

Aligning ontologies in the context of large size and number

Imen Megdiche ¹

In collaboration with:

A. Laadhar¹, F. Ghazzi ³, F. Ravat ¹, P. Roussille ¹, O. Teste ¹,
C. Trojhan ²

¹SIG/IRIT

²MELODI/IRIT

³MIRACL-SFAX

MAY 27, 2019

Outline

- 1 Context and Motivation
- 2 Alignment of Large Ontologies
 - A Local Matching Learning Approach
 - Evaluation Results
- 3 Alignment of Numerous Ontologies
 - Holistic Matching
 - Holistic Reference Alignments
 - LPHOM
- 4 Future directions and Open Questions

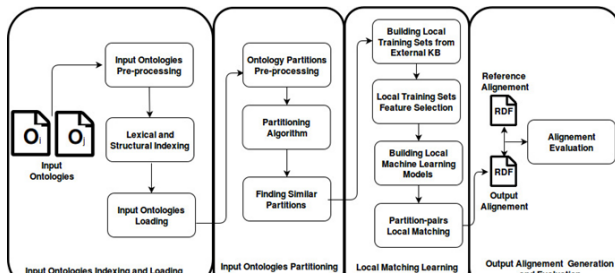
Context and Motivation

- Ontology matching is a key task in the management of semantically heterogeneous sources.
- Match large size resources
 - Large size ontologies refer to ontologies composed of ten of thousands of concepts. In the Biomedical domain : SNOMED (122,464), FMA (78,989) and NCI (66,724) classes.
 - State-of-the-art Ontology Matching Systems lack the automation of aligning large biomedical ontologies, e.g., Similarity measures or thresholds.
- Match more than 2 resources
 - A very rich state of the art approaches design to deal with pair of ontologies : *pairwise ontology matching*.
 - **BUT**, the increasing number of data producers (like IOT sensors, linked open data sources, etc) generating massive and heterogeneous data: need solutions able to handle more than two resources.

Aligning Large Biomedical Ontologies

Proposed Solution

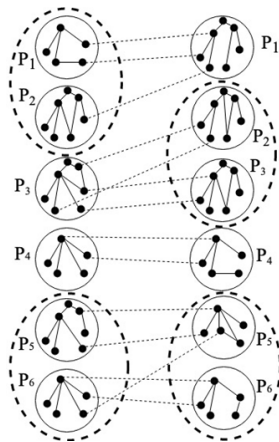
- **Large size ontologies** → 'Divide and conquer' strategy: generate small partitions and resolve matching on pairwise small partitions
- **Automate the settings for pairwise small partitions** → Local matching learning : use machine learning to find best local tuning; thresholds, choice of similarity measures and their weights



Proposed Solution: Ontology Partitionning (Laadhar et al., 2019a)

- 1 Apply Hierarchical Agglomerative clustering for each input ontology : generate a set of partitions
- 2 Using Cross-referencing (external knowledge resources : Uberon) + cross-searching to generate anchors
- 3 Merge the set of partitions based on anchors : not large + isolated-partitions

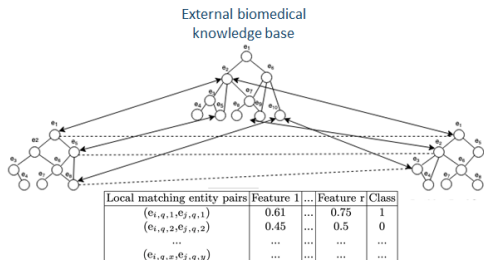
⇒ Each pair of partitions is a local matching task.



Proposed Solution : Matching Learning (Laadhar et al. 2019b)

Automate the generation and settings of matching pairwise partitions

- Automatic generation of a local training set for each local matching task without reference alignments
 - Using EKB resources (uberon)
 - Element and structural level features
- Resampling \Rightarrow balanced training data
- Apply feature selection



Evaluation results on LargeBio Dataset OAEI'2018

Participation in OAEI'2017 (PopMAP), OAEI'2018 (PopMAP++)
: third rank in Anatomy track.

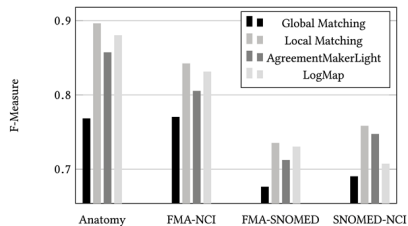
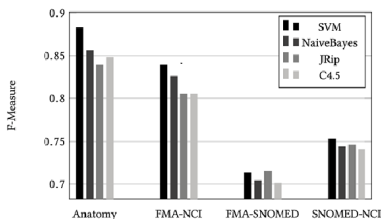


Figure: Local Matching Learning :
Comparing Different ML Algorithms

Figure: Comparing Local Matching
to the state-of-the-art matching
systems

Amir Laadhar (2019a), Amir Laadhar (2019b)

Aligning Large number of Ontologies

Holistic Matching

- Holistic data integration approaches able to integrate many data sources together Rahm (2016).
 - Clustering-based approaches : Toni Gruetze (2012) Saleem et al. (2008)
- From our point of view:
 - Holistic matching should deal simultaneously with the matching of different data sources (ontologies)
 - The matching process extends the ontology pairwise matching using a set $\Omega = \{O_1, \dots, O_N\}$ of ontologies with $N \geq 2$.



Holistic Matching

- Some developed tools: LPHOM (Megdiche et al. (2016)) OAEI'2016, Holontology (Roussille et al. 2018b) OAEI'2018.
- Some approaches: PORSCHE (Saleem et al. (2007)), HCM (Grütze et al. (2012)).

⇒ few works compared to pairwise matching

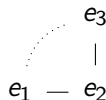
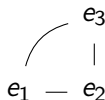
Bottlenecks: How to boost researches without validating these approaches What about holistic reference alignments ?



Holistic Matching: Alignments

For $N = 3$, each correspondence c_i is defined as a triple correspondence $\langle \{e_1, e_2, e_3\}, \equiv, n \rangle$ where $e_1 \in O_1$, $e_2 \in O_2$ and $e_3 \in O_3$:

- clique correspondence is $\langle \{e_1, e_2, e_3\}, \equiv, 1 \rangle$
- clique-relaxed correspondence is $\langle \{e_1, e_2, e_3\}, \equiv, \frac{2}{3} \rangle$



Our proposal (Roussille et al. 2018)

A *pseudo-holistic* approach to build holistic alignments from available pairwise alignments in two steps:

- 1 build a graph G_H of all combinations of correspondences in existing pairwise alignments;
- 2 build the holistic alignments according to different levels of relaxation with respect to complete graphs (cliques):
 - clique-strict level : search complete sub-graphs from N input ontologies.
 - clique-relaxed level: two proposed methods

Our proposal Philippe Roussille (2018)

Two methods for clique-relaxed level:

- Method 1 is a direct relaxation of cliques with N nodes.
- Method 2 take all combinations of input graphs $\in [0, N]$. We select only one tuple of nodes based on the intra-ontology relations and the best confidence value of *clique_likeness*. *clique_likeness* computes the confidence of the clique-relaxed subgraph, which is the geometric distance of a subgraph compared to a clique

$$\text{clique_likeness}(G_i) = \frac{2 * |E_i|}{|V_i| * (|V_i| - 1)}$$



General Idea (DEXA 2015, RCIS 2015)

- Reduce ontology matching problem, for 1:1 alignments, to the maximum weighted graph matching problem (MWGM)
- Extend MWGM problem with a set of linear constraints expressing OM requirements.

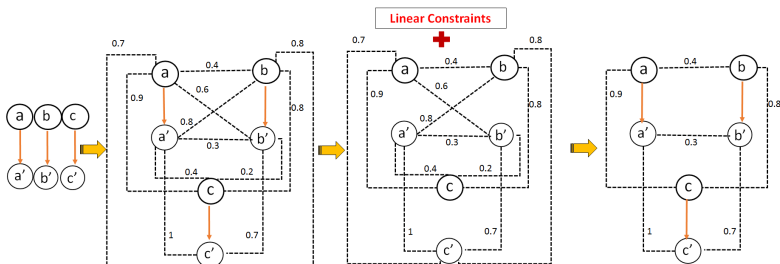


Figure: LPHOM to MWGM

LPHOM Model

$$\max \sum_{i=1}^{N-1} \sum_{j=i+1}^N \sum_{k=1}^{nbC_i, nbC_j} \sum_{l=1}^{nbOP_i, nbOP_j} sim_{ik,jl} x_{ik,jl} + \sum_{m=1}^{nbOP_i, nbOP_j} \sum_{n=1}^{nbDP_i, nbDP_j} sim_{im,jn} y_{im,jn} + \sum_{q=1}^{nbDP_i, nbDP_j} \sum_{r=1}^{nbDP_i, nbDP_j} sim_{iq,jr} z_{iq,jr}$$

$$s.t. \sum_{l=1}^{nbC_j} x_{ik,jl} \leq 1, \forall k \in [1, nbC_i] \quad (C1 \text{ Classes})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\sum_{n=1}^{nbOP_j} y_{im,jn} \leq 1, \forall m \in [1, nbOP_i] \quad (C1 \text{ Object Properties})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\sum_{r=1}^{nbDP_j} z_{iq,jr} \leq 1, \forall q \in [1, nbDP_i] \quad (C1 \text{ Data Properties})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$x_{ik,jl} + x_{i'k',jl} \leq 1 \quad (C2 \text{ Classes})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\forall k, k' \in [1, nbC_i], \forall l \in [1, nbC_j]$$

$$y_{im,jn} + x_{i'k',jl} \leq 1 \quad (C2 \text{ Object Properties})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\forall m, m' \in [1, nbOP_i], \forall n \in [1, nbOP_j]$$

$$z_{iq,jr} + x_{i'k',jl} \leq 1 \quad (C2 \text{ Data Properties})$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\forall q, q' \in [1, nbDP_i], \forall r \in [1, nbDP_j]$$

$$y_{im,jn} \leq x_{i'k',jl} + x_{i'k',jr} \quad (C3)$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\forall m \in [1, nbOP_i], \forall n \in [1, nbOP_j]$$

$$\forall k', k'' \in [1, nbC_i], \forall l', l'' \in [1, nbC_j]$$

$$z_{iq,jr} \leq x_{i'k',jl} \quad (C4)$$

$$\forall i \in [1, N-1], j \in [i+1, N]$$

$$\forall q \in [1, nbDP_i], \forall r \in [1, nbDP_j]$$

- **Objectif Function**
maximizes the profit (similarities) of the selected alignments.
- **Binary Decision Variables**
Classes $x_{ik,jl}$, Object properties $y_{im,jn}$ and Data Properties $z_{iq,jr}$
- Four types of **Linear Constraints**.



Conclusions and Open Questions

- Aligning Large Ontologies : Experiment with different domain datasets, which KBR to use ? Automatic selection of KBR ?
- Holistic matching on large ontologies and how to generate reference alignments?
- Holistic Matching Learning ?
- Holistic Matching : how to gather expert reference alignments to complement the holistic reference alignment ?

Bibliography I

- Amir Laadhar, I. M. F. R. O. T. F. G., F. Ghozzi. (2019a). The impact of imbalanced training data on local matching learning of ontologies. In *International conference on business information systems (bis 2019), seville, spain* (pp. –).
- Amir Laadhar, I. M. F. R. O. T. F. G., F. Ghozzi. (2019b). Partitioning and local matching learning of large biomedical ontologies. In *SAC 2019, limassol, cyprus* (pp. 2285–2292).
- Grütze, T., Böhm, C., & Naumann, F. (2012). Holistic and scalable ontology alignment for linked open data. In *WWW2012*.

Bibliography II

- Megdiche, I., Teste, O., & dos Santos, C. T. (2016). LPHOM results for OAEI 2016. In *(ISWC2016), kobe, japan* (pp. 190–195).
- Philippe Roussille, O. T. C. T., Imen Megdiche. (2018). Boosting holistic ontology matching: Generating graph clique-based relaxed reference alignments for holistic evaluation. In *EKAW 2018, proceedings* (pp. 355–369).
- Rahm, E. (2016). The case for holistic data integration. In *ADBIS 2016* (pp. 11–27).
- Saleem, K., Bellahsene, Z., & Hunt, E. (2007). Performance oriented schema matching. In *DEXA 2007* (pp. 844–853).

Bibliography III

- Saleem, K., Bellahsene, Z., & Hunt, E. (2008). PORSCHE: Performance ORiented SCHEMA mediation. *Information Systems*, 33(7-8), 637-657.
- Toni Gruetze, F. N., Christoph Böhm. (2012, 4). Holistic and scalable ontology alignment for linked open data. In *Proceedings of the 5th linked data on the web workshop at the 21th international world wide web conference*.