

AROA Results for OAEI 2020*

Lu Zhou and Pascal Hitzler

DaSe Lab, Kansas State University, Manhattan KS 66506, USA
{luzhou, hitzler}@ksu.edu

Abstract. This paper introduces the results of an ontology alignment system named Association Rule-based Ontology Alignment (AROA) in the Ontology Alignment Evaluation Initiative (OAEI) 2020 campaign. This ontology alignment system focuses on producing simple and complex alignment between ontologies that are populated with instance data. This is the second participation of AROA in the OAEI campaign, and it produces the best performance in terms of relaxed F-measure on two benchmarks in complex track, which are populated GeoLink and populated Enslaved.

1 Presentation of the system

1.1 State, purpose, general statement

AROA (Association Rule-based Ontology Alignment) system aims to automatically generate simple and complex alignment between two and more ontologies. These ontologies are required to have shared common instance data because AROA relies on association rule mining and requires these instances as input to discover interesting relations. After generating a set of association rules, AROA utilizes the simple and complex correspondence patterns that have been widely accepted in the Ontology Matching community [4, 5] to further narrow a large number of rules down to more meaningful ones and finally establishes the alignments.

1.2 Specific techniques used

Figure 1 illustrates the overview of AROA alignment system. In this section, we introduce each step of AROA alignment system along with some concepts that we frequently use in the AROA system, such as association rule mining, FP-growth algorithm, and complex alignment generation.

Clean Triple. First, AROA extracts all triples as the format of ⟨Subject, Predicate, Object⟩ from the source and target ontologies. Each item in a triple is expressed as a web URI. After collecting all of the triples, we clean the data based on the following criteria: we only keep the triples that contain at least one entity under the source or the target ontology namespace or the triples contain `rdf:type` information, as our algorithm relies on this information.

*Copyright © 2020 for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

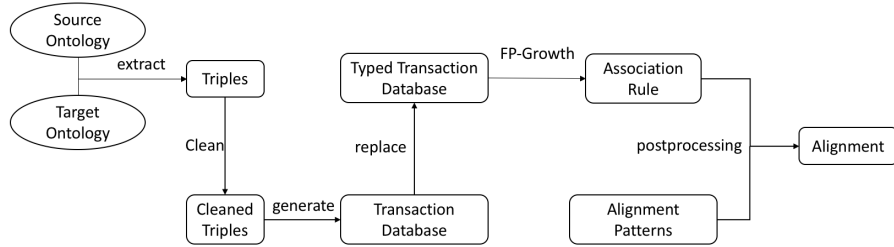


Fig. 1. Overview of AROA Alignment System

Generate Transaction Database. After the filtering process, we generate the transaction database as the input for the FP-growth algorithm. Let $I = \{i_1, i_2, \dots, i_n\}$ be a set of distinct attributes called items. Let $D = \{t_1, t_2, \dots, t_m\}$ be a set of transactions where each transaction in D has a unique transaction ID and contains a subset of the items in I . Table 1 shows a list of transactions corresponding to a list of triples. Instance data can be displayed as a set of triples, each consisting of subject, predicate, and object. Here, subjects represent the identifiers and the set of corresponding properties with the objects represent transactions, which are separated by the symbol “|”. I.e., a transaction is a set $T = (s, Z)$ such that s is a subject, and each member of Z is a pair (p, o) of a property and an object such that (s, p, o) is an instance triple.

Generate Typed Transaction Database. Then we replace the object in the triples with its `rdf:type`¹ because we focus on generating schema-level (rather than instance-level) mapping rules between two ontologies, and the type

¹If there are multiple types of the object, it can also combine the subject and predicate as additional information to determine the correct type, or keep both types as two triples.

Table 1. Triples and Corresponding Transactions

s_1	p_1	o_1		
s_1	p_2	o_2		
s_1	p_4	o_4		
s_2	p_1	o_1	s_1	$p_1 o_1, p_2 o_2, p_4 o_4$
s_2	p_2	o_2	s_2	$p_1 o_1, p_2 o_2, p_3 o_3, p_4 o_4$
s_2	p_3	o_3	s_3	$p_1 o_1, p_2 o_2$
s_2	p_4	o_4		
s_3	p_1	o_1		
s_3	p_2	o_2		

Table 2. Original Transaction Database

TID	Itemsets
x_1	<code>gbo:hasAward y₁, gmo:fundedBy y₂</code>
x_2	<code>gbo:hasFullName y₃, gmo:hasPersonName y₄</code>
x_3	<code>rdf:type gbo:Cruise, rdf:type gmo:Cruise</code>

Table 3. Typed Transaction Database

TID	Itemsets
x_1	<code>gbo:hasAward gbo:Award, gmo:fundedBy gmo:FundingAward</code>
x_2	<code>gbo:hasFullName xsd:string, gmo:hasPersonName gmo:PersonName</code>
x_3	<code>rdf:type gbo:Cruise, rdf:type gmo:Cruise</code>

information of the object is more meaningful than the original URI. If an object in a triple has `rdf:type` of a class in ontology, we replace the URI of the object with its class. If the object is a data value, the URI of the object is replaced with the datatype. If the object already is a class in ontology, it remains unchanged. Tables 2 and 3 show some examples of the conversion.

Generate Association Rules. Our alignment system mainly depends on a data mining algorithm called association rule mining, which is a rule-based machine learning method for discovering interesting relations between variables in large databases [3]. Many algorithms for generating association rules have been proposed, like Apriori [1] and FP-growth algorithm [2]. In this paper, we use FP-growth to generate association rules between ontologies, since the FP-growth algorithm has been proven superior to other algorithms [2]. The FP-growth algorithm is run on the transaction database in order to determine which combinations of items co-occur frequently. The algorithm first counts the number of occurrences of all individual items in the database. Next, it builds an FP-tree structure by inserting these instances. Items in each instance are sorted by descending order of their frequency in the dataset so that the tree can be processed quickly. Items in each instance that do not meet the predefined thresholds, such as minimum support and minimum confidence (see below for these terms), are discarded. Once all large itemsets have been found, the association rule creation begins. Every association rule is composed of two sides. The left-hand-side is called the antecedent, and the right-hand-side is the consequent. These rules indicate that whenever the antecedent is present, the consequent is likely to be

Table 4. Examples of Association Rules

Antecedent	Consequent
$p_4 o_4, p_1 o_1$	$p_2 o_2$
$p_2 o_2$	$p_1 o_1$
$p_4 o_4$	$p_1 o_1$

Table 5. The Alignment Pattern Types Covered in AROA System

Pattern	Category
Class Equivalence	1:1
Class Subsumption	1:1
Property Equivalence	1:1
Property Subsumption	1:1
Class by Attribute Type	1:n
Class by Attribute Value	1:n
Property Typecasting Equivalence	1:n
Property Typecasting Subsumption	1:n
Typed Property Chain Equivalence	m:n
Typed Property Chain Subsumption	m:n

as well. Table 4 shows some examples of association rules generated from the transaction database in Table 1.

Generate Alignment. AROA utilizes some simple and complex correspondences that have been widely accepted in Ontology Matching community to further filter rules [4, 5] and finally generate the alignments. There are a total of 10 different types of correspondences that AROA covers this year. Table 5 lists all the simple and complex alignment correspondences and corresponding categories. Since the association rule mining might generate a large number of rules, in order to narrow the association rules down to a smaller set, AROA follows these patterns to generate corresponding alignments. For example, Class by Attribute Type (CAT) is a classic complex alignment pattern. This type of pattern was first introduced in [4]. It states that a class in the source ontology is in some relationship to a complex construction in the target ontology. This complex construction may comprise an object property and its range. Class C_1 is from ontology O_1 , and object property op_1 and its range t_1 are from ontology O_2 .

Association Rule format: $\text{rdf:type}|C_1 \rightarrow \text{op}_1|t_1$

Example: $\text{rdf:type}|gbo:PortCall \rightarrow gmo:atPort|gmo:Place$

Generated Alignment: $gbo:PortCall(x) \rightarrow gmo:atPort(x, y) \wedge gmo:Place(y)$

In this example, this association rule implies that if the subject x is an individual of class $gbo:PortCall$, then x is subsumed by the domain of $gmo:atPort$ with its range $gmo:Place$. The equivalence relationship can be generated by combining another association rule holding the reverse information. Other simple and complex alignments are also generated by following the same steps.

1.3 Adaptations made for the evaluation

AROA is an instance-based ontology alignment system. Therefore, AROA embeds Apache Jena Fuseki server in the system. The ontologies are first downloaded from the SEALS repository. And then, AROA uploads and stores the

Table 6. The Number of Alignments Found on Populated GeoLink Benchmark

Alignment Patterns	Category	Reference Alignment	AROA	
			# of Correct Entities	# of Correct Relation
-	-	-		
Class Equiv.	1:1	10	10	10
Class Subsum.	1:1	2	1	0
Property Equiv.	1:1	7	5	5
Property Typecasting Subsum.	1:n	5	3	0
Property Chain Equiv.	m:n	26	15	13
Property Chain Subsum.	m:n	17	7	0

ontologies in the embedded Fuseki server, which might take some time for this step to load large-size ontology pairs.

2 Results

This year, AROA alignment system evaluates its performance on the populated GeoLink benchmark [5, 6] and populated Enslaved benchmark [7]. In the populated GeoLink benchmark, there are 19 simple mappings, including 10 class equivalence, 2 class subsumption, and 7 property equivalence. And there are 48 complex mappings, including 5 property subsumption, 26 property chain equivalence, and 17 property chain subsumption. In the populated Enslaved benchmark, 15 simple mappings are all class equivalences. And there are 83 complex mappings, including 68 property chain equivalence and 15 property chain subsumption. Table 6 and Table 7 list the alignment patterns and categories in the populated GeoLink and populated Enslaved Benchmark with the results of AROA system. We list the numbers of identified mappings for each pattern. There are two dimensions that we can look into the details to understand the performance. The first dimension is the entity identification, which means, given an entity in the source ontology, the system should be able to generate related entities in the target ontology. Another dimension is relationship identification, in which the system should detect the correct relationship between these entities, such as equivalence and subsumption. Therefore, we list the number of correct entities and the number of correct relationships in order to understand the strengths and weaknesses of the system. For example, In the Table 6, AROA correctly identifies all 1:1 class equivalence including entity and relationship. AROA also finds one class subsumption alignment, which is the class *PortCall* in the GeoLink Base Ontology (GBO) is related to the class *Fix* in the GeoLink Modular Ontology (GMO). However, it outputs the relationship between

Table 7. The Number of Alignments Found on Populated Enslaved Benchmark

Alignment Patterns	Category	Reference Alignment	AROA	
			# of Correct Entities	# of Correct Relation
-	-	-		
Class Equiv.	1:1	15	11	11
Property Chain Equiv.	m:n	68	29	29
Property Chain Subsum.	m:n	15	3	0

Table 8. The Performance Comparison on Populated GeoLink and Populated Enslaved Benchmarks

Matcher	Populated GeoLink						Populated Enslaved					
	(1:1)	(1:n)	m:n	Relaxed_Precision	Relaxed_F-measure	Relaxed_Recall	(1:1)	(1:n)	m:n	Relaxed_Precision	Relaxed_F-measure	Relaxed_Recall
Reference Alignment	19	5	43	-	-	-	15	0	83	-	-	-
AMLC	13	0	0	0.50	0.32	0.23	12	0	18	0.73	0.40	0.28
AROA	15	3	22	0.87	0.60	0.46	11	0	32	0.80	0.51	0.38
CANARD	15	2	17	0.89	0.54	0.39	3	0	16	0.42	0.19	0.13

PortCall and *Fix* as equivalence, which it should be subsumption. Therefore, we count the number of correct entities as 1 and the number of correct relations as 0. This criterion is also applied to other patterns. In the Table 7, AROA detects 73% (11 out of 15) of the simple class equivalences and 38% (32 out of 83) of the complex mappings in the populated Enslaved benchmark. In addition, we compare the performance of AROA against other complex alignment systems in Table 8. AMLC, AROA, and CANARD are only three systems can produce complex relations on the complex benchmarks. AROA found the highest number of complex alignments and achieved the best performance in terms of relaxed recall and relaxed f-measure on both benchmarks.²

3 General comments

From the performance comparison, AMLC, AROA, and CANARD can generate almost correct complex alignment, which means some alignments found by these two systems may not be completely correct, but it can be easily improved by semi-automated fashion. For example, the system can produce correct entities that should be involved in a complex alignment, but it doesn't output the correct relationship. Another possible situation is that the system can detect the correct relationship but fails to find all the entities. Based on these situations, we will investigate the incorrect alignments and improve the algorithm to find the relationship and entities as accurately as possible.

4 Conclusions

This paper introduces the AROA ontology alignment system and its preliminary results in the OAEI 2020 campaign. This year, AROA evaluates its performance on populated GeoLink and populated Enslaved benchmarks and achieves the best performance in terms of relaxed recall and relaxed f-measure among the three complex alignment systems. We will continue to evaluate AROA on other benchmarks and improve the algorithm in the near future.

5 Acknowledgement

This work has been supported by the National Science Foundation under Grant No. 2033521, KnowWhereGraph: Enriching and Linking Cross-Domain Knowl-

²<http://oaei.ontologymatching.org/2020/results/complex/geolink/index.html>

edge Graphs using Spatially-Explicit AI Technologies and the Andrew W. Mellon Foundation through the Enslaved project (identifiers 1708-04732 and 1902-06575).

References

1. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules in large databases. In: Bocca, J.B., Jarke, M., Zaniolo, C. (eds.) VLDB'94, Proc. of 20th International Conference on Very Large Data Bases, September 12-15, 1994, Santiago de Chile, Chile. pp. 487–499. Morgan Kaufmann (1994), <http://www.vldb.org/conf/1994/P487.PDF>
2. Han, J., et al.: Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Min. Knowl. Discov.* **8**(1), 53–87 (2004). <https://doi.org/10.1023/B:DAMI.0000005258.31418.83>, <https://doi.org/10.1023/B:DAMI.0000005258.31418.83>
3. Piatetsky-Shapiro, G.: Discovery, analysis, and presentation of strong rules. In: *Knowledge Discovery in Databases*, pp. 229–248. AAAI/MIT Press (1991)
4. Ritzke, D., Meilicke, C., Sváb-Zamazal, O., Stuckenschmidt, H.: A pattern-based ontology matching approach for detecting complex correspondences. In: Shvaiko, P., Euzenat, J., Giunchiglia, F., Stuckenschmidt, H., Noy, N.F., Rosenthal, A. (eds.) *Proceedings of the 4th International Workshop on Ontology Matching (OM-2009) collocated with the 8th International Semantic Web Conference (ISWC-2009)* Chantilly, USA, October 25, 2009. CEUR Workshop Proceedings, vol. 551. CEUR-WS.org (2009), <http://ceur-ws.org/Vol-551/om2009.Tpaper3.pdf>
5. Zhou, L., et al.: A complex alignment benchmark: Geolink dataset. In: Vrandečić, D., Bontcheva, K., Suárez-Figueroa, M.C., Presutti, V., Celino, I., Sabou, M., Kaffee, L., Simperl, E. (eds.) *The Semantic Web - ISWC 2018 - 17th International Semantic Web Conference, Monterey, CA, USA, October 8-12, 2018, Proceedings, Part II*. Lecture Notes in Computer Science, vol. 11137, pp. 273–288. Springer (2018). https://doi.org/10.1007/978-3-030-00668-6_17, https://doi.org/10.1007/978-3-030-00668-6_17
6. Zhou, L., Cheatham, M., Krisnadhi, A., Hitzler, P.: Geolink data set: A complex alignment benchmark from real-world ontology. *Data Intell.* **2**(3), 353–378 (2020). https://doi.org/10.1162/dint_a_00054, https://doi.org/10.1162/dint_a_00054
7. Zhou, L., Shimizu, C., Hitzler, P., Sheill, A.M., Estrecha, S.G., Foley, C., Tarr, D., Rehberger, D.: The enslaved dataset: A real-world complex ontology alignment benchmark using wikibase. In: d’Aquin, M., Dietze, S., Hauff, C., Curry, E., Cudré-Mauroux, P. (eds.) *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*. pp. 3197–3204. ACM (2020). <https://doi.org/10.1145/3340531.3412768>, <https://doi.org/10.1145/3340531.3412768>