

CogMap: A Cognitive Support Approach to Property and Instance Alignment

Jan Nöbner, David Martin, Peter Z. Yeh, and Peter F. Patel-Schneider

Nuance Communications, AI Research Group, Sunnyvale, CA 94085, USA
jan.noessner@gmail.com, <first>.<last>@nuance.com

Abstract. The iterative user interaction approach for data integration proposed by Falconer and Noy can be generalized to consider interactions between integration tools (generators) that generate potential schema mappings and users or analysis tools (analyzers) that select the best mapping. Each such selection then provides high-confidence guidance for the next iteration of the integration tool. We have implemented this generalized approach in COGMAP, a matching system for both property and instance alignments between heterogeneous data. The generator in COGMAP uses the instance alignment from the previous iteration to create high-quality property alignments and presents these alignments and their consequences to the analyzer. Our experiments show that multiple iterations as well as the interplay between instance and property alignment serve to improve the final alignments.

1 Introduction

In recent years, companies have spent more and more effort in building knowledge graphs based on light-weight ontologies, which incorporate data from multiple heterogeneous sources (which we will henceforth call “information stores”). A key challenge of these efforts is determining the best alignment of the schema of a new store to the ontology of the knowledge graph, while minimizing the “manual” analytical effort required of a human knowledge engineer.

Most of the current ontology alignment systems, such as those evaluated recently in the annual ontology alignment evaluation initiative [1], have several limitations. Most of these alignment algorithms solve one integration problem (deriving a mapping between two ontologies) using a fully-automated, “one-shot” approach. Thus, they are often not able to improve by iterating over previous alignments. Partly for this reason, the results of fully automated algorithms are often error prone [33] and cannot be reliably used for high-quality data integration.

Currently much information to be integrated is obtained from non-ontological sources such as relational databases or XML documents. Classical ontology alignment systems are often not able to process this data [1]. To address this need, systems like ONTODB [21] and standards like D2RQ [4] have emerged. However, these solutions do not include semi-automated alignment algorithms which take instance information into account.

Our approach, implemented in the COGMAP system, follows a cognitively-inspired, iterative approach. With multiple iterations the system is able to improve over time,

since it builds on the results of previous iterations (or, in the case of the first iteration, seed queries given by the user). At each iteration, the results are augmented with new information that has been verified by a user or automated verification capability.

COGMAP uses instance information to perform property alignment. While most state-of-the-art schema alignment algorithms do not take instance information into account, focusing exclusively on the alignment of classes and properties and mainly considering their labels or structural information [8, 30], using instance matching has attracted more and more attention over the last years [17].

COGMAP explores instances by not only focusing on data properties but also taking object properties into account. In the case of databases, it follows foreign keys; with RDF information stores it explores sub-tags. To the best of our knowledge, there exists no other approach which explores the space of potential mappings between information stores as we do.

COGMAP is not restricted to the alignment of information based on formal ontologies. It also supports relational databases and XML documents, which can serve either as the source or the target of an alignment. In addition, COGMAP allows support for other data formats to be added in a modular fashion.

2 Related Work

Many schema alignment systems have been developed in ontology matching. The development of these systems has largely been driven by the available benchmark datasets of the ontology alignment evaluation initiative. An overview of the current systems and their evaluation is given by Grau *et al.* [15]. The most important datasets, however, cover only a small problem space.

Although most ontology matching systems ignore instances, there exists a strand of literature which combines schema alignment and instance alignment [17]. Bilke and Naumann [3] developed an approach that first aligns instances and uses this information for schema alignment. Their evaluation is based on artificially populated data whereas we employ real-world data information stores like FREEBASE and DBPEDIA. Bilke *et al.* [2, 27], Thor *et al.* [35], Gal [13], and Leme *et al.* [25] use instances to align schema and resolve conflicts. Another fully automated system that integrates both schema and instance alignment is PARIS [34]. Its algorithms are, however, resource intensive, in some cases taking days to produce a solution. In contrast, the COGMAP algorithms are much less resource intensive and can be run on a typical desktop computer. Wang *et al.* investigates the problem of having only a few non-overlapping instances by approaching the mapping problem as a classification problem. However, this approach is limited to mapping concepts and ignores properties. Duan *et al.* [6] use hashing techniques to speed up instance-based matching. Nunes *et al.* [31] present an instance-based algorithm for complex data property matching. A prominent example is the system RIMOM, which dynamically combines several alignment strategies including instance alignment [26]. Due to its recent excellent achievements at the ontology alignment evaluation initiative, we chose this system for our evaluation.

To the best to our knowledge, none of these approaches is exploring object properties with an iterative cognitive support approach. QUICKMIG [5] is a migration tool for

database systems which follows a semi-automated approach. However, it considers only exact value matches and their results are not used to improve the ongoing iterations.

A smaller number of systems utilize learning. A prominent example is SILK [20, 19, 18] which learns expressive linking rules by using genetic programming. However, its target user is a technical expert who can, e.g., analyse complex matching trees while COGMAP focuses on domain experts by hiding technical complexity. LIMES [28] focuses on runtime improvements by using the triangle inequality. However, it does not allow a user-centric iterative approach. Furthermore, neither system is able to map data properties to object properties, which is required by the real-world datasets we examined. (See the algorithm section for details.)

Recently, the ontology alignment evaluation initiative initiated an interactive track which simulates interactive matching [32], where a human expert is involved to validate mappings found by the matching system. The client was modified to allow interactive matchers to ask an oracle, which emulates a perfect user. The interactive matcher can present a correspondence to the oracle, which then tells the user whether the correspondence is right or wrong. However, the initiative uses a dataset which does not contain any instance data and thus is not suitable for evaluating our approach. The two most successful participating systems 2014 were AML[12] with respect to gained f-measure due to the interactive approach and LOGMAP [23] with respect to efficiency (number of interactions required). We have included both systems in our evaluation.

Tools have been developed to support the alignment of databases to ontologies. One example is ONTOP (`ontop.inf.unibz.it`), which provides a PROTÉGÉ plug-in to facilitate the creation of integration rules. ONTOP focuses on fast execution of already existing data integration rules, but not on the (semi-)automated construction of them. Furthermore, its target ontology is assumed to be small and to contain only schema information but no instance information. There have also been attempts to build graphical tools for supporting the user in data integration. KARMA [24], for example, loads data from different information stores and uses instance information for schema alignment. However, its approach is different from our algorithm. KARMA learns the general structure of fields based on previous alignments made whereas COGMAP operates on instance information. Two disadvantages of KARMA's approach are that it generally assumes that fields (e.g., ids) have similar structures in different datasets and its algorithms require a large amount of training data.

3 The Cognitive Support Approach

Researchers in ontology and schema matching have recently recognized the need for various types of cognitive support in aligning complex conceptual models [9, 11]. Most approaches are based on advanced visualization of the models to be integrated and the mappings created by the user [14]. While the appropriate use of visualizations is known to be a key aspect for successful manual data integration, visualizations quickly reach their limits in the presence of very complex or very large models.

As a result, recent work has tried to go beyond pure visualization support to include cognitively efficient interaction strategies to support the user [10]. Falconer [9] proposed an interactive strategy for data integration where the integration task is dis-

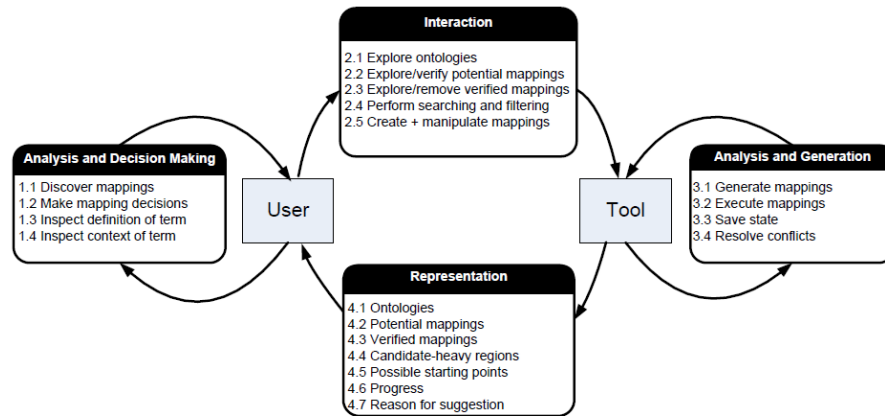


Fig. 1. The cognitive support model for data integration by Falconer [10].

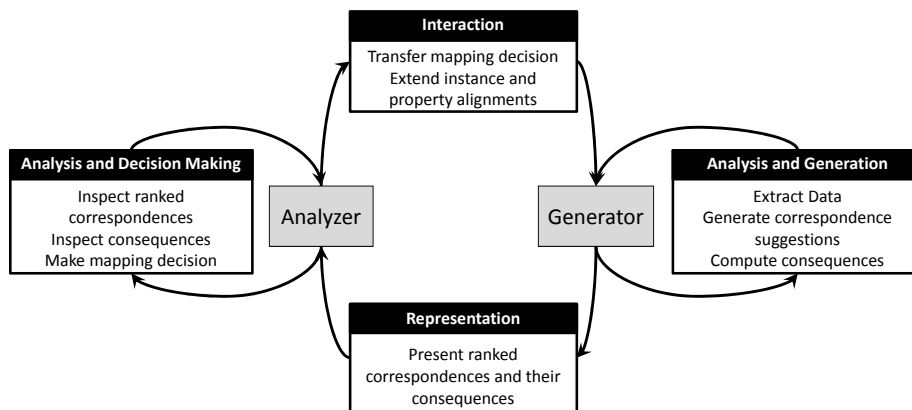


Fig. 2. Modified cognitive support Model as implemented in COGMAP.

tributed between the user and the tool (Figure 1). The MappingAssistant [33] project used a modified cognitive support model for data integration, focusing on detecting and correcting incorrect data integration rules.

In our implementation of the cognitive support model, which we call COGMAP, we go one step further and allow the “user” to be either a human user or an intelligent automated agent. Thus, in our implementation of the cognitive support model (Figure 2), we distinguish between an *analyzer* and a *generator*, instead of a user and a tool.¹

¹ As stated, the analyzer can be a human. As this paper is about the effectiveness of the overall approach, we only use simple agents. Sophisticated automated agents or humans can utilize world knowledge or judgements to select better alignments instead of just picking the highest-scoring ones.

In each iteration, COGMAP extracts data based on the results of previous iterations (or, in the first iteration, based on given seed queries), generates property correspondence suggestions, and computes the consequences for the top suggestions. These consequences are the instances that would be aligned if the analyzer selects (verifies) this property correspondence. Next, COGMAP sends the ranked correspondences and their consequences to the analyzer. The analyzer then inspects this information and selects a correspondence. The selected property correspondence, and the resulting instance alignment, are added to the evolving results sets, which allows the system to improve its suggestions in subsequent iterations. The algorithm terminates when no properties remain to generate new correspondences, or no correspondence is selected by the analyzer.

COGMAP is designed to cope with many different types of data stores, in many different formats. We currently have implemented support for RDF accessed via SPARQL, relational databases, and general XML based files. This list is easily extended by implementing our `Connector` interface.

4 Algorithm

The primary focus of the COGMAP algorithm (Algorithm 1) is to construct property correspondences and instance alignments. An alignment (or mapping) consists of a set of correspondences. According to Euzenat *et al.* [7], a correspondence is a 4-tuple $\langle e_s, e_t, r, c \rangle$, where e_s and e_t are source and target entities, r is a semantic relation, and c is a confidence value (usually, $c \in [0, 1]$). Like most ontology alignment systems [1], we focus on equivalent relations $\langle e_s, e_t, \equiv, c \rangle$.

The algorithm can be split into three phases. The *data extraction* (lines 4-6) and *data exploration* (lines 13-16) phases are only executed in the first iteration ($i = 0$). The *alignment generation and selection* (lines 7-12) phase is repeated until no more correspondences are found. This phase includes the decision making of the analyzer. The following subsections will explain the phases in more depth.

4.1 Data Extraction and Exploration

We adapt the terms data property, object property, and instances from the semantic web literature, extending them to databases and XML documents in an obvious fashion. For example, instances include database rows and XML nodes.

In the *data extraction* phase, we extract all data property names and their corresponding values for M instances into a source table T_s and a target table T_t (line 5). The left block (2nd column) of Table 1 illustrates the general form of T_s and T_t after extraction.

In the *data exploration* phase, we explore the search space by following object properties. In other words, for each object property op of an instance i , we examine the object which is the value of that property. Then, for each data property of that object, we add its value to the row for i . Thus, the right blocks of Table 1 (op_a, op_b, \dots) are added during this phase. The reason this exploration happens at the end of the first iteration ($i = 0$, line 13-16) is that there may exist many object properties to follow. This often leads to a large amount of data. Thus, the idea is to restrict the exploration to the smaller instance

Algorithm 1 High-level algorithm of COGMAP.

Input: S_s, S_t : Seed queries for source and target**Input:** M : number of extracted instances of each information store (default: 5000)**Input:** k : number of suggestions (default: 5)**Output:** \mathcal{X}, \mathcal{Y} : Set of user-verified property correspondences and instance correspondences

GETALIGNMENTS

```
1:  $\mathcal{X}, \mathcal{Y} \leftarrow \emptyset$ 
2:  $i \leftarrow 0$ 
3: repeat
   $\triangleright$  Data Extraction
4: if  $i=0$  then
5:    $T_s, T_t \leftarrow$  Extract  $M$  instances and their data properties and values based on seeds  $S_s$ 
   and  $S_t$ .
6: end if
   $\triangleright$  Alignment Generation and Selection
7:  $X_i \leftarrow$  Compute top- $k$  property correspondence suggestions based on  $T_s, T_t$  and  $\mathcal{Y}$  (if not
  empty).
8: for every  $x \in X_i$  do
9:    $Y_x \leftarrow$  Compute instance alignment consequences for  $x$  based on  $T_s, T_t$  and  $\mathcal{Y}$  (if not
  empty).
10: end for
11: Analyzer selects the optimal  $x \in X_i$  based on  $X_i$  and  $\{Y_x \mid x \in X_i\}$ .
12: add  $x$  to  $\mathcal{X}, \mathcal{Y} \leftarrow Y_x$ .
   $\triangleright$  Data Exploration
13: if  $i=0$  then
14:    $I_s, I_t \leftarrow$  Extract source and target instance sets from instance alignment  $\mathcal{Y}$ .
15:    $T_s, T_t \leftarrow$  Extend tables by following the object-property assertions of  $I_s$  and  $I_t$ .
16: end if
17:  $i \leftarrow i + 1$ 
18: until No more suggestions found
```

Table 1. General form of source and target table. Initially, the direct data properties and the corresponding data are imported (left block). Second, and subsequent, steps further explore the data by including object properties (right blocks). (dp =data property, op = object property, i = instance, and v = value).

	dp_1	\dots	dp_n	$dp_{1,a}$	op_a	$dp_{n,a}$	op_b
i_1	$v_{1,1}$	\dots	$v_{1,n}$	$v_{1,1,a}$	\dots	$v_{1,n,a}$	\dots
\dots	\dots	\dots	\dots	\dots	\dots	\dots	\dots
i_m	$v_{m,1}$	\dots	$v_{m,n}$	$v_{m,1,a}$	\dots	$v_{m,n,a}$	\dots

sets I_s and I_t . These instance sets are extracted from the first instance alignment \mathcal{Y} (line 15) such that $I_s = \{e_s \mid \langle e_s, e_t, \equiv, c \rangle \in \mathcal{Y}\}$ and $I_t = \{e_t \mid \langle e_s, e_t, \equiv, c \rangle \in \mathcal{Y}\}$. Then, COGMAP only follows the object properties for the instance sets I_s and I_t which are usually much smaller.

For RDF repositories we utilize SPARQL queries to access data. We do not rely on the completeness of domain and range restrictions for extracting properties, since they are often poorly defined (e.g., in DBPEDIA). Instead we take the distinct set of all properties of the relevant instances (those retrieved by S for extraction or those identified as values of object properties for expansion) as the relevant data properties. For relational data extraction, we just add the limit M to the seed SQL query S , execute the query, and store the result in table T . For exploration, we follow the foreign keys according to the definitions in the database schema. For XML files, we extract every attribute and every direct child node that has a primitive value from the initial XPath expression S . We store both attribute and child node values as data properties in Table 1. For the exploration phase, we inspect the children x of all non-primitive nodes. From these nodes, we again store the values of every attribute and every direct child node that has primitive values.

For some properties, a given instance may have multiple values. For simplicity, we consider only single values in this presentation. In practice, we have found the concatenation of multiple values to be effective. More sophisticated strategies will be developed in future work. On the other hand, there may exist properties and/or instances which have almost no assertions, especially in large RDF knowledge bases and XML documents. To ensure the effectiveness of the approach, those assertions might need to be ignored. To cope with that issue, COGMAP has an optional parameter ϕ to filter instances and data properties with sparse value assertions.

4.2 Alignment Generation and Selection

The goal of COGMAP is to establish instance correspondences and correspondences between data and object properties. In doing so, the space of ontology elements (in RDF stores) or schema elements (in relational and XML stores) that COGMAP considers is constrained by the seed information S . In addition, instance alignments are constrained by the domains and ranges of property alignments. For example, if we align a property $e_s = \text{id}$ to a property $e_t = \text{movieName}$, then the resulting instance alignment is bounded by $S_s = \text{movie}$ and $S_t = \text{film}$ as domains. Thus, there is no need to consider possible correspondences involving other instances in the information stores.

In addition to the standard data-property to data-property, object-property to object-property, and instance to instance correspondences, we also support object- and data-property to data-property correspondences (Example: $e_t = \text{film/country/. /name/}$ and $e_s = \text{movie/language}$).

Line 7 of Algorithm 1 first computes the top- k property correspondence suggestions X_i . In the first iteration ($i = 0$), COGMAP uses all available instance data since no instance alignment exists yet ($\mathcal{Y} = \emptyset$). For computational reasons, we use an implicit cutoff at this initial stage. In the following iterations, we can improve the suggestions by considering the instance correspondences from the previous iteration and comparing the property values only for the instance pairs in \mathcal{Y} . In these iterations, we do not apply any threshold but rank the correspondences at the end.

Then, we compute the consequences Y_x for the top- k suggestions X_i (line 8-10). That is, for each of those suggested property correspondences, we compute the instance alignment that will result if this correspondence is selected by the analyzer. Initially

Table 2. Selected implemented components.

Aggregators		
Unions or joins of sets of correspondences. Average, maximum, or multiples of confidence values if correspondences share the same source entity e_s and target entity e_t .		
Name	Description	Filters
TopKFilter	Returns the top- k correspondences with the highest confidence value c .	
OneToOneFilter	Returns a functional one-to-one alignment. We implemented a greedy strategy. First, it orders the correspondences in descending order. Then, it traverses through the list and drops all correspondences whose entities e_s or e_t have been already matched.	
Property Matchers		
Name	Description	
PropertyNameMatcher	Matches properties according to their name.	
ValueLengthMatcher	Matches properties p_1 and p_2 with close average value length /	
DistinctValueMatcher	close percentages of distinct row entries l_1 and l_2 . The similarity is computed with $Min(l_1, l_2) / (Max(l_1, l_2))$.	
InstanceBasedMatcher	If instance alignments $\mathcal{Y} = \emptyset$, we align properties by concatenating all values of all instances for each property and compute the string similarity. If $\mathcal{Y} \neq \emptyset$, we compute the property similarities for every instance pair in \mathcal{Y} separately and average over the results.	
Instance Matchers		
Align instances by concatenating every value for every property and computing their string similarity. If a specific property p is given, consider only values of that property. If an instance alignment \mathcal{Y} is given, traverse through that alignment and update the similarities based on the string values of all the given property value(s).		

($i = 0$), all source instances are compared against all target instances. In following iterations the instance alignment \mathcal{Y} from the previous iteration is used to compute the new alignment. The threshold applied for the instance alignment equals the confidence value c of $\langle e_1, e_2, \equiv, c \rangle \in X_i$.

The value of k is relatively unimportant here. As long as a correct correspondence is in the top- k suggestions, the results of the approach will not be significantly affected. We have found that $k = 5$ is generally adequate to achieve this condition, and results in a reasonable load on human analysts. In an automated setting it would be easy to use a larger k , which might produce slightly better results at the price of of somewhat longer run times (to score the extra suggestions).

Finally, the analyzer selects the optimal $x \in X_i$ based on the suggestions X_i and the consequences $\{Y_x | x \in X_i\}$. As noted above, this selection can either be made by a human user or by an automated selection function that takes the confidence value and the suggestions into account. In this paper, we use only a simple automated agent that selects the best-scoring alignment. Employing humans or more-sophisticated agents would presumably produce better results, but then any advantage of the approach might only come from the intelligence in the human or agent—using a simple agent means that the benefits come from the overall alignment philosophy. (We plan to address elsewhere the user interface issues associated with supporting selections by a human.) After se-

Table 3. Benchmark Statistics.

	Benchmark (1)		Benchmark (2)	
	DBPEDIA	EPG	FREEBASE	FANDANGO
	People	Cast	Film	Movie
Format	RDF	RDB	RDF	XML
Data Properties	6	18	233	26
Object Properties	0	10	234	14
Instances	1,045,474	6,857	247,608	100,959

lection of x , we update the seed property alignment \mathcal{X} and the seed instance alignment \mathcal{Y} for the next iteration.

COGMAP supports many different components to match instances and properties. Every `Filter`, `Aggregator`, and `Matcher` is a component. Each component has an `execute()` method, which returns a set of alignments.

The components are organized as a tree. The `Matchers` form the leaves. They take a source table T_s , a target table T_t , a set of previously verified property alignments \mathcal{X} and a set of instance alignments \mathcal{Y} from the previous iteration as input. An `Aggregator` executes every component in the list `cs` and aggregate the results. It might, for example, just take the maximum confidence value c of all correspondences with equal entities e_s and e_t . A `Filter` reduces the size of the alignment of its component after executing it. A simple filter might, for example, only return the correspondences for which confidence values c are above a certain threshold.

Table 2 lists a selection of implemented components and a short explanation of their functionality. COGMAP incorporates mechanisms to deal with different date and number formats, which are omitted here for brevity, and easy interfaces to facilitate new component development. Figure 3 provides example trees built from these components.

5 Evaluation

We have selected two natural alignment tasks using real-world data, assessed the performance of COGMAP on them benchmarks, and compared its performance with that of AML, LOGMAP, and RIMOM.

5.1 Benchmarks

The first benchmark aligns all people from DBPEDIA in FOAF format (wiki.dbpedia.org/Downloads39#persondata) with cast information of all programs playing on TV in the U.S. over a two week window from a commercial Electronic Program Guide (EPG) database. The second benchmark aligns FREEBASE films (www.freebase.com/film/film) with movie data from FANDANGO (www.fandango.com). Table 3 provides details on the number of instances and properties of each benchmark.

We designed and selected these benchmarks because existing state-of-the-art ontology alignment benchmarks, e.g., in the Ontology Alignment Evaluation Initiative

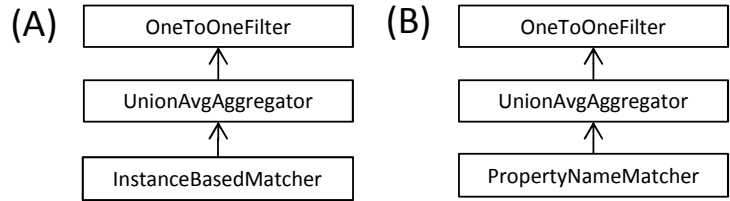


Fig. 3. Experiment Configurations.

campaigns², lack sufficient instance data, which are required by COGMAP. Because of the size of these benchmarks, it was not possible to prepare in advance an official gold-standard. Instead, a human judge was employed to grade the correctness of the alignment results, which we discuss below.

5.2 Experiment Setup

We used the following experiment setup to answer three key questions:

- What is the impact of implementing a cognitive support model for alignment?
- What is the impact of using instance data for alignment?
- How general is our solution?

We first setup our solution, COGMAP, using configuration (A) in Figure 3. A description of each component used in configuration (A) can be found in Table 2. COGMAP analyzes the property correspondences, and selects the one with the highest confidence value to iterate on (see lines 11 and 12 of Algorithm 1).

We then created two variants of COGMAP to answer the first two experimental questions above. We first created a variant—called INSTMAP—by ablating the cognitive support model used by COGMAP. INSTMAP still uses instance data but does not iterate on the results to further improve alignment.

We also created a second variant—called Baseline—by ablating both the cognitive support model and the use of instance data. Baseline performs alignment using a property name matcher, but the other configuration components are the same (see configuration (B) in Figure 3).

Moreover, we selected three state-of-the-art ontology alignment systems [16] to compare COGMAP against, in order to assess its practical impact. The three systems are AML [12], LOGMAP [23], and RIMOM [26]. AML is focused on computational efficiency and designed to handle very large ontologies. It is the leading system in the conference and anatomy tracks of the 2014 ontology alignment evaluation, in terms of f-measure. LOGMAP provides a scalable logical ontology alignment framework. RIMOM automatically combines multiple alignment strategies with the goal of finding the optimal alignment results. We selected these systems because they are the most established systems in the 2014 ontology alignment evaluation, and an executable version is available to the public.

² See oaei.ontologymatching.org.

Table 4. Benchmark 1 results.

	Baseline	AML	LOGMAP	RIMOM	INSTMAP	COGMAP
$nDCG@3$	0.38	0.76	0.38	0.38	0.76	1.00
$nDCG@6$	0.35	0.51	0.25	0.48	0.74	0.89
$P@3$	0.33	0.67	0.33	0.33	0.67	1.00
$P@6$	0.33	0.33	0.17	0.50	0.67	0.83
Runtime in sec	0.4	2.7	9.1	5.8	1.9	3.0

Finally, we applied all systems above to both Benchmarks 1 and 2 to assess their generality, and hence answer the third experimental question. Unless otherwise noted, we set the number of instances to use from each benchmark to $M = 5000$, and the fraction of non-null values required for each property to $\phi = 0.1$. We also converted each benchmark into the RDF OWL syntax because many of the ontology matching systems compared cannot directly consume databases or XML files.

All experiments were run on a desktop PC with 4GB of RAM and an Intel i5 duo-core processor. We used FAST-JOIN [38] as the underlying matching algorithm for instances. FAST-JOIN combines both token-based similarity (Jaccard, Cosine, or Dice) and string edit distance. Moreover, it is currently the fastest matching algorithm (see [22]), by implementing efficient pruning and hashing techniques, with soundness and completeness guarantees. This efficiency is required because of our large benchmarks, which make it infeasible to compare every source instance with every target instance.

The output of each system was graded by a human judge familiar with the data sources in each benchmark³ using the metrics of *Precision at n* ($P@n$) and the *normalized (logarithmic) Discounted Cumulative Gain at n* ($nDCG@n$) [39] where n denotes that the top- n results. Precision P is defined as:

$$P = \frac{|\text{correct correspondences}|}{|\text{retrieved correspondences}|}$$

and $nDCG$ is defined as:

$$nDCG = \frac{rel_1 + \sum_{i=2}^n \frac{rel_i}{\log_2 i}}{(1 + \sum_{i=2}^n \frac{1}{\log_2 i})}$$

where rel_i is 1 if the correspondence at position i is correct and 0 else. $nDCG@n$ gives more weight to correct correspondences that are ranked higher.

5.3 Results and Discussions

Tables 4 and 5 show the results for benchmarks 1 and 2, respectively. From these results, we observed that COGMAP outperformed INSTMAP in most cases. We attribute

³ Determining the correctness of the correspondences produced by each system was simple for the human judge. We thus believe that the use of a human judge in this manner did not introduce any biases and did not affect the comparison.

Table 5. Benchmark 2 results.

	Baseline	AML	LOGMAP	RIMOM	INSTMAP	COGMAP
<i>nDCG@3</i>	0.38	0.76	0.38	0.38	1.00	1.00
<i>nDCG@6</i>	0.25	0.60	0.49	0.38	0.90	1.00
<i>nDCG@9</i>	0.20	0.68	0.39	0.30	0.79	1.00
<i>nDCG@12</i>	0.17	0.68	0.38	0.25	0.69	0.85
<i>P@3</i>	0.33	0.67	0.33	0.33	1.00	1.00
<i>P@6</i>	0.17	0.50	0.50	0.33	0.83	1.00
<i>P@9</i>	0.11	0.67	0.33	0.22	0.67	1.00
<i>P@12</i>	0.08	0.67	0.33	0.17	0.50	0.75
Runtime in sec	2.6	7.0	21.7	33.2	20.5	29.1

this improvement to the only difference between the two systems: COGMAP uses a cognitive support model while INSTMAP does not. Hence, the use of a cognitive support model has a positive impact on alignment results.

We also observed that INSTMAP outperformed Baseline in all cases. We attribute this improvement to the only difference between the two systems: the use of instance data. For example, Baseline could not correctly align the following data properties in benchmark 1 by matching just the names of these properties.

```

first_name ⇔ givenName
last_name ⇔ surName
full_name ⇔ name

```

However, INSTMAP correctly found these alignments because of the overlap between the instances of these properties. Hence, these results show that the use of instance data also has a positive impact on performance.

Finally, we observed that COGMAP outperformed all three state-of-the-art ontology matching systems compared, i.e. AML, LOGMAP, and RIMOM. We attribute this improvement to the following factors:

- COGMAP uses instance data for alignment.
- COGMAP uses an iterative cognitive model for alignment.
- COGMAP can ignore rarely used properties by using the ϕ parameter.

Given the different characteristics of these two benchmark, the results above suggest the general utility of an alignment system like COGMAP that combines a cognitive support model with the use of instance data. Moreover, the additional computation does not contribute to a significant increase in runtime. Across both benchmarks, COGMAP had comparable (or better) runtime than the other state-of-the-art systems compared.

Figures 4 and 5 show the impact of varying ϕ (the fraction of non-null values required for each property) and M (the number of instances used) for COGMAP and INSTMAP on both benchmarks. These results demonstrate the relative robustness of COGMAP to these parameter settings compared to INSTMAP, and further demonstrate the positive impact of using a cognitive support model. For example, we observed on both benchmarks that the performance of COGMAP only became negatively impacted for larger values of ϕ , which was in contrast to INSTMAP. Similarly, the performance of COGMAP increased at a faster rate compared to INSTMAP as M was increased, and plateaued sooner than INSTMAP.

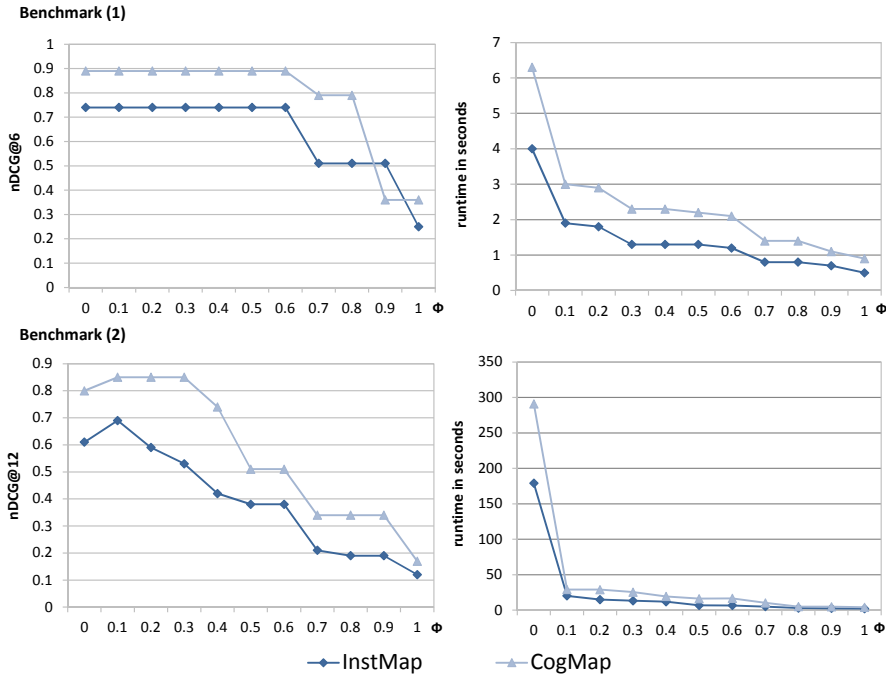


Fig. 4. Results for varying ϕ (number of non-null values) for COGMAP and INSTMAP for both benchmarks. For high ϕ , $nDGC$ and runtime decrease because fewer alignment candidates remain.

6 Conclusion and Future Work

This paper presents a cognitive based approach for aligning properties by taking instance information into account. The approach is implemented in the system COGMAP which iteratively suggests property correspondences and their consequences in terms of instance alignments. In each round, the system is able to improve these alignments based on the user verifications of the previous round. Experiments show that the cognitive based approach outperforms both a baseline approach and the purely instance-based approach.

Currently, the system is restricted to aligning instances and properties. In future work, we will enable class alignments and complex matchings [37]. These complex matchings will be described using the R2RML standard (www.w3.org/TR/r2rml).

We will extend exploration of the knowledge sources. First, we will integrate object properties that are more than one hop away. This will require efficient pruning techniques to avoid an intolerable blowup of both data size and processing requirements. Second, we will use the organization of the knowledge structure (ontologies and schemas, when they are specified) to widen the search space by, e.g., exploring the data of the superclasses.

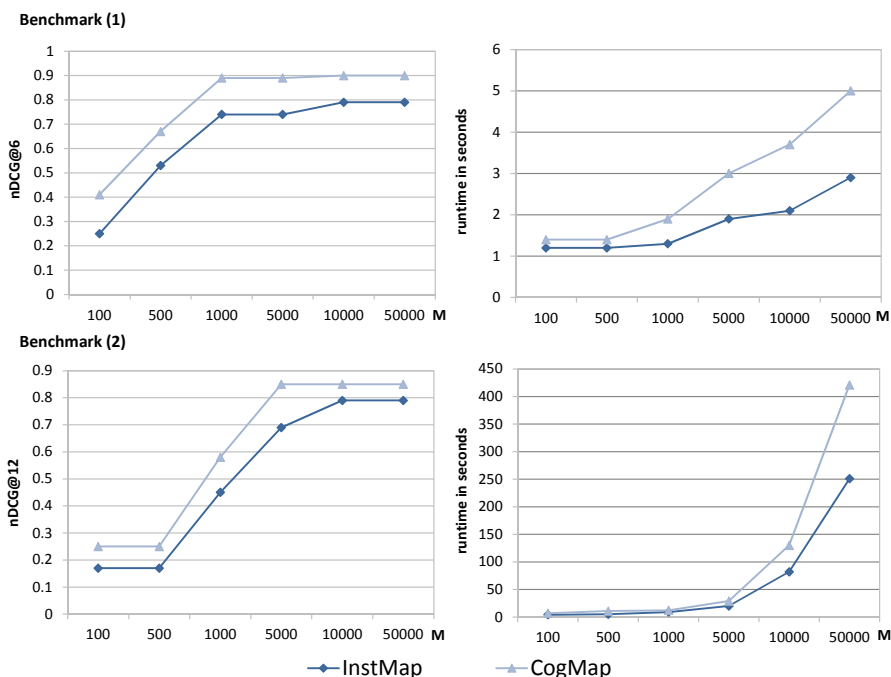


Fig. 5. Results for varying M (number of instances) for COGMAP and INSTMAP for both benchmarks. The more instances included, the higher the overlap and hence better results ($nDCG$).

The knowledge structures will also help to improve the alignment itself by including ideas from [30, 29]). Additionally, tree structure learning algorithms, inspired by [36], will be used to learn the optimal composition of matching trees.

Finally, we plan to explore the possibility of integration into KARMA [24], which we believe would provide a suitable graphical user interface.

References

1. J. L. Aguirre, B. C. Grau, K. Eckert, J. Euzenat, A. Ferrara, R. W. van Hague, L. Hollink, E. Jiménez-Ruiz, C. Meilicke, A. Nikolov, et al. Results of the ontology alignment evaluation initiative 2012. In *Proc. 7th ISWC workshop on ontology matching (OM)*, pages 73–115, 2012.
2. A. Bilke, J. Bleiholder, F. Naumann, C. Böhm, K. Draba, and M. Weis. Automatic data fusion with hummer. In *Proceedings of the 31st international conference on Very large data bases*, pages 1251–1254. VLDB Endowment, 2005.
3. A. Bilke and F. Naumann. Schema matching using duplicates. In *Data Engineering, 2005. ICDE 2005. Proceedings. 21st International Conference on*, pages 69–80. IEEE, 2005.
4. C. Bizer and A. Seaborne. D2rq-treating non-rdf databases as virtual rdf graphs. In *Proceedings of the 3rd international semantic web conference (ISWC2004)*, volume 2004, 2004.

5. C. Drumm, M. Schmitt, H.-H. Do, and E. Rahm. Quickmig: automatic schema matching for data migration projects. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 107–116. ACM, 2007.
6. S. Duan, A. Fokoue, O. Hassanzadeh, A. Kementsietsidis, K. Srinivas, and M. J. Ward. Instance-based matching of large ontologies using locality-sensitive hashing. In *The Semantic Web–ISWC 2012*, pages 49–64. Springer, 2012.
7. J. Euzenat, P. Shvaiko, et al. *Ontology matching*, volume 18. Springer, 2007.
8. J. Euzenat, P. Valtchev, et al. Similarity-based ontology alignment in owl-lite. In *ECAI*, volume 16, page 333, 2004.
9. S. Falconer. Cognitive support for semi-automatic ontology mapping. *PhD Thesis, University of Victoria*, 2009.
10. S. Falconer and N. Noy. Interactive techniques to support ontology matching. *Schema Matching and Mapping*, pages 29–51, 2011.
11. S. Falconer and M. Storey. A cognitive support framework for ontology mapping. *The Semantic Web*, pages 114–127, 2007.
12. D. Faria, C. Pesquita, E. Santos, M. Palmonari, I. F. Cruz, and F. M. Couto. The agreement-makerlight ontology matching system. In *On the Move to Meaningful Internet Systems: OTM 2013 Conferences*, pages 527–541. Springer, 2013.
13. A. Gal. Interpreting similarity measures: Bridging the gap between schema matching and data integration. In *Data Engineering Workshop, 2008. ICDEW 2008. IEEE 24th International Conference on*, pages 278–285. IEEE, 2008.
14. M. Granitzer, V. Sabol, K. W. Onn, D. Lukose, and K. Tochtermann. Ontology alignment - a survey with focus on visually supported semi-automatic techniques. *Future Internet*, 2(3):238–258, 2010.
15. B. C. Grau, Z. Dragisic, K. Eckert, J. Euzenat, A. Ferrara, R. Granada, V. Ivanova, E. Jiménez-Ruiz, A. O. Kempf, P. Lambrix, et al. Results of the ontology alignment evaluation initiative 2013. In *Proc. 8th ISWC workshop on ontology matching (OM)*, pages 61–100, 2013.
16. B. C. Grau, Z. Dragisic, K. Eckert, J. Euzenat, A. Ferrara, R. Granada, V. Ivanova, E. Jiménez-Ruiz, A. O. Kempf, P. Lambrix, et al. Results of the ontology alignment evaluation initiative 2014. In *Proc. 8th ISWC workshop on ontology matching (OM)*, pages 61–100, 2013.
17. A. Isaac, L. Van Der Meij, S. Schlobach, and S. Wang. *An empirical study of instance-based ontology matching*. Springer, 2007.
18. R. Isele and C. Bizer. Learning linkage rules using genetic programming. In *Proceedings of the Sixth International Workshop on Ontology Matching*, pages 13–24, 2011.
19. R. Isele and C. Bizer. Learning expressive linkage rules using genetic programming. *Proceedings of the VLDB Endowment*, 5(11):1638–1649, 2012.
20. R. Isele and C. Bizer. Active learning of expressive linkage rules using genetic programming. *Web Semantics: Science, Services and Agents on the World Wide Web*, 23:2–15, 2013.
21. S. Jean, H. Dehainsala, D. N. Xuan, G. Pierra, L. Bellatreche, and Y. Ait-Ameur. Ontodb: It is time to embed your domain ontology in your database. In *Advances in Databases: Concepts, Systems and Applications*, pages 1119–1122. Springer, 2007.
22. Y. Jiang, G. Li, J. Feng, and W.-S. Li. String similarity joins: An experimental evaluation. *Proceedings of the VLDB Endowment*, 7(8), 2014.
23. E. Jiménez-Ruiz and B. C. Grau. Logmap: Logic-based and scalable ontology matching. In *The Semantic Web–ISWC 2011*, pages 273–288. Springer, 2011.
24. C. A. Knoblock, P. Szekely, J. L. Ambite, A. Goel, S. Gupta, K. Lerman, M. Muslea, M. Taheriyani, and P. Mallick. Semi-automatically mapping structured sources into the semantic web. In *The Semantic Web: Research and Applications*, pages 375–390. Springer, 2012.

25. L. A. P. P. Leme, M. A. Casanova, K. K. Breitman, and A. L. Furtado. Instance-based owl schema matching. In *Enterprise Information Systems*, pages 14–26. Springer, 2009.
26. J. Li, J. Tang, Y. Li, and Q. Luo. Rimom: A dynamic multistrategy ontology alignment framework. *Knowledge and Data Engineering, IEEE Transactions on*, 21(8):1218–1232, 2009.
27. F. Naumann, A. Bilke, J. Bleiholder, and M. Weis. Data fusion in three steps: Resolving schema, tuple, and value inconsistencies. *IEEE Data Eng. Bull.*, 29(2):21–31, 2006.
28. A.-C. N. Ngomo and S. Auer. Limes-a time-efficient approach for large-scale link discovery on the web of data. *integration*, 15:3, 2011.
29. M. Niepert, J. Noessner, C. Meilicke, and H. Stuckenschmidt. Probabilistic-logical web data integration. In *Reasoning Web. Semantic Technologies for the Web of Data*, pages 504–533. Springer, 2011.
30. J. Noessner, M. Niepert, C. Meilicke, and H. Stuckenschmidt. Leveraging terminological structure for object reconciliation. In *The Semantic Web: Research and Applications*, pages 334–348. Springer, 2010.
31. B. P. Nunes, A. Mera, M. A. Casanova, B. Fetahu, L. A. P. P. Leme, and S. Dietze. Complex matching of rdf datatype properties. In *Database and Expert Systems Applications*, pages 195–208. Springer, 2013.
32. H. Paulheim, S. Hertling, and D. Ritze. Towards evaluating interactive ontology matching tools. In *The Semantic Web: Semantics and Big Data*, pages 31–45. Springer, 2013.
33. H. Stuckenschmidt, J. Noessner, and F. Fallahi. User-centric data integration with the mappingassistant. In *Enterprise Information Systems*, pages 323–339. Springer, 2013.
34. F. M. Suchanek, S. Abiteboul, and P. Senellart. Paris: Probabilistic alignment of relations, instances, and schema. *Proceedings of the VLDB Endowment*, 5(3):157–168, 2011.
35. A. Thor, T. Kirsten, and E. Rahm. Instance-based matching of hierarchical ontologies. In *BTW*, volume 103, pages 436–448, 2007.
36. J. Volz, C. Bizer, M. Gaedke, and G. Kobilarov. Silk-a link discovery framework for the web of data. *LDOW*, 538, 2009.
37. B. Walshe, R. Brennan, and D. O’Sullivan. A comparison of complex correspondence detection techniques. In *OM*, 2012.
38. J. Wang, G. Li, and J. Fe. Fast-join: An efficient method for fuzzy token matching based string similarity join. In *Data Engineering (ICDE), 2011 IEEE 27th International Conference on*, pages 458–469. IEEE, 2011.
39. Y. Wang, W. Liwei, Y. Li, D. He, W. Chen, and T.-Y. Liu. A theoretical analysis of ndcg ranking measures. In *26th Annual Conference on Learning Theory*, 2013.