

OAEI 2017 results of KEPLER

Marouen KACHROUDI¹, Gayo DIALLO², and Sadok BEN YAHIA¹

¹ Université de Tunis El Manar, Faculté des Sciences de Tunis
Informatique Programmation Algorithmique et Heuristique
LIPAH-LR 1 ES14, 2092, Tunis, Tunisie

{marouen.kachroudi, sadok.benyahia}@fst.rnu.tn

² BPH Center - INSERM U1219, Team ERIAS & LaBRI UMR5800,
Univ. Bordeaux

gayo.diallo@u-bordeaux.fr

Abstract. This paper presents and discusses the results produced by KEPLER for the 2017 Ontology Alignment Evaluation Initiative (OAEI 2017). This method is based on the exploitation of three different strategy levels. The proposed alignment method KEPLER is enhanced by the integration of powerful treatments inherited from other related domains, such as Information Retrieval (IR) [1]. For scaling, the method is equipped with a partitioning module. For the management of multilingualism, KEPLER develops a well-defined strategy based on the use of a translator, and this provides very encouraging results.

1 Presentation of the system

Given the substantial growth of the semantic Web users that create and update knowledge all over the world in a multitude of conceptualizations. This process has been accelerated due to a few initiatives which encourage all the active participants to make their data available to the public. These actors often publish their data sources in their own respective languages, in order to make this information interoperable and accessible to members of all communities [2]. As a solution, the ontology alignment process aims to provide semantic interoperable bridges between heterogeneous and distributed information systems. Indeed, the informative volume reachable via the semantic Web stresses needs of techniques guaranteeing the share, reuse and interaction of all resources [3]. The explication of the associated concepts related to a particular domain of interest resorts to ontologies, considered as the kernel of the semantic Web. In this register, KEPLER is an ontology alignment system dealing with the key challenges related to heterogeneous ontologies on the semantic Web, and it uses several hybrid alignment strategies. KEPLER is designed to discover alignments for both normal size and large scale ontologies. In addition, the proposed alignment approach has the ability to treat multilingual ontologies as well as monolingual ones.

1.1 State, purpose, general statement

The proposed method, KEPLER, exploits besides the classic techniques, an external resource, *i.e.*, a translator to deal with multilingualism. KEPLER implements an alignment strategy which aims at exploiting all the wealth of the used ontologies.

1.2 Specific techniques used

The main idea of KEPLER is to exploit the expressiveness of the OWL language to detect and compute the similarity between entities of two given ontologies through 6 complementary modules as presented in figure 1.

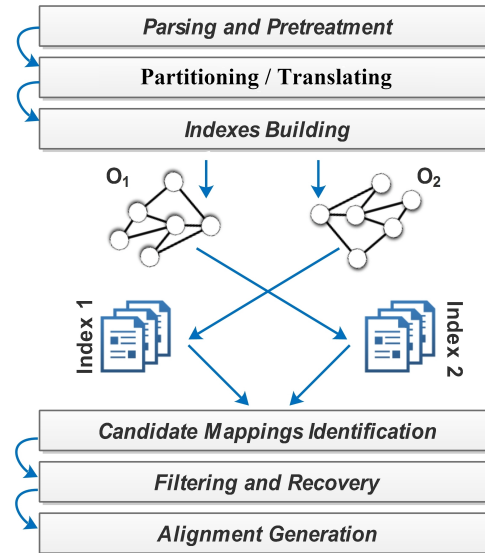


Fig. 1. KEPLER workflow.

Entities are described using OWL primitives with their semantics. We can then consider ontology as a semantic graph where entities are nodes connected by links which are OWL primitives. These links have specified semantic primitives. Consequently, if two ontologies in the same domain are similar, their semantic graphs are also the same.

Parsing and pretreatment This module allows to extract the ontological entities initially represented by a primary form of lists. In other words, at the parsing stage, we seek primarily to transform an OWL ontology in a well defined structure that preserves and highlight all the information contained in this ontology. Furthermore, in the resulting informative format, it has a considerable impact on the results of the similarity computation thereafter. Thus, we get couples formed by the entity name and its associated labels.

Partitioning This module aims at splitting ontologies into smaller parts to support the alignment task [4]. Consequently, partitioning a set $\mathcal{B}(\mathcal{C})$ is to find subsets $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n$, encompassing semantically close elements bound by a relevant set of relationships, *i.e.*, $\mathcal{O} = \bigcup \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n\}$, where \mathcal{B}_i is an ontological block, and n is the resulting number of extracted blocks. Hence, we can define an ontological portion as a reduced ontology that could be extracted from another larger one by splitting up the latter according to its both constituents : structures and semantics. One way to obtain such a

partitioning, can be to maximize the relationships inside a block and minimize the relationship between the blocks themselves. The partitioning quality result can be evaluated using different criteria:

- *The size of the generated blocks*: that must have a reasonable size, *i.e.*, a number of elements that can be handled by an alignment tool;
- *The number of the generated blocks*: this number should be as small as possible to limit the number of block pairs to be aligned;
- *The compactness degree of a block*: a block is said to be substantially compact if relations (lexical and structural ones) are stronger inside the block and low outside.

Translation : An originality of our system, is to solve the heterogeneity problem mainly due to multilingualism, given the importance of this research area [5, 6]. This challenge brings us to choose between two alternatives, either we consider the translation path to one of the languages according to the two input ontologies, or we consider the translation path to a chosen pivot language. At this stage, we must have a foreseeable vision for the rest of our approach. Specifically, at the semantic alignment stage we use an external resource, *i.e.*, WordNet³. The latter is a lexical database for the English language. Therefore, the choice is governed by the use of WordNet, and we will prepare a translation of the two ontologies to the pivot language, which is English. To perform the translation phase we chose Bing Microsoft⁴ tool.

Indexation : Indexing is one of the novelties of our approach. It consists in reducing the search space through the use of effective search strategy on the built indexes which represent the input ontologies components. To enable faster searching, the driving idea that was previously used in some works [1] is to execute the analysis in advance and store it in an optimized format for the search.

Candidate Mappings Identification : The role of this module is to find the entities in common between the indexes. Once the indexes are set up, the querying step of the latter is activated. Thus, the query implementation satisfies the terminology search and semantic aspects at once as we are querying documents in a vector representation that contain a given ontological entity and its synonyms obtained via WordNet. It is worthy to mention that indexes querying is done in both senses.

Filtering and Recovery : The filtering module consists of two complementary sub-modules, each one is responsible of a specific task in order to refine the set of primarily aligned candidates. At this stage, once the list of candidates is ready, the alignment method uses the first filter. We should note that indexes querying may includes a set of redundant mappings. Doing so, this filter eliminates the redundancy. It goes through the list of candidates and for each candidate, it checks if there are duplicates. If this is the case, it removes the redundant element(s). At the end of filtering phase, we have a candidates list without redundancy, however, there is always the concern of *false positives*,

³ <https://wordnet.princeton.edu/>

⁴ <https://www.bing.com/translator>

in fact, there was the need to establish a second filter. Once the redundant candidates are deleted, the system uses the second filter that eliminates *false positives*. This filter is applied to what we call *partially* redundant entities. An entity is considered as *partially* redundant if it belongs to two different mappings (*i.e.*, being given three ontological entities e_1 , e_2 and e_3 . If on the one hand, e_1 is aligned to e_2 , and secondly, e_1 is aligned to e_3 , this last alignment is qualified as doubtful. We note that our method generates (1 : 1) alignments. To overcome this challenge, the alignment method compares the topology of the two suspicious entities (e_3 neighbors with e_1 neighbors, e_2 neighbors with e_1 neighbors) with respect to the redundant entity e_1 , and retains the couple having the highest topological proximity value. All candidates are subject of this filter, and as output we have the final alignment file.

Alignment Generation : The result of the alignment process provides a set of mappings, which are serialized in the RDF format.

2 Results

In this section, we present the results obtained by KEPLER in the OAEI 2017.

2.1 Anatomy

This track consists of two real world ontologies to be matched, the source ontology describing the Adult Mouse Anatomy (with 2744 classes) and the target ontology is the NCI Thesaurus describing the Human Anatomy (with 3304 classes). For this track, and according to figure 2, KEPLER succeeded to extract 74% of correct mappings with a precision about 95%. Figure 2 summarizes the evaluation metrics values for Anatomy track. To this end, it is important to mention that KEPLER has managed to support the ontologies of the Anatomy database thanks to the *Ontopart* module [7, 4].

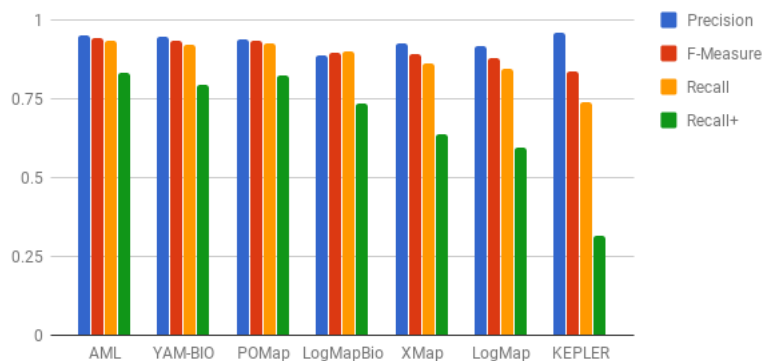


Fig. 2. KEPLER evaluation metrics among other pioneering systems for Anatomy track.

2.2 Conference

The conference track consists of 15 ontologies from the conference organization domain and each ontology must be matched against every other ontologies. The dataset describes the domain of organizing conferences from different perspectives. Precision values varies between 76% and 58%. Recall values varies between 48% and 68%. The metrics are obtained according to several evaluation scenarios.

2.3 Multifarm

This dataset is composed of a subset of the Conference track, translated in nine different languages (*i.e.*, Chinese, Czech, Dutch, French, German, Portuguese, Russian, Spanish and Arabic). With a special focus on multilingualism, it is possible to evaluate and compare the performance of alignment approaches through these test cases. Based on several previous contributions [8–13], the designed main goal of the MultiFarm track is to evaluate the ability of the alignment systems to deal with multilingual ontologies. It serves the purpose of evaluating the strength and weakness of a given system across languages.

KEPLER uses a specific technique to determine the equivalence between ontology entities described in different natural languages. We chose to use the English as a pivot language. The use of a pivot language ensures greater consistency of obtained translations since it starts from the same text. In the *different ontologies* case, the method is ranked fourth with a recall value of 0.31 as depicted by figure 3.

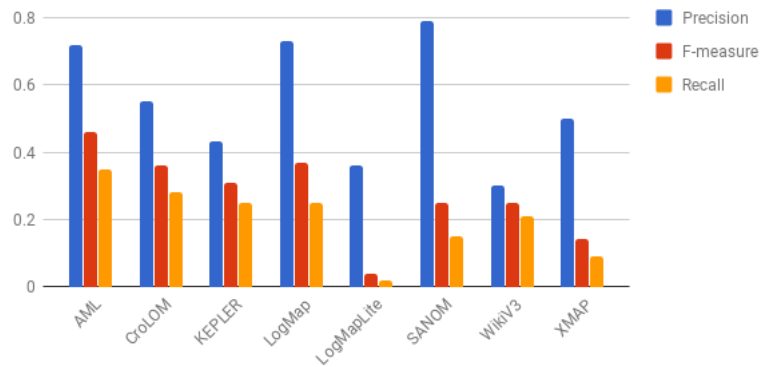


Fig. 3. KEPLER evaluation metrics among other pioneering systems for Multifarm track (*different ontologies*).

Whereas in the *same ontologies* case, the method occupies the first place with a recall value of 0.52 as flagged by figure 4.

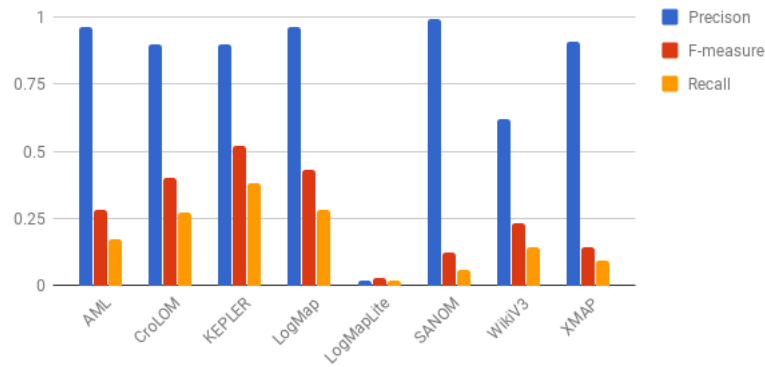


Fig. 4. KEPLER evaluation metrics among other pioneering systems for Multifarm track (*same ontologies*).

2.4 Large Biomedical Ontologies and Phenotype

In the scalability register, this track consists of finding alignments between the Foundational Model of Anatomy (FMA), SNOMED CT, and the National Cancer Institute Thesaurus (NCI). These ontologies are semantically rich and contain tens of thousands of classes. The Large BioMed Track consists of three matching problems, *i.e.*, (1) FMA-NCI matching problem, (2) FMA-SNOMED matching problem and (3) SNOMED-NCI matching problem. KEPLER handles large ontologies in two phases: the first phase consists on partitioning the ontologies into a set of blocks and the second phase selects two suitable blocks giving the highest value of similarity to be aligned. KEPLER treated (*Task 1: FMA-NCI small fragments*) [Precision : 0.96 / Recall : 0.83] according to figure 5.

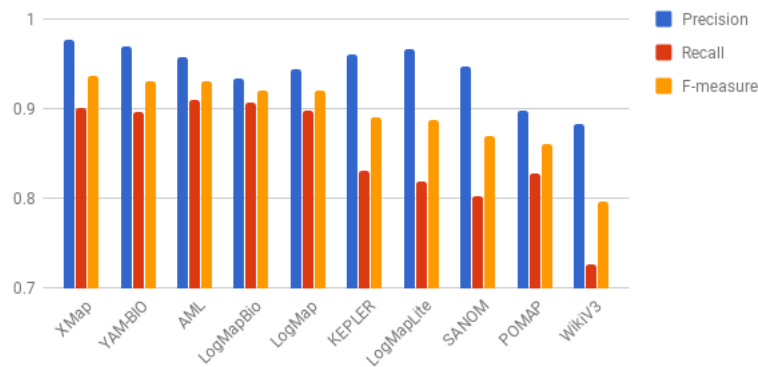


Fig. 5. KEPLER evaluation metrics among other pioneering systems for LargeBio track.

As depicted by figure 6, KEPLER processed also the task 3 of the LargeBio dataset (*FMA-SNOMED small fragments*) with a Precision of 0.82 and Recall of 0.55. In the Phenotype track, our method succeeds in processing only the DOID-ORDO sub-case by identifying 1824 matches for 1237 expected ones.

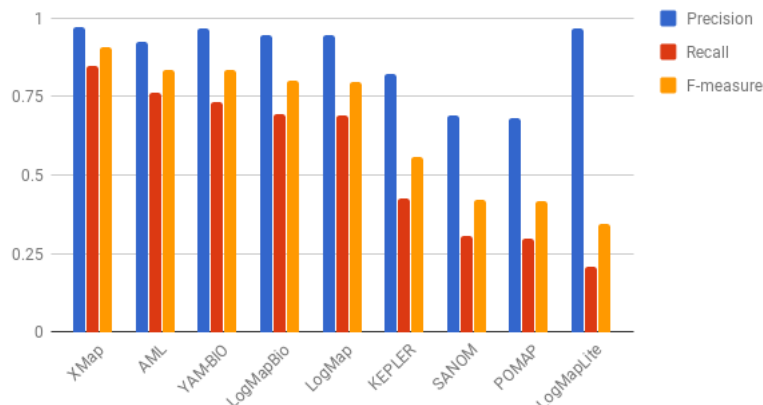


Fig. 6. KEPLER evaluation metrics among other pioneering systems for LargeBio track.

3 Conclusion

In this paper, we briefly presented the alignment system KEPLER with comments of the results obtained according to the OAEI 2017 tracks, corresponding to the SEALS platform evaluation modalities. Several observations regarding these results were highlighted, in particular the impact of the elimination of any ontological resource on the similarity values. KEPLER is an ongoing work which borrows its idea from two previous systems, CLONA [12] and SERVOMAP [1]. It showed promising results for its first participation. As future work, we plan to consolidate our system to more support the instance based ontology alignment in a wider range and context. We have dealt with this issue before [14, 15], but the test base update imposes other challenges, in terms of the used ontological languages and the evolutive semantic description formalisms.

References

1. Diallo, G.: An effective method of large scale ontology matching. *Journal of Biomedical Semantics* **5(44)** doi:10.1186/2041-1480-5-44 (2014)
2. Berners-Lee, T.: Designing the web for an open society. In: *Proceedings of the 20th International Conference on World Wide Web (WWW2011)*, Hyderabad, India (2011) 3–4
3. Suchanek, F.M., Varde, A.S., Nayak, R., Senellart, P.: The hidden web, xml and semantic web: A scientific data management perspective. *Computing Research Repository* (2011) 534–537

4. Kachroudi, M., Zghal, S., Ben Yahia, S.: Ontopart: at the cross-roads of ontology partitioning and scalable ontology alignment systems. *International Journal of Metadata, Semantics and Ontologies* **8**(3) (2013) 215–225
5. Diallo, G.: Efficient building of local repository of distributed ontologies. In: *Proceedings of the Seventh International Conference on Signal-Image Technology and Internet-Based Systems, SITIS 2011, Dijon, France, November 28 - December 1, 2011.* (2011) 159–166
6. Dramé, K., Diallo, G., Delva, F., Dartigues, J., Mouillet, E., Salamon, R., Mouglin, F.: Reuse of termino-ontological resources and text corpora for building a multilingual domain ontology: An application to alzheimer's disease. *Journal of Biomedical Informatics* **48** (2014) 171–182
7. Kachroudi, M., Hassen, W., Zghal, S., Ben Yahia, S.: Large ontologies partitioning for alignment techniques scaling. In: *Proceedings of the 9th International Conference on Web Information Systems and Technologies (WEBIST)*, 8-10 May, Aachen, Germany (2013) 165–168
8. Kachroudi, M., Ben Yahia, S., Zghal, S.: Damo - direct alignment for multilingual ontologies. In: *Proceedings of the 3rd International Conference on Knowledge Engineering and Ontology Development (KEOD)*, 26-29 October, Paris, France (2011) 110–117
9. Kachroudi, M., Zghal, S., Ben Yahia, S.: When external linguistic resource supports cross-lingual ontology alignment. In: *Proceedings of the 5th International Conference on Web and Information Technologies (ICWIT 2013)*, 9-12, May, Hammamet, Tunisia (2013) 327–336
10. Kachroudi, M., Zghal, S., Ben Yahia, S.: Using linguistic resource for cross-lingual ontology alignment. *International Journal of Recent Contributions from Engineering* **1**(1) (2013) 21–27
11. Kachroudi, M., Zghal, S., Ben Yahia, S.: Bridging the multilingualism gap in ontology alignment. *International Journal of Metadata, Semantics and Ontologies* **9**(3) (2014) 252–262
12. El Abdi, M., Souid, H., Kachroudi, M., Ben Yahia, S.: Clona results for oaei 2015. In: *Proceedings of the 12th International Workshop on Ontology Matching (OM-2015) Colocated with the 14th International Semantic Web Conference (ISWC-2015)*. Volume 1545 of CEUR-WS., Bethlehem (PA US) (2015) 124–129
13. Kachroudi, M., Diallo, G., Ben Yahia, S.: Initiating cross-lingual ontology alignment with information retrieval techniques. In: *Actes de la 6^{ème} Edition des Journées Francophones sur les Ontologies (JFO'2016)*, Bordeaux, France (2016) 57–68
14. Damak, S., Souid, H., Kachroudi, M., Zghal, S.: Exona results for oaei 2015. In: *Proceedings of the 12th International Workshop on Ontology Matching (OM-2015) Colocated with the 14th International Semantic Web Conference (ISWC-2015)*. Volume 1545 of CEUR-WS., Bethlehem (PA US) (2015) 145–149
15. Zghal, S., Kachroudi, M., Damak, S.: Alignement d'ontologies à base d'instances indexées. In: *Actes de la 6^{ème} Edition des Journées Francophones sur les Ontologies (JFO'2016)*, Bordeaux, France (2016) 69–74