

# Multilingual Ontology Matching Evaluation – A First Report on using MultiFarm

Christian Meilicke<sup>1</sup>, Cássia Trojahn<sup>2</sup>,  
Ondřej Šváb-Zamazal<sup>3</sup>, Dominique Ritze<sup>1</sup>

<sup>1</sup> University of Mannheim

<sup>2</sup> INRIA & LIG, Grenoble

<sup>3</sup> University of Economics, Prague

**Abstract.** This paper reports on the first usage of the MultiFarm dataset for evaluating ontology matching systems. This dataset has been designed as a comprehensive benchmark for multilingual ontology matching. In this first set of experiments, we analyze how state-of-the-art matching systems – not particularly designed for the task of multilingual ontology matching – perform on this dataset. Our experiments show the hardness of MultiFarm and result in baselines for any algorithm specifically designed for multilingual ontology matching. Moreover, this first reporting allows us to draw relevant conclusions for both multilingual ontology matching and ontology matching evaluation in general.

## 1 Introduction

Ontology matching is the task of finding correspondences that link concepts, properties or instances between two ontologies. Different approaches have been proposed for performing this task. They can be classified along the features in the ontologies (labels, structures, instances, semantics) they take into account or with regard to the kind of disciplines they belong to (e.g., statistics, combinatorial, semantics, linguistics, machine learning, or data analysis) [11, 9, 4].

With the aim of establishing a systematic evaluation of matching systems, the Ontology Alignment Evaluation Initiative (OAEI)<sup>4</sup> [3] has been carried out over the last years. It is an annual evaluation campaign that offers datasets, from different domains, organized by different groups of researchers. However, most of the OAEI datasets have been focused on monolingual tasks. A detailed definition on multilingual and cross-lingual ontology matching tasks can be found in [14]. The multilingual datasets so far available contain single pairs of languages, as the MLDirectory dataset,<sup>5</sup> which consists of website directories in English and Japanese; and the VLCR dataset,<sup>6</sup> that aims at matching the thesaurus of the Netherlands Institute for Sound and Vision, written in Dutch, to the WordNet and DBpedia, in English. Furthermore, these datasets contain only partial

<sup>4</sup> <http://oaei.ontologymatching.org/>

<sup>5</sup> <http://oaei.ontologymatching.org/2008/mldirectory/>

<sup>6</sup> <http://www.cs.vu.nl/~laurah/oaei/2009/>

reference alignments or are not fully open. Thus, they are not suitable for an extensive evaluation.

For overcoming the lack of a comprehensive benchmark for multilingual ontology, the MultiFarm<sup>7</sup> dataset has been designed. This dataset is based on the OntoFarm [16] dataset, which has been used successfully in OAEI in the Conference track. MultiFarm is composed of a set of seven ontologies translated in eight different languages and the complete corresponding alignments between these ontologies.

In this paper, we report on the first usage of MultiFarm for multilingual ontology matching evaluation. In [10], we have deeply discussed the design of MultiFarm, focusing on its multilingual features and the specificities of the translation process, with a very preliminary report on its evaluation. Here, we extend this preliminary evaluation and provide a deep discussion on the performance of matching systems. Our evaluation is based on a representative subset of MultiFarm<sup>8</sup> and a set of state-of-the-art matching systems participating in OAEI campaigns. These systems have not particularly been designed for matching ontologies described in different languages. The choice for these settings was caused by the fact that – to our knowledge – there exists no multilingual ontology matching system that is executable out of the box. For example, an implementation of a multilingual matching system is described in [5], however, it is not available for download.

We expect that the results of these systems can be topped by specific methods, however, with our experiments we establish a first non-trivial baseline specific for the MultiFarm dataset. To our knowledge, such a comprehensive evaluation has not yet been conducted so far in the field of multilingual ontology matching.

The rest of the paper is organised as follows. In Section 2, we first introduce the OntoFarm dataset and then we present its multilingual counterpart. We shortly discuss the hardness of MultiFarm and present the results that have been gathered in previous OAEI campaigns on the OntoFarm. In Section 3, we present the evaluation setting used to carry out our experiments and list the tools we have evaluated. In particular, we discuss why and how we applied specific configurations to some of the tools. In Section 4, we finally describe the results of our experiments. We mainly focus on highly aggregated results due to the enormous amount of generated data. In Section 5, we conclude the paper and discuss directions for future work.

## 2 Background on MultiFarm

In the following, we shortly describe the OntoFarm dataset, explain how MultiFarm has been constructed, and roughly report about evaluation results of the OAEI Conference track.

---

<sup>7</sup> The dataset has been thoroughly described in [10] and is available at <http://web.informatik.uni-mannheim.de/multifarm/>

<sup>8</sup> We have been discarding Russian and Chinese languages.

## 2.1 OntoFarm

The OntoFarm dataset is based on a set of 16 ontologies from conference organisation domain. All contained ontologies differ in numbers of classes, properties, and in their DL expressivity. They are very suitable for ontology matching tasks since they were independently designed by different people who used various kinds of resources for ontology design:

- actual conferences and their web pages,
- actual software tools for conference organisation support, and
- experience of people with personal participation in organisation of actual conferences

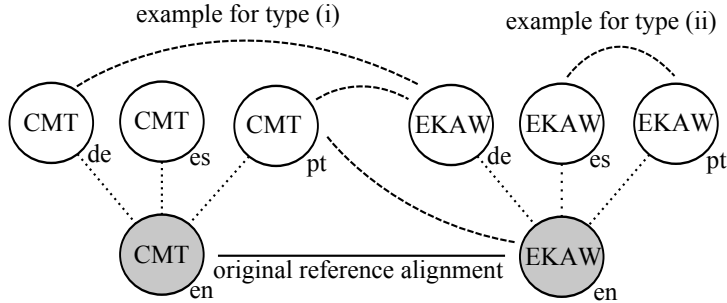
Thus, the OntoFarm dataset describes a quite realistic matching scenario and has been successfully applied in the OAEI within the Conference track since 2006. In 2008, a first version of the reference alignments was created and then annually enriched and updated up to current 21 reference alignments built between seven (out of 16) ontologies. Each of them has between four to 25 correspondences. The relatively small number of correspondences in the reference alignments are based on the fact that the reference alignments contain only simple equivalence correspondences. Due to different modeling styles of the ontologies, for many concepts and properties thus no equivalent counterparts exist.

## 2.2 MultiFarm

For generating the MultiFarm dataset, those seven OntoFarm ontologies, for which reference alignments are available, were manually translated into eight different languages (Chinese, Czech, Dutch, French, German, Portuguese, Russian, and Spanish). Since native speakers with a certain knowledge about the ontologies translated them, we do not expect any serious errors but of course they can never be excluded at all. Based on these translations, it is possible to re-create cross-lingual variants of the original test cases from the OntoFarm dataset as well as to exploit the translations more directly. Thus, the MultiFarm dataset contains two types of cross-lingual reference alignments.

We have depicted a small subset of the dataset shown in Figure 1. This figure indicates the cross-lingual reference alignments between different ontologies, derived from original alignments and translations (type (i)), and cross-lingual reference alignments between the same ontologies, which are directly based on the translations or on exploiting transitivity of translations (type (ii)). Reference alignments of type (i) cover only a small subset of all concepts and properties. We have explained this above for the original test cases of the OntoFarm dataset. In contrast, for test cases of type (ii) there are (translated) counterparts for each concept and property.

Overall, the MultiFarm dataset has  $36 \times 49$  test cases. 36 is a number of pairs of languages – each English ontology has its 8 language variants. 49 is the number of all reference alignments for each language pair. This is implied from



**Fig. 1.** Constructing MultiFarm from OntoFarm. Small subset that covers two ontologies and three translations. The solid line refers to a reference alignment of the OntoFarm dataset; dotted lines refer to translations; dashed lines refer to new cross-lingual reference alignments.

the number of original reference alignments (21) which is doubled (42) due to the fact that there is a difference between  $CMT_{en}-EKAW_{de}$  and  $CMT_{de}-EKAW_{en}$  in comparison with the original test cases where the test cases  $CMT-EKAW$  and  $EKAW-CMT$  are not distinguished. Additionally, we can also construct new reference alignments for matching each ontology on its translation which gives us seven additional reference alignments for each pair.

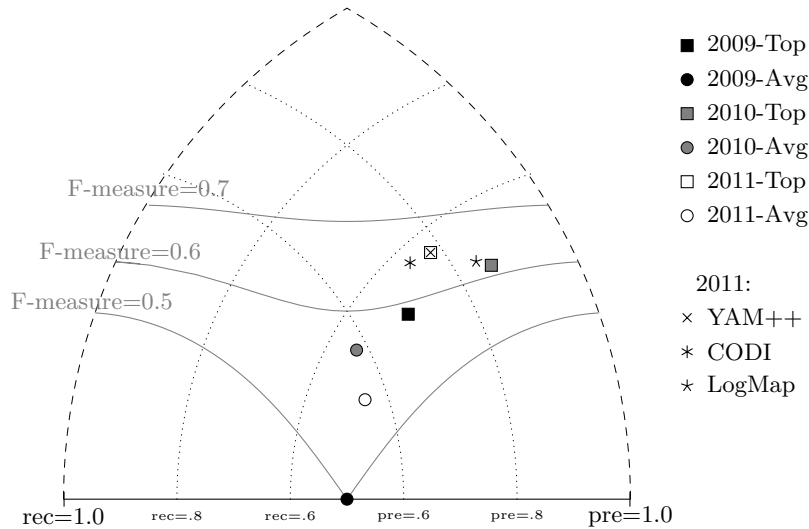
The main motivation for creating the MultiFarm dataset has been the ability to create a comprehensive set of test cases of type (i). We have especially argued in [10] that type (ii) test cases are not well suited for evaluating multilingual ontology matching systems, because they can be solved with very specific methods that are not related to the multilingual matching task.

### 2.3 Test Hardness

The OntoFarm dataset has a very heterogeneous character due to different modeling styles by various people. This leads to a high hardness of the resulting test cases. For example, the object property `writtenBy` occurs in several OntoFarm ontologies. When only considering the labels, one would expect that a correspondence like `writtenBy = writtenBy` correctly describes that these object properties are equivalent. However, in ontology  $\mathcal{O}_1$  the property indicates that a paper (domain) is written by an author (range), while in  $\mathcal{O}_2$  the property describes that a review (domain) is written by a reviewer (range). Therefore, this correspondence is not contained in the reference alignment between  $\mathcal{O}_1$  and  $\mathcal{O}_2$ . Similarly, comparing the English against the Spanish variant, there are the object properties `writtenBy` and `escrito por`. Pure translation would, similarly to the monolingual example, not result in detecting a correct correspondence. For that reason, the MultiFarm type (i) test cases go far beyond being a simple translation task.

The cross-lingual test cases of MultiFarm are probably much harder than the monolingual test cases of the OntoFarm. It is thus important to know how

matching systems perform on the OntoFarm dataset. These results can be understood as an upper bound that will be hard to top by results achieved for MultiFarm. In Figure 2, we have depicted some results of previous OAEI campaigns in a precision/recall triangular graph. This graph shows precision, recall, and F-measure in a single plot. It includes the best (squares) and average (circles) results of the 2009, 2010 and 2011 Conference track as well as results of the three best ontology matching systems (triangles) from 2011. Best results are considered according to the highest F-measure which corresponds to exactly one ontology matching system for each year. In 2011, YAM++ achieved the highest F-measure that is why its triangle overlaps with the white square depicting the best result of 2011.



**Fig. 2.** Precision/recall triangular graph for the last three Conference tracks. Horizontal line depicts level of precision/recall while values of F-measure are depicted by areas bordered by corresponding lines F-measure=0.[5|6|7].

On the one hand, Figure 2 shows that there is an improvement every year, except the average results of the last year. A reason might be the availability of the complete dataset over several years. Since the MultiFarm dataset has not been used in the past, we expect that evaluation results also improve over the years. On the other hand, we can see that recall is not very high (.63 in 2010 and .60 in 2011 for the best matching systems). This indicates that test cases of the OntoFarm dataset are especially difficult regarding recall measure.

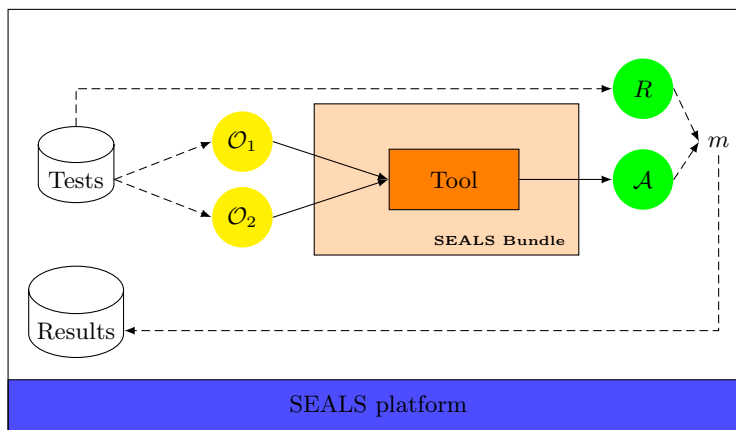
### 3 Evaluation Settings

In the following, we explain how we executed our evaluation experiments and list the matching systems that have been evaluated.

### 3.1 Evaluation Workflow

Following a general definition, *matching* is the process that determines an *alignment*  $\mathcal{A}$  for a pair of ontologies  $\mathcal{O}_1$  and  $\mathcal{O}_2$ . Besides the ontologies, there are other input parameters that are relevant for the matching process, namely: (i) the use of an input alignment  $\mathcal{A}'$ , which is to be extended or completed by the process; (ii) parameters that affect the matching process, for instance, weights and thresholds; and (iii) external resources used by the matching process, for instance, common knowledge and domain specific thesauri.

In this paper, we focus on evaluating a standard matching task. (i) In most of our experiments, we do not modify the parameters that affect the matching process. For two systems, we made an exception from this rule and report very briefly on the results. (ii) We do not use an additional input alignment at all. Note that most systems do not support such a functionality. (iii) We put no restriction on the external resources that are taken into account by the evaluated systems. Thus, we use the system standard settings for our evaluation. However, we obviously focus on the matching process where labels and annotations of  $\mathcal{O}_1$  and  $\mathcal{O}_2$  are described in different languages.



**Fig. 3.** Execution of tools.

The most common way to evaluate the quality of a matching process is to evaluate the correctness (precision) and completeness (recall) of its outcome  $\mathcal{A}$  by comparing  $\mathcal{A}$  against a reference alignment  $\mathcal{R}$ . Since 2010, in the context of OAEI campaigns, the process of evaluating matching systems has been automated thanks to the SEALS platform (Figure 3). For OAEI 2011, participants have been invited to wrap their tools into a format that can be executed by the platform, i.e. the matching process is not conducted by the tool developer but by the organisers of an evaluation using the platform. For the purpose of this paper, we benefit from the large number of matching tools that become available for our evaluation. Furthermore, evaluation test cases are available in the SEALS

repositories and can be used by everyone. Thus, all of our experiments can be completely reproduced.

### 3.2 Evaluated Matching Systems

As stated before, a large set of matching systems has already been uploaded to the platform in the context of OAEI 2011. We apply most of these tools to the MultiFarm dataset. In particular, we evaluated the tools AROMA [2], CIDER [6], CODI [7], CSA [15], LogMap and LogMapLt [8], MaasMatch [12], MapSSS [1], YAM++ [13] and Lily [17]. For most of these tools, we used the version submitted to OAEI 2011. However, some tool developers have already submitted a new version with some modifications. This is the case for CODI, LogMap and MapSSS. Moreover, the developer of LogMap has additionally uploaded a lite version of their matching systems called LogMapLt.

There have also been some systems participating in OAEI 2011 that are not listed here. We have not added them to the evaluation for different reasons. Some of these systems cannot finish the MultiFarm matching process in less than several weeks while others generate empty alignments for nearly all matching tasks or terminate with an error. We have to emphasise that none of these systems has originally been designed to solve the multilingual matching task.

## 4 Results

In the following, we discuss the evaluation results on different perspectives: first, aggregating the results obtained for all pairs of test cases (and languages) per matcher and then focusing on the different pairs of languages.

### 4.1 Differences in Test Cases

As explained in Section 2, the dataset can be divided in (i) those test cases where the ontologies to be matched are translations of different ontologies and (ii) those test cases where the same original ontology has been translated into two different languages and the translated ontologies have to be matched. We display the results for test cases of type (i) on the left and those for type (ii) on the right of Table 1. We have ordered the systems according to the F-measure for the test cases of type (i). The best results, in terms of F-measure, are achieved by CIDER (18%) followed by CODI (13%), LogMap (11%) and MapSSS (10%). CIDER has both better precision and recall scores than any other system. Compared to the top-results that have been reported for the original Conference dataset (F-measure > 60%), the test cases of the MultiFarm dataset are obviously much harder. However, an F-measure of 18% is already a remarkable result given the fact that we executed CIDER in its default setting.

The outcomes for test cases of type (ii) differ significantly. In particular, the results of MapSSS (67% F-measure) are surprisingly compared to the results presented for test cases of type (i). This system can leverage the specifics of type

| matcher   | (i) different ontologies |           |        |           | (ii) same ontologies |           |        |           |
|-----------|--------------------------|-----------|--------|-----------|----------------------|-----------|--------|-----------|
|           | size                     | precision | recall | F-measure | size                 | precision | recall | F-measure |
| CIDER     | 1433                     | 0.42      | 0.12   | 0.18      | 1090                 | 0.66      | 0.06   | 0.12      |
| CODI      | 923                      | 0.43      | 0.08   | 0.13      | 7056                 | 0.77      | 0.48   | 0.59*     |
| LogMap    | 826                      | 0.39      | 0.06   | 0.11      | 469                  | 0.71      | 0.03   | 0.06      |
| MapSSS    | 2513                     | 0.16      | 0.08   | 0.10      | 6008                 | 0.97      | 0.51   | 0.67*     |
| LogMapLt  | 826                      | 0.26      | 0.04   | 0.07      | 387                  | 0.56      | 0.02   | 0.04      |
| MaasMatch | 558                      | 0.24      | 0.03   | 0.05      | 290                  | 0.56      | 0.01   | 0.03      |
| CSA       | 17923                    | 0.02      | 0.06   | 0.03      | 8348                 | 0.49      | 0.36   | 0.42*     |
| YAM++     | 7050                     | 0.02      | 0.03   | 0.03      | 4779                 | 0.22      | 0.09   | 0.13*     |
| Aroma-    | 0                        | -         | 0.00   | -         | 207                  | 0.54      | 0.01   | 0.02      |
| Lily      | 0                        | -         | 0.00   | -         | 11                   | 1.00      | 0.00   | 0.00      |

**Table 1.** Results aggregated per matching system.

(ii) test cases to cope with the problem of matching labels expressed in different languages. Similar to MapSSS, we also observe a higher F-measure for CODI, CSA, and YAM++. We have marked those systems with an asterisk. Note that all these systems have an F-measure of at least five times higher than the F-measure for test cases of type (i). For all other systems, we observe a slightly decreased F-measure comparing test cases of type (i) with type (ii).

Again, we have to highlight the differences between both types of test cases. Reference alignments of type (i) cover only a small fraction of all concepts and properties described in the ontologies. This is not the case for test cases of type (ii). Here, we have complete alignments that connect each concept and property with an equivalent counterpart in the other ontology. There seems to be a clear distinction between systems that are specialised or configured to generate complete alignments in the absence of (easy) usable label description, and other systems that focus on generating good results for test cases of type (i).

Comparing these results with the results for the OAEI 2011 Benchmark track, it turns out that all systems marked with an asterisk have been among the top five systems of this track. All Benchmark test cases have a similar property, namely, their reference alignments contain for each entity of the smaller ontology exactly one counterpart in the larger ontology. An explanation for this can be that these systems have been developed or at least configured to score well for the Benchmark track. For that reason, they generate good results for test cases of type (ii), while their results for test cases of type (i) are less good. MapSSS and CODI are an exception. These systems generate good results for both test cases of type (i) and (ii).

## 4.2 Differences in Languages

Besides aggregating the results per matcher, we have analysed the results per pair of languages (Table 2), for the case where different ontologies are matched (type (i) in Table 1). As stated before, the multilingual aspect of the matching process can be negligible for matchers that are able to adapt their strategies to match structurally similar ontologies. We have also compared the matchers with a simple edit distance algorithm on labels (edna).



| pairs   | edna | Aroma | CIDER | CODI | CSA  | Lily | LogMap | LogLt | MaasMatch | MapSSS | YAM++ | average     |
|---------|------|-------|-------|------|------|------|--------|-------|-----------|--------|-------|-------------|
| cz-de   | 0.01 | -     | 0.12  | 0.11 | 0.03 | -    | 0.09   | 0.09  | 0.02      | 0.07   | 0.03  | 0.06        |
| cz-en   | 0.01 | -     | 0.20  | 0.12 | 0.04 | -    | 0.06   | 0.04  | 0.03      | 0.08   | -     | 0.07        |
| cz-es   | 0.00 | -     | 0.14  | 0.13 | 0.02 | -    | 0.11   | 0.11  | -         | 0.11   | 0.03  | 0.08        |
| cz-fr   | 0.01 | -     | 0.08  | 0.01 | 0.01 | -    | 0.01   | 0.01  | 0.01      | 0.01   | 0.03  | 0.02        |
| cz-nl   | 0.00 | -     | 0.09  | 0.09 | 0.04 | -    | 0.04   | 0.04  | 0.04      | 0.05   | 0.03  | 0.05        |
| cz-pt   | 0.00 | -     | 0.15  | 0.15 | 0.04 | -    | 0.13   | 0.13  | 0.02      | 0.12   | 0.04  | 0.09        |
| de-en   | 0.01 | -     | 0.31  | 0.22 | 0.03 | -    | 0.22   | 0.20  | 0.20      | 0.16   | -     | <b>0.17</b> |
| de-es   | 0.01 | -     | 0.25  | 0.20 | 0.02 | -    | 0.19   | 0.06  | -         | 0.15   | 0.03  | <b>0.11</b> |
| de-fr   | 0.00 | -     | 0.18  | 0.18 | 0.01 | -    | 0.17   | 0.04  | 0.04      | 0.13   | 0.03  | 0.09        |
| de-nl   | 0.01 | -     | 0.22  | 0.08 | 0.03 | -    | 0.05   | 0.04  | 0.04      | 0.15   | 0.03  | 0.07        |
| de-pt   | 0.01 | -     | 0.10  | 0.09 | 0.03 | -    | 0.07   | 0.07  | 0.01      | 0.06   | 0.04  | 0.05        |
| en-es   | 0.00 | -     | 0.25  | 0.24 | 0.03 | -    | 0.18   | 0.04  | 0.04      | 0.18   | -     | <b>0.12</b> |
| en-fr   | 0.01 | -     | 0.20  | 0.24 | 0.03 | -    | 0.19   | 0.04  | 0.04      | 0.13   | -     | <b>0.11</b> |
| en-nl   | 0.01 | -     | 0.22  | 0.10 | 0.04 | -    | 0.07   | 0.10  | 0.07      | 0.15   | -     | <b>0.10</b> |
| en-pt   | 0.00 | -     | 0.15  | 0.11 | 0.06 | -    | 0.06   | 0.06  | 0.06      | 0.07   | -     | 0.07        |
| es-fr   | 0.01 | -     | 0.29  | 0.07 | 0.02 | -    | 0.06   | 0.01  | 0.04      | 0.06   | 0.03  | 0.07        |
| es-nl   | 0.01 | -     | 0.07  | 0.01 | 0.02 | -    | -      | -     | -         | 0.01   | 0.02  | 0.02        |
| es-pt   | 0.01 | -     | 0.29  | 0.26 | 0.06 | -    | 0.27   | 0.23  | 0.09      | 0.23   | 0.03  | <b>0.16</b> |
| fr-nl   | 0.01 | -     | 0.23  | 0.14 | 0.02 | -    | 0.13   | 0.12  | 0.13      | 0.11   | 0.03  | <b>0.10</b> |
| fr-pt   | 0.00 | -     | 0.11  | 0.06 | 0.02 | -    | 0.06   | -     | 0.04      | 0.02   | 0.03  | 0.04        |
| nl-pt   | 0.00 | -     | 0.02  | 0.04 | 0.03 | -    | 0.01   | 0.01  | 0.02      | 0.02   | 0.04  | 0.02        |
| average | 0.01 |       | 0.17  | 0.13 | 0.03 |      | 0.11   | 0.08  | 0.05      | 0.10   | 0.03  | 0.08        |

**Table 2.** Results per pairs of languages for different ontologies.

With exception of Aroma and Lily, which are not able to deal with the complexity of the matching task, for most of the test cases no matcher has lower F-measure than edna. For some of them, however, LogMap, LogMapLt, MaasMatch and YAM++, respectively, have not provided any alignment. YAM++ has a specific behaviour and is not able to match the English ontologies to any other languages. For the other matchers, it (incidentally) happens mostly for the pairs of languages that do not share the same root language (e.g. es-nl or de-es). The exception is LogMapLt, which is not able to identify any correspondence between fr-pt, even if these languages have the same root language (e.g. Latin) and thus have a similar vocabulary. It could be expected that matchers should be able to find a higher number of correspondences for the pairs of languages where there is an overlap in their vocabularies because most of the matcher apply some label similarity strategy. However, it is not exactly the case in MultiFarm. The dataset contains many complex correspondences that cannot be found by a single translation process or by string comparison. This can be partially corroborated by the very low performance of edna in all test cases.

Looking at the results for each pair of languages, per matcher, the best five F-measures are obtained for de-en (31%), es-fr/es-pt (29%), de-es/en-es (25%), all for CIDER, en-es/en-fr (24%), for CODI, and fr-nl (23%) again for CIDER. We could observe that 3 ahead pairs contain languages with some degree of overlap

in their vocabularies (i.e., de-en, es-fr, es-pt). For each individual matcher, seven out of eight matchers have their best scores for these pairs (exception is YAM++ that scores better for cz-pt and de-pt), with worst scores in cz-fr, es-nl, which have very different vocabularies.

When aggregating the results per pair of languages, that order is mostly preserved (highly affected by CIDER): de-en (17%), es-pt (16%), en-es (12%), de-es/en-fr (11%), followed by fr-nl/en-nl (10%). The exception is for the pair es-fr, where the aggregated F-measure decreases to 7%. Again, the worst scores are obtained for cz-fr, nl-pt and es-nl. We can observe that, for most of the cases, the features of the languages (i.e., their overlapping vocabularies) have an impact in the matching results. However, there is no universal pattern and we have cases with similar languages where systems score very low (fr-pt, for instance). This has to be further analysed with a deep analysis of the individual pairs of ontologies.

## 5 Discussion

Some of the reported results are relevant for multilingual ontology matching in general, while others help us to understand the characteristics of the MultiFarm dataset. The latter ones are relevant for any further evaluation that builds on the dataset. Moreover, we can also draw some conclusions that might be important for the use of datasets in the general context of ontology matching evaluation.

*Exploiting structural information* Very good results for test cases of type (ii) can be achieved by methods non-specific to multilingual ontology matching. The result of MapSSS is an interesting example. This was also one of the main reasons why the MultiFarm dataset has been constructed as a comprehensive collection for test cases of type (i) and (ii). We suggest to put a stronger focus on test cases of type (i) in the context of evaluating multilingual ontology matching techniques. Otherwise, it remains unclear whether the measured results are based on multilingual techniques or on exploiting that the matched ontologies can be interpreted as versions of the same ontology.

*Finding a good configuration* The results for test cases of type (i) show that state-of-the-art matching systems are not very well suited for the tasks of matching ontologies described in different languages, especially when executed in their default setting. We started another set of experiments by running some tools (CODI, LogMap, Lily) in a manually configured setting better suited for the matching task. A first glimpse, the results shows that it is possible to increase the average F-measure up to a value of 26%. Thus, we are planning to further investigate the influence of configurations on multilingual matching tasks within more extensive experiments.<sup>9</sup>

---

<sup>9</sup> We would like to thank Ernesto Jimenez-Ruiz (LogMap [8]) and Peng Wang (Lily [17]) for supporting us with a quick hint about a good, manually modified configuration for running their systems on MultiFarm.

*The role of language features* We cannot neglect certain language features (like their overlapping vocabularies) in the matching process. Once most of the matchers take advantage of label similarities it is likely that it may be harder to find correspondences between Czech and Portuguese ontologies than Spanish and Portuguese ones. In our evaluation, for most of the systems, the better performance where incidentally observed for the pairs of languages that have some degree of overlap in their vocabularies. This is somehow expected, however, we could find exceptions to this behavior. In fact, MultiFarm requires systems exploiting more sophisticated matching strategies than label similarity and for many ontologies in MultiFarm it is the case.

*Implications on analyzing OAEI results* Aside from the topic of multilingual ontology matching, the results implicitly emphasise the different characteristics of test cases of type (i) and (ii). An example can be found when comparing results for the OAEI Benchmark and Conference track. The Benchmark track is about matching different versions (some slightly modified, some heavily modified) of the same ontology. The Conference dataset is about matching different ontologies describing the same domain. This difference finds its counterparts in the distinction between type (i) and type (ii) ontologies in the MultiFarm dataset. Without taking this distinction into account, it is not easy to draw valid conclusions on the generality of measured results.

## 6 Future Work

Even though we reported about diverse aspects, we could not analyse or evaluate all interesting issues. The following listing shows possible extensions and improvements for further evaluations.

- Executing (all or many) matching systems with a specifically tailored configuration;
- Including Chinese and Russian versions of the ontologies in the evaluation setting;
- Exploiting automatic translation strategies and evaluate their impact on the matching process;
- Analysing the role of diacritics : in some languages, the same word written with or without accent can have a different meaning, e.g., in French ‘où’ (where) is different from ‘ou’ (or).

There are many things left to do, however, we have shown that the MultiFarm dataset is a useful, comprehensive, and a difficult dataset for evaluating ontology matching systems. We strongly recommend to apply this resource and to compare measured results against the results presented in this paper. In particular, we encourage developers of ontology matching systems, specifically designed to match ontologies described in different languages, to make use of the dataset and to report about achieved results.

## Acknowledgements

Some of the authors are partially supported by the SEALS project (IST-2009-238975). Ondřej Šváb-Zamazal has been partially supported by the CSF grant no. P202/10/0761.

## References

1. M. Cheatham. MapSSS results for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 184–189, 2011.
2. J. David. Aroma results for OAEI 2009. In *Proc. 4th ISWC workshop on Ontology Matching (OM)*, pages 147–152, 2009.
3. J. Euzenat, C. Meilicke, H. Stuckenschmidt, P. Shvaiko, and C. T. dos Santos. Ontology alignment evaluation initiative: Six years of experience. *Journal on Data Semantics*, 15:158–192, 2011.
4. J. Euzenat and P. Shvaiko. *Ontology matching*. Springer, Heidelberg (DE), 2007.
5. B. Fu, R. Brennan, and D. O’Sullivan. Using pseudo feedback to improve cross-lingual ontology mapping. In *Proc. 8th Extended Semantic Web Conference (ESWC)*, pages 336–351, 2011.
6. J. Gracia, J. Bernad, and E. Mena. Ontology matching with CIDER: evaluation report for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 126–133, 2011.
7. J. Huber, T. Sztylek, J. Noessner, and C. Meilicke. Codi: Combinatorial optimization for data integration: results for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 134–141, 2011.
8. E. Jimenez-Ruiz, A. Morant, and B. C. Grau. LogMap results for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 163–170, 2011.
9. Y. Kalfoglou and M. Schorlemmer. Ontology mapping: the state of the art. *The Knowledge Engineering Review*, 18(1):1–31, 2003.
10. C. Meilicke, R. G. Castro, F. Freitas, W. R. van Hage, E. Montiel-Ponsoda, R. R. de Azevedo, H. Stuckenschmidt, O. Šváb-Zamazal, V. Svátek, A. Taminin, C. Trojahn, and S. Wang. Multifarm: A benchmark for multilingual ontology matching. *Journal of Web Semantics*. Submitted 2011.
11. E. Rahm and P. Bernstein. A survey of approaches to automatic schema matching. *The VLDB Journal*, 10(4):334–350, 2001.
12. F. Schadd and N. Roos. Maasmatch results for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 171–178, 2011.
13. F. Schadd and N. Roos. YAM++ – results for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 228–235, 2011.
14. D. Spohr, L. Hollink, and P. Cimiano. A machine learning approach to multilingual and cross-lingual ontology matching. In *Proc. 10th International Semantic Web Conference (ISWC)*, pages 665–680, 2011.
15. Q.-V. Tran, R. Ichise, and B.-Q. Ho. Cluster-based similarity aggregation for ontology matching. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 142–147, 2011.
16. O. Šváb, V. Svátek, P. Berka, D. Rak, and P. Tomášek. Ontofarm: Towards an experimental collection of parallel ontologies. In *Poster Track of ISWC 2005*, 2005.
17. P. Wang. Lily results on SEALS platform for OAEI 2011. In *Proc. 6th ISWC workshop on Ontology Matching (OM)*, pages 156–162, 2011.