



UNIVERSITY
OF TRENTO

DIPARTIMENTO DI INGEGNERIA E SCIENZA DELL'INFORMAZIONE

38050 Povo – Trento (Italy), Via Sommarive 14
<http://www.disi.unitn.it>

COMPUTING MINIMAL MAPPINGS

Fausto Giunchiglia, Vincenzo Maltese and Aliaksandr Autayeu

December 2008

Technical Report # DISI-08-078

Also: at the ISWC Ontology Matching Workshop (OM 2009), 25th
October 2009, Washington DC, USA

Computing minimal mappings

Fausto Giunchiglia, Vincenzo Maltese, Aliaksandr Autayeu

Dipartimento di Ingegneria e Scienza dell'Informazione (DISI)
Università di Trento
{fausto, maltese, autayeu}@disi.unitn.it

Abstract. Given two classifications, or lightweight ontologies, we compute the minimal mapping, namely the subset of all possible correspondences, called mapping elements, between them such that i) all the others can be computed from them in time linear in the size of the input ontologies, and ii) none of them can be dropped without losing property i). In this paper we provide a formal definition of minimal mappings and define a time efficient computation algorithm which minimizes the number of comparisons between the nodes of the two input ontologies. The experimental results show a substantial improvement both in the computation time and in the number of mapping elements which need to be handled.

Keywords: Knowledge Organization Systems, ontology matching, minimal mappings

1 Introduction

Given any two graph-like structures, e.g., database and XML schemas, classifications, thesauri and ontologies, matching is usually identified as the problem of finding those nodes in the two structures which semantically correspond to one another. Any such pair of nodes, along with the semantic relationship holding between the two, is what we informally call a *mapping element*.

In the last few years a lot of work has been done on this topic both in the digital libraries [15, 20, 21, 22] and the computer science [2, 3, 4, 5, 6, 8, 9] communities, but the problem is still largely unsolved even if substantial progress has been made [6]. This difficulty is not unexpected because a solution to the matching problem would amount to solving the semantic heterogeneity problem at the level of schematic metadata (e.g., ontologies).

In this paper we focus on a specific sub-problem of matching, namely that of finding *minimal mappings*, that is, the subset of all possible correspondences, called *mapping elements*, such that i) all the others can be computed from them in time linear in the size of the input graphs, and ii) none of them can be dropped without losing property i). We concentrate on lightweight ontologies, as formally defined in [7], namely on acyclic graph structures where each node is labelled by a natural language sentence which can be used to compute the meaning of that node as a Description Logic (DL) formula, and where the formula associated to each node is subsumed by the formula of the node above. We assume that each mapping element is associated with one of

the following semantic relations: disjointness (\perp), equivalence (\equiv), more specific (\sqsupseteq) and less specific (\sqsubseteq), as computed for instance by semantic matching [5].

The main advantage of minimal mappings is that they are the minimal amount of information that needs to be dealt with. Notice that this is a rather important feature as the number of possible mapping elements can grow up to $n*m$ with n and m being the size of the two input ontologies. In particular, minimal mappings become crucial with large ontologies, e.g., DMOZ, with 10^5 - 10^6 nodes, where even relatively small subsets of the number (10^{12}) of possible mapping elements are unmanageable. Minimal mappings provide clear usability advantages. Many systems and corresponding interfaces, mostly graphical, have been provided for the management of mappings but all of them hardly scale with the increasing number of nodes, and the resulting visualizations are rather messy [3]. Furthermore, the maintenance of smaller sets makes the work of the user much easier, faster and less error prone [11].

The main contributions of this paper are a formal definition of *minimal* and, dually, *redundant mappings*, evidence of the fact that the minimal mapping always exists and it is unique and an algorithm for computing it. This algorithm has the following main features:

1. It can be proved to be correct and complete, in the sense that it always computes the minimal mapping;
2. It is very efficient as it minimizes the number of calls to the node matching function, namely to the function which computes the relation between two nodes. Notice that node matching in the general case amounts to logical reasoning (i.e., SAT reasoning) [5], and it may require exponential time;
3. As also described in the evaluation section (Section 5), our algorithm computes the mapping of maximum size (including redundant elements) as it maximally exploits the information codified in the graph of the lightweight ontologies in input. This, in turn, avoids missing mapping elements due to pitfalls in the node matching functions, for instance as a consequence of missing background knowledge [8].

As far as we know very little work has been done on the problem of minimal mappings. In the context of digital libraries, Marshall and Madhusudan [17] use a syntactic element level matcher (based on string matching) to match nodes of two concept maps. This technique is proposed as a valid tool to support educational processes (e.g. to match student maps against an instructor one). The solution proposed is an adaptation of the Similarity Flooding (SF) algorithm which is known to be very accurate in the results. Unfortunately, experiments conducted on it [18] show that SF is very poor in recall and results one of the worst in terms of F-measure. Kovacs and Micsik [16] propose semantic matching as a valid support for information retrieval. In their approach the query is matched against a collection of documents which are both described in form of logical expressions (DNF clauses). Using a set of heuristics, their element level matcher identifies relevant documents in a reasonable amount of time. Unfortunately no evaluation is provided in terms of precision and recall. Hence, it is not possible to appreciate the quality of the proposed solution. In general the computation of minimal mappings can be seen as a specific instance of the mapping inference problem [4]. Closer to our work, in [9, 10, 11] the authors use Distributed Description Logics (DDL) [12] to represent and reason about existing ontology

mappings. They introduce a few debugging heuristics which remove mapping elements which are redundant or generate inconsistencies from a given set [10]. The main problem of this approach, as also recognized by the authors, is the complexity of DDL reasoning [11]. In our approach, instead of pruning redundant elements, we directly compute the minimal set. Among other things, our approach allows us to minimize the number of calls to node matching.

The paper is organized as follows. Section 2 provides a motivating example. Section 3 provides the definition for redundant and minimal mappings, and it shows that the minimal set always exists and it is unique. Section 4 describes the algorithm while Section 5 evaluates it. Finally, Section 6 draws some conclusions and outlines the future work.

2 A motivating example

Classifications are perhaps the most natural tool humans use to organize information content. Information items are hierarchical arranged under topic nodes moving from general ones to more specific ones as long as we go deep in the hierarchy. This attitude is well known in Knowledge Organization as the principle of organizing from the general to the specific [19], called synthetically the *get-specific principle* in [1, 7, 23]. Consider the two fragments of classifications depicted in Fig. 1. They are designed to arrange more or less the same content, but from different perspectives. The second is a fragment taken from the Yahoo web directory¹ (category Computers and Internet).

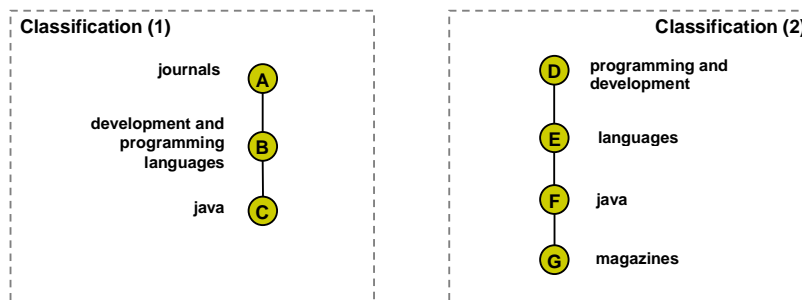


Fig. 1. Two classifications

Following the approach described in [7] and exploiting dedicated NLP techniques tuned to short phrases (for instance, as described in [13]), each node label can be translated in an unambiguous, propositional DL expression. The resulting formulas are reported in Fig. 2. Here each string denotes a concept (e.g., journals#1) and the numbers at the end of the strings denote a specific concept constructed from a WordNet sense. Notice that the formula associated to each node contains the formula of the node above to capture the fact that the meaning of each node is contextualized by the meaning of its ancestor nodes. As a consequence, the backbone structure of the resulting lightweight ontologies is represented by subsumption relations between nodes.

¹<http://dir.yahoo.com/>

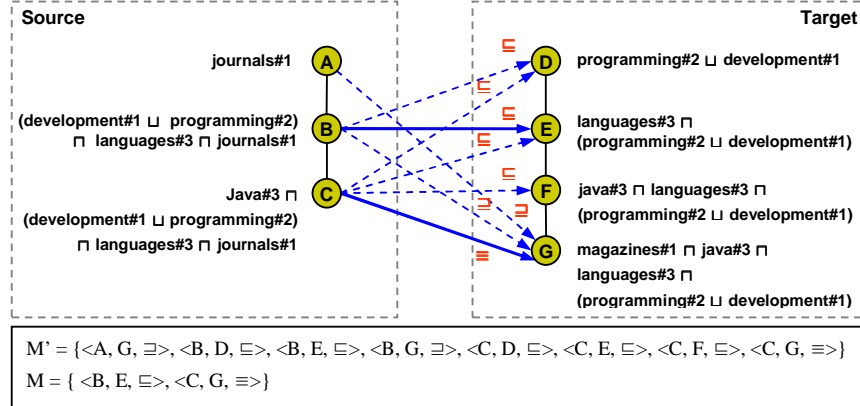


Fig. 2. The minimal and redundant mapping between two lightweight ontologies

Fig. 2 also reports the resulting mapping elements. Notice however that not all the mapping elements have the same semantic valence. For instance, $B \supset D$ is a trivial consequence of $B \supset E$ and $E \supset D$, and similarly for $C \supset F$ and $C \supset G$. We represent the elements in the minimal mapping using solid lines and redundant elements using dashed lines. M' is the set of maximum size (including the maximum number of redundant elements) while M is the minimal. The problem is how to compute the minimal set in the most efficient way.

3 Redundant and minimal mappings

Adapting the definition in [7] we define a lightweight ontology as follows:

Definition 1 (Lightweight ontology). A lightweight ontology O is a rooted tree $\langle N, E, L^F \rangle$ where:

- N is a finite set of nodes;
- E is a set of edges on N ;
- L^F is a finite set of labels expressed in a Propositional DL language such that for any node $n_i \in N$, there is one and only one label $l_i^F \in L^F$;
- $l_{i+1}^F \sqsubseteq l_i^F$ with n_i being the parent of n_{i+1} .

The superscript F is used to emphasize that labels are in a formal language. Fig. 2 above provides an example of (a fragment of) two lightweight ontologies. We then define mapping elements as follows:

Definition 2 (Mapping element). Given two lightweight ontologies O_1 and O_2 , a mapping element m between them is a triple $\langle n_1, n_2, R \rangle$, where:

- $n_1 \in N_1$ is a node in O_1 , called the source node;
- $n_2 \in N_2$ is a node in O_2 , called the target node;
- $R \in \{ \equiv, \sqsubseteq, \supset, \perp \}$ is the strongest semantic relation holding between n_1 and n_2 .

The partial order is such that disjointness is stronger than equivalence which, in turn, is stronger than subsumption (in both directions), and such that the two subsumption symbols are unordered. This in order to return subsumption only when equivalence does not hold or one of the two nodes being inconsistent (this latter case generating at the same time both a disjointness and a subsumption relation), and similarly for the order between disjointness and equivalence. Notice that, under this ordering, there can be at most one mapping between two nodes.

The next step is to define the notion of redundancy. The key idea is that, given a mapping element $\langle n_1, n_2, R \rangle$, a new mapping element $\langle n_1', n_2', R' \rangle$ is redundant with respect to the first if the existence of the second can be asserted simply by looking at the relative positions of n_1 with n_1' , and n_2 with n_2' . In algorithmic terms, this means that the second can be computed without running the time expensive node matching functions. We have identified four basic redundancy patterns as follows:

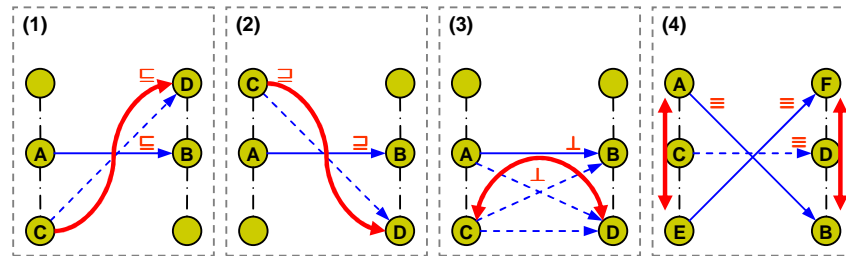


Fig. 3. Redundancy detection patterns

In Fig. 3, the blue dashed mappings are redundant w.r.t. the solid blue ones. The bold red solid lines show how a semantic relation propagates. Let us discuss the rationale for each of the patterns:

- **Pattern (1):** each mapping element $\langle C, D, \sqsupseteq \rangle$ is redundant w.r.t. $\langle A, B, \sqsupseteq \rangle$. In fact, C is more specific than A which is more specific than B which is more specific than D. As a consequence, by transitivity C is more specific than D.
- **Pattern (2):** dual argument as in pattern (1).
- **Pattern (3):** each mapping element $\langle C, D, \perp \rangle$ is redundant w.r.t. $\langle A, B, \perp \rangle$. In fact, we know that A and B are disjoint, that C is more specific than A and that D is more specific than B. This implies that C and D are also disjoint.
- **Pattern (4):** Pattern 4 is the combinations of patterns (1) and (2).

Notice that patterns (1) and (2) are still valid in case we substitute subsumption with equivalence. However, in this case we cannot exclude the possibility that a stronger relation holds between C and D. A trivial example of where this is not the case is provided in Fig. 4.

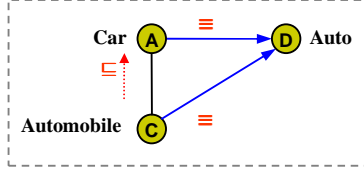


Fig. 4. Two non redundant mappings

On the basis of the patterns and the considerations above we can define redundant elements as follows. Here $\text{path}(n)$ is the path from the root to the node n .

Definition 3 (Redundant mapping element). Given two lightweight ontologies O_1 and O_2 , a mapping M and a mapping element $m' \in M$ with $m' = \langle C, D, R' \rangle$ between them, we say that m' is redundant in M iff one of the following holds:

- (1) If R' is \sqsubseteq , $\exists m \in M$ with $m = \langle A, B, R \rangle$ and $m \neq m'$ such that $R \in \{\sqsubseteq, \equiv\}$, $A \in \text{path}(C)$ and $D \in \text{path}(B)$;
- (2) If R' is \supseteq , $\exists m \in M$ with $m = \langle A, B, R \rangle$ and $m \neq m'$ such that $R \in \{\supseteq, \equiv\}$, $C \in \text{path}(A)$ and $B \in \text{path}(D)$;
- (3) If R' is \perp , $\exists m \in M$ with $m = \langle A, B, \perp \rangle$ and $m \neq m'$ such that $A \in \text{path}(C)$ and $B \in \text{path}(D)$;
- (4) If R' is \equiv , conditions (1) and (2) must be satisfied.

See how Definition 3 maps to the four patterns in Fig. 3. Fig. 2 in Section 2 provides examples of redundant elements. Definition 3 can be proved to capture all and only the cases of redundancy.

Theorem 1 (Redundancy, soundness and completeness). Given a mapping M between two lightweight ontologies O_1 and O_2 , a mapping element $m' \in M$ is redundant if and only if it satisfies one of the conditions of Definition 3.

The soundness argument is the rationale described for the patterns above. Completeness can be shown by constructing the counterargument that we cannot have redundancy in the remaining cases. We can proceed by enumeration, negating each of the patterns, encoded one by one in the conditions appearing in the Definition 3. The complete proof is given in the appendix. Fig. 5 provides some examples of non redundancy. The first, based on pattern (1), tells us that the existence of a link between two nodes does not necessarily propagate to the two nodes below. For example we cannot derive that $\text{Canine} \sqsubseteq \text{Dog}$ from the set of axioms $\{\text{Canine} \sqsubseteq \text{Mammal}, \text{Mammal} \sqsubseteq \text{Animal}, \text{Dog} \sqsubseteq \text{Animal}\}$, and it would be wrong to do so. The second, based on pattern (3), shows that disjointness cannot propagate to the target (or to the source) one level up. For example we cannot derive that $\text{Dog} \perp \text{Animal}$ only from $\{\text{Dog} \sqsubseteq \text{Canine}, \text{Cat} \sqsubseteq \text{Animal}, \text{Canine} \perp \text{Cat}\}$. The third example, based on pattern (4), tells us that we cannot derive equivalence if the source node C or target D is not between

the source and target nodes of the two equivalence mappings. Notice that, by chance, the other equivalence mapping holds.

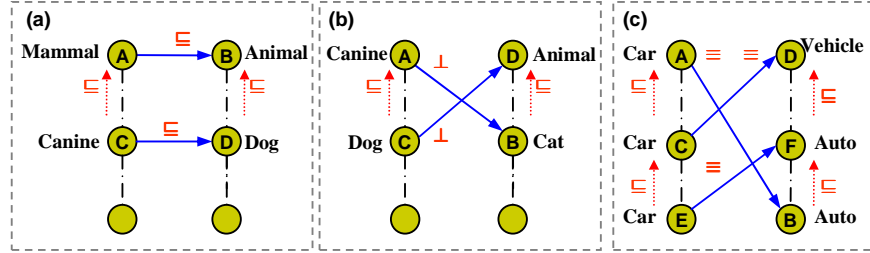


Fig. 5. Some examples of non redundant mapping elements

The notion of redundancy allows us to formalize the notion of minimal mapping as follows:

Definition 4 (Minimal mapping). Given two lightweight ontologies O_1 and O_2 , we say that a mapping M between them is minimal iff:

- a) $\nexists m \in M$ such that m is redundant (minimality condition);
- b) $\nexists M' \supset M$ satisfying condition a) above (maximality condition).

A mapping element is minimal if it belongs to the minimal mapping.

Note that conditions (a) and (b) ensure that the minimal set is the set of maximum size with no redundant elements. As an example, the set M in Fig. 2 is minimal. Comparing this mapping with M' we can observe that all elements in the set $M' - M$ are redundant and that, therefore, there are no other supersets of M with the same properties. In effect, $\langle A, G, \supseteq \rangle$ and $\langle B, G, \supseteq \rangle$ are redundant w.r.t. $\langle C, G, \equiv \rangle$ for pattern (2); $\langle C, D, \equiv \rangle$, $\langle C, E, \equiv \rangle$ and $\langle C, F, \equiv \rangle$ are redundant w.r.t. $\langle C, G, \equiv \rangle$ for pattern (1); $\langle B, D, \equiv \rangle$ is redundant w.r.t. $\langle B, E, \equiv \rangle$ for pattern (1). Note that M contains far less mapping elements w.r.t. M' .

As last observation, for any two given lightweight ontologies, the minimal mapping always exists and it is unique.

Theorem 2 (Minimal mapping, existence and uniqueness). Given two lightweight ontologies O_1 and O_2 , there is always one and only one minimal mapping between them.

A proof is given in the appendix. Keeping in mind the patterns in Fig. 3, the minimal set can be efficiently computed using the following key intuitions:

1. Equivalence can be “opened” into two subsumption mapping elements;
2. Taking any two paths in the two ontologies, a minimal subsumption mapping element (in both directions of subsumption) is an element with the lowest node in one path whose formula is subsumed by the other node and the highest node in the other path which subsumes the formula in the other node.

3. Taking any two paths in the two ontologies, a minimal disjointness mapping element is the one with the highest nodes in both paths such that their formulas satisfy disjointness.

4 Computing minimal and redundant mappings

The patterns described in the previous section allow us not only to identify minimal and redundant mapping elements, but they also suggest how to significantly reduce the amount of calls to the node matchers. By looking for instance at pattern (2) in Fig. 3, given a mapping element $m = \langle A, B, \exists \rangle$ we know in advance that it is not necessary to compute the semantic relation holding between A and any descendant C in the sub-tree of B since we know in advance that it is \exists . At the top level the algorithm is organized as follows:

- **Step 1:** based on the ideas described in the previous section, compute the set of disjointness and subsumption mapping elements which are *minimal modulo equivalence*. By this we mean that they are minimal modulo collapsing, whenever possible, two subsumption relations of opposite direction into a single equivalence mapping element;
- **Step 2:** eliminate the redundant subsumption mapping elements. In particular, collapse all the pairs of subsumption elements (of opposite direction) between the same two nodes into a single equivalence element. This will result into the *minimal mapping*;
- **Step 3:** Compute the mapping of maximum size (including minimal and redundant mapping elements) or, similarly, given any two nodes return the mapping element existing between the two nodes, or the fact that such element does not exist. During this step the existence of a (redundant) element is computed as the result of the propagation of the elements in the minimal mapping. Notice that redundant equivalence mapping elements can be computed due to the propagation of minimal equivalence elements or of two minimal subsumption elements of opposite direction. However, it can be easily proved that in the latter case they correspond to two partially redundant equivalence elements, where a partially redundant equivalence element is an equivalence element where one direction is a redundant subsumption mapping element while the other is not.

The first two steps are performed at matching time, while the third is activated whenever the user wants to exploit the pre-computed mapping elements for instance for their visualization. The following three subsections analyze the three steps above in detail.

4.1 Step 1: Computing the minimal mapping modulo equivalence

The minimal mapping is computed by a function **TreeMatch** whose pseudo-code is described in Fig. 6. M is the minimal set while $T1$ and $T2$ are the input lightweight ontologies. **TreeMatch** is called on the root nodes of $T1$ and $T2$. It is crucially dependent on the node matching functions **NodeDisjoint** (Fig. 7) and **NodeSubsumedBy** (Fig. 8) which take two nodes $n1$ and $n2$ and return a positive answer in case of disjointness or subsumption, or a negative answer if it is not the case or they are not

able to establish it. Notice that these two functions hide the heaviest computational costs; in particular their computation time is exponential when the relation holds and, exponential in the worst case, but possibly much faster, when the relation does not hold. The main motivation for this is that the node matching problem, in the general case, should be translated into disjointness or subsumption problem in propositional DL (see [5] for a detailed description).

```

10 node: struct of {cnode: wff; children: node[];}
20 T1,T2: tree of (node);
30 relation in { $\sqsubseteq$ ,  $\sqsupseteq$ ,  $\equiv$ ,  $\perp$ };
40 element: struct of {source: node; target: node; rel: relation;};
50 M: list of (element);
60 boolean direction;

70 function TreeMatch(tree T1, tree T2)
80   {TreeDisjoint(root(T1),root(T2));
90   direction := true;
100  TreeSubsumedBy(root(T1),root(T2));
110  direction := false;
120  TreeSubsumedBy(root(T2),root(T1));
130  TreeEquiv();
140  };

```

Fig. 6. Pseudo-code for the tree matching function

The goal, therefore, is to compute the minimal mapping by minimizing the calls to the node matching functions and, in particular minimizing the calls where the relation will turn out to hold. We achieve this purpose by processing both trees top down. To maximize the performance of the system, **TreeMatch** has therefore been built as the sequence of three function calls: the first call to **TreeDisjoint** (line 80) computes the minimal set of disjointness mapping elements, while the second and the third call to **TreeSubsumedBy** compute the minimal set of subsumption mapping elements in the two directions modulo equivalence (lines 90-120). Notice that in the second call, **TreeSubsumedBy** is called with the input ontologies with swapped roles. These three calls correspond to Step 1 above. Line 130 in the pseudo code of **TreeMatch** implements Step 2 and it will be described in the next subsection.

TreeDisjoint (Fig. 7) is a recursive function which finds all disjointness minimal elements between the two sub-trees rooted in n_1 and n_2 . Following the definition of redundancy, it basically searches for the first disjointness element along any pair of paths in the two input trees. Exploiting the nested recursion of **NodeTreeDisjoint** inside **TreeDisjoint**, for any node n_1 in T_1 (traversed top down, depth first) **NodeTreeDisjoint** visits all of T_2 , again top down, depth first. **NodeTreeDisjoint** (called at line 30, starting at line 60) keeps fixed the source node n_1 and iterates on the whole target sub-tree below n_2 till, for each path, the highest disjointness element, if any, is found. Any such disjoint element is added only if minimal (lines 90-120). The condition at line 80 is necessary and sufficient for redundancy. The idea here is to exploit the fact that any two nodes below two nodes involved in a disjointness mapping element are part of a redundant element and, therefore, to stop the recursion thus saving a lot of time expensive calls ($n*m$ calls with n and m the number of the nodes in the two trees). Notice that this check needs to be performed on the full path. **NodeDis-**

joint checks whether the formula obtained by the conjunction of the formulas associated to the nodes $n1$ and $n2$ is unsatisfiable (lines 150-170).

```

10 function TreeDisjoint(node n1, node n2)
20   {c1: node;
30   NodeTreeDisjoint(n1, n2);
40   foreach c1 in GetChildren(n1) do TreeDisjoint(c1,n2);
50   };

60 function NodeTreeDisjoint(node n1, node n2)
70   {n,c2: node;
80   foreach n in Path(Parent(n1)) do if (<n,n2,⊥> ∈ M) then return;
90   if (NodeDisjoint(n1, n2)) then
100    {AddMappingElement(<n1,n2,⊥>);
110    return;
120    };
130   foreach c2 in GetChildren(n2) do NodeTreeDisjoint(n1,c2);
140   };

150 function boolean NodeDisjoint(node n1, node n2)
160 {if (Unsatisfiable(mkConjunction(n1.cnode,n2.cnode))) then
    return true;
170 else return false; };

```

Fig. 7. Pseudo-code for the **TreeDisjoint** function

TreeSubsumedBy (Fig. 8) recursively finds all minimal mapping elements where the strongest relation between the nodes is \sqsubseteq (or dually, \sqsupseteq in the second call; in the following we will concentrate only on the first call).

```

10 function boolean TreeSubsumedBy(node n1, node n2)
20   {c1,c2: node; LastNodeFound: boolean;
30   if (<n1,n2,⊥> ∈ M) then return false;
40   if (!NodeSubsumedBy(n1, n2)) then
50     foreach c1 in GetChildren(n1) do TreeSubsumedBy(c1,n2);
60   else
70     {LastNodeFound := false;
80     foreach c2 in GetChildren(n2) do
90       if (TreeSubsumedBy(n1,c2)) then LastNodeFound := true;
100    if (!LastNodeFound) then AddSubsumptionMappingElement(n1,n2);
120    return true;
140    };
150   return false;
160   };

170 function boolean NodeSubsumedBy(node n1, node n2)
180 {if (Unsatisfiable(mkConjunction(n1.cnode,negate(n2.cnode)))) then
    return true;
190 else return false; };

200 function AddSubsumptionMappingElement(node n1, node n2)
210 {if (direction) then AddMappingElement(<n1,n2,⊆>);
220 else AddMappingElement(<n2,n1,⊇>); };

```

Fig. 8. Pseudo-code for the **TreeSubsumedBy** function

Notice that **TreeSubsumedBy** assumes that the minimal disjointness elements are already computed; as a consequence, at line 30 it checks whether the mapping element between the nodes $n1$ and $n2$ is already in the minimal set. If this is the case it stops the recursion. This allows computing the stronger disjointness relation rather than subsumption when both hold (namely with an inconsistent node). Given $n2$, lines 40-50 implement a depth first recursion in the first tree till a subsumption is found. The test for subsumption is performed by function **NodeSubsumedBy** that checks whether the formula obtained by the conjunction of the formulas associated to the node $n1$ and the negation of the formula for $n2$ is unsatisfiable (lines 170-190). Lines 60-140 implement what happens after the first subsumption is found. The key idea is that, after finding the first subsumption, **TreeSubsumedBy** keeps recursing down the second tree till it finds the last subsumption. When this happens, the resulting mapping element is added to the minimal mapping (line 100). Notice that both **NodeDisjoint** and **NodeSubsumedBy** call the function **Unsatisfiable** which embeds a call to a SAT solver.

To fully understand **TreeSubsumedBy**, the reader should check what happens in the four situations in Fig. 9. In case (a) the first iteration of the **TreeSubsumedBy** finds a subsumption between A and C. Since C has no children, it skips lines 80-90 and directly adds the mapping element $\langle A, C, \sqsupseteq \rangle$ to the minimal set (line 100). In case (b), since there is a child D of C the algorithm iterates on the pair A-D (lines 80-90) finding a subsumption between them. Since there are no other nodes under D, it adds the mapping element $\langle A, D, \sqsupseteq \rangle$ to the minimal set and returns true. Therefore **LastNodeFound** is set to true (line 90) and the mapping element between the pair A-C is recognized as redundant. Case (c) is similar. The difference is that **TreeSubsumedBy** will return false when checking the pair A-D (line 30), thanks to previous computation of minimal disjointness mapping elements, and therefore the mapping element $\langle A, C, \sqsupseteq \rangle$ is recognized as minimal. In case (d) the algorithm iterates after the second subsumption mapping element is identified. It first checks the pair A-C and iterates on A-D concluding that subsumption does not hold between them (line 40). Therefore, it recursively calls **TreeSubsumedBy** between B and D. In fact, since $\langle A, C, \sqsupseteq \rangle$ will be recognized as minimal, it is not worth checking $\langle B, C, \sqsupseteq \rangle$ for pattern (1). As a consequence $\langle B, D, \sqsupseteq \rangle$ is recognized as minimal together with $\langle A, C, \sqsupseteq \rangle$.

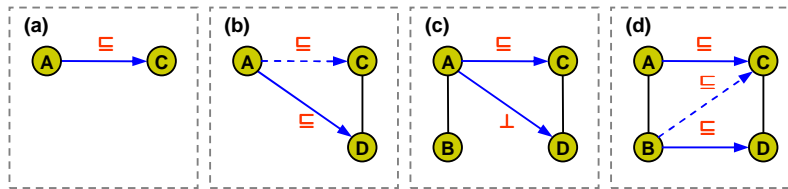


Fig. 9. Examples of applications of the **TreeSubsumedBy**

Five observations. The first is that, even if, overall, **TreeMatch** implements three loops instead of one, the wasted (linear) time is largely counterbalanced by the exponential time saved by avoiding a lot of useless calls to the SAT solver. The second is that, when the input trees $T1$ and $T2$ are two nodes, **TreeMatch** behaves as a node matching function which returns the semantic relation holding between the input

nodes. The third is that the call to **TreeDisjoint** before the two calls to **TreeSubsumedBy** allows us to implement the partial order on relations defined in the previous section. In particular it allows returning only a disjointness mapping element when both disjointness and subsumption hold (see Definition 2 of mapping). The fourth is that, in the body of **TreeDisjoint**, the fact that the two sub-trees where disjointness holds are skipped is what allows not only implementing the partial order (see the previous observation) but also saving a lot of useless calls to the node matching functions (line 2). The fifth and last observation is that the implementation of **TreeMatch** crucially depends on the fact that the minimal elements of the two directions of subsumption and disjointness can be computed independently (modulo inconsistencies).

4.2 Step 2: Computing the minimal mapping

The output of Step 1 is the set of all disjointness and subsumption mapping elements which are minimal modulo equivalence. The final step towards computing the minimal mapping is that of collapsing any two subsumption relations, in the two directions, holding between the same two nodes into a single equivalence relation. The tricky part here is that equivalence is in the minimal set only if both subsumptions are in the minimal set. We have three possible situations:

1. None of the two subsumptions is minimal (in the sense that it has not been computed as minimal in Step 1): nothing changes and neither subsumption nor equivalence is memorized as minimal;
2. Only one of the two subsumptions is minimal while the other is not minimal (again according to Step 1): this case is solved by keeping only the subsumption mapping as minimal. Of course, during Step 3 (see below) the necessary computations will have to be done in order to show to the user the existence of an equivalence relation between the two nodes;
3. Both subsumptions are minimal (according to Step 1): in this case the two subsumptions can be deleted and substituted with a single equivalence element.

Notice that Step 3 can be computed very easily in time linear with the number of mapping elements output of Step 1: it is sufficient to check for all the subsumption elements of opposite direction between the same two nodes and to substitute them with an equivalence element. This is performed by function **TreeEquiv** in Fig. 6.

4.3 Step 3: Computing the mapping of maximum size

For brevity we concentrate on the following problem: given two lightweight ontologies T1 and T2 and the of minimal mapping M compute the mapping element between two nodes n1 in T1 and n2 in T2 or the fact that no element can be computed given the current available background knowledge. Corresponding pseudo-code is given in Fig. 10. **ComputeMappingElement** is structurally very similar to the **NodeMatch** function described in [5], modulo the key difference that no calls to SAT are needed. **ComputeMappingElement** always returns the strongest mapping element. The test for redundancy performed by **IsRedundant** reflects the definition of redundancy provided in Section 3 above. For lack of space, we provide below only the code which does the check for the first pattern; the others are analogous. Given for example a mapping element $\langle n1, n2, \exists \rangle$, condition 1 is verified by checking whether

in M there is an element $\langle c1, c2, \sqsupseteq \rangle$ or $\langle c1, c2, \equiv \rangle$ with $c1$ ancestor of $n1$ and $c2$ descendant of $n2$. Notice that **ComputeMappingElement** calls **IsRedundant** at most three times and, therefore, its computation time is linear with the number of mapping elements in M .

```

10 function mapping ComputeMappingElement(node n1, node n2)
20   {isLG, isMG: boolean;
30   if ((⟨n1,n2,⊥⟩ ∈ M) || IsRedundant(⟨n1,n2,⊥⟩)) then
       return ⟨n1,n2,⊥⟩;
40   if (⟨n1,n2,≡⟩ ∈ M) then return ⟨n1,n2,≡⟩;
50   if ((⟨n1,n2,⊃⟩ ∈ M) || IsRedundant(⟨n1,n2,⊃⟩)) then isLG := true;
60   if ((⟨n1,n2,⊃⟩ ∈ M) || IsRedundant(⟨n1,n2,⊃⟩)) then isMG := true;
70   if (isLG && isMG) then return ⟨n1,n2,≡⟩;
80   if (isLG) then return ⟨n1,n2,⊃⟩;
90   if (isMG) then return ⟨n1,n2,⊃⟩;
100  return NULL;
110 };

120 function boolean IsRedundant(mapping ⟨n1,n2,R⟩)
130   {switch (R)
140     {case ⊃: if (VerifyCondition1(n1,n2)) then return true; break;
150     case ⊃: if (VerifyCondition2(n1,n2)) then return true; break;
160     case ⊥: if (VerifyCondition3(n1,n2)) then return true; break;
170     case ≡: if (VerifyCondition1(n1,n2) &&
                  VerifyCondition2(n1,n2)) then return true;
180     };
190   return false;
200   };

210 function boolean VerifyCondition1(node n1, node n2)
220   {c1,c2: node;
230   foreach c1 in Path(n1) do
240     foreach c2 in SubTree(n2) do
250       if ((⟨c1,c2,⊃⟩ ∈ M) || (⟨c1,c2,≡⟩ ∈ M)) then return true;
260   return false;
270   };

```

Fig. 10. Pseudo-code to compute a mapping element

5 Evaluation

The algorithm presented in the previous sections, let us call it MinSMatch, has been implemented by taking the node matching routines of the state of the art matcher S-Match [5] and by changing the way the tree structure is matched. The evaluation has been performed by directly comparing the results of MinSMatch and S-Match on several real-world datasets. All tests have been performed on a Pentium D 3.40GHz with 2G of RAM running Windows XP SP3 operating system with no additional applications running except the matching system. Both systems were limited to allocating no more than 1G of RAM. The tuning parameters were set to the default values. The selected datasets had been already used in previous evaluations, see [14]. Some

of these datasets can be found at Ontology Alignment Evaluation Initiative web site². The first two datasets describe courses and will be called **Cornell** and **Washington**, respectively. The second two come from the arts domain and will be referred to as **Topia** and **Icon**, respectively. The third two datasets have been extracted from the Looksmart, Google and Yahoo! directories and will be referred to as **Source** and **Target**. The fourth two datasets contain portions of the two business directories eCI@ss³ and UNSPSC⁴ and will be referred to as **Eclass** and **Unspsc**. Table 1 describes some indicators of the complexity of these datasets.

#	Dataset pair	Node count	Max depth	Average branching factor
1	Cornell/Washington	34/39	3/3	5.50/4.75
2	Topia/Icon	542/999	2/9	8.19/3.66
3	Source/Target	2857/6628	11/15	2.04/1.94
4	Eclass/Unspsc	3358/5293	4/4	3.18/9.09

Table 1. Complexity of the datasets

Consider Table 2. The reduction in the last column is calculated as $(1-m/t)$, where m is the number of elements in the minimal set and t is the total number of elements in the mapping of maximum size, as computed by MinSMATCH. As it can be easily noticed, we have a significant reduction, in the range 68-96%.

#	S-Match	MinSMATCH		
	Total mapping elements (t)	Total mapping elements (t)	Minimal mapping elements (m)	Reduction, %
1	223	223	36	83.86
2	5491	5491	243	95.57
3	282638	282648	30956	89.05
4	39590	39818	12754	67.97

Table 2. Mapping sizes.

The second interesting observation is that in Table 2, in the last two experiments, the number of total mapping elements computed by MinSMATCH is higher (compare the second and the third column). This is due to the fact that in the presence of one of the patterns, MinSMATCH directly infers the existence of a mapping element without testing it. This allows MinSMATCH, differently from S-Match, to avoid missing elements because of failures of the node matching functions (because of lack of background knowledge [8]). One such example from our experiments is reported below (directories Source and Target):

```
\Top\Computers\Internet\Broadcasting\Video Shows
\Top\Computing\Internet\Fun & Games\Audio & Video\Movies
```

We have a minimal mapping element which states that Video Shows \sqsubseteq Movies. The element generated by this minimal one, which is captured by MinSMATCH and

² <http://oei.ontologymatching.org/2006/directory/>

³ <http://www.eclass-online.com/>

⁴ <http://www.unspsc.org/>

missed by S-Match (because of the lack of background knowledge about the relation between ‘Broadcasting’ and ‘Movies’) states that Broadcasting \sqsupseteq Movies.

To conclude our analysis, Table 3 shows the reduction in computation time and calls to SAT. As it can be noticed the time reductions are substantial, in the range 16% - 59%, but where the smallest savings are for very small ontologies. The interested reader can refer to [5, 14] for a detailed qualitative and performance evaluation of SMatch w.r.t. other state of the art matching algorithms.

#	Run Time, ms			SAT calls		
	S-Match	MinSMatch	Reduction, %	S-Match	MinSMatch	Reduction, %
1	472	397	15.88	3978	2273	42.86
2	141040	67125	52.40	1624374	616371	62.05
3	3593058	1847252	48.58	56808588	19246095	66.12
4	6440952	2642064	58.98	53321682	17961866	66.31

Table 3. Run time and SAT problems

6 Conclusion

In this paper we have provided a definition and a very fast algorithm for the computation of the minimal mapping between two lightweight ontologies and for the follow-up computation of the mapping of maximum size upon user request (for instance in order to visualize them). We have evaluated the resulting system with respect to the state-of-the-art matching system S-Match [5]. The results show a substantial improvement in the (much lower) computation time, in the (much lower) number of elements which need to be stored and handled and in the (higher) total number of mapping elements which are computed. The last phenomenon is a consequence of the fact that, by minimizing the number of calls to the node matching functions and by maximally exploiting the information codified in the tree structure of the two input ontologies, our algorithm minimizes the impact of the lack of background knowledge.

The future work includes the development of a suitable user interface which exploits minimal mappings thus avoiding the messy visualizations which are generated whenever the number of mapping elements grows, but also the experimentation with large scale mapping tasks. At the moment we are considering applying the system to various large Knowledge Organization Systems (e.g., NALT, AGROVOC, LCSH).

References

1. F. Giunchiglia, M. Marchese, I. Zaihrayeu, 2006. Encoding Classifications into Lightweight Ontologies. *Journal of Data Semantics* 8, pp. 57-81.
2. P. Shvaiko, J. Euzenat, 2007. *Ontology Matching*. Springer-Verlag New York, Inc. Secaucus, NJ, USA.
3. P. Shvaiko, J. Euzenat, 2008. Ten Challenges for Ontology Matching. In *Proceedings of the 7th International Conference on Ontologies, DataBases, and Applications of Semantics (ODBASE 2008)*.
4. J. Madhavan, P. A. Bernstein, P. Domingos, A. Y. Halevy, 2002. Representing and Reasoning about Mappings between Domain Models. At the 18th National Conference on Artificial Intelligence (AAAI 2002).

5. F. Giunchiglia, M. Yatskevich, P. Shvaiko, 2007. Semantic Matching: algorithms and implementation. *Journal on Data Semantics*, IX, 2007.
6. C. Caracciolo, J. Euzenat, L. Hollink, R. Ichise, A. Isaac, V. Malaisé, C. Meilicke, J. Pane, P. Shvaiko, 2008. First results of the Ontology Alignment Evaluation Initiative 2008.
7. F. Giunchiglia, I. Zaihrayeu, 2007. Lightweight Ontologies. In *The Encyclopedia of Database Systems*, to appear. Springer, 2008.
8. F. Giunchiglia, P. Shvaiko, M. Yatskevich, 2006. Discovering missing background knowledge in ontology matching. *Proceedings of the 17th European Conference on Artificial Intelligence (ECAI 2006)*, pp. 382–386.
9. H. Stuckenschmidt, L. Serafini, H. Wache, 2006. Reasoning about Ontology Mappings. *Proceedings of the ECAI-06 Workshop on Contextual Representation and Reasoning, 2006*.
10. C. Meilicke, H. Stuckenschmidt, A. Tamilin, 2006. Improving automatically created mappings using logical reasoning. In the proceedings of the 1st International Workshop on Ontology Matching OM-2006, *CEUR Workshop Proceedings Vol. 225*.
11. C. Meilicke, H. Stuckenschmidt, A. Tamilin, 2008. Reasoning support for mapping revision. *Journal of Logic and Computation*, 2008.
12. A. Borgida, L. Serafini. Distributed Description Logics: Assimilating Information from Peer Sources. *Journal on Data Semantics* pp. 153-184.
13. I. Zaihrayeu, L. Sun, F. Giunchiglia, W. Pan, Q. Ju, M. Chi, and X. Huang, 2007. From web directories to ontologies: Natural language processing challenges. In 6th International Semantic Web Conference (ISWC 2007).
14. P. Avesani, F. Giunchiglia and M. Yatskevich, 2005. A Large Scale Taxonomy Mapping Evaluation. In *Proceedings of International Semantic Web Conference (ISWC 2005)*, pp. 67-81.
15. M. L. Zeng, L. M. Chan, 2004. Trends and Issues in Establishing Interoperability Among Knowledge Organization Systems. *Journal of the American Society for Information Science and Technology*, 55(5) pp. 377–395.
16. L. Kovács. A. Micsik, 2007. Extending Semantic Matching Towards Digital Library Contexts. *Proceedings of the 11th European Conference on Digital Libraries (ECDL 2007)*, pp. 285-296.
17. B. Marshall, T. Madhusudan, 2004. Element matching in concept maps. *Proceedings of the 4th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL 2004)*, pp.186-187.
18. M Yatskevich, 2003. Preliminary evaluation of schema matching systems. Technical Report, DIT-03-028, University of Trento, (2003).
19. B. Hjørland, 2008. What is Knowledge Organization (KO)?. *Knowledge Organization. International Journal devoted to Concept Theory, Classification, Indexing and Knowledge Representation* 35(2/3) pp. 86-101.
20. D. Soergel, 1972. A Universal Source Thesaurus as a Classification Generator. *Journal of the American Society for Information Science* 23(5), pp. 299–305.
21. D. Vizine-Goetz, C. Hickey, A. Houghton, and R. Thompson. 2004. Vocabulary Mapping for Terminology Services. *Journal of Digital Information*, Volume 4, Issue 4.
22. M. Doerr, 2001. Semantic Problems of Thesaurus Mapping. *Journal of Digital Information*, Volume 1, Issue 8.
23. F. Giunchiglia, I. Zaihrayeu, U. Kharkevich, 2007. Formalizing the get-specific document classification algorithm. In 11th European Conference on Research and Advanced Technology for Digital Libraries (ECDL2007). LNCS Springer Verlag.

Appendix: proofs of the theorems

Theorem 1 (Redundancy, soundness and completeness). Given a mapping M between two lightweight ontologies O_1 and O_2 , a mapping element $m' \in M$ is redundant if and only if it satisfies one of the conditions of Definition 3.

Proof:

Soundness: The argumentation provided in section 3 as a rationale for the patterns already provides a full demonstration for soundness.

Completeness: We can demonstrate the completeness by showing that we cannot have redundancy in the cases which do not fall in the conditions listed in Definition 3. We proceed by enumeration, negating each of the conditions. There are some trivial cases we can exclude in advance:

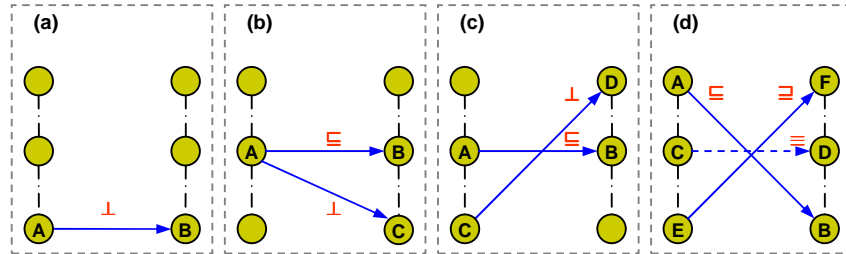


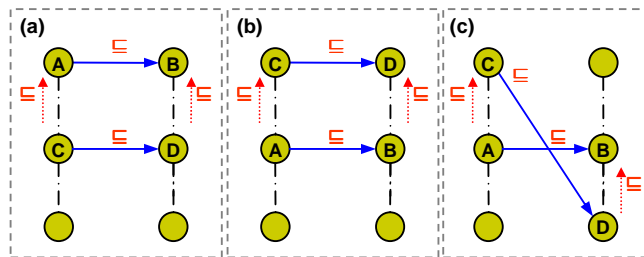
Fig. 11. Some trivial cases which do not fall in the redundancy patterns

- The trivial case in which m' is the only mapping element between the lightweight ontologies. See Fig. 11 (a);
- Incomparable symbols. The only cases of dependency across symbols are captured by conditions (1) and (2) in Definition 3, where equivalence can be used to derive the redundancy of a more or less specific mapping element. This is due to the fact that equivalence is exactly the combination of more and less specific. No other symbols can be expressed in terms of the others. This means for instance that we cannot establish implications between an element with more specific and one with disjointness. In Fig. 11 (b) the two elements do not influence each other;
- All the cases of inconsistent nodes. See for instance Fig. 11 (c). If we assume the element $\langle A, B, \sqsubseteq \rangle$ to be correct, then according to pattern (1) the mapping element between C and D should be $\langle C, D, \sqsubseteq \rangle$. However, in case of inconsistent nodes the stronger semantic relation \perp holds. The algorithm presented in section 4 correctly returns \perp in these cases;
- Cases of underestimated strength not covered by the previous cases, namely the cases in which equivalence holds instead of the (weaker) subsumption. Look for instance at Fig. 11 (d). The two subsumptions in $\langle A, B, \sqsubseteq \rangle$ and $\langle E, F, \sqsubseteq \rangle$ must be equivalences. As a consequence, $\langle C, D, \equiv \rangle$ is redundant for pattern (4). In fact, the chain of subsumptions $E \sqsubseteq \dots \sqsubseteq C \sqsubseteq \dots \sqsubseteq A \sqsubseteq B \sqsubseteq \dots \sqsubseteq D \sqsubseteq \dots \sqsubseteq F$ allows to conclude that $E \sqsubseteq F$ holds and therefore $E \equiv F$. Symmetrically, we can con-

clude that $A \equiv B$. Note that the mapping elements $\langle A, B, \sqsubseteq \rangle$ and $\langle E, F, \supseteq \rangle$ are minimal. We identify the strongest relations by propagation (at step 3 of the proposed algorithm, as described at the beginning of section 4).

We refer to all the other cases as the *meaningful cases*.

Condition (1): its negation is when $R \neq \sqsubseteq$ or $A \notin \text{path}(C)$ or $D \notin \text{path}(B)$. The cases in which $R = \sqsubseteq$ are shown in Fig. 12. For each case, the provided rationale shows that available axioms cannot be used to derive $C \sqsubseteq D$ from $A \sqsubseteq B$. The remaining meaningful cases, namely only when $R = \equiv$, are similar.

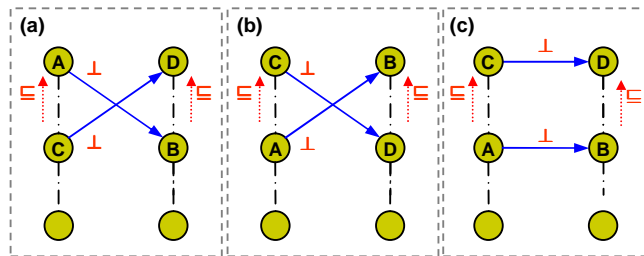


	$A \notin \text{path}(C)$	$D \notin \text{path}(B)$	Rationale
(a)	NO	YES	$C \sqsubseteq \dots \sqsubseteq A, D \sqsubseteq \dots \sqsubseteq B, A \sqsubseteq B$ cannot derive $C \sqsubseteq D$
(b)	YES	NO	$A \sqsubseteq \dots \sqsubseteq C, B \sqsubseteq \dots \sqsubseteq D, A \sqsubseteq B$ cannot derive $C \sqsubseteq D$
(c)	YES	YES	$A \sqsubseteq \dots \sqsubseteq C, D \sqsubseteq \dots \sqsubseteq B, A \sqsubseteq B$ cannot derive $C \sqsubseteq D$

Fig. 12. Completeness of condition (1)

Condition (2): it is the dual of condition (1).

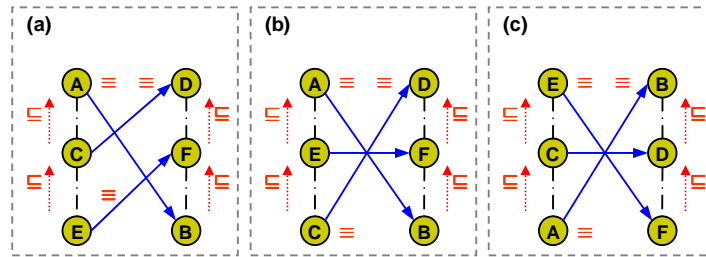
Condition (3): its negation is when $R \neq \perp$ or $A \notin \text{path}(C)$ or $B \notin \text{path}(D)$. The cases in which $R = \perp$ are shown in Fig. 13. For each case, the provided rationale shows that available axioms cannot be used to derive $C \perp D$ from $A \perp B$. There are no meaningful cases for $R \neq \perp$.



	$A \notin \text{path}(C)$	$B \notin \text{path}(D)$	Rationale
(a)	NO	YES	$C \sqsubseteq \dots \sqsubseteq A, B \sqsubseteq \dots \sqsubseteq D, A \perp B$ cannot derive $C \perp D$
(b)	YES	NO	$A \sqsubseteq \dots \sqsubseteq C, D \sqsubseteq \dots \sqsubseteq B, A \perp B$ cannot derive $C \perp D$
(c)	YES	YES	$A \sqsubseteq \dots \sqsubseteq C, D \sqsubseteq \dots \sqsubseteq B, A \perp B$ cannot derive $C \perp D$

Fig. 13. Completeness of condition (3)

Condition (4): it can be easily noted from Fig. 3 that the redundant elements identified by pattern (4) are exactly all the mapping elements $m' = \langle C, D, \equiv \rangle$ with source C and target D respectively between (or the same of) the source node and target node of two different mapping elements $m = \langle A, B, \equiv \rangle$ and $m'' = \langle E, F, \equiv \rangle$. This configuration allows to derive from m and m'' the subsumptions in the two directions which amount to the equivalence. The negation of condition 4 is when $R \neq \equiv$ in m or m'' or $A \notin \text{path}(C)$ or $D \notin \text{path}(B)$ or $C \notin \text{path}(E)$ or $F \notin \text{path}(D)$. In almost all the cases (14 over 15) in which $R = \equiv$ we just move the source C or the target D outside these ranges. For sake of space we show only some of such cases in Fig. 14. The rationale provided for cases (a) and (b) shows that we cannot derive $C \equiv D$ from $A \equiv B$ and $E \equiv F$. The only exception (the remaining 1 case over 15), represented by case (c), is when $A \notin \text{path}(C)$ and $D \notin \text{path}(B)$ and $C \notin \text{path}(E)$ and $F \notin \text{path}(D)$. This case however is covered by condition 4 by inverting the role of m and m'' . The remaining cases, namely when $R \neq \equiv$ in m or m'' , are not meaningful.



	$A \notin \text{path}(C)$	$D \notin \text{path}(B)$	$C \notin \text{path}(E)$	$F \notin \text{path}(D)$	Rationale
(a)	NO	NO	NO	YES	$E \equiv \dots \equiv C, C \equiv \dots \equiv A,$ $B \equiv \dots \equiv F, F \equiv \dots \equiv D,$ $A \equiv B$ and $E \equiv F$ cannot derive $C \equiv D$ (we can only derive $C \equiv D$).
(b)	NO	NO	YES	YES	$C \equiv \dots \equiv E, E \equiv \dots \equiv A,$ $B \equiv \dots \equiv F, F \equiv \dots \equiv D,$ $A \equiv B$ and $E \equiv F$ cannot derive $C \equiv D$ (we can only derive $C \equiv D$).
...					
(c)	YES	YES	YES	YES	Covered by condition (4) inverting the roles of m and m''

Fig. 14. Completeness of condition (4)

This completes the demonstration.□

Theorem 2 (Minimal mapping, existence and uniqueness). Given two lightweight ontologies O_1 and O_2 , there is always one and only one minimal mapping between them.

Proof:

The proof is based on two main observations:

Observation 1: A redundant mapping element $m' \in M'$ can be caused by one and only one of the redundancy conditions in Definition 3. In other words, redundancy conditions are mutually exclusive. In particular, apart from the cases of inconsistent nodes, subsumption and disjointness mutually exclude themselves. For equivalence, note that for condition 4 in Definition 3 a mapping element m cannot satisfy both condition 1 and 2. In fact, since $m \neq m'$, it cannot be $A \in \text{path}(C)$ and $D \in \text{path}(B)$ (see condition 1) and $C \in \text{path}(A)$ and $B \in \text{path}(D)$ at the same time (see condition 2). Conditions 1 and 2 are both needed to build the redundant equivalence.

Observation 2: We can define a strict partial order over mapping elements for each relation in $\{\subseteq, \supseteq, \perp\}$: given two lightweight ontologies O_1 and O_2 , a mapping M between them and two distinct mapping elements $m, m' \in M$ with $m = \langle A, B, R \rangle$ and $m' = \langle C, D, R' \rangle$, we say that $m' < m$ iff one of the following holds:

- (1) If R' is \subseteq , $R \in \{\subseteq, \equiv\}$, $A \in \text{path}(C)$ and $D \in \text{path}(B)$;
- (2) If R' is \supseteq , $R \in \{\supseteq, \equiv\}$, $C \in \text{path}(A)$ and $B \in \text{path}(D)$;
- (3) If R' and R are \perp , $A \in \text{path}(C)$ and $B \in \text{path}(D)$;

It can be easily noticed that in all the three cases above $<$ enforces a partial order, since, given two mapping elements with source and target on the same paths, they are ordered when their structural configuration is like in the corresponding patterns in Fig. 3, while they are not ordered when the configuration is different or are on different paths. The fact that the ordering defined is partial is a direct consequence of the tree structure of lightweight ontologies.

It is clear that under the strict partial order above, if we “open” equivalence relationships in the two subsumptions of opposite direction, the minimal mapping is the set of all the maximal elements of the partially ordered set, where subsumptions of opposite direction involving the same nodes are collapsed into a (minimal) equivalence mapping element. For the properties of partial orders, this set always exists and it is unique. \square
