

ONTECTAS: Bridging the Gap between Collaborative Tagging Systems and Structured Data

Ali Moosavi, Tianyu Li, Laks V.S. Lakshmanan, and Rachel Pottinger

University of British Columbia, Vancouver, BC, Canada
{amoosavi, lty419, laks, rap}@cs.ubc.ca

Abstract. Ontologies define a set of terms and the relationships (e.g., IS-A and HAS-A) between them; they are the building block of the emerging semantic web. An ontology relating the tags in a collaborative tagging system (CTS) makes the CTS easier to understand. We propose an algorithm to automatically construct an ontology from CTS data and conduct a detailed empirical comparison with previous related work on four real data sets – Del.icio.us, LibraryThing, CiteULike, and IMDb. We also verify the effectiveness of our algorithm in detecting IS-A and HAS-A relationships.

Keywords: ontology, taxonomy, tag, collaborative tagging systems.

1 Introduction

Ontologies organize information in content management systems and are the core building blocks of the emerging Semantic Web. Substantial work has been done in extracting ontologies automatically from large repositories like text corpora, databases, and the web. This paper focuses on collaborative social tagging systems (CTSs) such as Del.icio.us (for tagging bookmarks), Flickr (for tagging photos), IMDb (for tagging movies), LibraryThing (for tagging books) and CiteULike (for tagging publications). These systems permit users to tag and share resources (documents, photos, videos, etc.). Our goal is to create a generic ontology of the tags from a CTS. By ontology, we mean a set of concepts from a domain, represented by the tags, and their (IS-A and HAS-A) relationships.

Learning an ontology from a CTS can help make the CTS more useful. For example, browsing an ontology of tags from a CTS can help users better refine their queries, either to find more items by using a more general term or to find fewer items by using a more specific term. This is especially important in a CTS since the resources are typically labeled by a small, sparse, set of tags — so discovering content in CTSs by simple keyword search is much harder than in document and web search. Another application of domain specific ontology builders is to enhance search engines with ontologies. E.g., the prototype Clever Search system [15] merges words and their word senses in the general ontology, WordNet¹, and returns more relevant result items to the user.

¹ <http://wordnet.princeton.edu>

In principle, we could use a general purpose ontology such as WordNet to browse a CTS; there are two disadvantages. First, tags in CTSs are not based on a fixed vocabulary but constantly evolve. Thus, one cannot expect WordNet (or similar systems) to capture the vocabulary in a dynamic CTS, e.g., “Mac OS X”. Secondly, as we demonstrate in Section 7, even when terms corresponding to tags in a CTS *are* present in WordNet, in many cases, valid IS-A relationships between them that are found by our algorithm are missing in WordNet. This mirrors a similar finding for the ontology extracted from Wikipedia using YAGO [25]; using a combination of WordNet and Wikipedia found significantly more ontological relationships (including IS-A) that were absent in WordNet.

This paper studies the following problem: given a collaborative tagging system consisting of users, resources (also called items), and tags assigned by users to items, extract an ontology consisting of tags in the CTS and IS-A and HAS-A relationships between the tags. We consider HAS-A relationships in addition to IS-A: indeed, IS-A and HAS-A relationships are among those most used in ontologies with rich relationships, such as WordNet.

Our algorithm for ontology extraction from CTSs is predicated on the hypothesis that tags assigned to a resource by a group of users tend to contain both child and parent tags. We have conducted experiments to validate this assumption in the full version of our paper[20]. A possible explanation for this phenomenon is that different users may use tags at different levels of abstraction (from an underlying ontology in their mind); thus tags for the same item may include more abstract or more specific terms as an aggregation effect of various tagging behaviors. We leverage this hypothesis using association rules [1] and lexico-syntactic patterns to find relationships between tags. Our approach accounts for bi-grams (which can affect the precision of detected relationships), multi-word tags, and also infer non-trivial IS-A relationships from detected ones. We make the following contributions:

- We propose (Sections 4 through 6) an algorithm for ontology extraction from a CTS, called ONTECTAS (for ONTOlogy Extraction from Collaborative Tagging Systems). The highlights of the algorithm include:
 - Candidate IS-A relationships are mined using association rules, making use of both forward and reverse confidence (Section 4.1).
 - Invalid tuples are pruned based on discovering bi-grams (Section 4.3).
 - Headword detection is leveraged for discovering relationships between multi-word tags (Section 4.4).
 - Lexico-syntactic patterns are used for detecting IS-A and HAS-A relationships. To our knowledge, we are the first to explicitly extract HAS-A relationships from CTSs (Section 5).
 - Based on items in the ontology having a common (IS-A) child, additional IS-A relationships are inferred (Section 6).
- We demonstrate via a comprehensive set of experiments on four real datasets that our algorithm outperforms previous algorithms w.r.t. quality and richness of the extracted ontology (Section 7).

Section 2 discusses related work. Section 3 formalizes the problem studied in this paper. Section 8 concludes and discusses future work.

2 Related Work

Some other works have studied extracting ontologies from CTSs. Some approaches [16,18,2] match CTS tags to concepts in general purpose ontologies such as WordNet, resulting in a graph of tags. However, because CTSs are ad-hoc and use terms dynamically, general purpose ontologies miss many terms as well as edges (i.e., relationships). For example, our experiments show that WordNet misses more than 25% of correct edges between concepts extracted from Del.icio.us, *even when both parent and child concepts are in WordNet*.

Schmitz [23] constructs weighted graphs based on conditional probabilities between pairs of tags. His algorithm cannot identify the exact relationship (e.g., IS-A and HAS-A) between terms — it simply says they are related, not how. By contrast, our algorithm pinpoints IS-A and HAS-A relationships between terms.

Heymann and Garcia-Molina [10] create an ontology by vectorizing the tags and finding the cosine similarity between tags. However, their method puts every tag from the similarity matrix into the taxonomy which causes many erroneous edges. Their work lacks an evaluation.

Schmitz et al. [22] use association rule mining to build a tree of related tags from a CTS; however, they do not explain how the edges are built or what types of relationships they model. We explain this in depth and also use lexico-syntactic patterns and a search engine to detect accurate IS-A and HAS-A relationships. [24] extends [22] and [10] by considering the tag's context. Barla and Bielíková [3] consider tag context similarly to [24].

The DAG algorithm [5] distinguishes between subjective and objective tags. After calculating feature vectors for each objective tag, DAG places tags with higher entropy in higher levels of abstraction. Like many other previous works, DAG does not determine the type of relationship between concepts.

Lin et al. [19] build a subsumption graph from the folksonomy and use a random walk to sort tags by generality ranking. They put tags in the taxonomy based on support and confidence between candidate nodes from the graph. They only consider a single sense for each tag, which leads to missed relationships. The authors claim building transactions for tags associated to items by specific users will lead to the best taxonomy because it preserves most of the information. In contrast, we found that user information does not improve taxonomy quality.

Körner et al. [14] categorize users by the kind of tags they use. They show that excluding some users can reduce noise and improve precision. This improvement is orthogonal to the contribution we make in this paper and is applicable in our context as well. We leave adapting ONTECTAS to this as future work.

Hearst [9] defines a set of patterns that indicate IS-A relationships between words in text documents. [4,6] find patterns for detecting HAS-A relationships from text corpora. To our knowledge, our work is the first to extend the lexico-syntactic patterns to find relationships of any type between tags in CTSs.

In sum, in contrast to previous works on ontology extraction from CTSs, our method is capable of detecting both HAS-A and IS-A relationships and explicitly identifying each. Our multi-stage algorithm also extracts high quality relationships between multi-word tags.

3 Problem Statement

A *collaborative tagging system* [22] is a 4-tuple $C = (U, T, I, Y)$ where U is a set of users, T is the set of tags used by the users, I is the set of items (resources) to which tags are assigned by users, and Y , the set of tag assignments, is a ternary relation on tags, users, and items, i.e., $Y \subseteq U \times T \times I$.

Specific CTSs may vary in detail from our definition above, e.g., IMDb does not have user information. We can model such CTSs by dropping U and defining $Y \subseteq T \times I$ as a binary relation. CTSs such as [11] allow users to declare their own IS-A relationships. User-supplied IS-A relationships can augment those automatically extracted but cannot supplant them because of the scale.

This paper studies how to efficiently extract IS-A and HAS-A relationships between tags in a given CTS. The output ontology consists $\langle \text{tag1}, \text{tag2}, \text{label} \rangle$ tuples where tag1 is the super class and label is either IS-A or HAS-A.² E.g., the tuple $\langle \text{OS}, \text{Windows}, \text{IS-A} \rangle$ indicates that Windows a kind of OS.

4 Ontology Extraction from Collaborative TAGging Systems (ONTECTAS) Algorithm

Algorithm 1. ONTECTAS

Input: (D) A set of $\langle \text{item}, \text{tag} \rangle$ 2-tuples or $\langle \text{user}, \text{item}, \text{tag} \rangle$ 3-tuples
Output: (O) Ontology of tags with IS-A and HAS-A relationships

- 1: $D' \leftarrow$ Preprocess D. /* D' is a set of $\langle \text{item}, \text{tag} \rangle$ tuples*/
- 2: $\langle T_{\text{basic}}, F \rangle \leftarrow$ Association_Rule_Tuple_Detection(D') /*Algorithm 2*/
- 3: $T_{\text{pruned}} \leftarrow$ Bigram_Filtering(T_{basic}) /*Algorithm 3*/
- 4: $\langle T_{\text{headword}}, O \rangle \leftarrow$ Headword_Detection(T_{pruned}) /*Algorithm 4*/
- 5: $O \leftarrow O \cup$ IS-A_Relationship_Detection(T_{headword} , IS-A-patterns,
IS-A - threshold) /*Algorithm 5*/
- 6: $O \leftarrow O \cup$ HAS-A_Relationship_Detection(T_{headword} , HAS-A-patterns,
HAS-A - threshold)
- 7: $T_{\text{co_parent}} \leftarrow$ Co-Parent_Pruning(T_{headword}, F) /*Algorithm 6*/
- 8: Return $O \cup$ IS-A_Relationship_Detection($T_{\text{co_parent}}$, IS-A - patterns,
IS-A - threshold) /*Algorithm 5*/

Our ONTECTAS algorithm for ontology extraction (Algorithm 1) consists of six phases. First, data is preprocessed and cleaned. Next, we extract candidate tag tuples via association rule mining using forward and reverse confidence. We then

² In both relationships tag2 IS-A tag1 and tag1 HAS-A tag2, we refer to tag1 as the super class label or the parent label for convenience, by abusing terminology.

remove tuples corresponding to bigrams. Next, we detect headwords of multi-word tags and use this to infer additional IS-A relationships. We then use lexico-syntactic patterns to extract additional IS-A and HAS-A relationships. Finally, we leverage pairs of tags sharing a common child in the extracted ontology to infer additional IS-A relationships. The next three sections describe the phases.

4.1 Preprocessing

The preprocessing step is primarily a cleaning step. It takes as input a CTS and performs the following tasks: (1) Any user information is projected away; we found looking at transactions at the level of group of users was most effective in ontology extraction. (2) Words with non-English characters are removed from the input data using the same method as in [5]. This adequately removed non-English words from all of our datasets. (3) Basic stemming: singular nouns are substituted for their plural forms. (4) Since tags occurring very infrequently are not statistically reliable, we removed tags or items that occurred fewer than 5 times. This threshold was determined empirically. (5) Verbs and verb phrases are removed by applying the Stanford parser³ to each tag. This prunes tags that are used for organizing but convey no meaning about the item being tagged [7].

4.2 Detecting Potential Relationships Using Association Rules

Adapting tagged data to market basket analysis requires defining how to build transactions from tags, which in turn requires defining “co-occurrence”. We explored three different definitions of co-occurrence. Empirically, we determined that the most effective co-occurrence definition is the following: Tags t and t' co-occur if both were used to tag the same item (by possibly different users). The frequency of $\{t, t'\}$ equals the number of distinct items which were assigned both tags t and t' . Our careful study of the best definition of co-occurrence [20] allows us to more optimally use association rules than previous approaches, e.g., [23].

We use the FP-tree association rule mining algorithm [8] to extract frequent tag sets⁴ and interesting rules from the set of transactions. The *support* of a tag set X is the proportion of transactions containing tag set X and the *confidence* of a rule is defined as $\text{confidence}(X \Rightarrow Y) = \text{support}(X \cup Y) / \text{support}(X)$ — i.e., the proportion of transactions in which X and Y occur together among those in which X appears. In this paper, we refer to the well-known definition of confidence as forward confidence (FC). We also introduce a new notion, reverse confidence (RC) as follows: $\text{reverse_confidence}(X \Rightarrow Y) := \text{support}(X \cup Y) / \text{support}(Y)$.

We assume tags assigned by users tend to contain both a term in the ontology and another term that has relationship with it. Therefore, if two keywords co-occur frequently, they are likely to be related. We use *support* to filter sets of tags with a cardinality of two. However, popular unrelated terms may occur together frequently, so we use *confidence* to remove tuples containing unrelated tags. Because terms which co-occur with high confidence are sometimes synonyms

³ <http://nlp.stanford.edu/software/lex-parser.shtml>

⁴ Tag sets correspond to itemsets in the context of frequent itemset mining.

Algorithm 2. Association_Rule_Tuple_Detection

Input: (D) A set of 2-tuples in form of $\langle item, tag \rangle$ **Output:** (T) Preliminary tag tuples, (F) Set of frequent itemsets

```

1: Group  $D$  by item. /*create:  $\langle item, \{tag_1, \dots, tag_k\} \rangle$ */
2:  $S \leftarrow$  Union of tags associated with each item (i.e.,  $S$  is set of transactions)
3:  $F \leftarrow$  Frequent itemsets of size two from  $S$  where support  $> min\_support$ 
   /* $FC_i$  and  $RC_i$  are forward and reverse confidence respectively*/
4: for all  $F_i \in F$  do
5:   if  $((FC_i \geq min\_conf.)$  and  $(RC_i \leq 1 - min\_conf.)$ ) OR  $((RC_i \geq min\_conf.)$ 
     and  $(FC_i \leq 1 - min\_conf.)$ ) then
6:     Add  $F_i$  to  $T$ 
7:   end if
8: end for
9: Return  $\langle T, F \rangle$ 

```

(e.g., “os” and “operating system”), we use confidence in the reverse direction to ensure that terms are related with IS-A or HAS-A relationships. Different values for $min_support$ and $min_conf.$ can drastically change the size of the ontology; in our experiments these values were chosen empirically. At the end of this step, we have not yet classified the relationships into IS-A and HAS-A.

4.3 Pruning Edges between Bi-gram Elements

In this phase, bi-gram tuples which are common phrases are automatically pruned using a search engine. Usually bi-grams are compound nouns in the form of “adjective + noun” (e.g., free software) or “noun + noun” (e.g., web browser). Bi-grams do not contain IS-A or HAS-A relationships but sometimes are incorrectly detected as edges of an ontology since they co-occur frequently.

Finding bigrams by using a search engine [26,12,17] has not previously been applied to extracting relationships between CTS tags. ONTECTAS sends two keyword queries to a search engine for each relationship tuple (Algorithm 3). The queries are the quoted permutations of the terms in the tuple. If the ratio of the number of results returned for the two queries is larger than a threshold, the terms in the relationship tuple are regarded as bi-grams. E.g., if the relationship tuple is $\langle software, free \rangle$, the queries are “free software” and “software free”. Since the ratio is higher than the threshold for this tuple, it is detected as a bi-gram and pruned. We experimentally found that the optimal threshold for detecting bi-grams is between 50 and 100. Because words in text documents have Zipfian distribution, [12] suggests using a logarithmic transformation of returned result counts. We found that the logarithmic transformation is also more accurate in detecting bi-grams.

4.4 Detecting Headwords in Multi-word Tags

Since many CTS tags are multi-word tags in form of compound phrases such as “science-fiction” and “object-oriented-data-model”, we use headword detection

Algorithm 3. Bi-gram_Filtering

Input: (T) A set of 2-tuples of the form $\langle tag1, tag2 \rangle$ **Output:** (T') A reduced set of 2-tuples

```

1:  $T' \leftarrow T$ 
2: for all  $T'_i \in T'$  do
3:    $ratio1 \leftarrow$  # of hits of querying “tag1 tag2” as a phrase
4:    $ratio2 \leftarrow$  # of hits of querying “tag2 tag1” as a phrase
5:    $ratio \leftarrow \frac{\log(\max(ratio1, ratio2))}{\log(\min(ratio1, ratio2))}$ 
6:   if  $ratio \geq bi - gram\_threshold$  then
7:     remove  $T'_i$  from  $T'$ 
8:   end if
9: end for
10: Return  $T'$ 

```

to extract additional IS-A relationships (Algorithm 4). First, the Stanford parser detects the *headwords* for each phrase. A headword is a phrase’s grammatically most important word; it determines the phrase’s syntactic type. We then extract an IS-A relationship for each multi-word tag by putting the headword as the parent of the whole phrase. E.g., we can infer “object-oriented data model” IS-A “model”. In this phase, more candidate tuples are produced by using either whole phrases or their headwords as the tags in tuples.

Algorithm 4. Headword_Detection

Input: (T) A set of 2-tuples of the form $\langle tag1, tag2 \rangle$ **Output:** (T') A set of enhanced 2-tuples, (O) Ontology with IS-A relationships

```

1:  $T' \leftarrow T$ 
2: for all  $T_i \in T$  do
3:   if  $T_i$  contains multi-tags then
4:      $head1 \leftarrow$  headword in  $tag1$ 
5:      $head2 \leftarrow$  headword in  $tag2$ 
6:      $O \leftarrow O \cup \{\langle head1, tag1, IS-A \rangle\}$ 
7:      $O \leftarrow O \cup \{\langle head2, tag2, IS-A \rangle\}$ 
8:      $T' \leftarrow T' \cup \{\langle head1, tag2 \rangle, \langle head2, tag1 \rangle, \langle head1, head2 \rangle\}$ 
9:   end if
10: end for
11: Return  $\langle T', O \rangle$ 

```

5 Using Lexico-Syntactic Patterns

Finally, we analyze occurrences of lexico-syntactic patterns to detecting IS-A and HAS-A relationships. Due to data sparsity, lexico-syntactic patterns do not occur frequently enough to accurately detect relationships between terms [21]. Hence, we build on [13] and query the web for more occurrences of the patterns.

The core of our lexico-syntactic search is shown in lines 3-6 of Algorithm 5: given two tags and a pattern, we generate two keyword queries by considering

Algorithm 5. IS-A_Relationship_Detection

Input: (T, P , threshold) Where T is a set of $\langle tag1, tag2 \rangle$ tuples, P is a set of patterns

Output: (I) A set of 2-tuples in form of $\langle parent_tag, child_tag \rangle$

```

1: for all  $t_i \in T$  do
2:   for all  $p_j \in P$  do
3:      $hits1 \leftarrow \#$  of hits of querying " $t_i.tag1 p_j t_i.tag2$ " as a phrase
4:      $hits2 \leftarrow \#$  of hits of querying " $t_i.tag2 p_j t_i.tag1$ " as a phrase
5:      $ratio_j.F \leftarrow \frac{hits1}{hits2}$ 
6:      $ratio_j.R \leftarrow \frac{hits2}{hits1}$ 
7:   end for
8:    $maximum_F \leftarrow \max(ratio_j.F)$  over all  $j$ 
9:    $maximum_R \leftarrow \max(ratio_j.R)$  over all  $j$ 
10:   $maximum \leftarrow \max(maximum_F, maximum_R)$ 
11:  if  $((maximum = maximum_F)$  and  $(maximum_F \geq threshold))$  then
12:     $I \leftarrow I \cup \{ \langle tag1, tag2, IS-A \rangle \}$ 
13:  else
14:    if  $((maximum = maximum_R)$  and  $(maximum_R \geq threshold))$  then
15:       $I \leftarrow I \cup \{ \langle tag2, tag1, IS-A \rangle \}$ 
16:    end if
17:  end if
18: end for
19: Return  $I$ 

```

the two possible permutations of the tags in the pattern. E.g., given (“human”, “body”, “s”), the two generated queries will be “human’s body” and “body’s human”. Then, the ratios for both forward and reverse occurrences direction are calculated. It is clear that given any set of patterns for any relationship, this algorithm can be applied. We use the following patterns from [9] to identify IS-A relationships: (1) Pattern 1: NP₁ such as NP₂; (2) Pattern 2: NP₁ including NP₂; (3) Pattern 3: NP₁ especially NP₂.

Our HAS-A relationships are supersets of meronymy (part-of relationships), and are not limited to the physical perspective. We consider two noun phrases NP₁ and NP₂ to have a HAS-A relationship (with NP₁ as the parent) if one of the following statements is true: (1) NP₂ is a part of NP₁. E.g., “body” is a part of “human”; or (2) NP₁ has/have NP₂. E.g., “human” has “mind” and “google” has “googleMaps”; or (3) NP₁ may have NP₂. E.g., “human” may have “disease”.

From the existing lexico-syntactic patterns mentioned in the literature such as [4,6], we use three following patterns to detect HAS-A relationships:(1) Pattern 1: NP₁’s NP₂; (2) Pattern 2: NP₂ of the NP₁; (3) Pattern 3: NP₂ of NP₁.

While patterns 1 and 2 are among the most common English patterns [6], pattern 3 is not. However, pattern 3 can be used to detect HAS-A relationship between tags such as the tuple $\langle \text{Coffee}, \text{Caffeine} \rangle$.

All patterns for a relationship are fed into a search engine. If the largest ratio of a pattern is above a threshold, that tuple is labeled with the corresponding relationship and added to the ontology. Algorithm 5 shows the IS-A detection

algorithm. The HAS-A algorithm is similar, but requires that pattern 1 *and* one of patterns 2 and 3 are above the threshold. Both thresholds were found experimentally. In our experiments, the IS-A threshold was 7 and HAS-A threshold ranged from 20 to 50.

6 Exploiting Co-parents to Find More IS-A Relationships

Examining the ontology built thus far reveals an interesting property when pairs of tags share the same child. Consider the following example: the ontology may contain “fiction \rightarrow urban-fantasy” and “fantasy \rightarrow urban-fantasy”, where “fiction” and “fantasy” are both parents for “urban-fantasy” w.r.t. the IS-A relationship.⁵ However, the IS-A relationship between “fiction” and “fantasy” may be missing. One possible reason for this is that people tend to use the more specific tags leading to “fiction \rightarrow urban-fantasy” and “fantasy \rightarrow urban-fantasy”, so that “fiction \rightarrow fantasy” does not occur above the relatively high threshold needed to avoid noise.

Hence we have the following hypothesis: in a co-parent structure it is more likely than usual that the two parents are in an IS-A relationship. Hence, we include the following additional step (Algorithm 6) to ONTECTAS: for such co-parent pairs, we re-examine the pair’s confidences under a lower threshold and extract candidate tuples for an IS-A relationship.

Algorithm 6. Co_Parent_Pruning

Input: (T) A set of tuples with IS-A relationships in form $\langle parentTag, childTag \rangle$; (F) A set of frequent itemsets

Output: (T') An enhanced set of tuples with IS-A relationships

1: $T' \leftarrow T$

2: $G \leftarrow$ A graph where each tuple in T corresponds to an edge from $parentTag$ to $childTag$.

3: $S \leftarrow$ All tuples of tags $\langle parent_1, parent_2, child \rangle$

s.t. (1) $edge(parent_1 \rightarrow child) \in G$ and (2) $edge(parent_2 \rightarrow child) \in G$ and
(3) $edge(parent_1 \rightarrow parent_2) \notin G$ and (4) $edge(parent_2 \rightarrow parent_1) \notin G$.

4: **for all** $\langle parent_1, parent_2, child \rangle \in S$ **do**

5: **if** $\{parent_1, parent_2\}$ is frequent and if it satisfies lower forward and reverse confidence thresholds **then**

6: Add $\langle parent_1, parent_2 \rangle$ to T' with the more frequent tag as the parent.

7: **end if**

8: **end for**

9: Return T'

As a final step of the ONTECTAS algorithm, following standard practice in ontology extraction algorithms, if the graph of relationships is disconnected, we add a generic “Entity” root node and make it the parent of all orphan nodes.

⁵ Here, —fiction” \rightarrow “urban-fantasy” means “urban-fantasy” IS-A “fiction”.

7 Experiments

7.1 Datasets and Assumptions

Our experiments used four real datasets: Del.icio.us (a social bookmarking web service), IMDb (the Internet Movie Database), LibraryThing (for tagging books) and CiteULike (a service for storing, organizing, and sharing scholarly papers). Table 1 shows the characteristics of the datasets. User information is not available in the IMDb dataset, so competing algorithms were unable to create ontologies from it.

Table 1. Corpus Details in Some Collaborative Tagging Systems

	Del.icio.us (Dec. 2007)	CiteULike (Jan. 2010)	IMDb (Nov. 2009)	LibraryThing (corpus from Delft*)
Number of Tags	6,933,179	431,160	2,593,747	10,469
Number of Items	54,401,067	2,081,799	356,162	37,232
Number of Users	978,979	60,220	N/A	7,279
Number of Tag Assignments	450,113,886	7,922,454	2,625,237	2,415,517

* <http://homepage.tudelft.nl/5q88p/LT>

To show that general purpose ontologies are insufficient, we validated that WordNet misses many relationships between terms *even when it contains both terms*. To show this, we evaluated a sample ontology (from Del.icio.us) both manually and by using all parent-child senses (meanings) in WordNet. We limited our experiments to relationships where both parent and child term exist in WordNet. This gives WordNet an advantage since many tags do not appear in WordNet at all. In this case, we found WordNet is missing 26.9% of manually validated relationships discovered by ONTECTAS. For example, WordNet contains 3 senses for “python”, but none of these senses is related to programming; as a result, “programming \rightarrow python” is missing in WordNet.

Since our approach is successful, it is clear that our hypothesis that *a group of users tend to tag items with both parent and child tags* is validated. The full version of this paper [20] shows detailed experiments which validate this empirically. We discuss our results, beginning with HAS-A relationship detection.

7.2 Evaluation of ONTECTAS in Detecting HAS-A Relationships

Table 2 shows the precision of ONTECTAS in detecting HAS-A relationships. *None of the other competing algorithms address HAS-A relationships from CTSs*. Table 2 only reports precision for ONTECTAS, the first algorithm to detect HAS-A from CTS data.

One challenge in detecting HAS-A relationships was that pattern-based search engine queries such as “human’s middle” and “middle of human” are frequently part of phrases such as “human’s middle finger” and “middle of human history”. Clearly, there is room for improvement in ONTECTAS’ precision in HAS-A detection, which we plan to address in future work.

Table 2. Precision in detecting HAS-A relationships

	Del.icio.us	CiteULike	LibraryThing	IMDb
Precision	51.6%	61.9%	55.5%	33.3%

7.3 Evaluation of ONTECTAS in Detecting IS-A Relationships

In the following, we focus on IS-A relationships. All competing algorithms do not distinguish between IS-A and other relationships such as synonyms, whereas we clearly isolate IS-A relationships. We lump all other relationships into ANY and compare the performance of ONTECTAS on IS-A with that of other algorithms on IS-A and ANY, giving them an advantage, since in this evaluation, we do not give credit to ONTECTAS for correctly finding HAS-A relationships. We use the following standard performance measures: (1) Precision: We consider the precision of ONTECTAS on IS-A with that of other algorithms on IS-A+ ANY. Precision for both is the number of correct edges over the number of all edges. (2) Maximum depth and average depth of the IS-A taxonomy. (3) Average number of children. A higher value of the last two measures implies richer ontology is extracted. In addition, following [19], we compare all algorithms with a gold standard to see how they fare in trying to recreate manually-curated ontologies.

For depth and breadth metrics, we calculate these metrics on an ontology with only correct relationships to ensure algorithms cannot earn an artificially and unfairly high score on these by finding many incorrect relationships!

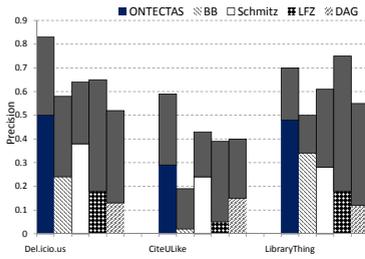
Absolute recall for ontology extraction from a large CTS is very hard to measure. Instead, we propose a new metric: *relative recall*. Relative recall for an algorithm is the number of valid IS-A relationships found by the algorithm divided by the total number of valid IS-A relationships found by all algorithms.

7.4 Comparing ONTECTAS to Other Algorithms

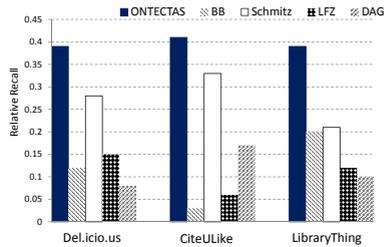
We compare ONTECTAS with the four algorithms from Section 2: 1) the algorithm from [19] (abbreviated “LFZ”) 2) the DAG algorithm [5] (“DAG-ALG”) 3) Schmitz’s algorithm [23] (“Schmitz”), and 4) Barla and Bieliková’s algorithm [3] (“BB”). Since these algorithms cannot process the IMDb dataset due to the lack of user information, we only compare them on Del.icio.us, LibraryThing, and CiteULike.

To have a fair comparison, we implemented the above algorithms as closely as possible to the way their authors had implemented them; we used the parameters that were described in the papers and contacted the authors for additional information about how to make their algorithms as competitive as possible.

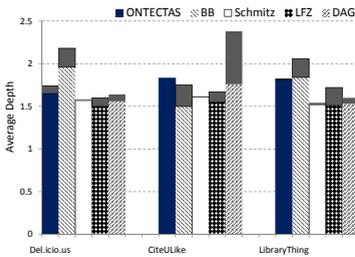
Validating the edges manually required that each algorithm output a small number of edges. To do so, we put another threshold on the number of times a tag, an item, or a user must occur in order to be considered. To be fair, we used the same threshold to ensure that each algorithm output fewer than 150 edges.



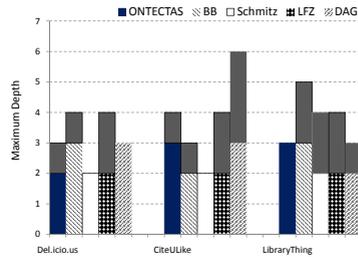
(a) Precision



(b) Relative Recall — IS-A



(c) Average Depth



(d) Maximum Depth

Fig. 1. Comparison of ONTECTAS to other algorithms for different metrics. Lower bars show IS-A relationships and higher bars show “any” relationships.

Figure 1(a) shows the algorithms’ precision for both IS-A relationships (the lower bars) and ANY relationships (the higher bars); for IS-A relationships, the precision of ONTECTAS is 0.50 for Del.icio.us, 0.48 for LibraryThing and 0.29 for CiteULike. ONTECTAS outperforms the precision of all other algorithms on all datasets. We also compare our precision on IS-A with that on IS-A+ ANY for the other algorithms since they do not distinguish IS-A from non-IS-A. Even then ONTECTAS outperforms the other algorithms in del.icio.us and CiteULike. On LibraryThing, the performance is close to the winner.

Figure 1(b) compares the algorithms’ relative recall for IS-A relationships. ONTECTAS is the best performer for all three datasets. One reason for DAG-ALG’s bad relative recall is that it detected many popular tags such as “web” and “software” as subjective tags, and pruned them before discovering the edges. BB had relatively low precision and recall in CiteULike because it detected many relationships with the tag “no-tag”, which is a popular tag rather than an ontological tag. ONTECTAS performs the best for relative recall for ANY relationships [20].

Figures 1(c) and 1(d) measure the depth of the validated ontology detected by each algorithm for both IS-A (lower bars) and ANY relationships (higher bars). These measures quantify the richness of the ontology. If there are multiple paths from the root to a node n , the depth is the longest path. Because the other algorithms find just ANY relationship between elements in an ontology, rather determining the types of relationships, like ONTECTAS does, we measure both the IS-A relationships and ANY relationships found. We do not consider HAS-A since no other algorithms detect it. Notice that this gives an advantage to the competing algorithms. For the depth metrics, other algorithms usually find a long chain with combination of synonyms and IS-A relationships. Since ONTECTAS detects mostly IS-A or HAS-A (and not synonyms), maximum depth for ANY relationship in ONTECTAS is close to maximum depth of IS-A relationship because in general of chains containing IS-A and HAS-A are rare.

For IS-A relationships, ONTECTAS has the highest maximum depth for two out of three datasets. In the full version of the paper [20], we show that the average number of children is similar to the average depth. For the average number of children, ONTECTAS has the best performance for CiteULike, is roughly tied for Library thing, and is second best for Del.icio.us.

Even when competing algorithms are given credit for ANY relationships and ONTECTAS only for finding IS-A, ONTECTAS performs fairly well. This is because there are so many IS-A relationships detected as compared to the other relationship types.

For all of the depth/children metrics, we note that *all algorithms perform markedly better using our preprocessing step of removing verb phrases*. This step helped a lot in removing non-ontological tags such as “to-read” in the Del.icio.us dataset. By applying this to all algorithms, we have improved all algorithms’ performance, not just ONTECTAS’s. Figure 1 also shows that most of the algorithms performed better on most measures for the Del.icio.us and Library-Thing datasets than on CiteULike. This validates the fact that the tags in these datasets are of better quality than the ones in CiteULike. This shows that we can compare different CTSs on the quality of tagging actions, using an ontology creation algorithm.

In summary, ONTECTAS outperforms the four other algorithms on precision and relative recall for IS-A relationships, and does well on the structural metrics of maximum depth, average depth and average number of children.

7.5 Comparing with a Gold Standard

Following [19], we compared how the algorithms extracted IS-A relationships against a “gold standard” ontology — the concept hierarchy from the Open Directory Project (ODP)⁶. To judge precision, recall, and F-measure, we use the lexical and taxonomic metrics from [19]. The lexical metrics measure how well the algorithms did in recreating the *concepts*, and the taxonomic metrics show how well the algorithms did in recreating the *structure*. Notice that comparing with a static ontology considered as gold standard has its problems since it

⁶ <http://dmoz.org>

Table 3. Gold standard based lexical and taxonomic comparison

	Lexical					Taxonomic				
	ONT	LFZ	BB	DAG	Schmitz	ONT	LFZ	BB	DAG	Schmitz
Precision	0.261	0.743	0.183	0.745	0.128	0.480	0.077	0.434	0.123	0.329
Recall	0.240	0.006	0.244	0.025	0.007	0.723	0.023	0.711	0.783	0.256
F-Measure	0.044	0.011	0.043	0.049	0.014	0.577	0.035	0.539	0.212	0.288

may miss important concepts and relationships and a good algorithm that finds concepts and relationships manually verified to be correct may get penalized unfairly. We will return to this point. The full version of this paper [20] shows the formal definitions of the measures and the detailed results. Due to space limitations, we only cover the highlights in this paper.

We looked at the 25 highest-level concepts common across the five algorithms. Table 3 shows the results. Bolded entries represent the best performance.

ONTECTAS has the second highest overall lexical recall and f-measure, which shows that it did well at finding the desired concepts. While DAG had the highest lexical precision and f-measure, and BB had the highest lexical recall, they both did very poorly on taxonomic precision, leading to a low taxonomic f-measure.

LFZ had a very good lexical precision; however, this is achieved by reporting a very small number of correct concepts. ONTECTAS is superior to LFZ in terms of all three taxonomic measures.

Because the 25 highest level common concepts were very uneven in size, we performed an analysis of the 6 largest subtrees — otherwise algorithms would be testing against subtrees that were only one or two concepts large. When we considered only the 6 largest subtrees, ONTECTAS had the best lexical and taxonomic f-measure.

Comparing to a gold standard shows how well algorithms do against a manually created ontology. But since a gold standard ontology is static, this metric may unfairly penalize algorithms that genuinely find correct concepts and relationships. E.g., “dialect” and “software IS-A technology” is incorrect according to this standard. Thus, comparing algorithms should take into account other components discussed above as well.

8 Conclusion and Future Work

We proposed an algorithm (ONTECTAS) for building ontologies of keywords from collaborative tagging systems. ONTECTAS uses association rule mining, bi-gram pruning, exploiting pairs of tags with the same child, and lexico-syntactic patterns to detect relationships between tags. We also provided a thorough analysis of ONTECTAS and how it compares to other algorithms. Some of the important open problems include detecting spam users, improving accuracy of ontology extraction via supervised learning and by means of incorporation of part-of-speech detection. Our ongoing work addresses some of these.

References

1. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules in large databases. In: VLDB, pp. 487–499 (1994)
2. An, Y.J., Geller, J., Wu, Y.-T., Chun, S.A.: Automatic generation of ontology from the deep web. In: Database and Expert Systems Applications (2007)
3. Barla, M., Bieliková, M.: On deriving tagsonomies: Keyword relations coming from crowd. In: Conference on Computational Collective Intelligence (2009)
4. Berland, M., Charniak, E.: Finding parts in very large corpora. In: Annual Meeting of the Association for Computational Linguistics, pp. 57–64 (1999)
5. Eda, T., Yoshikawa, M., Uchiyama, T.: The effectiveness of latent semantic analysis for building up a bottom-up taxonomy from folksonomy tags. *World Wide Web* 12(4), 421–440 (2009)
6. Girju, R., Badulescu, A., Moldovan, D.: Automatic discovery of part-whole relations. *Comput. Linguist.* 32(1), 83–135 (2006)
7. Golder, S., Huberman, B.A.: The structure of collaborative tagging systems. *Journal of Information Science* 32(2), 198–208 (2005)
8. Han, J., Pei, J., Yin, Y., Mao, R.: Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data Mining and Knowledge Discovery* 8(1), 53–87 (2004)
9. Hearst, M.A.: Automatic acquisition of hyponyms from large text corpora. In: Conference on Computational linguistics, pp. 539–545 (1992)
10. Heymann, Garcia-Molina.: Collaborative creation of communal hierarchical taxonomies in social tagging systems. Technical Report 2006-10, Stanford (2006)
11. Hotho, A., Jäschke, R., Schmitz, C., Stumme, G.: BibSonomy: A social bookmark and publication sharing system. In: Conceptual Structures Tool Interoperability Workshop at the International Conference on Conceptual Structures (2006)
12. Keller, F., Lapata, M.: Using the web to obtain frequencies for unseen bigrams. *Computational Linguistics* 29(3), 459–484 (2003)
13. Keller, F., Lapata, M., Ourioupina, O.: Using the web to overcome data sparseness. In: ACL Conference on Empirical Methods in NLP, pp. 230–237 (2002)
14. Körner, C., Benz, D., Hotho, A., Strohmaier, M., Stumme, G.: Stop thinking, start tagging: tag semantics emerge from collaborative verbosity. In: WWW (2010)
15. Kruse, P.M., Naujoks, A., Rsnier, D., Kunze, M.: Clever search: A wordnet based wrapper for internet search engines. In: Proceedings of the 2nd GermaNet Workshop (2005)
16. Laniado, D., Eynard, D., Colombetti, M.: Using wordnet to turn a folksonomy into a hierarchy of concepts. In: Semantic Web Application and Perspectives - Fourth Italian Semantic Web Workshop, pp. 192–201 (December 2007)
17. Lapata, M., Keller, F.: Web-based models for natural language processing. *ACM Transactions on Speech and Language Processing* 2, 1–31 (2005)
18. Lin, H., Davis, J., Zhou, Y.: An integrated approach to extracting ontological structures from folksonomies. In: Aroyo, L., Traverso, P., Ciravegna, F., Cimiano, P., Heath, T., Hyvönen, E., Mizoguchi, R., Oren, E., Sabou, M., Simperl, E. (eds.) *ESWC 2009. LNCS*, vol. 5554, pp. 654–668. Springer, Heidelberg (2009)
19. Liu, K., Fang, B., Zhang, W.: Ontology emergence from folksonomies. In: CIKM, pp. 1109–1118 (2010)
20. Moosavi, A., Li, T., Lakshmanan, L.V., Pottinger, R.: ONTECTAS: Bridging the gap between collaborative tagging systems and structured data (full version), <http://www.cs.ubc.ca/~rap/ontectas.pdf>

21. Sánchez, D., Moreno, A.: Learning non-taxonomic relationships from web documents for domain ontology construction. *DKE* 64(3), 600–623 (2008)
22. Schmitz, C., Hotho, A., Jäschke, R., Stumme, G.: Mining association rules in folksonomies. In: *Classification, Data Analysis, and Knowledge Organization* (2006)
23. Schmitz, P.: Inducing ontology from flickr tags. In: *Collaborative Web Tagging Workshop at WWW* (2006)
24. Schwarzkopf, E., Heckmann, D., Dengler, D., Kroner, E.: Mining the structure of tag spaces for user modeling. In: *Wkshp. on Data Mining for User Model.* (2007)
25. Suchanek, F.M., Kasneci, G., Weikum, G.: Yago: a core of semantic knowledge. In: *WWW*, pp. 697–706 (2007)
26. Zhu, X., Rosenfeld, R.: Improving trigram language modeling with the world wide web. In: *Acoustics, Speech, and Signal Processing*, pp. 533–536 (2001)