterministic by nature.

For the *second set of experiments*, we implemented an evaluation protocol based on the assumptions made at the beginning of this paper. Each input dataset was split into 10 parts of the same size. Consequently, we used 3 parts (30%) of the data as training data and the

**Table 2.** 10-fold cross validation F-Measure results.

Dataset	WOMBAT		EUCLID			EAGLE
	Simple	Complete	Linear	Conj.	Disj.	
Person 1	<b>1.00</b>	<b>1.00</b>	0.64	0.97	<b>1.00</b>	0.99
Person 2	<b>1.00</b>	0.99	0.22	0.78	0.96	0.94
Restaurants	<b>0.98</b>	0.97	0.97	0.97	0.97	0.97
DBLP-ACM	0.97	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>	<b>0.98</b>
Abt-Buy	0.60	0.61	0.06	0.06	0.52	<b>0.65</b>
Amazon-GP	0.70	0.67	0.59	0.71	<b>0.73</b>	0.71
DBP-LMDB	0.99	<b>1.00</b>	0.99	0.99	0.99	0.99
DBLP-GS	<b>0.94</b>	<b>0.94</b>	0.90	0.91	0.91	0.93
Average	<b>0.90</b>	<b>0.90</b>	0.67	0.80	0.88	<b>0.90</b>

rest 7 parts (70%) for testing. This was to implement the idea of the dataset growing and the specification (and therewith the links) for the new version of the dataset having to be derived by learning from the old dataset. During the learning process, the score function was the F-measure achieved by each refinement of the portion of the training data related to  $S \times T$  selected for training (dubbed  $S' \times T'$  previously). The F-measures reported are those achieved by LS on the test dataset. We used the same settings for EAGLE and EUCLID as in the experiments before. The results (see Table 3) show clearly that our simple operator outperforms all other approaches in this setting. Moreover, the complete version of WOMBAT reaches the best F-measure on 2 datasets and the second-best F-measure on 3 datasets. This result of central importance as it shows that WOMBAT is well suited for the task for which it was designed. Interestingly, our approach also outperforms the approaches that rely on negative examples (i.e. EUCLID and EAGLE). The complete version of WOMBAT seems to perform worse than the simple version because it can only explore a tree of depth 3. However, this limitation was necessary to test both implementations using the same hardware.

**Table 3.** A comparison of WOMBAT F-Measure against 4 state-of-the-art approaches on 8 different benchmark datasets using 30% of the original data as training data.

Dataset	WOMBAT		EUCLID			EAGLE
	Simple	Complete	Linear	Conj.	Disj.	
Person 1	<b>1.00</b>	<b>1.00</b>	0.95	0.96	0.99	0.92
Person 2	<b>0.99</b>	0.79	0.80	0.82	0.88	0.69
Restaurants	<b>0.97</b>	0.88	0.87	0.84	0.89	0.88
DBLP-ACM	<b>0.95</b>	0.91	0.88	0.89	0.91	0.85
Abt-Buy	<b>0.44</b>	0.40	0.29	0.29	0.29	0.27
Amazon-GP	<b>0.54</b>	0.41	0.31	0.30	0.32	0.32
DBP-LMDB	<b>0.98</b>	<b>0.98</b>	0.97	0.96	0.97	0.89
DBLP-GS	<b>0.91</b>	0.74	0.83	0.76	0.74	0.69
Average	<b>0.85</b>	0.76	0.74	0.73	0.75	0.69

In the *third set of experiments*, we measured the effect of increasing the amount of training data on the precision, recall and F-score achieved by both simple and complete versions of WOMBAT. The results are presented in Figure 3. Our results suggest that the complete version of WOMBAT is partly more stable in its results (see ABT-Buy and DBLP-Google Scholar) and converges faster towards the best solution that it can find. This suggests that once trained on a dataset, our approach can be used on subsequent versions of real datasets, where a small number of novel resources is added in each new version, which is the problem setup considered in this paper. On the other hand, the simple version is able to find better LS as it can

<sup>5</sup> <http://oaei.ontologymatching.org/2010/>

<sup>6</sup> <https://github.com/AKSW/LIMES/tree/master/evaluationsResults/wombat>

explore longer sequences of mappings.

In the *fourth set of experiments*, we measured the learning time for each of the benchmark datasets. The results are also presented in Figure 3. As expected, the simple approach is time-efficient to run even without any optimization. While the complete version of WOMBAT without pruning is significantly slower (up to 1 order of magnitude), the effect of pruning can be clearly seen as it reduces the runtime of the algorithm while also improving the total space that the complete version of WOMBAT can explore. These results are corroborated by our *fifth set of experiments*, in which we evaluated the pruning technique of the complete version of WOMBAT. In those experiments, for each of aforementioned benchmark datasets we computed what we dubbed as *pruning factor*. The pruning factor is the number of searched nodes (search tree size plus pruned nodes) divided by the maximum size of the search tree (which we set to 2000 nodes in this set of experiments). The results are presented in Table 5. Our average *pruning factor* of 2.55 shows that we can discard more than 3000 nodes while learning specifications.

In a *final set of experiments*, we compared the two versions of WOMBAT against the 2 systems proposed in [8]. To be comparable, we used the same evaluation protocol in [8], where 2% of the gold standard was used as training data and the remaining 98% of the gold standard as test data. The results (presented in Table 4) suggests that WOMBAT is capable of achieving better or equal performance in 4 out of the 6 evaluation data sets. While WOMBAT achieved inferior F-measures for the other 2 data sets, it should be noted that the competing systems are optimised for a low number of examples and they also get negative examples as input. Overall, these results can thus be regarded as positive as they suggest that our approach can generalise a small number of examples to a sensible LS.

Overall, our results show that  $\psi$  and  $\varphi$  are able to learn high-quality LS using only positive examples. When combined with our pruning algorithm, the complete version of  $\psi$  achieves runtimes that are comparable to those of  $\varphi$ . Given its completeness,  $\psi$  can reach specifications that simply cannot be learned by  $\varphi$  (see Figure 4 for an example of such a LS). Hence, the complete operator has the potential to be able to detect more specifications that cover corner cases in difficult datasets. This can be seen in the results on the Abt-Buy dataset (see Table 2), where the complete operator is slightly better than the simple one. However, for practical applications,  $\varphi$  seems to be a good choice.

**Table 4.** Comparison of WOMBAT F-Measure against the approaches proposed in [8] on 6 benchmarks using 2% of the original data as training data.

Dataset	Pessimistic	Re-weighted	Simple	Complete
Persons 1	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>	<b>1.00</b>
Persons 2	0.97	<b>1.00</b>	0.80	0.84
Restaurants	0.95	0.94	<b>0.98</b>	0.88
DBLP-ACM	0.93	<b>0.95</b>	0.94	0.94
Amazon-GP	0.39	0.43	<b>0.53</b>	0.45
Abt-Buy	0.36	<b>0.37</b>	<b>0.37</b>	0.36
Average	0.77	<b>0.78</b>	0.77	0.74

## 6 Related Work

There is a significant body of related work on *positive only learning*, which we can only briefly cover here. The work presented by [13] showed that logic programs are *learnable* with arbitrarily low expected error from positive examples only. [18] introduced an algorithm for learning from labeled and unlabeled documents based

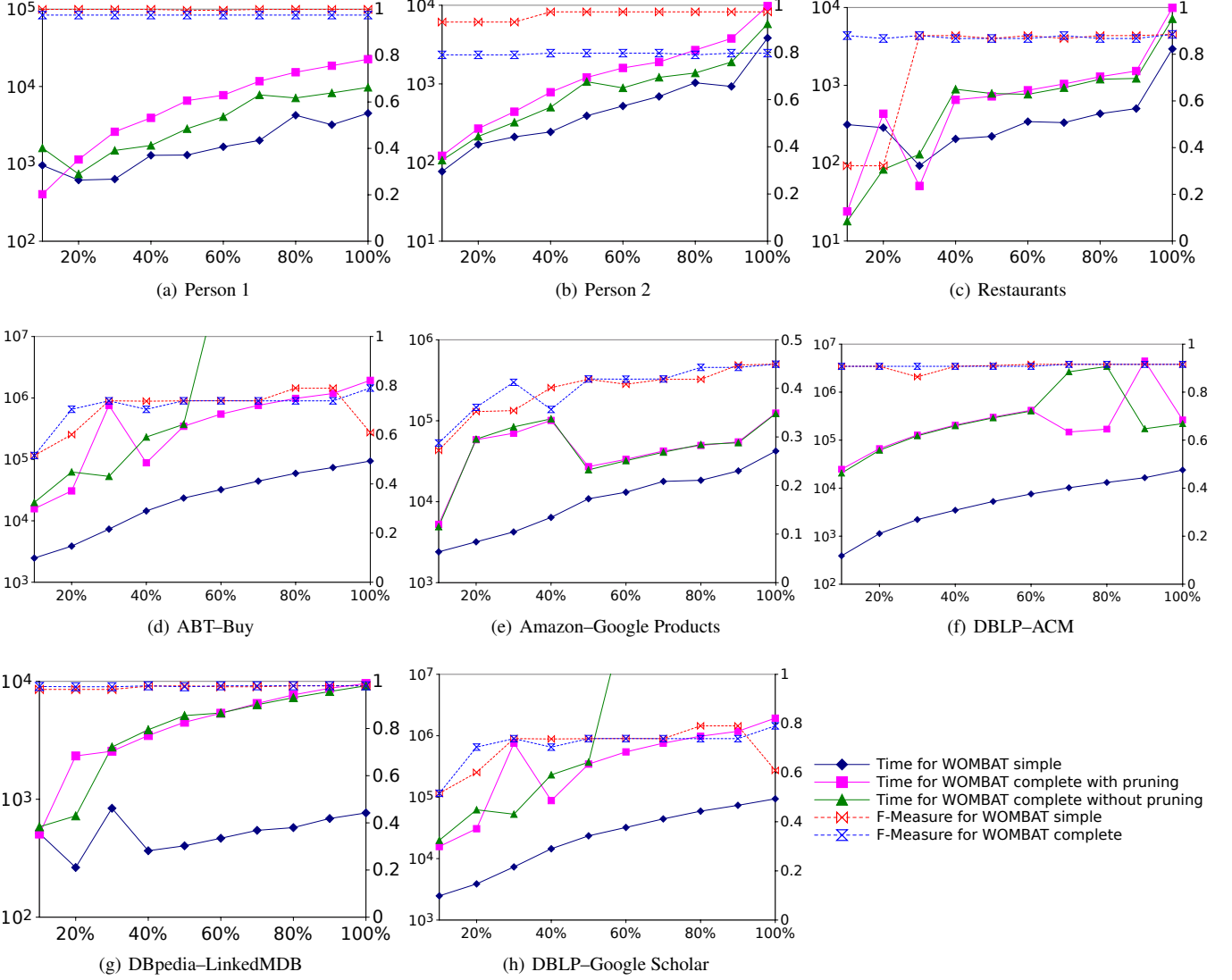
**Table 5.** The *pruning factor* of the benchmark datasets.

Dataset	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
Person 1	1.57	2.13	1.85	2.13	2.13	2.13	2.13	2.13	2.13	2.13
Person 2	1.29	1.29	1.57	1.57	1.57	1.57	1.57	1.57	1.57	1.57
Restaurant	1.17	1.45	1.17	1.45	1.45	1.45	1.45	1.45	1.45	1.45
DBLP-ACM	6.23	5.58	6.79	6.85	6.85	6.85	6.79	6.79	6.93	6.79
Abt-Buy	3.38	3.00	3.00	3.39	3.39	3.39	1.79	3.39	3.39	3.39
Amazon-GP	1.14	1.38	1.33	1.37	1.38	1.45	1.54	1.59	1.60	1.60
DBP-LMDB	1.00	1.86	2.86	1.86	1.86	2.33	2.36	2.36	2.36	2.36
DBLP-GS	1.79	1.93	2.01	2.36	2.45	1.66	2.44	2.26	1.97	2.05

on the combination of Expectation Maximization (EM) and a naive Bayes classifier. [2] provides an algorithm for learning from positive and unlabeled examples for statistical queries. The pLSA algorithm [23] extends the original probabilistic latent semantic analysis, which is a purely unsupervised framework, by injecting a small amount of supervision information from the user.

For learning with *refinement operators*, significant previous work exists in the area of Inductive Logic Programming and more generally concept learning which we only briefly sketch here. A milestone was the Model Inference System in [20]. Shapiro describes how refinement operators can be used to adapt a hypothesis to a sequence of examples. Afterwards, refinement operators became widely used as a learning method. In [22] some general results regarding refinement operators in quasi-ordered spaces were published. Nonexistence conditions for ideal refinement operators relating to infinite ascending and descending refinement chains and covers have been developed. This has been used to show that ideal refinement operators for clauses ordered by  $\theta$ -subsumption do not exist. Unfortunately, we could not make use of these results directly, because proving properties of covers in description logics without using a specific language is likely to be harder than directly proving the results. [?] discussed refinement for different versions of subsumption, in particular weakenings of logical implication. A few years later, it was shown in [?] how to extend refinement operators to learn general prenex conjunctive normal form. Perfect operators, i.e. operators which are weakly complete, locally finite, non-redundant, and minimal, were discussed in [?]. Because such operators do not exist for clauses ordered by  $\theta$ -subsumption, as previously shown in [22], weaker versions of subsumption were considered. This was later extended to theories, i.e. sets of clauses [?]. A less widely used property of refinement operators, called flexibility, was discussed in [?]. Flexibility essentially means that previous refinements of an operator can influence the choice of the next refinement. The article discusses how flexibility interacts with other properties and how it influences the search process in a learning algorithm. For description logics, a significant body of work has been devoted to the study of refinement operators. In [3] and later [4], algorithms for learning in description logics (in particular for the language  $\mathcal{ALC}$ ) were created which also make use of refinement operators. Recent studies of refinement operators include [11, 12] which analysed properties of  $\mathcal{ALC}$  and more expressive description logics. A constructive existence proof for ideal (complete, proper and finite) operators in the lightweight  $\mathcal{EL}$  description logics has been shown in [10].

Most LD approaches for *learning LS* developed are supervised. One of the first approaches to target this goal was presented in [6]. While this approach achieves high F-measures, it also requires large amounts of training data. Hence, methods based on active learning have also been developed (see, e.g., [7, 17]). Still, these approaches



**Figure 3.** Runtime and F-measure results of WOMBAT. The x-axis represents the fraction of positive examples from the gold standard used for training. The left y-axis represents the learning time in milliseconds with time out limit of  $10^7$  ms, processes running above this upper limit were terminated, all time plots are in log scale. The right y-axis represents the F-measure values.

are not guaranteed to require a small amount of training data to converge [19, 16]. The main advantage of unsupervised learning techniques is that they do not require any training data to discover mappings. Moreover, the classifiers they generate can be used as initial classifiers for supervised LD approaches. In general, these approaches assume some knowledge about the type of links that are to be discovered. For example, unsupervised approaches for ontology alignment such as PARIS [21] aim to discover exclusively `owl:sameAs` links. Newer unsupervised techniques for learning LS include approaches based on probabilistic models [21] and genetic programming [19, 16], which all assume that a 1-to-1 mapping is to be discovered. The main drawback of this approach is that it is not deterministic, which was addressed by EUCLID [16]. To the best of our knowledge, this paper presents the first LD approach designed to learn from positive examples only.

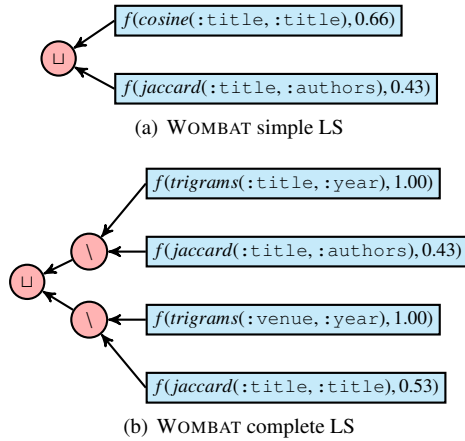
## 7 Conclusions and Future Work

We presented the (to the best of our knowledge) first approach to learn LS from positive examples via generalisation over the space of LS. We presented a simple operator  $\varphi$  that aims to achieve this goal as well as the complete operator  $\psi$ . We evaluated  $\varphi$  and  $\psi$  against state-of-the-art link discovery approaches and showed that we outperform them on benchmark datasets. We also considered scalability and showed that  $\psi$  can be brought to scale similarly to  $\varphi$  when combined with the pruning approach we developed. In future work, we aim to parallelize our approach as well as extend it by trying more aggressive pruning techniques for better scalability.

## REFERENCES

- [1] Sören Auer, Jens Lehmann, Axel-Cyrille Ngonga Ngomo, and Amrपालi Zaveri, ‘Introduction to linked data and its lifecycle on the web’, in *Reasoning Web*, pp. 1–90, (2013).





**Figure 4.** Best LS learned by WOMBAT for the DBLP-GoogleScholar data set.

- [2] François Denis, Rémi Gilleron, and Fabien Letouzey, ‘Learning from positive and unlabeled examples’, *Theoretical Computer Science*, **348**(1), 70–83, (2005). Algorithmic Learning Theory (ALT 2000) 11th International Conference, Algorithmic Learning Theory 2000.
- [3] Floriana Esposito, Nicola Fanizzi, Luigi Iannone, Ignazio Palmisano, and Giovanni Semeraro, ‘Knowledge-intensive induction of terminologies from metadata’, in *The Semantic Web - ISWC 2004: Third International Semantic Web Conference, Hiroshima, Japan, November 7-11, 2004. Proceedings*, pp. 441–455. Springer, (2004).
- [4] Luigi Iannone, Ignazio Palmisano, and Nicola Fanizzi, ‘An algorithm based on counterfactuals for concept learning in the semantic web’, *Applied Intelligence*, **26**(2), 139–159, (2007).
- [5] R. Isele, A. Jentzsch, and C. Bizer, ‘Efficient Multidimensional Blocking for Link Discovery without losing Recall’, in *WebDB*, (2011).
- [6] Robert Isele and Christian Bizer, ‘Learning Linkage Rules using Genetic Programming’, in *Sixth International Ontology Matching Workshop*, (2011).
- [7] Robert Isele, Anja Jentzsch, and Christian Bizer, ‘Active learning of expressive linkage rules for the web of data’, in *ICWE*, pp. 411–418, (2012).
- [8] Mayank Kejriwal and Daniel P Miranker, ‘Semi-supervised instance matching using boosted classifiers’, in *The Semantic Web. Latest Advances and New Domains*, 388–402, Springer, (2015).
- [9] Hanna Köpcke, Andreas Thor, and Erhard Rahm, ‘Evaluation of en-

- tity resolution approaches on real-world match problems’, *Proc. VLDB Endow.*, **3**(1-2), 484–493, (September 2010).
- [10] Jens Lehmann and Christoph Haase, ‘Ideal downward refinement in the EL description logic’, in *Inductive Logic Programming, 19th International Conference, ILP 2009, Leuven, Belgium*, (2009).
- [11] Jens Lehmann and Pascal Hitzler, ‘Foundations of refinement operators for description logics’, in *ILP*, volume 4894 of *Lecture Notes in Computer Science*, pp. 161–174. Springer, (2007).
- [12] Jens Lehmann and Pascal Hitzler, ‘Concept learning in description logics using refinement operators’, *Machine Learning journal*, **78**(1-2), 203–250, (2010).
- [13] Stephen Muggleton, ‘Learning from positive data’, in *Inductive logic programming*, 358–376, Springer, (1997).
- [14] Axel-Cyrille Ngonga Ngomo, ‘Link discovery with guaranteed reduction ratio in affine spaces with minkowski measures’, in *Proceedings of ISWC*, (2012).
- [15] Axel-Cyrille Ngonga Ngomo and Klaus Lyko, ‘Eagle: Efficient active learning of link specifications using genetic programming’, in *Proceedings of ESWC*, (2012).
- [16] Axel-Cyrille Ngonga Ngomo and Klaus Lyko, ‘Unsupervised learning of link specifications: deterministic vs. non-deterministic’, in *Proceedings of the Ontology Matching Workshop*, (2013).
- [17] Axel-Cyrille Ngonga Ngomo, Klaus Lyko, and Victor Christen, ‘Coala – correlation-aware active learning of link specifications’, in *Proceedings of ESWC*, (2013).
- [18] Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell, ‘Text classification from labeled and unlabeled documents using em’, *Machine learning*, **39**(2-3), 103–134, (2000).
- [19] Andriy Nikolov, Mathieu d’Aquin, and Enrico Motta, ‘Unsupervised learning of link discovery configuration’, in *The Semantic Web: Research and Applications*, 119–133, Springer, (2012).
- [20] E. Y. Shapiro, ‘Inductive inference of theories from facts’, in *Computational Logic: Essays in Honor of Alan Robinson*, eds., J. L. Lassez and G. D. Plotkin, 199–255, The MIT Press, (1991).
- [21] Fabian M. Suchanek, Serge Abiteboul, and Pierre Senellart, ‘PARIS: Probabilistic Alignment of Relations, Instances, and Schema’, *PVLDB*, **5**(3), 157–168, (2011).
- [22] P. R. J. van der Laag and S-H. Nienhuys-Cheng, ‘Existence and nonexistence of complete refinement operators’, in *ECML*, eds., F. Bergadano and L. De Raedt, volume 784 of *Lecture Notes in Artificial Intelligence*, pp. 307–322. Springer-Verlag, (1994).
- [23] Ke Zhou, Xue Gui-Rong, Qiang Yang, and Yong Yu, ‘Learning with positive and unlabeled examples using topic-sensitive pls’, *Knowledge and Data Engineering, IEEE Transactions on*, **22**(1), 46–58, (Jan 2010).