# Introducing the Disease and Phenotype OAEI Track\*

Ian Harrow<sup>1</sup>, Ernesto Jimenez-Ruiz<sup>2</sup>, Andrea Splendiani<sup>1</sup>, Martin Romacker<sup>1</sup>, Stefan Negru<sup>1</sup>, Peter Woollard<sup>1</sup>, Scott Markel<sup>1</sup>, Yasmin Alam-Faruque<sup>1</sup>, Martin Koch<sup>1</sup>, Erfan Younesi<sup>1</sup> and James Malone<sup>1</sup>

<sup>1</sup> Pistoia Alliance Ontologies Mapping Project, Pistoia Alliance Inc. USA <sup>2</sup> Department of Computer Science, University of Oxford, UK

## 1 Introduction

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 existing ontologies mapping tool providers to enable their evaluation.<sup>2</sup> A critical component of any Ontologies Mapping tool is the embedded ontology matching algorithm, therefore the project is supporting their development and evaluation through sponsorship and organisation of the new *Disease and Phenotype* track for the OAEI campaign<sup>3</sup> [1] which is described in this paper.

#### 2 Datasets

The *Disease and Phenotype* track<sup>4</sup> comprises two tasks that will involve the pairwise alignment of the HPO, MP, DOID and ORDO ontologies (Table 1 shows the metrics of these ontologies):

- Task 1: matching of the Human Phenotype Ontology (HPO) to the Mammalian Phenotype Ontology (MP).
- Task 2: matching of the Human Disease Ontology (DOID) to the Orphanet and Rare Diseases Ontology (ORDO).

The first task is important for translational science where HPO includes inherited diseases and MP originated from rodents as a model mammalian organism for many laboratory studies, including gene knock out. The second task includes representation of rare human diseases in both ontologies which are of fundamental importance for understanding how genetic variation can cause disease. Currently, such mappings are mostly curated by bioinformatics and disease experts who would benefit from automation supported by implementation of ontology matching algorithms into their workflows.

We have extracted a "baseline" reference alignments for the track based on the available BioPortal mappings [2] which are considered as a baseline since they are incomplete and may contain errors.

<sup>\*</sup> We have also submitted a 4-pages paper about the Pistoia Alliance Ontologies Mapping Project to the ISWC 2016 posters and demos track.

<sup>&</sup>lt;sup>1</sup> http://www.pistoiaalliance.org/projects/ontologies-mapping

<sup>&</sup>lt;sup>2</sup> https://pistoiaalliance.atlassian.net/wiki/display/PUB/Ontologies+ Mapping+Resources

<sup>&</sup>lt;sup>3</sup> http://oaei.ontologymatching.org/2016/

<sup>&</sup>lt;sup>4</sup> http://oaei.ontologymatching.org/2016/phenotype/description.html

Table 1. Metrics of the track ontologies. Source: NCBI BioPortal on 19th Aug 2016

| Ontology | Number of classes | Maximum depth | Avg. number of children |
|----------|-------------------|---------------|-------------------------|
| HPO      | 15,319            | 15            | 3                       |
| MP       | 11,720            | Undisclosed   | Undisclosed             |
| DOID     | 10,905            | 12            | 3                       |
| ORDO     | 13,105            | 11            | 16                      |

## **3** Evaluation process

The evaluation of the Disease and Phenotype Track will be run with support of the SEALS infrastructure.<sup>5</sup> Systems will be evaluated and ranked according to the following criteria:

- Precision and Recall with respect to a voted reference alignment that will be built automatically to generate consensus voting for the outputs of the participating systems.
- Recall with respect to manually generated mappings for three areas (carbohydrate, obesity and breast cancer).
- Manual assessment of a subset of the generated mappings, specially the ones that are not suggested by other systems.
- Performance in other tracks will also be taken into account, especially the OAEI interactive track [3] where the *Disease and Phenotype* dataset is also used.<sup>6</sup>

Additionally, systems able to discover complex logic relations in mappings beyond equivalence and subsumption will also be considered. The evaluation of these mappings will be in parallel to the evaluation of standard equivalence and subsumption mappings. Complex mappings should be provided in OWL 2 format.

## Acknowledgements

This work was partially funded by the Pistoia Alliance Ontology Mappings project, the EU project Optique (FP7-ICT-318338), and the EPSRC projects ED3 and DBOnto.

## References

- 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: Proceedings of the 10th International Workshop on Ontology Matching. (2015) 60–115
- 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)
- Dragisic, Z., Ivanova, V., Lambrix, P., Faria, D., Jimenez-Ruiz, E., Pesquita, C.: User validation in ontology alignment. In: International Semantic Web Conference. (2016)

<sup>&</sup>lt;sup>5</sup> http://oaei.ontologymatching.org/2016/seals-eval.html

<sup>&</sup>lt;sup>6</sup> http://oaei.ontologymatching.org/2016/interactive/index.html