

# Matching Disease and Phenotype Ontologies in the Ontology Alignment Evaluation Initiative

Ian Harrow<sup>1\*</sup>†, Ernesto Jiménez-Ruiz<sup>2†</sup>, Andrea Splendiani<sup>3</sup>, Martin Romacker<sup>4</sup>, Peter Woollard<sup>5</sup>, Scott Markel<sup>6</sup>, Yasmin Alam-Faruque<sup>7</sup>, Martin Koch<sup>8</sup>, James Malone<sup>9</sup> and Arild Waaler<sup>2</sup>

## Abstract

**Background:** The disease and phenotype track was designed to evaluate the relative performance of ontology matching systems that generate mappings between source ontologies. Disease and phenotype ontologies are important for applications such as data mining, data integration and knowledge management to support translational science in drug discovery and understanding the genetics of disease.

**Results:** 11 systems (out of 21 OAEI participating systems) were able to cope with at least one of the tasks in the *Disease and Phenotype* track. AML, FCA-Map, LogMap(Bio) and PhenoMF systems produced the top results for ontology matching in comparison to consensus alignments. The results against manually curated mappings proved to be more difficult most likely because these mapping sets comprised mostly subsumption relationships rather than equivalence. Manual assessment of unique equivalence mappings showed that AML, LogMap(Bio) and PhenoMF systems have the highest precision results.

**Conclusions:** Four systems gave the highest performance for matching disease and phenotype ontologies. These systems coped well with the detection of equivalence matches, but struggled to detect semantic similarity. This deserves more attention in the future development of ontology matching systems. The findings of this evaluation show that such systems could help to automate equivalence matching in the workflow of curators, who maintain ontology mapping services in numerous domains such as disease and phenotype.

**Keywords:** biomedical ontology; ontology alignment; OAEI; evaluation; phenotype; disease

## Background

The Pistoia Alliance Ontologies Mapping project<sup>[1]</sup> was set up to find or create better tools and services for mapping between ontologies (including controlled vocabularies) in the same domain and to establish best practices for ontology management in the Life Sciences. The project has developed a formal process to define and submit a request for information (RFI) from ontology matching system providers to enable their evaluation.<sup>[2]</sup> A critical component of any ontology alignment system is the embedded matching algorithm, therefore the Ontologies Mapping project is supporting their development and evaluation through sponsorship and organisation of the *Disease and Phenotype* track (added in 2016) for the OAEI campaign [1]. In this paper we describe the experiences

and results in the OAEI 2016 *Disease and Phenotype* track.<sup>[3]</sup>

The *Disease and Phenotype* track is based on a real use case where it is required to find two pairwise alignments between disease and phenotype ontologies: (i) Human Phenotype Ontology [3] (HP) to Mammalian Phenotype Ontology [4] (MP), and (ii) Human Disease Ontology [5] (DOID) to Orphanet Rare Disease Ontology<sup>[4]</sup> (ORDO). The first task maps between human and the more general mammalian phenotype ontologies. This is important for translational science in drug discovery, since mammalian models such as mice are widely used to study human diseases and their underlying genetics. Mapping human phenotypes to other mammalian phenotypes greatly facilitates the extrapolation from model animals to humans. The second task maps between two disease ontologies: the more generic DOID and the more specific ORDO, in the context of rare human diseases. These ontologies

\*Correspondence: [ian.harrow@pistoiaalliance.org](mailto:ian.harrow@pistoiaalliance.org)

<sup>1</sup>Pistoia Alliance Ontologies Mapping Project, Pistoia Alliance Inc, USA  
Full list of author information is available at the end of the article

†Equal contributor

<sup>[1]</sup><http://www.pistoiaalliance.org/projects/ontologies-mapping>

<sup>[2]</sup><https://pistoiaalliance.atlassian.net/wiki/display/PUB/Ontologies+Mapping+Resources>

Ontologies+Mapping+Resources

<sup>[3]</sup>The contents of this paper have been partially reported in the OAEI 2016 annual report [1], published within the “informal” proceedings of the Ontology Matching workshop [2].

<sup>[4]</sup>[http://www.orphadata.org/cgi-bin/inc/ordo\\_orphanet.inc.php](http://www.orphadata.org/cgi-bin/inc/ordo_orphanet.inc.php)

can support investigative studies to understand how genetic variation can cause or contribute to disease.

Currently, mappings between the aforementioned ontologies within the disease and phenotype domain are mostly generated manually by bioinformaticians and disease experts. Inclusion of automated ontology matching systems into such curation workflows is likely to improve the efficiency and scalability of this process to expand the coverage across many source ontologies. Automation of mappings is also important because the source ontologies are dynamic, often having more than ten versions per year which means the mappings must be maintained to remain useful and valid.

## Preliminaries

In this paper we assume that the ontologies are represented using the OWL 2 Web Ontology Language [6], which is a World Wide Web Consortium (W3C) recommendation.<sup>[5]</sup> Description Logics (DL) are the formal underpinning of OWL 2 [7].

An ontology *mapping* (also called *match* or *correspondence*) between entities of two ontologies  $\mathcal{O}_1, \mathcal{O}_2$  is typically represented as a 4-tuple  $\langle e, e', r, c \rangle$  where  $e$  and  $e'$  are entities of  $\mathcal{O}_1$  and  $\mathcal{O}_2$ , respectively;  $r \in \{\sqsubseteq, \supseteq, \equiv\}$  is a semantic relation; and  $c$  is a confidence value, usually, a real number within the interval  $(0 \dots 1]$ . Mapping confidence intuitively reflects how reliable a mapping is (*i.e.*, 1 = very reliable, 0 = not reliable).

An ontology *alignment*  $\mathcal{M}$  between two ontologies, namely  $\mathcal{O}_1, \mathcal{O}_2$ , is a set of mappings between  $\mathcal{O}_1$  and  $\mathcal{O}_2$ . In the ontology matching community, mappings are typically expressed using the RDF Alignment format [8]. In addition, mappings can also be represented through standard OWL 2 axioms (*e.g.*, [9]). This representation enables the use of the OWL 2 reasoning infrastructure that is currently available.

When mappings are translated into OWL 2 axioms, an *aligned ontology*  $\mathcal{O}^{\mathcal{M}} = \mathcal{O}_1 \cup \mathcal{O}_2 \cup \mathcal{M}$  is the result of merging the input ontologies and an alignment between them. The aligned ontology is also an OWL 2 ontology.

An ontology *matching system* is a program<sup>[6]</sup> that, given two input ontologies  $\mathcal{O}_1$  and  $\mathcal{O}_2$ , generates an ontology alignment  $\mathcal{M}^S$ .

An ontology *matching task* is typically composed by one or more pairs of ontologies with their correspondent *reference alignments*  $\mathcal{M}^{RA}$ . Reference alignments can be of different nature: gold standards, silver standards and baselines. *Gold standards* are typically (almost) complete mapping sets that have been manually curated by domain experts, while *silver standard*

mapping sets are not necessarily complete nor correct. Finally, *baseline* mappings typically represent a highly incomplete set of the total mappings. In this paper we use a type of *silver standard* that has been created by voting the mappings produced by several matching systems. In the remainder of the paper, we refer to this (silver standard) mapping set as *consensus alignments*.

The standard evaluation measures, for a system generated alignment  $\mathcal{M}^S$ , are *precision* (P), *recall* (R) and *f-measure* (F) computed against a reference alignment  $\mathcal{M}^{RA}$  as follows:

$$P = \frac{|\mathcal{M}^S \cap \mathcal{M}^{RA}|}{|\mathcal{M}^S|}, R = \frac{|\mathcal{M}^S \cap \mathcal{M}^{RA}|}{|\mathcal{M}^{RA}|}, F = 2 \cdot \frac{P \cdot R}{P + R} \quad (1)$$

Standard precision and recall have, however, limitations when considering the (OWL 2) semantics of the input ontologies and the mappings. Hence a mapping  $m$  such that  $m \in \mathcal{M}^S$  and  $m \notin \mathcal{M}^{RA}$  will penalise the standard precision value even though  $\mathcal{O}^{\mathcal{M}^{RA}} \models m$ , that is,  $m$  is inferred or entailed (using OWL 2 reasoning) by the union of the input ontologies  $\mathcal{O}_1$  and  $\mathcal{O}_2$  and the reference mappings  $\mathcal{M}^{RA}$ . Analogously, a mapping  $m$  such that  $m \notin \mathcal{M}^S$  and  $m \in \mathcal{M}^{RA}$  will penalise standard recall, even though the aligned ontology  $\mathcal{O}^{\mathcal{M}^S}$  can entail  $m$ . In this paper we adopt the notion of *semantic precision* and *recall* as defined in Equations 2 and 3 to mitigate the limitations of the standard measures (the interested reader please refer to [10, 11] for alternative definitions).

Semantic precision and recall, as presented in this paper, may still suffer from some limitations [12]. In order to reduce the impact of these limitations, when computing semantic precision and recall, equivalence mappings ( $\equiv$ ) are split into two subsumption mappings ( $\sqsubseteq$  and  $\supseteq$ ).

Note that when evaluating the mappings produced by a matching system against (incomplete) baseline mappings, only semantic recall should be taken into account.

$$P_{(sem)} = \frac{|\{m \in \mathcal{M}^S \mid \mathcal{O}^{\mathcal{M}^{RA}} \models m\}|}{|\mathcal{M}^S|} \quad (2)$$

$$R_{(sem)} = \frac{|\{m \in \mathcal{M}^{RA} \mid \mathcal{O}^{\mathcal{M}^S} \models m\}|}{|\mathcal{M}^{RA}|} \quad (3)$$

An ontology is *incoherent* [13] if it contains logical errors in the form of unsatisfiable concepts. If the union of the input ontologies  $\mathcal{O}_1$  and  $\mathcal{O}_2$  and the reference mappings  $\mathcal{M}^{RA}$  is incoherent, semantic precision and recall, as defined in Equations 2 and 3, may lead to

<sup>[5]</sup><https://www.w3.org/TR/owl2-overview/>

<sup>[6]</sup>Typically automatic, although there are systems that also allow human interaction

unexpected results. In this case, mapping repair (e.g., [13, 14, 15]) techniques should be applied before computing semantic precision and recall.

## Methodology

The Ontology Alignment Evaluation Initiative<sup>[7]</sup> (OAEI) is an annual campaign for the systematic evaluation of ontology matching systems [1, 16, 17, 18]. The main objective is the comparison of ontology matching systems on the same basis and to enable the reproducibility of the results. The OAEI included 9 different tracks organised by different research groups and involving different matching tasks.

The novel *Disease and Phenotype*<sup>[8]</sup> track was one of the new additions in the OAEI 2016 campaign. The track aims at evaluating the performance of systems in a real-world use case where pairwise alignments between disease and phenotype ontologies are required.

The *Disease and Phenotype* track closely followed the OAEI phases as summarised in Figure 1.

### Dataset

The *Disease and Phenotype* track comprises two matching tasks that involve the alignment of the Human Phenotype Ontology (HP), the Mammalian Phenotype Ontology (MP), the Human Disease Ontology (DOID), and the Orphanet Rare Disease Ontology (ORDO). Table 1 shows the metrics provided by BioPortal of these ontologies.

*Task 1:* pairwise alignment of the HP and the MP ontologies (HP-MP matching task).

*Task 2:* pairwise alignment of the DOID and the ORDO ontologies (DOID-ORDO matching task).

### Preparation phase

As specified by the OAEI the ontologies and (public) reference alignments were made available in advance during the first week of June 2016. The ontologies and mappings were downloaded from BioPortal [19] on June 2nd.

The mappings were obtained using a script that, given a pair of ontologies, uses BioPortal's REST API<sup>[9]</sup> to retrieve all mappings between those ontologies. We focused only on *skos:closeMatch* (BioPortal) mappings<sup>[10]</sup> as suggested in [20], and we

represented them as equivalence mappings.<sup>[11]</sup> The BioPortal-based alignment between HP and MP consisted in 639 equivalence mappings, while the alignment between DOID and ORDO included 1,018 mappings. Mappings were made available in both RDF Alignment and OWL 2 formats.

The preparatory phase gives the opportunity to both OAEI track organisers and participants to find and correct problems in the datasets. During this phase we noticed that the BioPortal mappings were highly incomplete.<sup>[12]</sup> Hence, the participants were notified that the BioPortal-based mappings were to be used as a *baseline* and not as a *gold standard* reference alignment. Given the limitations of the BioPortal mappings we were in need of creating a (blind) *consensus reference alignment* to perform the (automatic) evaluation (see details in the *Evaluation phase* section).

All (open) OAEI datasets were released on July 15th, 2016 and did not evolve after that.

### Execution phase

System developers had to implement a simple interface and to wrap their tools including all required libraries and resources in order to use the SEALS infrastructure.<sup>[13]</sup> The use of the SEALS infrastructure ensures that developers can perform a full evaluation locally and eases the reproducibility and comparability of the results.

This phase was conducted between July 15th and August 31st, 2016. During this time OAEI organisers attended technical issues reported by the developers. We also requested system developers to register their systems and their intention to participate in the different OAEI tracks by July 31st. Thirty systems were registered, from which 14 seemed potential participants of the *Disease and Phenotype* track.

### Evaluation phase

Participants were required to submit their wrapped tools by August 31st, 2016. From the 30 registered systems only 21 were finally submitted, and 13 were annotated (by the system developers) as participants of the *Disease and Phenotype* track. The final results were published on the OAEI website by October 15th.

The evaluation for the *Disease and Phenotype* track was semi-automatic with support of the SEALS infrastructure. Systems were evaluated according to the

<sup>[11]</sup>We did not consider mappings labelled as *skos:exactMatch* since they represent correspondences between entities with the same URI, and thus these mappings are redundant if translated into OWL 2 axioms.

<sup>[12]</sup>Our tests with last year participants revealed a large amount of missing valid mappings. The *Results* section quantifies this degree of incompleteness.

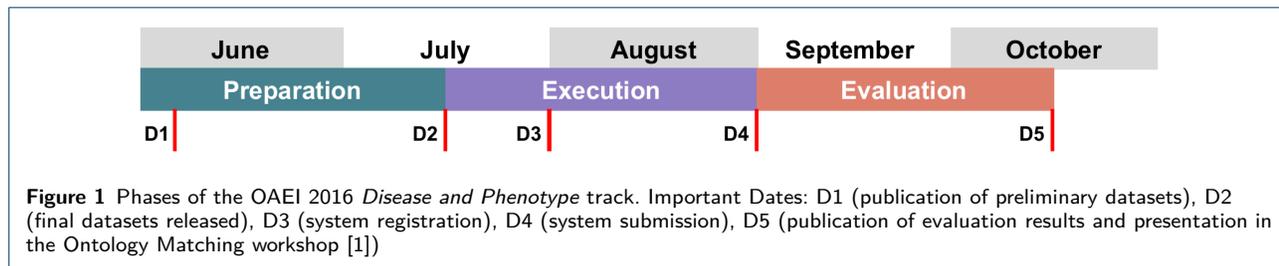
<sup>[13]</sup><http://oaei.ontologymatching.org/2016/seals-eval.html>

<sup>[7]</sup><http://oaei.ontologymatching.org/>

<sup>[8]</sup><http://oaei.ontologymatching.org/2016/phenotype/>

<sup>[9]</sup><http://data.bioontology.org/documentation#Mapping>

<sup>[10]</sup>[https://www.bioontology.org/wiki/index.php/BioPortal\\_Mappings](https://www.bioontology.org/wiki/index.php/BioPortal_Mappings)



**Table 1** Metrics of the track ontologies. Source: NCBI BioPortal on 2nd June 2016. Note that the metric “average number of children” excludes the leaf nodes.

Ontology	Number of axioms	Number of classes	Maximum depth	Avg. number of children
HP	137,289	11,786	15	3
MP	129,036	11,721	15	3
DOID	124,362	9,248	12	3
ORDO	188,991	12,936	11	16

### Algorithm 1 Steps followed in the evaluation

**Input:**  $\mathcal{O}_1, \mathcal{O}_2$ : ontologies in matching task;  $\mathcal{M}_m^{RA}$ : manually generated alignment; *Systems*: ontology matching systems participating in the task.

```

▷ Generation of system alignments with SEALS infrastructure:
1: for each Systemi in Systems do
2:    $\mathcal{M}_i^S \leftarrow \text{System}_i(\mathcal{O}_1, \mathcal{O}_2)$  ▷ Computes system alignment
3: end for
▷ Generation of consensus alignments:
4:  $\mathcal{M}_{c2}^{RA} \leftarrow \text{ConsensusAlignment}(\mathcal{M}_1^S \dots \mathcal{M}_n^S, 2)$  ▷ With vote 2
5:  $\mathcal{M}_{c3}^{RA} \leftarrow \text{ConsensusAlignment}(\mathcal{M}_1^S \dots \mathcal{M}_n^S, 3)$  ▷ With vote 3
▷ Aligned ontologies for consensus reference alignments:
6:  $\mathcal{O}^{\mathcal{M}_{c2}^{RA}} \leftarrow \mathcal{O}_1 \cup \mathcal{O}_2 \cup \mathcal{M}_{c2}^{RA}$  ▷ Repair  $\mathcal{M}_{c2}^{RA}$  if required
7:  $\mathcal{O}^{\mathcal{M}_{c3}^{RA}} \leftarrow \mathcal{O}_1 \cup \mathcal{O}_2 \cup \mathcal{M}_{c3}^{RA}$  ▷ Repair  $\mathcal{M}_{c3}^{RA}$  if required
▷ Evaluation for each system generated alignments:
8: for each  $\mathcal{M}_i^S$  in  $\mathcal{M}_1^S \dots \mathcal{M}_n^S$  do
▷ Aligned ontology for  $\mathcal{M}_i^S$ :
9:    $\mathcal{O}^{\mathcal{M}_i^S} \leftarrow \mathcal{O}_1 \cup \mathcal{O}_2 \cup \mathcal{M}_i^S$  ▷ Repair  $\mathcal{M}_i^S$  if required
▷ Results against consensus alignment with vote 2:
10:   $P_2 \leftarrow \text{SemanticPrecision}(\mathcal{M}_i^S, \mathcal{O}^{\mathcal{M}_{c2}^{RA}})$ 
11:   $R_2 \leftarrow \text{SemanticRecall}(\mathcal{M}_{c2}^{RA}, \mathcal{O}^{\mathcal{M}_i^S})$ 
▷ Results against consensus alignment with vote 3:
12:   $P_3 \leftarrow \text{SemanticPrecision}(\mathcal{M}_i^S, \mathcal{O}^{\mathcal{M}_{c3}^{RA}})$ 
13:   $R_3 \leftarrow \text{SemanticRecall}(\mathcal{M}_{c3}^{RA}, \mathcal{O}^{\mathcal{M}_i^S})$ 
▷ Results against manually generated alignment:
14:   $R_m \leftarrow \text{SemanticRecall}(\mathcal{M}_m^{RA}, \mathcal{O}^{\mathcal{M}_i^S})$ 
▷ Manual assessment of unique system mappings:
15:   $\mathcal{U}_i^S \leftarrow \text{UniqueMappings}(\mathcal{M}_i^S, \mathcal{O}^{\mathcal{M}_{c2}^{RA}})$ 
16:   $\{P_m, PC, NC\} \leftarrow \text{ManualAssessment}(\mathcal{U}_i^S)$ 
17: end for

```

following criteria for each of the matching tasks of the *Disease and Phenotype* track:

- Semantic precision and recall with respect to the *consensus alignments*.
- Semantic recall with respect to manually generated mappings.
- Manual assessment of unique mappings produced by a participant system.

Algorithm 1 formalizes the steps followed in the evaluation for each of the *Disease and Phenotype* matching tasks. The following subsections below comment on the main points of the evaluation process.

*Consensus alignments.* The *consensus alignments* are automatically generated based on the alignments produced by the participating systems in each of the matching tasks of the track. For the evaluation we have selected the *consensus alignments* of vote=2 (*i.e.*, mappings suggested by two or more systems) and vote=3 (*i.e.*, mappings suggested by three or more systems). In the case where both an equivalence and a subsumption mapping contribute to the consensus, the equivalence relationship prevails over the subsumption. The use of vote=2 and vote=3 was motivated by our experience in the creation of consensus alignments [21]. Consensus alignments with vote $\geq 4$  are typically highly precise but also very incomplete unless the number of contributing systems is significant.<sup>[14]</sup> Note that, when there are several systems of the same family (*i.e.*, systems participating with several variants), their (voted) mappings are only counted once in order to reduce bias.<sup>[15]</sup>

Note that consensus alignments have numerous limitations. It allows us to compare how the participating systems perform only in relation to each other. Some of the mappings in the consensus alignments may be erroneous (false positives), as it only requires 2 or 3 systems to agree on the erroneous mappings they find. Furthermore, the consensus alignments may not be complete, as there will likely be correct mappings that

<sup>[14]</sup>We may consider vote $\geq 4$  in future editions of the *Disease and Phenotype* track as the contributing participants increase.

<sup>[15]</sup>There could still be some bias through systems exploiting the same resource, e.g., UMLS.

no or only one system is able to find. Nevertheless, consensus alignments help to provide some insights into the performance of a matching system.

*Semantic precision and recall.* As introduced in the *preliminaries* section, the semantic precision and recall take into account the implicit knowledge derived from the ontologies and the mappings via OWL 2 reasoning.<sup>[16]</sup> Hence, the methods `SemanticPrecision` and `SemanticRecall` in Algorithm 1 receive as input a set of mappings  $\mathcal{M}$  and a coherent ontology  $\mathcal{O}$ . Both methods return as output the value of  $\frac{|\mathcal{M}'|}{|\mathcal{M}|}$  where  $\mathcal{M}'$  is a subset of  $\mathcal{M}$  such that the mappings  $m \in \mathcal{M}'$  are entailed by  $\mathcal{O}$  (*i.e.*,  $\mathcal{O} \models m$ ).

*Manually generated mappings.* These reference mappings were created through manual curation by eight disease informatics experts, who are authors of this paper, all working within or for the pharmaceutical industry for three areas of phenotype and disease; namely carbohydrate and glucose metabolism, obesity and breast cancer. These sets of reference mappings comprised of 29 pairwise mappings between HP and MP and 60 pairwise mappings between DOID and ORDO across the three areas. They included some relationships of equivalence, but most of them represented subsumption relationships. The three areas were selected as representative samples which were known already to be present across the four source ontologies. Inclusion of these manually defined mappings enabled a real-world evaluation of recall for the two matching tasks. The future editions of the track will increase the number of manual mappings through inclusion of additional areas relevant to the phenotype and disease domain.

*Unique mappings and manual assessment.* Unique mappings are mappings generated by an ontology matching system that have not been (explicitly) suggested by any of the other participating systems, nor entailed by the aligned ontology using the consensus alignment with `vote=2` ( $\mathcal{O}^{\mathcal{M}_{c2}^{RA}}$ ). The method `UniqueMappings` in Algorithm 1 receives as input a set of mappings  $\mathcal{M}$  and the (coherent) ontology  $\mathcal{O}^{\mathcal{M}_{c2}^{RA}}$  and returns as output  $\mathcal{M}'$  where  $\mathcal{M}' \subseteq \mathcal{M}$  such that the mappings  $m \in \mathcal{M}'$  are not entailed by  $\mathcal{O}^{\mathcal{M}_{c2}^{RA}}$  (*i.e.*,  $\mathcal{O}^{\mathcal{M}_{c2}^{RA}} \not\models m$ ).

Manual assessment over unique mappings has been performed by an expert in disease informatics from the pharmaceutical industry. This assessment aims at complementing the evaluation against the consensus alignments of those mappings that, although being

suggested or voted by only one matching system, may still be correct. We have focused the assessment on unique “equivalence” mappings and we have manually evaluated up to 30 mappings for each system in order to (roughly) estimate the percentage of correct mappings (*i.e.*, precision,  $P_m$  in Algorithm 1) and the positive/negative contribution to the total number of unique mappings ( $PC$  and  $NC$  in Algorithm 1), that is, the weight of the correct (*i.e.*, true positives) and incorrect (*i.e.*, false positives) mappings. Intuitively, the positive contribution (see Equation 4) of a system producing a small set of unique mappings will most likely be smaller than a system producing a larger set of unique (and mostly correct) mappings. The negative contribution (see Equation 5) will weight the number of incorrect unique mappings with respect to the total. Negative and positive contributions, for a set of unique mappings  $\mathcal{U}_i^S$  computed by a system  $i$ , are defined as follows:

$$\text{PositiveContribution}(\mathcal{U}_i^S) = \frac{|\mathcal{U}_i^S| \cdot \text{Precision}(\mathcal{U}_i^S)}{\sum_{j=1}^n |\mathcal{U}_j^S|} \quad (4)$$

$$\text{NegativeContribution}(\mathcal{U}_i^S) = \frac{|\mathcal{U}_i^S| \cdot (1 - \text{Precision}(\mathcal{U}_i^S))}{\sum_{j=1}^n |\mathcal{U}_j^S|} \quad (5)$$

## Results

We have run the evaluation of the *Disease and Phenotype* track in a Ubuntu Laptop with an Intel Core i7-4600U CPU @ 2.10GHz x 4 and allocating 15Gb of RAM. From the 13 systems registered to the track (out of 21 OAEI participants), 11 systems have been able to cope with at least one of the *Disease and Phenotype* matching tasks within a 24 hour time frame. Results for all OAEI tracks have been reported in [1].

### Participating systems

*AML* [23, 24] is an ontology matching system originally developed to tackle the challenges of matching biomedical ontologies. While its scope has since expanded, biomedical ontologies have remained one of the main drives behind its continued development. *AML* relies on the use of background knowledge and it also includes mapping repair capabilities.

*DiSMATCH* [25] estimates the similarity among concepts through textual semantic relatedness. *DiSMATCH* relies on a biomedical domain-adapted variant of a state-of-the-art semantic relatedness measure [26], which is based on Explicit Semantic Analysis.

<sup>[16]</sup>We rely on the OWL 2 reasoner *Hermit* [22].

*FCA-Map* [27] is an ontology matching system based on Formal Concept Analysis (FCA). FCA-Map attempts to push the envelope of the FCA to cluster the commonalities among classes at various levels.

*LogMap* [28, 29] relies on lexical and structural indexes to enhance scalability. It also incorporates approximate reasoning and repair techniques to minimise the number of logical errors in the aligned ontology.

*LogMapBio* [30] extends *LogMap* to use BioPortal [19] as a (dynamic) provider of mediating ontologies, instead of relying on a few preselected ontologies. *LogMapBio* retrieves the most suitable top-10 ontologies for the matching task.

*LogMapLt* is a “lightweight” variant of *LogMap*, which essentially only applies (efficient) string matching techniques.

*LYAM++* [31] is a fully automatic ontology matching system based on the use of external sources. *LYAM++* applies a novel orchestration of the components of the matching workflow [32].

*PhenomeNET* [33] alignment system comes in three flavours, which rely on three different versions of the PhenomeNET ontology [34]. PhenomeNET-Plain (PhenoMP) relies on a plain ontology which only uses the axioms provided by the HP ontology and the MP ontology. PhenomeNET-Map (PhenoMM) utilizes additional lexical equivalence axioms between HP and MP provided by BioPortal. Finally, PhenomeNET-Full (PhenoMF) relies on an extended version of the PhenomeNET ontology with equivalence mappings to the DOID and ORDO ontologies obtained via BioPortal and the AML matching system [23].

*XMap* [35] is a scalable matcher that implements parallel processing techniques to enable the composition of basic ontology matchers. It also relies on the use of external resources such as the UMLS Metathesaurus [36].

#### Use of specialised background knowledge

The use of (specialised) background knowledge is allowed in the OAEI, but participants are required to specify which sources their systems rely on to enhance the matching process. AML has three sources of background knowledge which can be used as mediators between the input ontologies: the Uber Anatomy Ontology [37] (Uberon), the Human Disease Ontology [5] (DOID) and the Medical Subject Headings<sup>[17]</sup> (MeSH). *LYAM++* also makes use of the Uberon ontology [37].

*LogMapBio* uses BioPortal [19] as dynamic mediating ontology provider, while *LogMap* uses normalisations and spelling variants from the general (biomedical) purpose UMLS Lexicon.<sup>[18]</sup> *XMAP* uses synonyms provided by the UMLS Metathesaurus [36]. Finally, PhenoMM, PhenoMF and PhenoMP rely on different versions of the PhenomeNET<sup>[19]</sup> ontology [34] with variable complexity as described above.

#### Evaluation against BioPortal (baseline) mappings

Table 2 shows the results in terms of semantic recall against the baseline mappings extracted from BioPortal as described in the *Methodology* Section (Preparation phase). In the DOID-ORDO task, *LYAM++* failed to complete the task while PhenoMM and PhenoMP produced empty mapping sets.

BioPortal mappings mostly represent correspondences with a high degree of lexical similarity and, as expected, most of the systems managed to produce alignments with a very high recall. DiSMATCH, *LYAM++*, PhenoMM (in the DOID-ORDO task) and PhenoMP were the exception and produced very low results with respect to the baseline mappings. As mentioned in the *Methodology* Section, since the BioPortal mappings were highly incomplete, the results in terms of (semantic) precision were not significant. For this reason, we needed to create consensus alignments for each task.

#### Creation of consensus alignments

In the MP-HP matching task 11 systems were able to produce mappings. Mappings voted by *LogMap* and PhenomeNET families were only counted once, and hence there were 7 independent system groups contributing to the consensus alignment. In the DOID-ORDO matching task 8 systems generated mappings and there were 6 independent system groups contributing to the consensus alignment.

Table 3 (resp. Table 4) shows the size of the different consensus alignments from vote=1, *i.e.*, mappings suggested by one or more system groups, to vote=7 (resp. vote=6), *i.e.*, mappings suggested by all system groups, in the HP-MP matching task (resp. DOID-ORDO task). It is noticeable that in the HP-MP task there were 0 mappings where all systems agreed, while in the DOID-ORDO task there were only 36. The number of mappings suggested by one system or more is specially large because PhenomeNET systems produce a large number of subsumption mappings. If only equivalence mappings of PhenomeNET systems are taken into account, the number of mappings with vote=1 would be 3,433 in the HP-MP task and 2,708 in the DOID-ORDO task.

<sup>[18]</sup><https://www.nlm.nih.gov/pubs/factsheets/umlslex.html>

<sup>[19]</sup><http://aber-owl.net/ontology/PhenomeNET>

<sup>[17]</sup><http://bioportal.bioontology.org/ontologies/MESH>

**Table 2** Recall against BioPortal (baseline) mappings

System	AML	DiSMATCH	FCA-Map	LYAM++	LogMap	LogMapBio	LogMapLt	PhenoMF	PhenoMM	PhenoMP	XMap
HP-MP	1.0	0.25	0.998	0.014	0.997	1.0	0.994	1.0	1.0	0.412	0.995
DOID-ORDO	0.993	0.048	0.984	-	0.942	0.950	0.943	0.994	0.0	0.0	0.967

**Table 3** Consensus alignments for the HP-MP matching task. Seven (family) system groups contributing.

Min. Votes	1	2	3	4	5	6	7
Mappings	217,039	<b>2,308</b>	<b>1,588</b>	1,287	677	152	0

**Table 4** Consensus alignments for the DOID-ORDO matching task. Six (family) system groups contributing.

Min. Votes	1	2	3	4	5	6
Mappings	50,998	<b>1,883</b>	<b>1,617</b>	1,447	991	36

As described in the *Methodology* Section we have selected the consensus alignments of vote=2 and vote=3. These consensus alignments for HP-MP contain 2,308 and 1,588 mappings, respectively; while for DOID-ORDO they include 1,883 and 1,617 mappings, respectively. Table 5 shows some examples of mappings included with the consensus alignments of vote=2 and vote=3. Also shown are some examples of manually created mappings and (correct/incorrect) unique mappings from ontology matching systems.

#### Results against consensus alignments

The union of the input ontologies together with the consensus alignments or the mappings computed by each of the systems was coherent and thus, we did not require to repair any of the mapping sets to calculate the semantic precision and recall. Note that the downloaded ontology versions from BioPortal did not contain any explicit or implicit disjointness. Tables 6 and 7 show the results achieved by each of the participating systems against the consensus alignments with vote=2 and vote=3. In the DOID-ORDO task, LYAM++, PhenoMM and PhenoMP failed to produce mappings and they were not included in Table 7.

We deliberately did not rank the systems since, as mentioned in the *Methodology* section, the consensus alignments may be incorrect or incomplete. We have simply highlighted the systems producing results relatively close to the consensus alignments. For example, in the HP-MP task, LogMap is the system producing an alignment that is closer to the mappings voted by at least 2 systems, while FCA-MAP produces results very close to the consensus alignments with vote=3.

The use of semantic precision and recall allowed us to provide a fair comparison for the systems PhenoMF, PhenoMM and PhenoMP. These systems discover a large set of subsumption mappings that are not explicit in the reference alignments, but they are still valid (*i.e.*, they are entailed by the aligned ontology using the reference alignment). For example, the standard precision of PhenoMF in the HP-MP task is 0.01 while the semantic precision reaches the value of 0.76.

Tables 6 and 7 also include the results of BioPortal mappings against the consensus alignments. Precision values are perfect, but recall is very low, which confirms our intuitions (recall *Preparation phase* section) about the incompleteness of BioPortal mappings.

It is striking how XMap and LogMapLt produced results very similar to the ones obtained by the BioPortal mappings. Closer scrutiny of these results showed us that the computed mappings were indeed very similar to the BioPortal mappings (*i.e.*, the F-measure of XMap and LogMapLt against the *baseline* mappings provided by BioPortal is  $\geq 0.95$  in both tasks). This could be expected for LogMapLt, since it only relies on simple string matching techniques as the matching system underlying BioPortal [38]. However, the results for XMap are unexpected since it produced top-results in the other biomedical-themed tracks of the OAEI 2016 [1].

#### Results against manually created mappings

Table 8 shows the results in terms of semantic recall against the manually created alignments. The results obtained in the HP-MP are relatively large positive values in general, especially for PhenoMF and PhenoMM that achieve a semantic recall of 0.90. The numbers for the DOID-ORDO, however, are much smaller values and only LogMap, LogMapBio and DiSMATCH are able to discover a few of the manually curated mappings. LogMapBio obtained the best semantic recall value with 0.17, which is far from the top results in the HP-MP task. The aforementioned results are also reflected when considering the consensus alignments. In the HP-MP task, both the consensus alignments with vote 2 and 3 obtained reasonably good results. However the picture changes dramatically in the DOID-ORDO task where none of the manually curated mappings are covered by the mappings agreed by 2 or more systems. The most likely explanation for this result is that the manual mappings for DOID-ORDO represent more complex subsumption mappings which were not possible to (semantically) derive for the other

**Table 5** Example mappings in the *Disease and Phenotype* track.

Entity 1	Entity 2	Rel.	Source
x-linked chondrodysplasia punctata (DOID_0060292)	chondrodysplasia punctata (Orphanet_93442)	≡	(only) consensus alignment vote=2
meningeal melanomatosis (DOID_8243)	diffuse leptomeningeal melanocytosis (Orphanet_252031)	≡	consensus alignment vote=3
reactive arthritis (DOID_6196)	reactive arthritis (Orphanet_29207)	≡	consensus alignment vote=3
hypoplastic scapulae (HP_0000882)	short scapula (MP_0004340)	≡	(only) consensus alignment vote=2
macrocytic anemia (HP_0001972)	macrocytic anemia (MP_0002811)	≡	consensus alignment vote=3
unerupted tooth (HP_0000706)	failure of tooth eruption (MP_0000121)	≡	consensus alignment vote=3
breast leiomyosarcoma (DOID_5285)	rare malignant breast tumor (Orphanet_180257)	⊂	manually created
abnormality of body weight (HP_0004323)	abnormal body weight (MP_0001259)	≡	manually created
microcephaly (HP_0000252)	decreased brain size (MP_0000774)	≡	AML unique mapping (correct)
skeletal dysplasia (HP_0002652)	abnormal skeletal muscle morphology (MP_0000759)	≡	AML unique mapping (incorrect)
carbohydrate metabolism disease (DOID_0050013)	disorder of carbohydrate metabolism (Orphanet_79161)	≡	LogMapBio unique mapping (correct)
spinocerebellar ataxia type 35 (DOID_0050982)	transglutaminase 6 (Orphanet_279644)	≡	LogMapBio unique mapping (incorrect)
female hypogonadism (HP_0000134)	small ovary (MP_0001127)	≡	PhenoMF unique mapping (correct)

**Table 6** Results against consensus alignments with vote=2 and vote=3 in the HP-MP task. Precision and Recall represent their semantic variants.

System-Mappings	Mappings	Precision-2	F-Measure-2	Recall-2	Precision-3	F-Measure-3	Recall-3
<i>BioPortal</i> (baseline)	639	1.00	0.50	0.33	1.00	0.60	0.43
AML	1,755	0.93	<b>0.86</b>	0.80	0.85	<b>0.90</b>	0.94
DiSMatch	644	0.55	0.30	0.21	0.45	0.28	0.20
FCA-Map	1,590	0.98	<b>0.85</b>	0.75	0.94	<b>0.93</b>	0.92
LYAM++	381	0.41	0.12	0.07	0.17	0.06	0.04
LogMap	2,011	0.94	<b>0.92</b>	0.91	0.77	<b>0.86</b>	0.97
LogMapBio	2,151	0.92	<b>0.92</b>	0.93	0.75	<b>0.85</b>	0.98
LogMapLt	667	1.00	0.51	0.34	1.00	0.62	0.45
PhenoMF	204,089	0.76	0.83	0.92	0.63	0.76	0.95
PhenoMM	198,149	0.77	0.83	0.91	0.64	0.76	0.94
PhenoMP	169,660	0.78	0.67	0.58	0.64	0.57	0.51
XMap	650	1.00	0.50	0.33	1.00	0.61	0.44

**Table 7** Results against disease alignments with vote=2 and vote=3 in the DOID-ORDO task. Precision and Recall represent their semantic variants.

System-Mappings	Mappings	Precision-2	F-Measure-2	Recall-2	Precision-3	F-Measure-3	Recall-3
<i>BioPortal</i> (baseline)	1,018	0.99	0.71	0.55	0.99	0.76	0.62
AML	2,098	0.85	<b>0.91</b>	0.97	0.78	<b>0.87</b>	1.00
DiSMatch	335	0.23	0.08	0.05	0.19	0.07	0.04
FCA-Map	1,803	0.97	<b>0.96</b>	0.96	0.89	<b>0.94</b>	0.99
LogMap	1,667	0.95	<b>0.91</b>	0.88	0.91	<b>0.92</b>	0.94
LogMapBio	1,804	0.92	<b>0.91</b>	0.90	0.86	<b>0.90</b>	0.95
LogMapLt	1,000	0.99	0.72	0.56	0.99	0.76	0.62
PhenoMF	40,612	0.95	<b>0.89</b>	0.83	0.95	<b>0.94</b>	0.92
XMap	1,030	0.98	0.72	0.57	0.98	0.77	0.63

mappings. Table 8 also shows the results for the BioPortal mappings, which, as expected, have a coverage of curated mappings very similar to the obtained by LogMapLt and XMap systems.

The use of semantic recall together with the standard measure, as in previous section, allowed us to provide more realistic results and a fair comparison with the PhenomeNET family systems. As it can be observed in the HP-MP task (Table 8), the standard recall, unlike the semantic recall, obtained by the other participants was very low and not comparable to the PhenomeNET family systems.

While the top performing algorithms were able to detect equivalence matches across whole source ontologies for the two mapping tasks giving high F-measures (Tables 6 and 7), it is clear from detection of the curated alignments that these proved much more difficult with a trend for lower semantic recall across both tasks

(Table 8). This result was not surprising because the curated alignments mostly comprised of subsumption relationships rather than equivalence. Table 5 shows two examples of curated mappings; the equivalence mapping between *abnormality of body weight* and *abnormal body weight* was suggested by at least one the systems, while the subsumption mapping between *breast leiomyosarcoma* and *rare malignant breast tumor* was not discovered by any of the systems.

#### Results for manual assessment of unique mappings

Tables 9 and 10 show the results of the manual assessment of the unique mappings generated by the participating systems. As mentioned in the *Methodology* section we manually analysed up to 30 unique equivalence mappings for each system to estimate the precision of the generated mappings not agreed with other systems. Table 5 shows examples of unique mappings

**Table 8** Results against curated alignments.

System-Mappings	HP-MP task		DOID-ORDO task	
	Standard Recall	Semantic Recall	Standard Recall	Semantic Recall
<i>BioPortal (baseline)</i>	0.17	0.52	0.00	0.00
AML	0.28	<b>0.76</b>	0.00	0.00
DiSMatch	0.07	0.14	0.02	<b>0.03</b>
FCA-Map	0.21	0.62	0.00	0.00
LYAM++	0.00	0.00	-	-
LogMap	0.24	0.66	0.02	<b>0.12</b>
LogMapBio	0.28	0.69	0.03	<b>0.17</b>
LogMapLt	0.17	0.52	0.00	0.00
PhenoMF	0.90	<b>0.90</b>	0.00	0.00
PhenoMM	0.90	<b>0.90</b>	-	-
PhenoMP	0.83	<b>0.83</b>	-	-
XMap	0.17	0.52	0.00	0.00
Consensus vote=1	0.90	<b>0.90</b>	0.05	<b>0.20</b>
Consensus vote=2	0.31	<b>0.79</b>	0.00	0.00
Consensus vote=3	0.24	0.66	0.00	0.00

computed by AML, LogMapBio and PhenoMF. Note that, we focus on equivalence mappings since PhenomeNET systems produce a large amount of (unique) subsumption mappings.

BioPortal mappings, as expected, contains a very low number of unique mappings in the DOID-ORDO task and no unique mappings in the HP-MP task.

It is noticeable in the HP-MP task that, although DiSMatch and LYAM++ produced very low results with respect to the consensus alignments (see Table 3), the positive contribution of their unique mappings is one of the highest. Nevertheless, their negative contribution has also an important weight. PhenomeNET systems produced the most precise set of unique mappings although their positive contribution was lower than other systems.

In the DOID-ORDO matching task, AML's unique mappings contains the higher number of true positives with a reasonable number of false positives. LogMapBio provided the best trade-off between positive and negative contribution.

The last row in Tables 9 and 10 shows (excluding BioPortal mappings) the total number of unique mappings, its (average) precision, and the total (aggregated) positive and negative contribution.

#### Results in the OAEI interactive matching track

The OAEI interactive track<sup>[20]</sup> aims at offering a systematic and automated evaluation of matching systems with user interaction to compare the quality of interactive matching approaches in terms of F-measure and number of required interactions. The interactive track relies on the datasets of the OAEI tracks: *Conferrence*, *Anatomy*, *Largebio*, and *Disease and Phenotype*; and it uses the reference alignments of each track as *oracle* in order to simulate the interaction with a domain expert with variable error rate [1].

<sup>[20]</sup><http://oaei.ontologymatching.org/2016/interactive/>

In this section we briefly present the results with the *Disease and Phenotype* datasets in the OAEI 2016 interactive track, which represents a side contribution of the work presented in this paper. For more details and results, the interested reader please refer to state-of-the-art papers on *interactive ontology alignment* [39, 40, 41, 1].

The consensus alignment with vote=3 was used as oracle in the *Disease and Phenotype* interactive track. Table 11 shows the obtained F-measure by AML and LogMap when simulating an interaction with a perfect user (*i.e.*, always gives the correct answer when asked about the validity of a mapping).<sup>[21]</sup> Both systems increase the F-measure with respect to the non-interactive results (see Tables 6 and 7) with a gain between 0.03 and 0.11. It is noticeable that the number of required requests by LogMap is around 4-5 times larger than AML.

#### Discussion

The OAEI has been proven to be an effective campaign to improve ontology matching systems. As a result, available techniques are more mature and robust. Nevertheless, despite the impressive state-of-the-art technology in ontology alignment, new matching tasks like those presented in this paper are very important for the OAEI campaign since they introduce new challenges to ontology alignment systems. For example, our preliminary tests with the *Disease and Phenotype* dataset revealed that only the 2015 versions of AML and LogMap, among the systems participating in the OAEI 2015, were able to cope with the track ontologies.

<sup>[21]</sup>From the *Disease and Phenotype* track participating systems only AML, LogMap and XMap implement an interactive algorithm. We have discarded XMap from the results since its number of oracle/user requests was very low in the *Disease and Phenotype* track.

**Table 9** Manual assessment of unique mappings and estimated positive and negative contribution in the HP-MP task.

System-Mappings	Unique Mappings	Precision	Positive Contrib.	Negative Contrib.
<i>BioPortal (baseline)</i>	0	-		
AML	122	0.87	8.63%	1.33%
DiSMATCH	291	0.83	<b>19.80%</b>	3.96%
FCA-Map	26	0.96	2.04%	0.08%
LYAM++	226	0.70	<b>12.91%</b>	5.53%
LogMap	130	0.93	9.90%	0.71%
LogMapBio	176	0.93	<b>13.40%</b>	0.96%
LogMapLt	0	-	-	-
PhenoMF	89	1.00	7.27%	<b>0.00%</b>
PhenoMM	85	1.00	6.94%	<b>0.00%</b>
PhenoMP	80	1.00	6.53%	<b>0.00%</b>
XMap	0	-	-	-
<b>Total</b>	<b>1,225</b>	<b>0.91</b>	<b>87.42%</b>	<b>12.58%</b>

**Table 10** Manual assessment of unique mappings and estimated positive and negative contribution in the DOID-ORDO task.

System-Mappings	Unique Mappings	Precision	Positive Contrib.	Negative Contrib.
<i>BioPortal (baseline)</i>	5	0.40		
AML	308	0.87	<b>30.40%</b>	4.68%
DiSMATCH	259	0.40	<b>11.80%</b>	17.70%
FCA-Map	61	0.83	5.79%	<b>1.16%</b>
LogMap	80	0.90	8.20%	<b>0.91%</b>
LogMapBio	144	0.97	<b>15.85%</b>	<b>0.55%</b>
LogMapLt	7	0.50	0.40%	0.40%
PhenoMF	3	1.00	0.34%	0.00%
XMap	16	0.56	1.03%	0.80%
<b>Total</b>	<b>878</b>	<b>0.75</b>	<b>73.81%</b>	<b>26.19%</b>

**Table 11** Results in the OAEI interactive track.

Task	System	F-measure	Gain	Requests
HP-MP	AML	0.93	0.03	388
	LogMap	0.97	0.11	1,928
DOID-ORDO	AML	0.96	0.09	413
	LogMap	0.99	0.07	1,602

In the OAEI 2016 campaign there were 11 systems that were able to produce results in at least one of the *Disease and Phenotype* matching tasks. The four systems: AML, FCA-Map, LogMap (and its *Bio* variant) and PhenoMF produced alignments relatively close to the consensus alignments for the *Disease and Phenotype* evaluation tasks as described in this paper. The results against curated alignments proved to be more challenging since they go beyond equivalent matches to include matches of semantic similarity, especially subsumption relationships. This finding suggests that while the systems performed well enough for detection of equivalent mappings, in future it would be good to improve their performance for detection of semantic similarity matches. For example, PhenomeNET systems showed potential advantage though exploiting a specialised background knowledge embedded within the system. LYAM++ is also specialised in the use of background knowledge, but it did not perform well in the *Disease and Phenotype* track, unlike in the OAEI Anatomy track, probably due to the lack of a suitable source of background knowledge for this track.

The OAEI also includes two biomedical-themed tracks, namely *Anatomy* and *Largebio* [1]. The complexity of the matching tasks is similar to the Anatomy track in terms of ontology size and expressiveness, while the Largebio tasks represent a significant leap in complexity with respect to the other OAEI test cases. The main differences with respect to the evaluation in the *Disease and Phenotype* track are the following: (i) we constructed two consensus reference alignments, unlike the Anatomy track where there exist a curated reference alignment [42] and the Largebio track where the reference alignment has been extracted from the UMLS Metathesaurus [9]; (ii) we performed an evaluation with respect to manually created mappings and a manual assessment of unique mappings produced by participating systems; and (iii) we used semantic precision and recall together with the standard measures.

The findings of the *Disease and Phenotype* evaluation show the potential of the top performing ontology matching systems that could help to automate the workflow of curators, who maintain ontology mapping services in numerous domains such as the disease and phenotype domain. Furthermore, the constructed

consensus alignments substantially improve available mapping sets provided by BioPortal.

## Conclusions

We have presented the methodology followed in the novel *Disease and Phenotype* track and the results in the OAEI 2016. The top systems in the track coped well with the detection of equivalence matches, but struggled to detect subsumption matches. This deserves more attention in the future development of ontology matching systems.

The Pistoia Alliance Ontologies Mapping project has gained much value from participation in the 2016 OAEI campaign through sponsorship and design of this new track on *Disease and Phenotype*. We believe that there is a real need for ontology matching algorithm developers to collaborate with ontology curators to improve the scale and quality of workflows necessary to build and maintain ontology mapping resources.

We are in an exploding information age with increasing amounts of human biology and genetics data in particular from sequencing technology improvements, biobanks and smart portable devices. This drives the need for stronger ontological standards, tools and services for ontology mapping to enable more efficient application of all this information. We expect that the *Disease and Phenotype* track will evolve in future campaigns as a strong use case which is widely applicable in the life sciences and beyond.

## Evolution of the track

The OAEI 2017 will include a new edition of the track, which will be composed by the same tasks as in 2016 (with updated ontology versions) and two additional tasks requiring the pairwise alignment of:

- HP and MESH (Medical Subject Headings) ontologies; and
- HP and OMIM (Online Mendelian Inheritance in Man) ontologies.

The alignment between HP and MESH is a new requirement of the Pistoia Alliance Ontologies Mapping project, while the mapping between HP and OMIM is placed within the scope of the Research Council of Norway project *BigMed* to improve the suggested genes and potential diagnosis associated to a given phenotype in state of the art tools like *PhenoTips* [43].

In the future editions of the *Disease and Phenotype* track, apart from including new datasets and updated versions, we aim to enhance the evaluation in a number of ways. We will consider *new metrics* like the mapping incoherence [13], the functional coherence [44] or the redundancy (minimality) [45] to evaluate the computed alignments. We also intend to redefine the *notion of semantic precision and recall*, using the semantic closure of the (aligned) ontologies, in order to

include the cases where the aligned ontology is incoherence (*i.e.*, contains unsatisfiable classes).

We plan to increase the number of *manually generated mappings* considering additional areas relevant to the phenotype and disease domain. In addition, we will also work towards the semi-automatic creation of *gold standard* reference alignments for the tasks by combining the consensus alignments and the manually generated mappings.

### List of abbreviations

**HP**: human phenotype ontology; **MP**: mammalian phenotype ontology; **DOID**: disease ontology; **ORDO**: orphanet rare disease ontology; **OA**: ontology alignment; **OAEI**: ontology alignment evaluation initiative; **DL**: description logics; **OWL**: web ontology language; **RDF**: resource description framework **RA**: reference alignment; **P**: precision; **R**: recall; **F**: f-measure; **O**: ontology; **M**: mappings;

### Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Availability of data and materials

OAEI 2016 datasets available from:

<http://oaei.ontologymatching.org/2016/phenotype/>.

OAEI 2017 datasets available from:

<http://oaei.ontologymatching.org/2017/phenotype/>.

Main entry point for the *Disease and Phenotype* track:

<http://sws.ifi.uio.no/oaei/phenotype/>

Competing interests

The authors declare that they have no competing interests.

Funding

This work was partially funded by the Pistoia Alliance Ontology Mappings project, the BIGMED project (IKT 259055), the HealthInsight project (IKT 247784), the SIRIUS Centre for Scalable Data Access (Research Council of Norway, project no.: 237889), the EU project Optique (FP7-ICT-318338), and the EPSRC projects ED3 and DBOnto.

Author's contributions

IH, EJR and AS organised and designed the experiments of the track. EJR conducted the automatic evaluation. IH prepared the manually curated mappings and performed the manual assessment which was checked by AS, MR, PW, SM, YAF and JM. All authors contributed to the writing of the manuscript. All authors read and approved the final manuscript.

Acknowledgements

We would like to thank the organisers and participants of the OAEI campaign. We also thank the anonymous reviewers for their comments and suggestions to improve the paper.

Author details

<sup>1</sup>Pistoia Alliance Ontologies Mapping Project, Pistoia Alliance Inc, USA.

<sup>2</sup>Department of Informatics, University of Oslo, Oslo, Norway. <sup>3</sup>Novartis,

Basel, Switzerland. <sup>4</sup>Roche Pharma Research and Early Development, pRED Informatics, Roche Innovation Center, Basel, Switzerland.

<sup>5</sup>GlaxoSmithKline R&D, Stevenage, UK. <sup>6</sup>BIOVIA 3DS, San Diego, USA.

<sup>7</sup>Eagle Genomics, Cambridge, UK. <sup>8</sup>OSTHUS, Aachen, Germany.

<sup>9</sup>FactBio, Cambridge, UK.

## References

- Achichi, M., Cheatham, M., Dragisic, Z., Euzenat, J., Faria, D., Ferrara, A., Flouris, G., Fundulaki, I., Harrow, I., Ivanova, V., Jiménez-Ruiz, E., Kuss, E., Lambrix, P., Leopold, H., Li, H., Meilicke, C., Montanelli, S., Pesquita, C., Saveta, T., Shvaiko, P., Splendiani, A., Stuckenschmidt, H., Todorov, K., dos Santos, C.T., Zamazal, O.: Results of the Ontology Alignment Evaluation Initiative 2016. In: 11th International Workshop on Ontology Matching (OM), pp. 73–129 (2016)
- Shvaiko, P., Euzenat, J., Jiménez-Ruiz, E., Cheatham, M., Hassanzadeh, O., Ichise, R. (eds.): Proceedings of the 11th International Workshop on Ontology Matching. *CEUR Workshop Proceedings*, vol. 1766 (2016)
- Köhler, S., *et al.*: The human phenotype ontology in 2017. *Nucleic Acids Research* **45**(D1), 865 (2017)
- Smith, C.L., Goldsmith, C.-A.W., Eppig, J.T.: The mammalian phenotype ontology as a tool for annotating, analyzing and comparing phenotypic information. *Genome Biology* **6**(1) (2004)
- Kibbe, W.A., *et al.*: Disease Ontology 2015 update: an expanded and updated database of human diseases for linking biomedical knowledge through disease data. *Nucleic Acids Research* **43**(D1) (2015)
- Cuenca Grau, B., Horrocks, I., Motik, B., Parsia, B., Patel-Schneider, P., Sattler, U.: OWL 2: The next step for OWL. *J. Web Semantics* **6**(4), 309–322 (2008)
- Horrocks, I., Kutz, O., Sattler, U.: The even more irresistible SROIQ. In: Tenth International Conference on Principles of Knowledge Representation and Reasoning, pp. 57–67 (2006)
- David, J., Euzenat, J., Scharffe, F., dos Santos, C.T.: The alignment API 4.0. *Semantic Web* **2**(1), 3–10 (2011)
- Jiménez-Ruiz, E., Cuenca Grau, B., Horrocks, I., Llavori, R.B.: Logic-based assessment of the compatibility of UMLS ontology sources. *J. Biomedical Semantics* **2**(S-1), 2 (2011)
- Euzenat, J.: Semantic precision and recall for ontology alignment evaluation. In: 20th International Joint Conference on Artificial Intelligence (IJCAI), pp. 348–353 (2007)
- Fleischhacker, D., Stuckenschmidt, H.: A Practical Implementation of Semantic Precision and Recall. In: 4th International Conference on Complex, Intelligent and Software Intensive Systems, pp. 986–991 (2010)
- David, J., Euzenat, J.: On fixing semantic alignment evaluation measures. In: 3rd International Workshop on Ontology Matching (OM) (2008)
- Meilicke, C.: Alignment incoherence in ontology matching. PhD thesis, University of Mannheim (2011). <https://ub-madoc.bib.uni-mannheim.de/29351>
- Jiménez-Ruiz, E., Meilicke, C., Grau, B.C., Horrocks, I.: Evaluating mapping repair systems with large biomedical ontologies. In: 26th International Workshop on Description Logics, pp. 246–257 (2013)
- Santos, E., Faria, D., Pesquita, C., Couto, F.M.: Ontology alignment repair through modularization and confidence-based heuristics. *PLoS One*. **10** (2015)
- Cheatham, M., Dragisic, Z., Euzenat, J., Faria, D., Ferrara, A., Flouris, G., Fundulaki, I., Granada, R., Ivanova, V., Jiménez-Ruiz, E., Lambrix, P., Montanelli, S., Pesquita, C., Saveta, T., Shvaiko, P., Solimando, A., dos Santos, C.T., Zamazal, O.: Results of the Ontology Alignment Evaluation Initiative 2015. In: 10th International Workshop on Ontology Matching (OM), pp. 60–115 (2015)
- Dragisic, Z., Eckert, K., Euzenat, J., Faria, D., Ferrara, A., Granada, R., Ivanova, V., Jiménez-Ruiz, E., Kempf, A.O., Lambrix, P., Montanelli, S., Paulheim, H., Ritze, D., Shvaiko, P., Solimando, A., dos Santos, C.T., Zamazal, O., Cuenca Grau, B.: Results of the ontology alignment evaluation initiative 2014. In: 9th International Workshop on Ontology Matching (OM), pp. 61–104 (2014)
- Shvaiko, P., Euzenat, J.: Ontology Matching: State of the Art and Future Challenges. *IEEE Trans. Knowl. Data Eng.* **25**(1), 158–176 (2013)
- Fridman Noy, N., Shah, N.H., Whetzel, P.L., Dai, B., Dorf, M., Griffith, N., Jonquet, C., Rubin, D.L., Storey, M.-A.D., Chute, C.G., Musen, M.A.: BioPortal: ontologies and integrated data resources at the click of a mouse. *Nucleic Acids Research* **37**(Web-Server-Issue) (2009)
- Faria, D., Jiménez-Ruiz, E., Pesquita, C., Santos, E., Couto, F.M.: Towards Annotating Potential Incoherences in BioPortal Mappings. In: 13th International Semantic Web Conference (ISWC), pp. 17–32 (2014)
- Jiménez-Ruiz, E., Grau, B.C., Horrocks, I.: Is my ontology matching system similar to yours? In: In 8th International Workshop on Ontology Matching (OM), pp. 229–230 (2013)
- Glimm, B., Horrocks, I., Motik, B., Stoilos, G., Wang, Z.: HermiT: An OWL 2 reasoner. *Journal of Automated Reasoning* **53**(3), 245–269 (2014)
- Faria, D., Pesquita, C., Santos, E., Palmonari, M., Cruz, I.F., Couto, F.M.: The AgreementMakerLight Ontology Matching System. In: OTM Conferences, pp. 527–541 (2013)
- Faria, D., Pesquita, C., Balasubramani, B.S., Martins, C., Cardoso, J., Curado, H., Couto, F.M., Cruz, I.F.: OAEI 2016 results of AML. In: 11th International Workshop on Ontology Matching (OM), pp. 138–145 (2016)
- Rybinski, M., del Mar Roldán García, M., García-Nieto, J., Montes, J.F.A.: Dismatch results for OAEI 2016. In: 11th International Workshop on Ontology Matching (OM), pp. 161–165 (2016)
- Gabrilovich, E., Markovitch, S.: Computing semantic relatedness using wikipedia-based explicit semantic analysis. In: Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI), pp. 1606–1611 (2007)
- Zhao, M., Zhang, S.: Identifying and validating ontology mappings by formal concept analysis. In: 11th International Workshop on Ontology Matching (OM), pp. 61–72 (2016)
- Jiménez-Ruiz, E., Cuenca Grau, B.: LogMap: Logic-based and Scalable Ontology Matching. In: Int'l Sem. Web Conf. (ISWC), pp. 273–288 (2011)
- Jiménez-Ruiz, E., Cuenca Grau, B., Zhou, Y., Horrocks, I.: Large-scale Interactive Ontology Matching: Algorithms and Implementation. In: 20th European Conference on Artificial Intelligence (ECAI), pp. 444–449 (2012)
- Chen, X., Xia, W., Jiménez-Ruiz, E., Cross, V.V.: Extending an ontology alignment system with biportal: a preliminary analysis. In: Proceedings of the ISWC 2014 Posters & Demonstrations Track, pp. 313–316 (2014)
- Tigrine, A.N., Bellahsene, Z., Todorov, K.: LYAM++ results for OAEI 2016. In: 11th International Workshop on Ontology Matching (OM), pp. 196–200 (2016)
- Tigrine, A.N., Bellahsene, Z., Todorov, K.: Light-Weight Cross-Lingual Ontology Matching with LYAM++. In: On the Move to Meaningful Internet Systems: OTM 2015 Conferences, pp. 527–544 (2015)
- Rodríguez-García, M.Á., Gkoutos, G.V., Schofield, P.N., Hoehndorf, R.: Integrating phenotype ontologies with PhenomeNET. In: Proceedings of the 11th International Workshop on Ontology Matching (OM), pp. 201–209 (2016)
- Hoehndorf, R., Schofield, P.N., Gkoutos, G.V.: PhenomeNET: a whole-phenome approach to disease gene discovery. *Nucleic Acids Research* **39**(18), 119 (2011)
- Djeddi, W.E., Khadir, M.T., Yahia, S.B.: XMap results for OAEI 2016. In: 11th International Workshop on Ontology Matching (OM), pp. 222–226 (2016)
- Bodenreider, O.: The unified medical language system (UMLS): integrating biomedical terminology. *Nucleic acids research* **32**(Database issue), 267–270 (2004)
- Mungall, C.J., Torniai, C., Gkoutos, G.V., Lewis, S.E., Haendel, M.A.: Uberon, an integrative multi-species anatomy ontology. *Genome biology* **13**(1) (2012)
- Ghazvinian, A., Noy, N.F., Musen, M.A.: Creating mappings for ontologies in biomedicine: Simple methods work. In: American Medical Informatics Association Annual Symposium (AMIA) (2009)
- Paulheim, H., Hertling, S., Ritze, D.: Towards evaluating interactive ontology matching tools. In: 10th Extended Semantic Web Conference (ESWC), pp. 31–45 (2013)
- Ivanova, V., Lambrix, P., Åberg, J.: Requirements for and evaluation of user support for large-scale ontology alignment. In: 12th Extended Semantic Web Conference (ESWC), pp. 3–20 (2015)
- Dragisic, Z., Ivanova, V., Lambrix, P., Faria, D., Jiménez-Ruiz, E., Pesquita, C.: User validation in ontology alignment. In: 15th International Semantic Web Conference (ISWC), pp. 200–217 (2016)

42. Bodenreider, O., Hayamizu, T.F., Ringwald, M., de Coronado, S., Zhang, S.: Of Mice and Men: Aligning Mouse and Human Anatomies. In: American Medical Informatics Association Annual Symposium (AMIA) (2005)
43. Girdea, M., Dumitriu, S., Fiume, M., Bowdin, S., Boycott, K.M., Chénier, S., Chitayat, D., Faghfoury, H., Meyn, M.S., Ray, P.N., So, J., Stavropoulos, D.J., Brudno, M.: PhenoTips: Patient Phenotyping Software for Clinical and Research Use. *Human Mutation* **34**(8), 1057–1065 (2013)
44. Schlicker, A., Domingues, F.S., Rahnenführer, J., Lengauer, T.: A new measure for functional similarity of gene products based on gene ontology. *BMC Bioinformatics* **7**(1), 302 (2006)
45. Giunchiglia, F., Maltese, V., Autayeu, A.: Computing minimal mappings between lightweight ontologies. *Int. J. on Digital Libraries* **12**(4), 179–193 (2012). doi:10.1007/s00799-012-0083-2