

Neighbourhood-based Cross-Language Ontology Matching

Juliana M. Destro ^{a,*}, Gabriel de Oliveira dos Santos ^a, Julio Cesar dos Reis ^a,
Ariadne Maria Brito Rizzoni Carvalho ^a, Ivan Ricarte ^c and Ricardo da S. Torres ^b

^a *Institute of Computing, University of Campinas, SP, Brazil*

E-mails: juliana.destro@ic.unicamp.br, mailgabriel@ic.unicamp.br, jreis@ic.unicamp.br, ariadne@ic.unicamp.br

^b *Department of ICT and Natural Sciences, Norwegian University of Science and Technology, Ålesund, Norway*

E-mail: ricardo.torres@ntnu.no

^c *School of Technology, University of Campinas, Limeira-SP, Brazil*

E-mail: ricarte@unicamp.br

Abstract. Cross-language ontology alignments play a key role for the semantic integration of data described in different languages. The task of automatically identifying ontology mappings in this context requires exploring similarity measures as well as ontology structural information. Such measures compute the degree of relatedness between two given terms from ontology's entities. The structural information in the ontologies may provide valuable insights about the concept alignments. Although the literature has extensively studied these measures for monolingual ontology alignments, the use of similarity measures and structural information for the creation of cross-language ontology mappings still requires further research. In this article, we define a novel technique for automatic cross-language ontology matching based on the combination of a composed similarity approach with the analysis of neighbour concepts to improve the effectiveness of the alignment results. Our composed similarity considers lexical, semantic, and structural aspects based on background knowledge to calculate the degree of similarity between contents of ontology entities in different languages. Experimental results with MultiFarm indicate a good effectiveness of our approach including neighbour concepts for mapping identification.

Keywords: ontology mapping, ontology cross-language alignment

1. Introduction

Ontologies¹ are used on a multitude of applications in computer science in the role of a specification mechanism or definition of a common vocabulary. Mapping establishes correspondences between different ontology entities and are relevant for the integration of heterogeneous data sources. There is a growing number of ontologies described in different natural languages. **In the biomedical domain, for example, LOINC² is available in 19 languages and SNOMED CT**

³ is available in 5 languages. The challenge of generating correspondences between different ontologies, created for diversified purposes, is aggravated when concepts are labeled in different natural languages, even in the same domain. Although automatic monolingual ontology matching has been extensively investigated [1], cross-language ontology matching still demands further investigations aiming to automatically identify correspondences between ontologies described in different languages [2].

In this context, accurate automatic methods are essential for ensuring the quality of the generated mappings. Current ontologies have highly grown in size.

*Corresponding author. E-mail: juliana.destro@ic.unicamp.br.

¹What is an ontology? <http://www-ksl.stanford.edu/kst/what-is-an-ontology.html> (As of March 2022).

²<https://loinc.org/international/> (As of March 2022)

³<https://www.snomed.org/> (As of March 2022)

As differences between the used alphabets hamper the use of simple string comparison techniques, **semantic similarity measures** play a key role to obtain well-defined ontology mappings because they allow calculating the level of lexical and semantic similarity between concepts [3]. Cross-language ontology matching approaches in the literature have not yet thoroughly investigated the influence of similarity calculation neither have they analyzed the influence of neighbour concepts in the matching process.

In this article, we propose an original cross-language ontology alignment technique based on the analysis of neighbour concepts relying on composed similarity measure, by combining both syntactic and semantic similarity techniques. Syntactic similarity computes a score calculated based on string analysis (extracted from labels of entities), whereas the semantic similarity is computed taking into account background knowledge, such as synonyms and the context in which terms appear (*e.g.*, use of external dictionaries and vocabularies). Our investigation explores a *Weighted Overlap* measure [4] relying on the neutral-domain semantic network *BabelNet* [5] and computes a weighted mean of semantic and syntactic similarities. The proposed technique also takes into account the similarity of those concepts immediately related to a given entity (the neighbours), both on source and target ontologies. The method finds the highest value of similarities among these concepts. In this investigation, we name such value as neighbourhood similarity. The neighbourhood similarity is used to improve the correctness of mappings and it is thus combined with the composed similarity whenever the initial value of composed similarity is in a **uncertain** range, that is, between a default and minimum threshold (set as parameters before the processing begins).

We carried out a series of experiments to investigate the quality of mappings generated by our technique. Our experiments explored conference-domain ontologies in 45 language pairs from the *MultiFarm*⁴ dataset [6]. *MultiFarm* provides curated mappings between multilanguage ontologies. This dataset has been used to assess cross-language ontology matching methods since 2011, by the **Ontology Alignment Evaluation Initiative**⁵. The obtained results indicate that syntactic and semantic similarities may have different weights in order to obtain a good accuracy. **Our experiments**

suggest that the language in which the ontologies are described, and the translation tool play an important role in the quality of generated alignments.

The remaining of this paper is organized as follows: Section 2 describes the related work; Section 3 formalizes the fundamental concepts of our proposal; Section 4 reports on our proposed technique; Section 5 describes the experimental results whereas Section 6 discusses our findings; Section 7 provides the conclusion remarks.

2. Background

There has been a number of investigations on specific aspects of cross-language for ontology matching. Meilicke *et al.* [7] studied the effectiveness of a set of matching systems based on a dataset defined to evaluate ontology alignment. Their results indicated the difficulties of traditional ontology matching algorithms for carrying out multilingual ontology alignment. Trojahn *et al.* [8] described an extensive survey of matching systems and strategies for accomplishing multilingual and cross-language ontology matching. More recently, Ivanova [2] provided a classification of the available approaches and strategies used by current cross-language mapping systems.

Several approaches have explored the translation effects and the use of a third language in cross-language ontology alignment. In particular, Fu *et al.* [9] analyzed the impact of automatic translations on multilingual ontology alignment, highlighting the translation's relevance for achieving adequate matching quality. Spohr *et al.* [10] studied the translation of concept labels to a third language for matching two ontologies described in different languages.

Ontology alignment techniques have considered the use of similarity methods, which aim to calculate the degree of relatedness between concepts exploring different sources (*e.g.*, dictionary, thesauri). Annane *et al.* [11] proposed a method to build a customized background knowledge resource to improve recall of generated mappings without sacrificing precision. Stoutenburg [12] argued that the use of ontologies combined with linguistic resources as background knowledge might enhance ontology matching processes. This appears as an alternative to syntactic similarity measures relying only on string comparison to determine the similarity value.

The use of multiple similarity measures for the ontology alignment task has been investigated in the lit-

⁴<https://www.irit.fr/recherches/MELODI/multifarm> (As of April 2019).

⁵<http://oaei.ontologymatching.org/> (As of March 2022)

erature. Nguyen and Conrad [13] proposed an ontology matching method based on the combination of lexical-based, structure-based, and semantic-based techniques. After obtaining the structural correspondences among the concepts, the method explores a semantic similarity based on WordNet dictionary and the results are combined. Their approach was evaluated with monolingual ontology alignments. Further investigations are necessary to understand whether a combination and use of semantic similarity can be relevant for cross-language ontology alignment.

Experimental studies have analyzed the influence of syntactic and semantic similarity methods and the structure of terms denoting concepts in ontologies in the context of cross-language alignment [14]. These studies highlight the potential influence of similarity measures.

Structural information (*i.e.* neighbourhood) and multiple similarity measures were used by Essayeh and Abed [15] to generate a similarity matrix. Lin *et al.* [16] also used structural information combined with other similarity methods but neither tackled the cross-language ontology alignment problem.

We investigated different proposed methods, including methods in participants of OAEI (Ontology Alignment Evaluation Initiative)⁶ MultiFarm track. Overall, a translator is used to overcome the natural language barrier and to enable applying monolingual methods to perform cross-language matching. This approach consists of translating the elements of concepts (such as class name, label, and commentaries of class) to the same language of the other ontology, or to a pivot language. Another common method leverages information retrieval techniques, for instance, PageRank and indexing, to define matchings.

WeSeE [17] proposed method uses a web search engine (Microsoft Bing Search API⁷) for retrieving web documents relevant for concepts in the ontologies. The method uses labels, comments, and URI fragments as search terms. For getting search terms from ontology concepts (*i.e.*, classes and properties). The search results of all concepts are then compared to each other. The similarity score is based on the similarity of search results.

The system AUTOMSv2 implements alignment methods already provided within an Alignment API. In order to solve the multi-language problem, AUTOMSv2

uses a free Java API named *WebTranslator*⁸. The AUTOMSv2 translation method is performed by converting the labels of classes and properties that are found to be in a non-English language (only WebTranslator supported languages) and creates a copy of an English-labeled ontology file for each non-English ontology. This process is performed before AUTOMSv2 profiling, configuration, and matching methods are executed. Therefore, their input will consider only English-labeled copies of ontologies, rendering the problem a monolingual matching.

CLONA [18] is an alignment system aiming to identify correspondences between two ontologies defined in two different natural languages. It consists of a six-step approach: (i) Parsing and Pretreatment, (ii) Translation, (iii) Indexation, (iv) Candidate Mappings Identification, and (vi) Alignment Generation. The second phase uses translation provided by Microsoft Translator⁹ to translate the non-English concepts into the chosen pivot language, English. CLONA uses the Lucene search engine¹⁰ to index the pre-processed and translated ontology to determine matching candidates. The documents at the indexes represent the semantic information collected from an external resource (*i.e.*, WordNet) about the entity. A search query is set up in Lucene to return all the matching candidates.

The proposal of *CroLOM* is based on natural languages processing techniques (such as lemmatization, stopwords elimination and stemming) to normalize labels extracted from ontologies. These entities are translated into English, as a pivot language, and the technique computes a Cartesian product among the concepts that compose the ontologies. They apply semantic and syntactic similarity measures in a hybrid way to identify potential mappings. The syntactic similarity is calculated from the *Levenshtein distance* [19], whereas the semantic similarity considers the category of words in WordNet. At this stage, an initial filter is applied to select candidate correspondences containing the maximum similarity value. Then, a second filter is applied to identify the correspondences that contain similarity value greater than a given threshold.

The *SOCOM++* [20] approach considers several setups with different parameters. In contrast to *CroLOM*, it translates concept labels of the source ontology to the same language of the target ontology, thus no pivot languages are considered. Afterwards, both ontologies

⁶<http://oaei.ontologymatching.org> (As of March 2022).

⁷<https://www.bing.com/partners/developers> (As of March 2022)

⁸ <http://webtranslator.sourceforge.net/> (As of March 2022)

⁹<https://www.microsoft.com/en-us/translator/> (As of March 2022)

¹⁰<https://lucene.apache.org/> (As of March 2022)

are described in the same language and monolingual matching methods are applied. In this process, the context of a given concept is analyzed considering all immediate neighbour concepts to improve the quality of the obtained alignment. This approach was designed to support user's influence on adjustments in the translation of the selected labels, and thus users can analyze the resulting mappings and propose changes.

The *AML (AgreementMakerLight)* is a general purpose ontology matching system based on the design principles of *AgreementMaker* [21]. *AML* relies primarily on lexical matching and structural algorithms for both matching and filtering. It makes use of external biomedical ontologies and the WordNet as sources of background knowledge.

In *YAM++*, concept labels of both ontologies (source and target) are translated into the English language. The concepts are filtered in a stage named candidate filtering. In this stage, heuristic filters are applied to selected candidate correspondences, reducing the search space. In the following stage, the method analyzes the neighbourhood of previously selected concepts to discover as many as possible high accurate mappings. Finally, the selected mappings go through a process of semantic verification [22], in which those correspondences considered inconsistent are removed.

XMAP [23] uses semantic similarity based on WordNet, combined with automatic translation provided by Microsoft Translator e taxonomic hierarchy to establish the similarity between two concepts.

LogMap considers a Lexical indexing, which is an inverted index used to store the lexical information. It exploits ontology modularization techniques to reduce the size of the problem. The relevant modules in the input ontologies together with (a subset of) the candidate mappings are encoded in *LogMap* using a Horn clause propositional representation. This approach extends Dowling-Gallier's algorithm [24] to track all mappings that may be involved in the unsatisfiability of a class and performs a greedy local repair; that is, it repairs unsatisfiabilities on-the-fly and only looks for the first available repair plan. It considers a Semantic Indexing, which allows to answer many entailment queries as an index lookup operation over the input ontologies and the mappings computed. The semantic index complements the use of the propositional encoding to detect and repair unsatisfiable classes in the input ontologies.

The *KEPLER*'s approach relies on a divide-and-conquer strategy. First, it splits up the ontology into small blocks, maximizing the relationship inside the

block, and minimizing the relationship between the blocks themselves. On the following step, it translates the ontologies to English as the pivot language, and uses the indexing strategy to reduce the searching space. It considers Candidate Mappings Identification, which queries documents in a vector space that contains a set of ontological entities and their synonyms obtained via WordNet for each ontology. Finally, the algorithm filters the candidate mappings by using two filters: the first filter eliminates the redundancy between these candidates by eliminating possible duplicates; the second filter eliminates false positive candidates.

SANOM [25] uses simulated annealing (SA), which is a probabilistic technique for approximating the global optimum of a given function, as the principal technique to find the mappings between two given ontologies while no ground truth is available. Although *SANOM* was not intended for the cross lingual task, its approach still has produced some results due to the structural similarity between the ontologies in Multifarm [26].

Exploiting Wikipedia as external knowledge base, *WikiV3* [27] uses the MediaWiki API and searches pages corresponding to a given element in the ontology (class comments, labels, class names, etc). When exploring the interlanguage links of Wikipedia through Wikidata, the system is also able to find mapping between ontologies of different languages.

Wiktionary [28] (or *Wiktionary Matcher*) matches concepts by linking labels to entries in *Wiktionary*¹¹, and then checks whether the concepts are synonymous in the external data set. A correspondence is added to the final alignment only based on the synonym relation.

VeeAlign [29] is a supervised Deep Learning based ontology alignment system. The proposed approach is to compute a contextualized representation of concepts as a function, using concept labels and their relationship with neighbouring concepts, producing a contextual vector. The contextualised concept representation is used to discover alignments without the requirement for background knowledge.

The novelty of our approach resides in the weighted combination of semantic and syntactic similarities, leveraging the concept of neighbourhood to improve the correctness of the generated mappings. Whereas existing methods described in this section use struc-

¹¹<https://en.wiktionary.org/wiki/Wiktionary> (As of March 2022)

Table 1: Comparing different methods and techniques of cross-lingual ontology matching. Neighbourhood-based is our proposed technique.

| System | Syntactic | Semantic | Lexicon | Structural | Other | External Resources |
|-------------------------|-----------|----------|---------|------------|---------------------|---------------------------------------|
| WeSeE | | | X | | | MS Bing Search |
| AUTOMsv2 | X | X | X | X | | WebTranslator |
| CLONA | | | X | | | Microsoft Translator Apache Lucene |
| CroLOM | X | X | | | | WordNet |
| SOCOM | | X | | X | | |
| AML | | | X | X | | |
| YAM++ | | X | X | | | |
| XMAP | X | X | | X | | WordNet |
| LogMap | | | X | X | | |
| Kepler | | | X | X | | WordNet |
| SANOM | | | | X | Simulated Annealing | |
| WikiV3 | | | X | | | Wikipedia |
| Wiktionary | | | X | | | Wiktionary |
| VeeAlign | | | | | Deep learning | |
| Neighbourhood based (*) | X | X | | X | | BabelNet Google Translator |

tural information mostly as another component in their similarity calculations, our approach proposed using the neighbouring concepts as a means to confirm or refute a mapping candidate.

Table 1 presents a comparison between methods and techniques proposed in the literature to tackle the ontology matching problem (some techniques are applied on both the monolingual and cross-lingual problems).

3. Formalization

This section formalizes the fundamental concepts in this investigation.

3.1. Ontologies

Ontologies define a common vocabulary in a domain [30]. They are used for semantic representation in computational systems, describing the definition of concepts and the relationship among them.

Definition 3.1 (Ontology). *An ontology \mathcal{O} describes a domain in terms of concepts, attributes and relationships. Formally, an ontology $\mathcal{O} = (\mathcal{C}_{\mathcal{O}}, \mathcal{R}, \mathcal{A}_{\mathcal{O}})$ consists in a set of classes or concepts $\mathcal{C}_{\mathcal{O}}$ interrelated by a set of directed relations \mathcal{R} . Each concept $c \in \mathcal{C}_{\mathcal{O}}$ has a unique identifier and it is associated with a set of attributes $\mathcal{A}_{\mathcal{O}}(c) = \{a_1, a_2, \dots, a_p\}$. Concepts are ontology entities represented by owl:Class construct in OWL¹². Each relation $r(c_1, c_2) \in \mathcal{R}$ can be described as a tuple $(c_1, c_2, r(c_1, c_2))$, where $r(c_1, c_2)$ is a function returning the type of relationship between (c_1, c_2) (e.g., “ \equiv ”, “ \sqsubseteq ”, etc.). The sym-*

¹²The W3C Web Ontology Language (OWL) is a Semantic Web language <https://www.w3.org/OWL/> (As of April 2019).

bols “ \equiv ” and “ \sqsubseteq ” represent relationships “equivalence” and “is-a”, respectively. Furthermore, the relationships can express domain-related relations. For instance, considering the biomedical domain, the concepts c_1 : “Insulin” and c_2 : “Diabetes” may be related by the following function: $r(c_1, c_2) = \text{“Treats”}$. Relations are entities represented by owl:ObjectProperty or owl:DatatypeProperty constructs in OWL [31].

Definition 3.2 (Neighbour Concepts). *We define neighbour concepts of a given entity $e \in \mathcal{C}_{\mathcal{O}}$ or $e \in \mathcal{R}$ the set of concepts with a direct relation to e . Formally, the neighbourhood of e is the set $nbh = \{cpt | cpt \in \mathcal{C}_{\mathcal{O}} \wedge dist(e, cpt) = 1\}$, where $dist(e, cpt)$ is the distance (in terms of the number of edges) between ‘ e ’ and ‘ cpt ’.*

Figure 1 presents an illustrative example of neighbour concepts. The neighbourhood of “Pancreas” is composed of “Endocrine System”, “Digestive System”, “Insulin” and “Glucagon”, because all of them are directly related to “Pancreas”. Because the distance between “Kidney” and “Pancreas” is equal to two, it is not considered a neighbour concept of “Pancreas”.

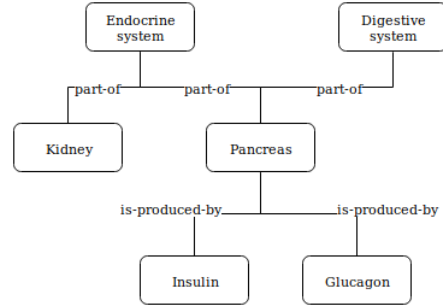


Figure 1. Example of neighbourhood.

3.2. Cross-Language Ontology Alignment

We formally characterize the cross-language ontology problem as follows. Let O_X and O_Y be ontologies described in different natural language “X” and “Y”, respectively; and entities $e_i \in O_X$ and $e_j \in O_Y$. The problem relies on automatically identifying the adequate set of 4-tuples $m_{e_i \rightarrow e_j} = (e_i, e_j, s_{i,j}, r(e_i, e_j))$, where $cf_{i,j}$ is the confidence value between (e_i, e_j) and falls under the interval $[0,1]$ and $r(e_i, e_j) \in \mathcal{R}$ is the relationship between these elements. For instance, considering the concepts $c_1 \in \mathcal{C}_{\text{Opt}}$ and $c_2 \in \mathcal{C}_{\text{En}}$, from ontologies described in Portuguese and English, respectively, such that $c_1 =$ with associated la-

bel “Cabeça” and $c_2 = \text{with}$ with associated label “Head”, the alignment between these concepts is $m_{c_1 \rightarrow c_2} = (c_1, c_2, 1, \equiv)$.

Definition 3.3 (Mappings). *The final result of the alignment process is a set containing the mappings found between the entities (classes, object properties and datatype properties) from two given ontologies. Formally, the mapping between the ontologies \mathcal{O}_X and \mathcal{O}_Y is given by each element of $\mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y}(\lambda) = \{m_{e_i \rightarrow e_j} | e_i \in \mathcal{O}_X \wedge e_j \in \mathcal{O}_Y \wedge s(e_i, e_j) \geq \lambda\}$, where “ λ ” is the threshold (minimum value) of confidence.*

3.3. Similarity Measures

Definition 3.4 (Similarity between entities). *Given two entities e_i and e_j from an ontology (or from different ontologies), the similarity value between them is defined as the maximum similarity value among the attributes (e.g. labels, synonyms, etc) of e_i and e_j . Formally:*

$$\text{syn}(e_i, e_j) = \arg \max \text{sim}(a_{ix}, a_{jy}) \quad (1)$$

where $\text{sim}(a_{ix}, a_{jy})$ is the similarity degree between the pair of attributes a_{ix} and a_{jy} from e_i and e_j , respectively. The similarity may be calculated in different linguistic levels, from string-based methods to semantic techniques [32].

Syntactic Similarity Measure. Levenshtein Distance [19] is an algorithm that computes a syntactic or string-based similarity, which can be understood as the minimum number of single-character editions (insertions, deletions or substitutions) needed to change a string s into s' . This algorithm has been chosen to compute the syntactic similarity in this investigation because Levenshtein Distance has been well-studied and has been extensively used to spelling correction, being considered a good alternative to syntactic analysis [33].

Semantic Similarity Measure. Semantic similarity between concepts is a metric to evaluate how similar two given concepts are, considering their meanings in a certain context. For instance, the words “lead” and “iron” are much more similar considering the metal context than “lead” and “leader”. On the other hand, when we consider the organizational context “lead” and “leader” may be more similar than “lead” and “iron”.

Table 2: Example of NASARI vector representation, where `synset_weight` represents dimensions $1 : n \wedge n \leq 300$.

| Babelnet SynsetId | Wikipedia PageTitle | synset1_weight1 | ... | synsetn_weightn |
|-------------------|---------------------|---------------------|-----|-------------------|
| bn:00000009n | 100 (number) | bn:00058285n_332.33 | ... | bn:00031261n_9.35 |
| bn:00000010n | 1000 (number) | bn:00058285n_347.11 | ... | bn:00024261n_2.11 |

There are algorithms to calculate semantic similarity. Usually, these algorithms explore an external resource such as vocabulary, dictionaries, and thesauri. In this work, we use Weighted Overlap applied to NASARI vectors, together with the neutral-domain semantic network *BabelNet* [5].

This choice relies on the studies of the influence of semantic similarity in neutral-domain context, using Weighted Overlap [14].

NASARI helps us to compute the similarity value in multilingual contexts because it uses vectors based on “*synsets*” (set of synonyms) used by *Babelnet* [34]. The vectors are created in two steps: first, for a given concept, it collects a set of Wikipedia pages where the concept is mentioned. The second step consists in processing the collected contextual information (*i.e.*, **information extracted from the Wikipedia pages**) using a statistical measure (lexical specificity [35]), aiming at finding the most relevant words and synsets appearing in the contextual information and assigning to each one of them a weight (based on the statistical measure). Each of these words and synsets are used as dimensions in the vector-based representation.

Table 2 shows the semantic vector-based representation of two *Babel synsets* (*i.e.*, the identification used in *BabelNet* to represent a given meaning of a word and all the synonyms expressing that meaning in a range of different languages). On each row of the NASARI vector table (exemplified by two rows in Table 2), the first column is the *Babel synsets* ID and the second column is the textual description of the synset (*e.g.*, the synsetID `bn:00000009n` represents the synset “*100 (number)*”). The vector dimensions are described from column three onwards, and are represented by a *Babelnet synset* ID and its correspondent weight (*e.g.*, vector dimension in column `synset1_weight1`, where `bn:00058285n` is the dimension and `332.33` is the weight). Vectors are truncated to the non-zero dimensions only (*i.e.*, all dimensions present weight above zero). Because vectors present *Babelnet synset* as their dimensions, they are comparable across languages.

NASARI leverages *Weighted Overlap (WO)* method applied to the semantic vectors representations [36] to

calculate the semantic similarity between two elements e_1 and e_2 (cf. Equation (2)):

$$sem(e_1, e_2) = WO(v_1, v_2) \quad (2)$$

Weighted Overlap calculates the similarity between the meanings of two given lexical items. Formally:

$$WO(v_1, v_2) = \frac{\sum_{i=1}^{|S|} (r_i^1 + r_i^2)^{-1}}{\sum_{i=1}^{|S|} (2i)^{-1}} \quad (3)$$

In Equation 3, S refers to the set of overlapping dimensions between the two vectors (*i.e.*, dimensions appearing on both vectors; in the example in Table 2, dimension bn:00058285n under column *synset1_weight1*). The r_q^j is the rank of dimension q in the vector v_j . Note that the weight is not used in *WO* equation; it is only used for ranking (*i.e.*, sorting) the dimensions.

Definition 3.5 (Composed Similarity). *We define the composed similarity by combining syntactic and semantic measures. Let $sem(e_1, e_2)$ (Equation (2)) be the semantic similarity and $syn(e_1, e_2)$ the syntactic one between the entities e_1 and e_2 , respectively. Formally:*

$$simC(e_1, e_2) = \frac{\alpha syn(e_1, e_2) + \beta sem(e_1, e_2)}{\alpha + \beta} \quad (4)$$

where α and β are constants.

Note that both semantic and syntactic similarities are a particular case of the composed similarity, when α and β are equal to zero, respectively.

We explore the composed similarity together with Neighbourhood Analysis (cf. Section 4) in our cross-language ontology alignment technique.

4. Cross-Language Ontology Alignment Relying on Neighbourhood Analysis

Our technique for cross-language ontology matching is based on a composed similarity measure relying on both syntactic and semantic similarity techniques, leveraging the similarity of local neighbour concepts (cf. Definition 3.2) to settle **uncertain** mappings.

Figure 2 presents the workflow of the proposed technique. The inputs are a source and target ontologies written in OWL (Web Ontology Language) format.

These ontologies are converted to an object (**object-oriented development**), preserving the relations and neighborhood relationship between concepts. Each entity of the source ontology is compared with all entities of the same type (*i.e.*, concepts are compared only with concepts, relations are compared only with relations) of the target ontology and their composed similarity is calculated. If the similarity surpasses the threshold, the pair is mapped. If not, the calculated similarity is compared with a minimum threshold to verify if the composed similarity is in an **uncertain** range. If the similarity value is above minimum a threshold, the similarity between the neighbour concepts is taken into account in a new validation against the threshold.

Algorithm 1 defines a cross-language alignment between two distinct ontologies \mathcal{O}_X and \mathcal{O}_Y expressed in different natural languages. The algorithm considers the following input arguments:

- Input ontologies $\mathcal{O}_X, \mathcal{O}_Y$
- $\lambda \in (0, 1]$ - default threshold
- $min_\lambda \in [0, \lambda)$ - minimum threshold
- α - Syntactic weight
- β - Semantic weight
- *pivot* - The pivot language

The algorithm starts with mapping set $\mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y} \leftarrow \emptyset$ (line 1) and the similarity variables with zero. It calculates the cartesian product from the set of entities $\mathcal{C}_{\mathcal{O}_X}$ and $\mathcal{C}_{\mathcal{O}_Y}$, and $\mathcal{R}_{\mathcal{O}_X}$ and $\mathcal{R}_{\mathcal{O}_Y}$ from ontologies \mathcal{O}_X and \mathcal{O}_Y , respectively. It considers automatic translation of labels of entities e_1 and e_2 to a pivot language, providing (w_1, w_2) , where $w_1 = translated(e_1)$ and $w_2 = translated(e_2)$ (line 9), leveraging Google Translate API during runtime. The algorithm computes the similarity value based on a syntactic measure (line 10). The syntactic similarity is calculated relying on the strings (w_1, w_2) . The semantic similarity value is also computed. To this end, for each tuple $(e_1, ling_{e_1}, e_2, ling_{e_2})$, composed of the entities e_1 and e_2 , and their respective natural languages $ling_{e_1}$ and $ling_{e_2}$, the algorithm calls the function $babelnet(e_1, ling_{e_1}, e_2, ling_{e_2})$ (line 13). This function uses *Babelnet synsets* and NASARI semantic vectors (cf. Section 3.3) to calculate the *Weighted Overlap* (Equation (3)).

The algorithm calculates the weighted average, assigning weights previously defined by α and β to the syntactic syn_{sim} and semantic sem_{sim} similarities, respectively. It results on the composed similarity $composed_{sim}$ (line 15). If the $composed_{sim}$ value is lower than the default threshold λ and is greater than

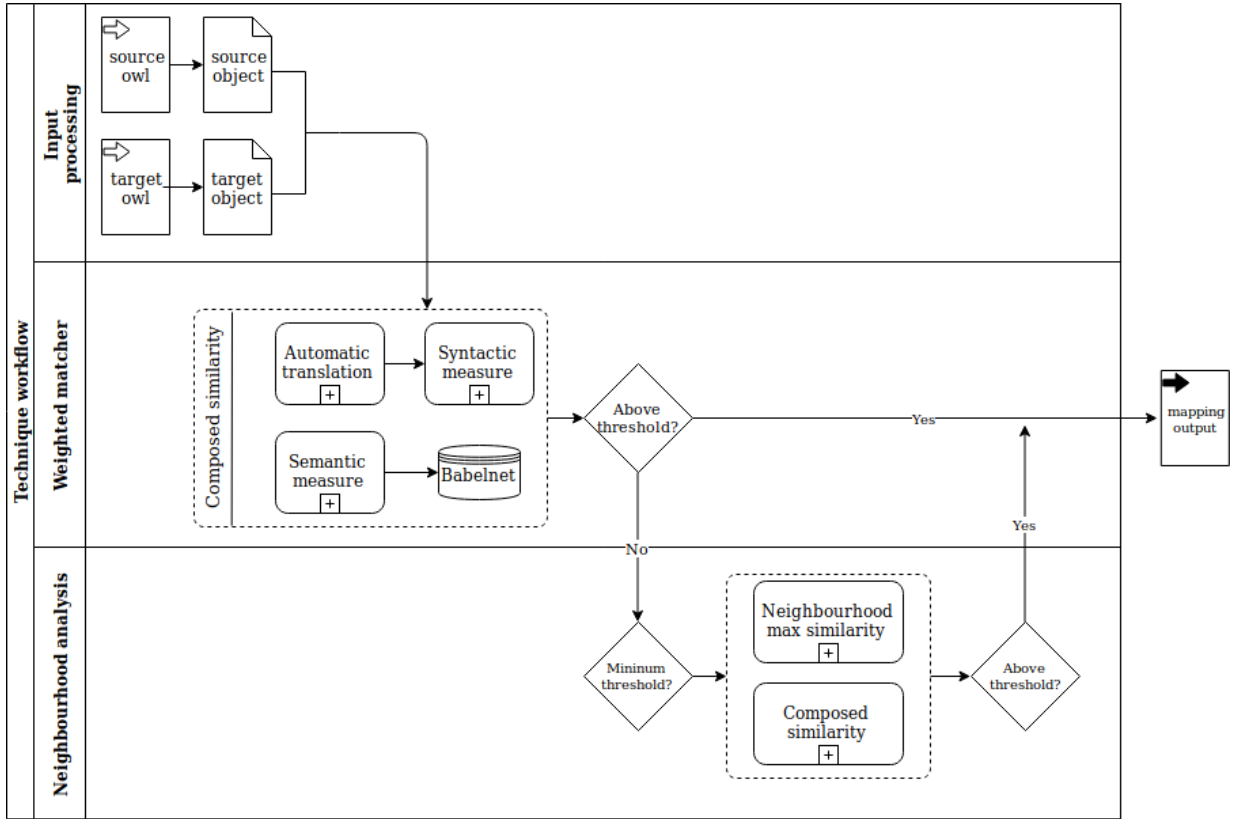


Figure 2. Schematic overview of the proposed cross-language ontology alignment based on neighbourhood analysis.

or equal to min_{λ} , the similarity is considered to be in an **uncertain** range and the algorithm verifies the neighbourhood of the involved concepts to ensure the quality of mappings.

The neighbourhood analysis computes the maximum similarity among the neighbour concepts of the considered entities (source and target). Algorithm 2 computes the similarity among the neighbours of the entities e_1 and e_2 (source and target entities given as input). Only concepts in the neighborhood are retrieved, for entities either in \mathcal{C}_O or \mathcal{R} . First, it extracts the concepts neighboring e_1 and e_2 to nbh_1 and nbh_2 , respectively (line 1 and 2 in Algorithm 2). The algorithm aims to find the pair of neighbour concepts (one from the source ontology and the other one from the target one) with the maximum similarity value based on the composed similarity measure.

Figure 3 presents an example to illustrate the technique of neighbourhood analysis. We consider two ontologies¹³, O_X and O_Y , where O_X is described in Portuguese language and O_Y is described in English language.

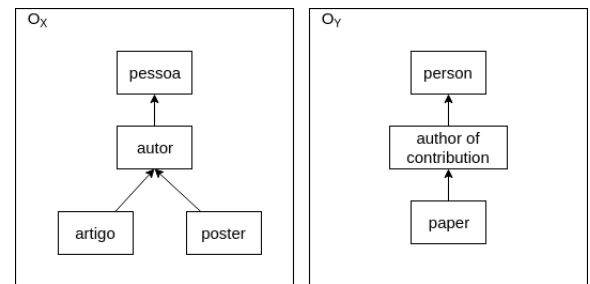


Figure 3. Ontologies O_X and O_Y under analysis.

¹³These ontologies are considered only for the purpose of this example. They were not extracted from real-world ontologies.

Algorithm 1: Cross-language ontology alignment based on composed similarity measure considering neighbourhood analysis

Require: $\mathcal{O}_X, \mathcal{O}_Y, \lambda, \min_\lambda \in [0, 1], \alpha, \beta, pivot$

- 1: $\mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y} \leftarrow \emptyset$ {Initialize the mapping as an empty set}
- 2: $syn_{sim} \leftarrow 0$
- 3: $sem_{sim} \leftarrow 0$
- 4: $composed_{sim} \leftarrow 0$
- 5: $nbh_{sim} \leftarrow 0$
- 6: **for all** $e_1 \in \mathcal{O}_X$ **do**
- 7: **for all** $e_2 \in \mathcal{O}_Y$ **do**
- 8: **if** $\alpha > 0$ **then**
- 9: $w_1 \leftarrow translate(e_1, pivot),$
- 10: $w_2 \leftarrow translate(e_2, pivot)$
- 11: $syn_{sim} \leftarrow syntactic_{sim}(w_1, w_2)$
- 12: **end if**
- 13: **if** $\beta > 0$ **then**
- 14: $sem_{sim} \leftarrow$
- 15: $semantic_{sim}(e_1, ling_{e_1}, e_2, ling_{e_2})$
- 16: **end if**
- 17: $composed_{sim} = \frac{\alpha syn_{sim} + \beta sem_{sim}}{\alpha + \beta}$ {Compute the composed similarity value}
- 18: {Analyze the neighbourhood of concepts if in an **uncertain** range}
- 19: **if** $composed_{sim} < \lambda$ **and**
- 20: $composed_{sim} \geq \min_\lambda$ **then**
- 21: $nbh_{sim} \leftarrow neighbourhood_{sim}(e_1, e_2)$
- 22: {Algorithm 2}
- 23: $similarity \leftarrow composed_{sim}^{(1-nbh_{sim})}$
- 24: **else**
- 25: $similarity \leftarrow composed_{sim}$
- 26: **end if**
- 27: **if** $similarity \geq \lambda$ **then**
- 28: $m_{e_1 \rightarrow e_2} \leftarrow (e_1, e_2, similarity, \equiv)$
- 29: $\mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y} \leftarrow \mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y} \cup \{m_{e_1 \rightarrow e_2}\}$
- 30: **end if**
- 31: **end for**
- 32: **end for**
- 33: **return** $\mathcal{M}_{\mathcal{O}_X \rightarrow \mathcal{O}_Y}$ {Generated mappings}

In the example of Figure 4, the entities “*autor*” and “*author of contribution*” are under analysis. We first calculate the syntactic and semantic similarity value between the two entities, and find a composed similarity value of 0.80.

Algorithm 2: Neighbourhood analysis.

Require: e_1, e_2 {Given the entities e_1 and e_2 from the source and target ontologies respectively}

{Extract the concept neighbourhood of entities e_1 and e_2 to nbh_1 and nbh_2 }

- 1: $nbh_1 \leftarrow neighbourhood(e_1)$
- 2: $nbh_2 \leftarrow neighbourhood(e_2)$
- 3: $maxSim \leftarrow 0$
- 4: **for all** $n_1 \in nbh_1$ **do**
- 5: **for all** $n_2 \in nbh_2$ **do**
- 6: $sim \leftarrow composed_{similarity}(n_1, n_2)$
- 7: {Compute the composed similarity value}
- 8: **if** $sim > maxSim$ **then**
- 9: $maxSim \leftarrow sim$
- 10: **end if**
- 11: **end for**
- 12: **end for**
- 13: **return** $maxSim$

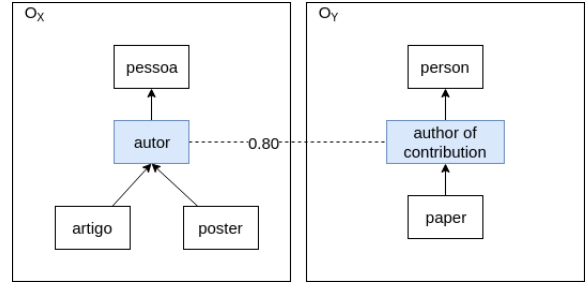


Figure 4. Pair of entities considered for mapping.

In our example, the minimum threshold to consider a mapping **uncertain** is 0.33 and default threshold is 0.95. Thus, by our Algorithm 1, it is necessary go through the neighbourhood analysis. The neighbours of “*autor*” are {“*pessoa*”, “*artigo*”, “*poster*”}, and the neighbours of “*author of contribution*” are {“*person*”, “*paper*”}. We apply a cartesian product to these sets to evaluate the composed similarity (*i.e.*, syntactic and semantic similarity combined) between the set of neighbours and retain maximum similarity value found. In this illustration, the maximum similarity is found for “*pessoa*” and “*person*”, with a 1.0 measure as depicted in Figure 5.

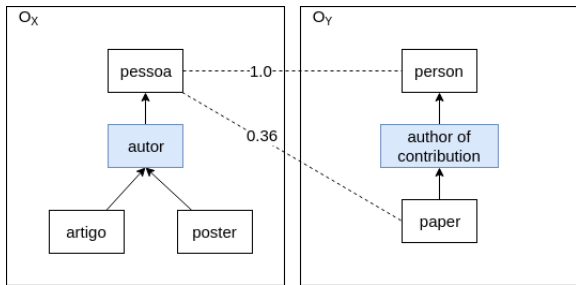


Figure 5. “*pessoa*” and “*person*” are the pair of neighbouring concepts with maximum similarity.

The neighbourhood similarity value returned by Algorithm 2 updates the similarity value considering $composed_{sim}^{(1-nbh_{sim})}$ in Algorithm 1 (line 19). Therefore, after the neighbourhood analysis process, the final similarity value is equal to $0.80^{(1-1.0)} = 0.80^{0.0} = 1.0$, thus these concepts under analysis are influenced by the neighbourhood analysis. Because the final similarity is greater than the default threshold then a mapping is created between “*autor*” and “*author of contribution*”. Note, when the neighbourhood similarity is high, close to 1, the resulting similarity also approaches 1, therefore it is likely to surpass the default threshold, and then be considered a candidate mapping.

Our method assumes as correct mappings when the neighbour concepts are quite similar even if the pair of concepts under analysis itself is not so similar. Finally, Algorithm 1 verifies whether the similarity value computed is greater than or equals to a beforehand input threshold λ (line 21 in Algorithm 1). If such condition is satisfied, the algorithm inserts the mapping $(e_1, e_2, 1, \equiv)$ into the set $\mathcal{M}_{\mathcal{O}_x \rightarrow \mathcal{O}_y}$ indicating a cross-language correspondence between the entities. The output mapping set file follows the general alignment format as the same used by the Alignment API¹⁴. The implementation of the defined algorithms can be obtained in our [institutional](#) project code repository¹⁵. [Access is granted per request.](#)

5. Experimental Evaluation

This evaluation aims to analyze the quality of mappings generated by our proposed technique which con-

Table 3: MultiFarm ontologies and statistics

| Languages | Ontologies | Classes | Object Properties | Datatype Properties | Total Entities |
|------------|---------------|---------|-------------------|---------------------|----------------|
| Arabic | conference-ar | 61 | 46 | 18 | 125 |
| Chinese | conference-cn | 61 | 46 | 18 | 125 |
| Czech | conference-cz | 61 | 46 | 18 | 125 |
| German | conference-de | 61 | 46 | 18 | 125 |
| English | conference-en | 61 | 46 | 18 | 125 |
| Spanish | conference-es | 61 | 46 | 18 | 125 |
| French | conference-fr | 61 | 46 | 18 | 125 |
| Dutch | conference-nl | 61 | 46 | 18 | 125 |
| Portuguese | conference-pt | 61 | 46 | 18 | 125 |
| Russian | conference-ru | 61 | 46 | 18 | 125 |

siders the structure of ontologies in the alignment of ontologies described in different natural languages. We conducted a series of 1260 experiments relying on a set of curated mappings manually established between ontologies described in distinct languages.

5.1. Datasets and Procedure

MultiFarm[6], version released in 2015, is the considered dataset in our experiments. This dataset is used in the *OAEI* and it is composed of a set of 5 ontologies of the Conference domain¹⁶, translated into 10 languages: Arabic (ar), English (en), Chinese (cn), Czech (cz), Dutch (nl), French (fr), German (de), Portuguese (pt), Russian (ru), Spanish (es), and the corresponding cross-language mappings between them. This dataset is based on the *OntoFarm* dataset, which has been successfully used for several years in the *OAEI* Conference track. Our experiments uses only Conference ontologies described in Table 3

This dataset was manually curated and may be used as a reference to assess algorithms that build automatic cross-lingual ontology mappings. For instance, the pair pt-es refers to the ontology mappings between Portuguese and Spanish conference ontologies. We consider 45 set of mappings with different pairs of language as follows: ar-cn, ar-cz, ar-de, ar-en, ar-es, ar-fr, ar-nl, ar-pt, ar-ru, cn-cz, cn-de, cn-en, cn-es, cn-fr, cn-nl, cn-pt, cn-ru, cz-de, cz-en, cz-es, cz-fr, cz-nl, cz-pt, cz-ru, de-en, de-es, de-fr, de-nl, de-pt, de-ru, en-es, en-fr, en-nl, en-pt, en-ru, es-fr, es-nl, es-pt, es-ru, fr-nl, fr-pt, fr-ru, nl-pt, nl-ru and pt-ru.

Our experiments built cross-language ontology mappings by using English as a pivot language for syntactic similarity measurement. [Selecting an appropriate pivot language is an important step in any task re-](#)

¹⁴<http://alignapi.gforge.inria.fr/index.html> (As of March 2022).

¹⁵<https://gitlab.ic.unicamp.br/jreis/evocros>

¹⁶Cmt, Conference, ConfOf, Iasted, Sigkdd

quiring this type of translation and there are resources available to help identify the best language-pair to be used for machine translation on both European languages [37] and Asian languages [38]. Our choice is due to the richness of available language resources between English and the other languages in the experiment. Babelnet is used for semantic similarity measurement and does not need a translation as it can retrieve the synsets used in NASARI vectors, by using the concepts original language. The results obtained by executing Algorithm 1 in different scenarios were compared with the reference mappings from the *MultiFarm dataset*, then metrics of precision, recall, and f-measure [39] were calculated.

We executed Algorithm 1 setting different weights and thresholds for similarity, but considering the minimum threshold equals to 0.33, based on fraction $\{\frac{1}{3}\}$ of the similarity spectrum analyzed $[0, 1]$. We used the mappings Conference-Conference of all 45 pair of languages in Multifarm as reference.

The weights of syntactic and semantic similarities in the composed similarity measure followed the fractions $\{\frac{1}{5}, \frac{1}{4}, \frac{1}{3}, \frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \frac{4}{5}\}$, considering the constraint $\alpha + \beta = 1$. We present results varying the threshold level to comprehend its role in the studied scenarios. We vary the threshold in $\{0.66, 0.75, 0.80, 0.95\}$, which were selected based on the fractions $\{\frac{2}{3}, \frac{3}{4}, \frac{4}{5}, \frac{19}{20}\}$. Table 4 shows the experiments configuration applied for each pair of language.

5.2. Experimental Results

Table 5 presents the achieved results with the highest f-measure obtained for each pair of language. Results vary by language pair, but they have an average threshold of 0.75, average syntactic weight of 0.75 and average semantic weight of 0.25.

The highest f-measure was found between West Germanic languages (English and German), 0.64591. The majority of the results with the lowest f-measure values involve the Arabic language, present in seven pairs of the lowest ten f-measures: ar-ru, de-ar, ar-fr, cn-ar, cz-ar, ar-pt, ar-es, cz-ru, de-ru, pt-ru. **This is mostly due to translation differences between the languages and the pivot. For example, for the pair pt-ru, one of the concept in Portuguese is “*contribuição de artigo completo*” and the translation in English using the Google Translate API is “*full article contribution*”. This concept is mapped in the gold standard to the Russian concept “”, which is translated to English as “*full labor*”, with a distance of 0.34, well below our**

Table 4: Experiments configurations. Different weights for syntactic and semantic similarity are applied to each threshold. Each configuration was applied to Conference-Conference mappings for each of the 45 pairs of languages obtaining a total of 1260 experiments.

| Similarity threshold | Syntactic measure | Semantic measure |
|----------------------|-------------------|------------------|
| 0.66 | 0.20 | 0.80 |
| | 0.25 | 0.75 |
| | 0.33 | 0.67 |
| | 0.50 | 0.50 |
| | 0.67 | 0.33 |
| | 0.75 | 0.25 |
| | 0.80 | 0.20 |
| 0.75 | 0.20 | 0.80 |
| | 0.25 | 0.75 |
| | 0.33 | 0.67 |
| | 0.50 | 0.50 |
| | 0.67 | 0.33 |
| | 0.75 | 0.25 |
| | 0.80 | 0.20 |
| 0.80 | 0.20 | 0.80 |
| | 0.25 | 0.75 |
| | 0.33 | 0.67 |
| | 0.50 | 0.50 |
| | 0.67 | 0.33 |
| | 0.75 | 0.25 |
| | 0.80 | 0.20 |
| 0.95 | 0.20 | 0.80 |
| | 0.25 | 0.75 |
| | 0.33 | 0.67 |
| | 0.50 | 0.50 |
| | 0.67 | 0.33 |
| | 0.75 | 0.25 |
| | 0.80 | 0.20 |

threshold. If we were going to translate the concept in Portuguese directly to Russian language, the result would be “”, with a distance of 0.40, an improvement from the previous 0.34 but still below the threshold.

An example of how our defined technique helps improving the results can be observed in conference-conference-en-pt set of mappings. The concept “*invited speaker*” in conference-en should be mapped to concept “*Palestrante convidado*” in conference-pt, according to the gold standard. Using similarity threshold equals to 0.66, $\alpha = 0.67$ and $\beta = 0.33$, the composed similarity value calculated between them was 0.4307. This similarity value is inside the **uncertain** range and thus the neighborhood similarity was verified to confirm the candidate mapping. The calculated neighborhood composed similarity value was 0.5847. The similarity value was then recalculated using the formula $composed_{sim}^{(1-nbh_{sim})}$, thus $0.4307^{(1-0.5847)} = 0.7048$, surpassing the 0.66 threshold. Therefore, the

mapping is included in the generated ontology set of mappings.

Another example is the concept “*chair of workshop track*” in conference-en and “*coordenador de trilha de workshop*” in conference-pt. Using a similarity threshold equals to 0.66, $\alpha = 0.67$ and $\beta = 0.33$, the composed similarity value calculated between them was 0.3828 and the neighborhood similarity was 1.0. The recalculated similarity value was 1.0 and then the mapping was included in the generated ontology alignment.

Table 5: Results with highest f-measure for each pair of language.

| Language pair | Similarity threshold | Syntactic measure | Semantic measure | Precision | Recall | F-Measure |
|---------------|----------------------|-------------------|------------------|-----------|---------|-----------|
| ar-es | 0.75 | 0.80 | 0.20 | 0.41071 | 0.35659 | 0.38174 |
| ar-fr | 0.75 | 0.80 | 0.20 | 0.44444 | 0.27273 | 0.33803 |
| ar-nl | 0.66 | 0.80 | 0.20 | 0.37989 | 0.47552 | 0.42236 |
| ar-pt | 0.66 | 0.75 | 0.25 | 0.35220 | 0.40288 | 0.37584 |
| ar-ru | 0.66 | 0.80 | 0.20 | 0.31356 | 0.28030 | 0.29600 |
| cn-ar | 0.66 | 0.80 | 0.20 | 0.52381 | 0.25191 | 0.34021 |
| cn-cz | 0.75 | 0.80 | 0.20 | 0.63492 | 0.30075 | 0.40816 |
| cn-de | 0.66 | 0.80 | 0.20 | 0.50476 | 0.37589 | 0.43089 |
| cn-en | 0.66 | 0.80 | 0.20 | 0.54717 | 0.40845 | 0.46774 |
| cn-es | 0.66 | 0.75 | 0.25 | 0.45082 | 0.39568 | 0.42146 |
| cn-fr | 0.75 | 0.80 | 0.20 | 0.58140 | 0.37037 | 0.45249 |
| cn-nl | 0.66 | 0.75 | 0.25 | 0.46667 | 0.48611 | 0.47619 |
| cn-pt | 0.66 | 0.80 | 0.20 | 0.41111 | 0.49664 | 0.44985 |
| cn-ru | 0.66 | 0.80 | 0.20 | 0.50980 | 0.37681 | 0.43333 |
| cz-ar | 0.66 | 0.80 | 0.20 | 0.30189 | 0.45070 | 0.36158 |
| cz-de | 0.75 | 0.80 | 0.20 | 0.46809 | 0.45205 | 0.45993 |
| cz-en | 0.80 | 0.67 | 0.33 | 0.68182 | 0.41667 | 0.51724 |
| cz-es | 0.80 | 0.75 | 0.25 | 0.64894 | 0.42657 | 0.51477 |
| cz-fr | 0.75 | 0.80 | 0.20 | 0.42636 | 0.39007 | 0.40741 |
| cz-nl | 0.95 | 0.80 | 0.20 | 0.70000 | 0.42282 | 0.52720 |
| cz-pt | 0.75 | 0.67 | 0.33 | 0.58654 | 0.42069 | 0.48996 |
| cz-ru | 0.66 | 0.75 | 0.25 | 0.41129 | 0.36429 | 0.38636 |
| de-ar | 0.66 | 0.75 | 0.25 | 0.28421 | 0.40602 | 0.33437 |
| de-en | 0.80 | 0.67 | 0.33 | 0.77570 | 0.55333 | 0.64591 |
| de-es | 0.80 | 0.67 | 0.33 | 0.82258 | 0.36429 | 0.50495 |
| de-fr | 0.66 | 0.80 | 0.20 | 0.36111 | 0.59091 | 0.44828 |
| de-nl | 0.80 | 0.75 | 0.25 | 0.59406 | 0.41667 | 0.48980 |
| de-pt | 0.75 | 0.80 | 0.20 | 0.45963 | 0.49664 | 0.47742 |
| de-ru | 0.66 | 0.80 | 0.20 | 0.40146 | 0.38732 | 0.39427 |
| en-ar | 0.75 | 0.80 | 0.20 | 0.42029 | 0.45313 | 0.43609 |
| en-es | 0.80 | 0.67 | 0.33 | 0.66667 | 0.40845 | 0.50655 |
| en-fr | 0.75 | 0.67 | 0.33 | 0.64516 | 0.30075 | 0.41026 |
| en-nl | 0.80 | 0.20 | 0.80 | 0.86111 | 0.40260 | 0.54867 |
| en-pt | 0.80 | 0.67 | 0.33 | 0.68293 | 0.38889 | 0.49558 |
| en-ru | 0.66 | 0.80 | 0.20 | 0.43972 | 0.42759 | 0.43357 |
| es-fr | 0.75 | 0.67 | 0.33 | 0.65672 | 0.33333 | 0.44221 |
| es-nl | 0.80 | 0.67 | 0.33 | 0.78351 | 0.50000 | 0.61044 |
| es-pt | 0.80 | 0.50 | 0.50 | 0.74545 | 0.51572 | 0.60967 |
| es-ru | 0.66 | 0.75 | 0.25 | 0.45600 | 0.39583 | 0.42379 |
| fr-nl | 0.75 | 0.67 | 0.33 | 0.62500 | 0.40441 | 0.49107 |
| fr-pt | 0.75 | 0.67 | 0.33 | 0.63636 | 0.36567 | 0.46445 |
| fr-ru | 0.66 | 0.80 | 0.20 | 0.36508 | 0.46309 | 0.40828 |
| nl-pt | 0.75 | 0.50 | 0.50 | 0.80682 | 0.46104 | 0.58678 |
| nl-ru | 0.75 | 0.80 | 0.20 | 0.56075 | 0.42254 | 0.48193 |
| pt-ru | 0.66 | 0.80 | 0.20 | 0.34597 | 0.48993 | 0.40556 |

6. Discussion

Cross-language ontology matching relies on several different approaches to obtain mappings that interrelate ontologies described in distinct languages. Cross-language ontology matching requires adequate techniques relying on similarity measures to overcome the matching task barrier. Ontological structure and similarity measures might help in matching algorithms to determinate the adequate mappings. Existing techniques can favor from the understanding of the benefits and limitations of syntactic and semantic similarity approaches to develop a better combination of them.

In this context, this investigation contributed with several experiments to determinate the impact of ontological structure and the similarity measures to be considered in the alignment of ontologies described in different languages. Our experiments were designed to help us understanding their effects in the matching process and the quality of the generated cross-lingual ontology mappings.

Our proposal concerns the influence of the ontological structure and similarity computation on cross-language ontology matching. Our goal was to understand how to combine them aiming to build accurate cross-lingual ontology mappings. To this end, we took into account the weighted average between syntactic and semantic similarities. Our approach considered neighbour concepts directed related to concepts under analysis in candidate mappings.

The choice of weights assigned to each similarity measure played an important role in the results. As we showed empirically, semantic and syntactic similarities might not have the same relevance, *i.e.*, the same weight. Considering the syntactic weight close to 0.70 generated the best mapping results, *i.e.*, it resulted in ontology mappings with the overall highest f-measure value. Thus, our technique may be understood as a good alternative to syntactic or semantic only methods. It might perform even better taking into account the correct parameters.

We found that the gain of effectiveness may vary according to the language describing the content of the ontologies. Comparing the results in Table 5, we observe that the results for the arabic language are generally in the lowest tier. A possible explanation for this behaviour might be the use of Google Translate for automatic translation to the pivot language. Although Google Translate has largely improved over the years[40], there are still some incorrect or mistypes in the translations that hinder syntactic measures. An

example is the word **أيستقبل**, incorrectly translated to “reception”. The correct translation is “reception”. Another example is **أليوم مَهَّت**, incorrectly translated to “today’s station” by Google Translate, when the correct translation would be “terminal date”.

The characteristics of the entity labels of the Conference ontology restrict the use of the semantic similarity measurement Weighted Overlap, because the entity labels are mostly complex sentences instead of words. Babelnet, the external source used in semantic measurement, is a dictionary, not a translation tool and therefore only able to identify synsets in words. This also explains the average semantic similarity being approximately $\frac{1}{4}$ of the overall weight. Thus, it might be useful considering algorithms such as stop-words elimination and stemming, etc. to break the complex sentences into simple structures.

The results showed an influence of threshold; as the threshold rises, the precision also increases. It may be explained by considering equivalence of only those concepts with a high level of similarity. However, the f-measure value reduces as the threshold increases. This happens because higher values assigned to threshold leads to the algorithm disregarding entities that are equivalent, but somehow were assigned a lower level of similarity than expected by the threshold. For instance, in en-es ontology mappings, the similarity between “strange” and “estranho” was equal to 0.89, but the given threshold is 0.95, thus “estranho” is not mapped to “strange”. As a result, the recall drops substantially, because many correct correspondences are ignored, and thus f-measure decreases. Empirically, we concluded that the thresholds generating the more accurate mappings were around $\lambda = 0.75$ and $\lambda = 0.80$.

Table 6 describes the results obtained by related work (ontology alignment systems) presented in OAEI 2018, 2019 and 2020. **The 2021 edition did not include results for same ontology alignments, focus of this study, only for different ontologies. This results use the blind dataset of same ontology (edas-edas) translated into two different languages. Although it is not possible to compare the results of our experiments directly to the OAEI results, it is still possible to see the opportunity for improvement, even on the task considered the easiest (same ontologies just translated into two different languages) and give us a sense of how our method could be a contribution to the field.**

Our obtained findings support the hypothesis that composing different types of similarity measures and

Table 6: Results obtained with existing ontology alignment systems in OAEI (Multifarm Track) in 2018, 2019, 2020 and 2021, considering the alignments between the same ontology edas-edas in different languages.

| Tool | Precision | Recall | F-measure |
|------------|-----------|--------|-----------|
| 2018 | | | |
| KEPLER | 0.85 | 0.36 | 0.49 |
| LogMap | 0.95 | 0.28 | 0.41 |
| AML | 0.96 | 0.16 | 0.27 |
| XMAP | 0.13 | 0.19 | 0.14 |
| 2019 | | | |
| AML | 0.93 | 0.17 | 0.27 |
| LogMap | 0.95 | 0.28 | 0.41 |
| Wiktionary | 0.94 | 0.07 | 0.12 |
| 2020 | | | |
| AML | 0.94 | 0.17 | 0.28 |
| LogMap | 0.95 | 0.28 | 0.41 |
| LogMapLt | 0.02 | 0.01 | 0.01 |
| VeeAlign | 0.91 | 0.08 | 0.14 |
| Wiktionary | 0.94 | 0.07 | 0.12 |

taking into account the ontology structure, by considering the similarity of neighbour concepts, **is a method that** can reveal satisfactory generated ontology mappings for cross-language ontology alignment.

7. Conclusion

Alignment of large ontologies described in different natural languages remains an open research challenge. In this investigation, we proposed an approach based on the weighted mean of syntactic and semantic similarities for this task. Our approach considered the influence of neighbour concepts on the cross-lingual alignment method, combining it with the composed similarity. The defined algorithms were implemented and we carried out a series of experiments to evaluate the effectiveness of this approach. Our findings based on experiments with standard datasets revealed the effectiveness of combining similarity measures and considering the neighbourhood of concepts in the cross-language ontology alignment problem. Future work involves to improve our cross-lingual alignment proposal by considering different combinations of background knowledge, such as specific-domain thesauri to evaluate the semantic similarity. In addition, we plan to investigate different ways of computing the syntactic

and semantic similarities considering additional stages in the pre-processing of entity labels.

Acknowledgements

This work has the financial support of CNPq (grant #307560/2016-3), São Paulo Research Foundations (FAPESP) (grants #2017/02325-5, #2014/12236-1, #2015/24494-8, #2016/50250-1, and #2017/20945-0) and the FAPESP-Microsoft Virtual Institute (grants #2013/50155-0 and #2014/50715-9). This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Finance Code 001¹⁷.

References

- [1] P. Shvaiko and J. Euzenat, in: *Ontology matching: State of the art and future challenges*, Vol. 25, 2013, pp. 158–176.
- [2] T. Ivanova, Cross-lingual and Multilingual Ontology Mapping - Survey, in: *Proceedings of the 19th International Conference on Computer Systems and Technologies*, CompSysTech'18, ACM, New York, NY, USA, 2018, pp. 50–57. ISBN 978-1-4503-6425-6. doi:10.1145/3274005.3274034. <http://doi.acm.org/10.1145/3274005.3274034>.
- [3] C. Pesquita, D. Faria, A.O. Falcão, P.W. Lord and F.M. Couto, Semantic Similarity in Biomedical Ontologies, *PLoS Computational Biology* 5(7) (2009). doi:10.1371/journal.pcbi.1000443. <https://doi.org/10.1371/journal.pcbi.1000443>.
- [4] M.T. Pilehvar, D. Jurgens and R. Navigli, Align, disambiguate and walk: A unified approach for measuring semantic similarity, in: *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Vol. 1, 2013, pp. 1341–1351.
- [5] R. Navigli and S.P. Ponzetto, BabelNet: The automatic construction, evaluation and application of a wide-coverage multilingual semantic network, *Artificial Intelligence* 193 (2012), 217–250.
- [6] C. Meilicke, R. Garcia-Castro, F. Freitas, W.R. van Hage, E. Montiel-Ponsoda, R.R. de Azevedo, H. Stuckenschmidt, O. Sváb-Zamazal, V. Svátek, A. Taminlin, C.T. dos Santos and S. Wang, MultiFarm: A Benchmark for Multilingual Ontology Matching, *Web Semantics: Science, Services and Agents on the World Wide Web* 15 (2012), 62–68.
- [7] C. Meilicke, C. Trojahn, O. Sváb-Zamazal and D. Ritze, Multilingual ontology matching evaluation—a first report on using multifarm, in: *Extended Semantic Web Conference*, Springer, 2012, pp. 132–147.
- [8] C. Trojahn, B. Fu, O. Zamazal and D. Ritze, State-of-the-art in multilingual and cross-lingual ontology matching, in: *Towards the Multilingual Semantic Web*, Springer, 2014, pp. 119–135.

¹⁷The opinions expressed in here are not necessarily shared by the financial support agencies.

- [9] B. Fu, R. Brennan and D. O’Sullivan, Cross-lingual Ontology Mapping - An Investigation of the Impact of Machine Translation, in: *Proceedings of the 4th Annual Asian Semantic Web Conference (ASWC 2009)*, Springer, 2009.
- [10] D. Spohr, L. Hollink and P. Cimiano, A machine learning approach to multilingual and cross-lingual ontology matching, in: *Proceedings of the 10th International Semantic Web Conference (ISWC 2011)*, Springer, 2011, pp. 665–680.
- [11] A. Annane, Z. Bellahsene, F. Azouaou and C. Jonquet, Building an effective and efficient background knowledge resource to enhance ontology matching, *Journal of Web Semantics* **51** (2018), 51–68.
- [12] S.K. Stoutenburg, Acquiring advanced properties in ontology mapping, in: *Proceedings of the 2nd PhD workshop on Information and knowledge management*, 2008, pp. 9–16.
- [13] T.T.A. Nguyen and S. Conrad, Ontology Matching using multiple similarity measures, in: *7th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management (IC3K)*, Vol. 01, 2015, pp. 603–611.
- [14] J.M. Destro, J.C. dos Reis, A.M. Brito, R. Carvalho and I.L.M. Ricarte, Influence of Semantic Similarity Measures on Ontology Cross-language Mappings, *Proceedings of the Symposium on Applied Computing* (2017), 323–329. ISBN 978-1-4503-4486-9. doi:10.1145/3019612.3019836. <http://doi.acm.org/10.1145/3019612.3019836>.
- [15] A. Essayeh and M. Abed, Towards ontology matching based system through terminological, structural and semantic level, *Procedia computer science* **60** (2015), 403–412.
- [16] H. Lin, Y. Wang, Y. Jia, J. Xiong, P. Zhang and X. Cheng, An ensemble matchers based rank aggregation method for taxonomy matching, in: *Asia-Pacific Web Conference*, Springer, 2015, pp. 190–202.
- [17] H. Paulheim, WeSeE-Match results for OEAI 2012., in: *Proc. 7th ISWC workshop on ontology matching (OM)*, 2012.
- [18] M. El Abdi, H. Souid, M. Kachroudi and S.B. Yahia, CLONA results for OEAI 2015., in: *OM*, 2015, pp. 124–129.
- [19] V.I. Levenshtein, Binary codes capable of correcting deletions, insertions and reversals., *Soviet Physics Doklady* **10**(8) (1966), 707–710, *Doklady Akademii Nauk SSSR*, V163 No4 845-848 1965.
- [20] B. Fu, R. Brennan and D. O’Sullivan, A configurable translation-based cross-lingual ontology mapping system to adjust mapping outcomes, *Web Semantics: Science, Services and Agents on the World Wide Web* **15** (2012), 15–36.
- [21] I.F. Cruz, F.P. Antonelli and C. Stroe, AgreementMaker: efficient matching for large real-world schemas and ontologies, *Proceedings of the VLDB Endowment* **2**(2) (2009), 1586–1589.
- [22] D. Ngo and Z. Bellahsene, Efficient Semantic Verification of Ontology Alignment, *2015 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT)* **1** (2015), 141–148.
- [23] W.E. Djeddi and M.T. Khadir, XMAP: a novel structural approach for alignment of OWL-full ontologies, in: *2010 International Conference on Machine and Web Intelligence*, IEEE, 2010, pp. 368–373.
- [24] W.F. Dowling and J.H. Gallier, Linear-time algorithms for testing the satisfiability of propositional Horn formulae, *The Journal of Logic Programming* **1**(3) (1984), 267–284.
- [25] M. Mohammadi, W. Hofman and Y.-H. Tan, Simulated annealing-based ontology matching, *ACM Transactions on Management Information Systems (TMIS)* **10**(1) (2019), 1–24.
- [26] M. Mohammadi, A.A. Atashin, W. Hofman and Y.-H. Tan, SANOM results for OEAI 2017., in: *OM@ ISWC*, 2017, pp. 185–189.
- [27] S. Hertling, WikiV3 results for OEAI 2017., in: *OM@ ISWC*, 2017, pp. 190–195.
- [28] J. Portisch, M. Hladik and H. Paulheim, Wiktionary matcher, in: *CEUR Workshop Proceedings*, Vol. 2536, RWTH, 2020, pp. 181–188.
- [29] V. Iyer, A. Agarwal and H. Kumar, VeeAlign: a supervised deep learning approach to ontology alignment., in: *OM@ ISWC*, 2020, pp. 216–224.
- [30] N.F. Noy and D.L. McGuinness, *Ontology Development 101: A Guide to Creating Your First Ontology*, (accessed: 04.Feb.2021).
- [31] J. Euzenat and P. Shvaiko, in: *Ontology Matching, 2nd Edition*, Springer, Heidelberg, 2013, pp. 85–90.
- [32] D. Dinh, J.C. Dos Reis, C. Pruski, M. Da Silveira and C. Reynaud-Delaître, Identifying relevant concept attributes to support mapping maintenance under ontology evolution, *Web semantics: Science, services and agents on the world wide web* **29** (2014), 53–66.
- [33] G. Navarro, A Guided Tour to Approximate String Matching, *ACM Comput. Surv.* **33**(1) (2001), 31–88, ISSN 0360-0300. doi:10.1145/375360.375365. <http://doi.acm.org/10.1145/375360.375365>.
- [34] J. Camacho-Collados, M.T. Pilehvar and R. Navigli, Nasari: Integrating explicit knowledge and corpus statistics for a multilingual representation of concepts and entities, *Artificial Intelligence* **240** (2016), 36–64.
- [35] P. Lafon, Sur la variabilité de la fréquence des formes dans un corpus, *Mots. Les langages du politique* **1**(1) (1980), 127–165.
- [36] J. Camacho-Collados, M.T. Pilehvar and R. Navigli, NASARI: a Novel Approach to a Semantically-Aware Representation of Items, in: *Proceedings of the North American Chapter of the Association of Computational Linguistics (NAACL 2015)*, Denver, USA, 2015, pp. 567–577.
- [37] A. Birch, M. Osborne and P. Koehn, Predicting success in machine translation, in: *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, 2008, pp. 745–754.
- [38] M. Paul, A. Finch and E. Sumita, How to choose the best pivot language for automatic translation of low-resource languages, *ACM Transactions on Asian Language Information Processing (TALIP)* **12**(4) (2013), 1–17.
- [39] D.M.W. Powers, Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation, *Journal of Machine Learning Technologies.* **2** (2011), 37–63.
- [40] Y. Wu, M. Schuster, Z. Chen, Q.V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey et al., Google’s neural machine translation system: Bridging the gap between human and machine translation, *arXiv preprint arXiv:1609.08144* (2016).