

# Enhancing the Navigability of Social Tagging Systems with Tag Taxonomies

Christoph Trattner<sup>†‡</sup>  
ctrattner@iicm.edu

Christian Körner<sup>†</sup>  
christian.koerner@tugraz.at

Denis Helic<sup>†</sup>  
dhelic@tugraz.at

<sup>†</sup> Institute for Information Systems and Computer Media

<sup>‡</sup> Knowledge Management Institute  
Graz University of Technology, Austria

## ABSTRACT

Tagging introduces an intuitive and easy method to organize resources in information systems. Although tags exhibit useful properties for e.g. personal organization of information, recent research has shown that the navigability of social tagging systems leaves much to be desired. When browsing social tagging systems users often have to navigate through huge lists of potential results before arriving at the desired resource. Thus, from a user point of view tagging systems are typically hard to navigate. To overcome this issue, we present in this paper a novel approach to supporting navigation in social tagging systems. We introduce tag-resource taxonomies that aim to support efficient navigation of tagging systems. To that end, we introduce an algorithm for the generation of these hierarchical structures. We evaluate the proposed algorithm and hierarchies from a theoretical, semantic and empirical point of view. With these evaluations we are able to show the high performance and usefulness of the proposed hierarchies.

## Categories and Subject Descriptors

H.4 [Information Systems Applications]: Miscellaneous

## General Terms

Tag Navigation

## Keywords

Tagging systems, navigation, tag taxonomy, resource taxonomy

## 1. INTRODUCTION

Tagging provides an easy and intuitive way to annotate, organize and retrieve resources on the web. For this reason the popularity of social tagging systems has increased tremendously in recent years. To give some examples: Delicious<sup>1</sup> enables the annotation of personal bookmarks with tags, Flickr<sup>2</sup> allows users to describe their

<sup>1</sup><http://www.delicious.com/>

<sup>2</sup><http://www.flickr.com/>

photos by tagging and Youtube<sup>3</sup> supports easier finding of videos via tags by content creators.

While there has been a lot of work on the structure of social tagging systems, little is known about the ways users use and navigate such systems. Previous work by Chi et al. [7] observed that the navigability leaves much to be desired. There, the authors showed that the number of new tags does not grow as quickly as the number of tagged resources in mature social tagging systems such as BibSonomy, CiteULike or Delicious. Therefore a lot of tags exist that refer to a large number of documents within such systems. To illustrate this problem from a user perspective: when users click on a popular Delicious tag such as “web” they retrieve 6.5 million resources in reverse chronological order – thus, rendering the system unusable from a navigational point of view.

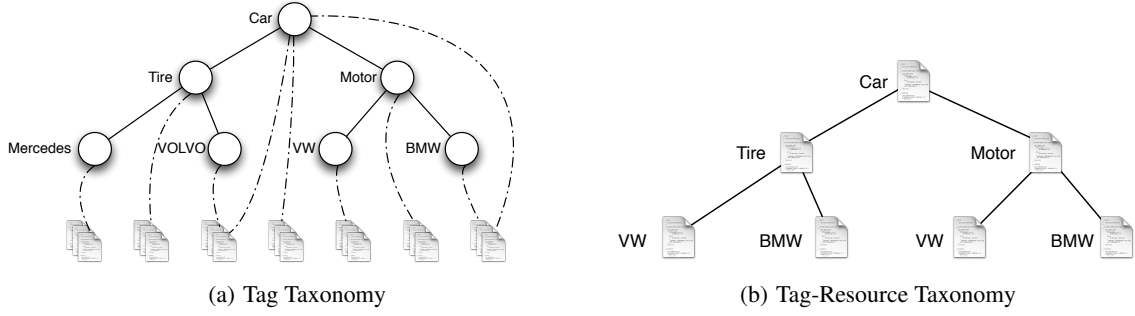
To overcome this issue recent research has investigated methods and strategies to make tagging systems more navigable. One prominent example of such endeavors are so-called tag taxonomies [18] – a method which allows the user to navigate to related concepts (tags) in a tagging system in a hierarchical and efficient manner (see also [16] for evaluation of several similar approaches). In this paper we introduce the notation of *tag-resource taxonomies*. Contrary to the idea of tag taxonomies, this approach enables the users not only to quickly navigate to related concepts but also to *resources* from a tagging system. With the approach of tag taxonomies, as it will be shown in this paper, efficient navigation to the resources of the tagging system is not possible. In a theoretical, semantic and empirical evaluation we show a high performance and usefulness of tag-resource taxonomies. To the best of our knowledge this is the first work that describes the notion of tag-resource taxonomies. Moreover, this approach significantly improves navigability of social tagging systems when compared with tag taxonomy approaches.

This paper is structured as follows: In Section 2 we introduce a novel approach to construct tag-resource hierarchies and illustrate the algorithms that were created for this purpose. This is followed by Section 3 explaining our evaluation. Section 4 gives an overview of related work. Finally Section 5 we conclude our findings and point to future work.

## 2. APPROACH

To tackle the issue of poor navigability in tagging systems, we introduce a novel approach to generate *tag-resource taxonomies*. The

<sup>3</sup><http://www.youtube.com/>



**Figure 1: Tag Taxonomy vs. Tag-Resource Taxonomy.**

goal of the approach is to offer the user a simple tool to navigate the tagging system in an efficient way. According to Kleinberg [20], efficient navigation in a network is possible if all resources are navigable in a polynomial of  $\log(n)$ , where  $n$  is the number of resources in the network. With the approach of tag-resource taxonomies, and as it is shown in Section 2.1, this prerequisite is fulfilled, i.e. it is possible to navigate a tagging system in a polynomial of  $\log(n)$ .

Basically, a tag-resource taxonomy is a hierarchy containing both resources and tags. The basis of a tag-resource taxonomy is the so-called *resource taxonomy*. A resource taxonomy is a hierarchy where the resources of a tagging system are arranged in a unique and taxonomic way, i.e. each resource of the tagging system occurs only once and parent nodes are more general than their child nodes.

Given such a resource taxonomy we construct the final tag-resource taxonomy by using a labeling algorithm that applies tag information to each resource in a descriptive and general manner. Hence, each resource in the resource taxonomy has one tag label attached to describe the underlying resource. The resulting tag-resource taxonomy presented to the user is then a tag hierarchy where the tags refer to a constant number of resources.

Figure 1 gives an example of a tag taxonomy compared to a tag-resource taxonomy. In a tag taxonomy tags appear only once in the hierarchy. However, resources can be referred by different tags. In a tag-resource taxonomy on the other hand resources occur only once while tags can appear on multiple and on different levels.

## 2.1 Why Usefulness of Tag Taxonomies for Navigation is Limited

### 2.1.1 Maximum Number of Clicks

A tag taxonomy allows the user to navigate to a designated tag (concept) efficiently, but navigation to a particular resource is still a problem due to the so-called pagination effect. As shown by [14] in tagging systems the tag-resource distribution follows a power-law function (see Figure 2), i.e. there are many tags that refer to a large number of resources. In BibSonomy or CiteULike for instance there are tags, which refer to hundreds or even thousands of resources. To make such frequently used tags still usable for the user, developers typically paginate the result list of such tags by a certain factor  $k$ . Hence, in the worst case the user has to click through the whole paginated result list to find the desired resource.

In detail, in the worst case the user would have to click

$$\max\{click(T_{tag})\} = \frac{|\max\{t\}|}{k} + \max\{depth(T_{tag})\} \quad (1)$$

times to reach a designated target resource with the approach of a tag taxonomy.

The term  $|\max\{t\}|$  in Equation 1 describes the size of the most frequently used tag in the tagging system. The term  $k$  stands for the pagination factor and  $\max\{depth(T_{tag})\}$  denotes the maximum depth of the tag taxonomy. As shown in [30] the size of the most frequently used tag can be estimated as  $|\max\{t\}| = c_1 \cdot |r|$ , where  $c_1$  is a constant typically ranging between  $[0.1, \dots, 0.2]$  and  $|r|$  is the number of unique resources in the tagging system.  $\max\{depth(T_{tag})\}$  can be estimated as,  $\log_{b/2} |t|$ , supposing that  $T_{tag}$  is a complete and fixed branched tree with branching factor  $b$ . The factor  $|t|$  describes the number of unique tags in the tagging system.  $|t|$  can be estimated as  $|t| = c_2 \cdot |r|$ , where  $c_2$  is a constant. Therefore, Equation 1 can be formalized as

$$\max\{click(T_{tag})\} = \frac{c_1 \cdot |r|}{k} + \log_{b/2}(c_2 \cdot |r|), \quad b \geq 2 \quad (2)$$

or

$$\max\{click(T_{tag})\} \approx \frac{c_1 \cdot |r|}{k} \quad (3)$$

supposing that  $\log_{b/2}(c_2 \cdot |r|) \ll \frac{c_1 \cdot |r|}{k}$ .

By generating a tag-resource taxonomy the worst case scenario is significantly better, especially for large numbers of  $|r|$ . Suppose the tag-resource taxonomy  $T_{res}$  is complete and has a fixed branching factor  $b$ , with  $b = k$ . A user would have to click

$$\max\{click(T_{res})\} = \max\{depth(T_{res})\} = \log_{k/2} |r|, \quad k \geq 2 \quad (4)$$

times in the worst case to reach a designated target resource. Then for large values of  $|r|$  we have:

$$\log_{k/2} |r| \ll \frac{c_1 \cdot |r|}{k} \quad (5)$$

Hence, according to the definition of Kleinberg [20] (see Section 2), and contrary to tag taxonomies, tag-resource taxonomies allow the user to navigate to the resources of a tagging system in an efficient manner, i.e. in a polynomial of  $\log(n)$ .

To give an example: Let us calculate the number of maximum clicks for the tag datasets as presented in Table 1 and compare the resulting tag taxonomy and tag-resource taxonomy for  $k = 10$ .

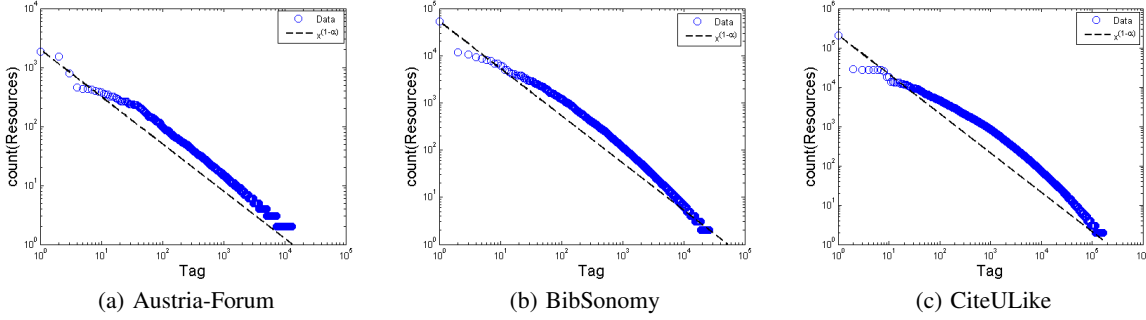


Figure 2: Tag distributions of the three data sets.

As shown in Table 2 (see  $\max\{click(T_{tag})\}$ ) with a tag taxonomy the user would have to click 184 times in the Austria-Forum tag dataset, respectively 5, 278 and 20, 799 times in the BibSonomy and CiteUlike tag dataset, to reach a desired resource in the worst case. Compared to this, with a tag-resource taxonomy, a user would have to click only 6.1 times in the Austria-Forum, 7.7 times in the BibSonomy and 8.5 times in the CiteULike tagging system to reach any designated target resource in the worst case (see  $\max\{click(T_{res})\}$ ).

	Austria-Forum [2]	BibSonomy [8]	CiteULike [5]
$ r $	19,430	233,712	949,851
$ t $	13,314	26,285	163,642
$\max\{t\}$	1,838	52,777	207,990
$\alpha$	2.2	1.9	2.0

Table 1: Statistics of Austria-Forum, BibSonomy and CiteU-Like tag dataset.

	Austria-Forum	BibSonomy	CiteULike
$\max\{click(T_{tag})\}$	184	5,278	20,799
$\max\{click(T_{res})\}$	6.1	7.7	8.5

Table 2: Tag Taxonomy vs. Tag-Resource Taxonomy: Maximum number of clicks for  $k = 10$ .

### 2.1.2 Number of Paginated Tags

Now, in order to calculate the number of tags suffering from the pagination effect we define the following equations: Since we know that the tag distribution (see Figure 2) has power-law qualities we approximate the number of paginated tags  $|t_p|$  as follows [9]

$$r_i = \frac{\alpha - 1}{t_{min}} \cdot \left(\frac{t_i}{t_{min}}\right)^{-\alpha}, \quad t_{min} > 0 \quad (6)$$

The parameter  $\alpha$  can be approximated with the method of maximum likelihood as

$$\alpha \cong 1 + |t| \left[ \sum_{i=1}^{|t|} \ln \frac{t_i}{t_{min}} \right]^{-1} \quad (7)$$

With  $r_i = k$  and  $t_{min} = 1$ , resolved by  $t_p$  the number of paginated tags  $|t_p|$  is then

$$|t_p| = |t| \cdot \left(\frac{\alpha}{k} - \frac{1}{k}\right)^{\left(\frac{1}{\alpha}\right)} \quad (8)$$

Example: Let us calculate the number of paginated tags for the tag datasets as shown in Table 1 for  $k = 10$ . Then, as shown in Table

3, within the Austria-Forum dataset 38% of all tags suffer from the pagination effect, respectively 28% for in the BibSonomy tag dataset and 32% in the CiteULike tag dataset. Or in other words, for a commonly used resource list of the length of  $k = 10$ , nearly 1/3 of all tags suffer from the pagination effect, i.e. the resources of such tags are not navigable in an efficient way!

	Austria-Forum	BibSonomy	CiteULike
$ t_p $ (%)	5079 (38%)	7401 (28%)	51748 (32%)

Table 3: Number of paginated tags for  $k = 10$ .

### 2.1.3 Mean Number of Clicks

Last but not least, we can approximate the mean number of clicks a user would need to reach a designated target resource in a tagging system navigating via a tag-resource taxonomy as follows:

$$\text{mean}\{click(T_{res})\} = \log_k(|r|) \quad (9)$$

The mean number of clicks with a tag taxonomy can be approximated as:

$$\text{mean}\{click(T_{tag})\} = \log_k(|t|) + \frac{1}{|t|} \sum_{i=1}^{|t|} \frac{r_i}{k} \quad (10)$$

In Table 4 example calculations for the mean number of clicks of tag taxonomies and tag-resource taxonomies with different branching factors  $k$  for different tag datasets are presented. As shown, on average, tag-resource taxonomies support the user with significantly less clicks (see  $\text{mean}\{click(T_{res})\}$ ) in navigating the resources of a tagging system than the approach of tag taxonomies (see  $\text{mean}\{click(T_{tag})\}$ ).

	k	Austria-Forum	BibSonomy	CiteULike
$\text{mean}\{click(T_{res})\}$	2	14.2	17.8	19.8
$\text{mean}\{click(T_{tag})\}$	2	29.5	22.4	30.7
$\text{mean}\{click(T_{res})\}$	5	6.1	7.6	8.5
$\text{mean}\{click(T_{tag})\}$	5	11.6	9.2	12.3
$\text{mean}\{click(T_{res})\}$	10	4.3	5.3	5.9
$\text{mean}\{click(T_{tag})\}$	10	6.4	5.6	7.3

Table 4: Tag Taxonomy vs. Tag-Resource Taxonomy: Mean number of clicks for different branching factors  $k$ .

## 2.2 Description of the Algorithm

### 2.2.1 Resource Taxonomy Generation Algorithm

As described in Section 2 the basis of the tag-resource taxonomy is the so-called resource taxonomy – a taxonomy where the resources of the tagging system are arranged in a taxonomic manner. In order to generate a resource taxonomy from tagging data we developed Algorithm 1. In words, the algorithm works as follows:

The algorithm takes a tag dataset and the desired taxonomy branching factor as input parameters. Since the algorithm should generate a resource taxonomy with the most general resource of the tagging system as root node and related and less general resources as children, the algorithm calculates in the first step degree centrality for all resource of the supplied tagging dataset and stores the centrality-resource pairs into a map  $C$ . Degree centrality was chosen since, on the one hand, it is computed fast, and on the other hand, it was observed in our previous research [4] that degree centrality in tagging systems is highly correlated to sophisticated centrality measures such as closeness or betweenness centrality. In the next step, the algorithm sorts the resources in  $C$  according to their centrality values in descending order.

Subsequently, the algorithm takes the first element of  $C$  (i.e. the most general resource) and sets that resource as the root node of the resource taxonomy. Thereafter, the algorithm starts iterating through the elements (resources) already present in resource taxonomy. For each resource in the resource taxonomy the algorithm calculates then the most similar resources (see *getMoreLikeThis* in Algorithm 1). Our algorithm calculates cosine similarity between all co-occurring resources taking also the  $tf \cdot idf$  values of the tag concepts into account. Additionally, the function returns only resources that are not already part of the constructed resource taxonomy. The results of this function are stored into a map  $SIM$ , with resources as key values and with the provided similarity values as corresponding map values. To account for resource generality we multiply resource similarity values with their corresponding centrality values. The final scores are normalized to fall into the range of  $[0...1]$ . After that, the resources in  $SIM$  are sorted by the scores in descending order. This procedure ensures that the resources in  $SIM$  are not only similar to the currently processed resource but also sorted by their generality values. Thereafter the algorithm appends a maximum of  $b$  resources to the currently processed resource. The algorithm stops, if no more similar resources can be found.

Note, due the fixed branching factor  $b$  the algorithm does not guarantee that all resources of the tagging dataset are contained in the resulting resource taxonomy. However, as it will be shown in Section 5 the probability that one or even more resources are missing is relatively small due to the high number of existing links between the resources of the resource-to-resource network of a given tag dataset. On the other hand, in a tag taxonomy the probability that one concept is missing is significantly higher. The reason for this behavior is the fact that the tag-to-tag network of a tagging system is typically substantially less connected.

### 2.2.2 Tag-Resource Taxonomy Generation Algorithm

To produce the final tag-resource taxonomy on the basis of generated resource taxonomy we developed Algorithm 2. In general it is a labeling algorithm taking a given resource taxonomy and a tagging dataset as input parameters. Tag information is used as label data. The algorithm tries to apply labels to the given resource taxonomy in such a way, that they are uniquely distinguishable and

---

### Algorithm 1 Resource Taxonomy Algorithm

---

```
INPUT: Tag Dataset  $D$ , Branching Factor  $b$ 
OUTPUT: Resource Taxonomy  $T$ 
 $C \leftarrow$  new HashMap[]
 $T \leftarrow$  new Tree[]
for each  $r_i \in F$  do
     $C[r_i] \leftarrow$  calculate degree centrality
end for
sortByValues( $C$ )
/*sort  $C$  by values in descending order*/
 $T[0] \leftarrow C[0]$ 
 $SIM \leftarrow$  new HashMap
for  $i = 0; i < \text{sizeof}(T); i++$  do
    /*get all similar resources of  $T[i]$  and store the resources as key values
    and the similarity values into  $SIM$ */
     $SIM \leftarrow$  getMoreLikeThis( $T[i]$ )
    for each  $r_i \in SIM$  do
         $T[r_i] \leftarrow T[r_i] \cdot C[r_i]$ 
    end for
    /*sort the resources in  $SIM$  by values in descending order*/
    sortByValues( $SIM$ )
    for  $j = 0; j < \text{sizeof}(SIM)$  and  $j < b; j++$  do
         $T[i].\text{append}(SIM[j])$ 
    end for
end for
return  $T$ 
```

---

---

### Algorithm 2 Tag-Resource Taxonomy Algorithm

---

```
INPUT: Resource Taxonomy  $T$ , Tag Dataset  $D$ 
OUTPUT: Tag-resource Taxonomy
 $COTAGS \leftarrow$  new HashMap[newArray][]
for  $i = 0; i < \text{sizeof}(T); i++$  do
     $Ts \leftarrow$  getTags( $T[i], D$ )
    for  $j = 0; j < \text{sizeof}(Ts); j++$  do
         $cotags \leftarrow$  getCoocTags( $Ts[j], D$ )
        sort( $cotags$ )
        remove all tags from  $cotags$  that are not contained in  $T[i]$ 
         $COTAGS[T[i]].\text{add}(cotags)$ 
    end for
end for
 $trails \leftarrow$  new HashSet[]
for each  $r_i \in T$  do
    /* $T$  is traversed in left-order*/
     $pl \leftarrow$  getParentLabels( $r_i$ )
    for each  $l_j \in COTAGS[r_i]$  do
        if  $!pl.\text{contains}(l_j)$  then
            if  $!(trails.\text{contains}(pl.\text{toString}() + l_j))$  then
                 $T[r_i].\text{applyLabel}(pl)$ 
                 $trails.\text{add}(pl.\text{toString}() + l_j)$ 
            end if
        end if
        if  $T[r_i]$  has no label then
             $T[r_i].\text{applyLabel}(\text{getTitle}(r_i))$ 
        end if
    end for
end for
return  $T$ 
```

---

the most descriptive for the given resource. The candidate tags are thereby ranked by the method of tag co-occurrence. However, since it can happen that resources in the resource taxonomy have the same tags in their parent tag trail, due to the lack of available tags in the tagging system, additional meta-data is taken into account. We use title information of the resources as an additional way for differentiation.

In words the algorithm works as follows: In the first step the algorithm calculates, for each resource in the resource taxonomy a list of co-occurring tags of all resource tags and stores this list sorted

Name	b	n	$\max\{click(T_{res})\}$	$\text{mean}\{click(T_{res})\}$
Res2	2	19,430	17	12.45
Res5	5	19,430	10	5.93
Res10	10	19,430	8	4.44

**Table 5:  $\max\{click(T_{res})\}$  and  $\text{mean}\{click(T_{res})\}$  for different branching factors  $b$ .**

in descending order into a map. After that, the algorithm traverses the resource taxonomy in left-order. In this loop the actual labeling procedure is performed. In detail, the labeling process looks as follows: For each resource in the resource taxonomy the corresponding co-occurrence vector is consulted and the first label, i.e. the most frequent tag, is tried to be applied to the currently processed resource. If the currently used candidate tag is already part of the tag trail of the currently processed resource (see variable *trails* in Algorithm 2) the next element, i.e. the next frequent tag label is chosen as candidate tag. If no uniquely distinguishable tag trail can be constructed, i.e. the candidate tag label from the co-occurrence vector is already present in the tag trail of the resource additional meta data is considered. We use title information of the currently processed resource for this purpose. Note, since tag and title information can be identical the proposed method is not completely free of collisions. However, to fix this issue one can include additional meta data information or other methods to generate a unique label such as appending an iterative number for each label that occurs more than once. The algorithm stops if all resources of the given resource taxonomy are labeled.

Figure 3 shows the branching factor distribution for a tag-resource taxonomy with branching  $b = 5$  generated from the Austria-Forum tag dataset. For branching factor  $b = 5$  the algorithm does not generate a complete  $b$ -tree (from levels 1 to 4 the resulting tree is complete, for levels  $> 4$  the tree is not complete). The reason for this behavior is the fact that in tag networks there are resources which are just connected to a few resources, i.e. if the branching factor  $b$  is beneath this threshold the resulting taxonomy becomes incomplete.

### 3. EVALUATION

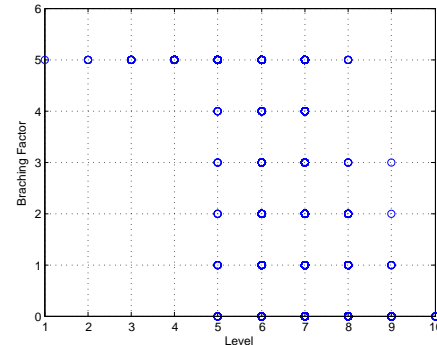
Now, since we have shown in theory that the approach of the so-called tag-taxonomies allows the user to navigate to the resource of a tagging system in an efficient way, we will provide in the following section results of four different experiment to show the usefulness of the proposed approach also in a practical setting.

#### 3.1 Dataset

For the experiments, we used the tag dataset from the *Austria-Forum* [28]. The Austria-Forum is a large online encyclopedia similar to Wikipedia providing the user with around 180,000 resources on topics related to Austria. In contrast to Wikipedia, Austria-Forum offers an integrated tagging system [29], which allows users to assign tags to resources and to navigate to related resources via tag clouds. As of October 16<sup>th</sup>, 2010 the Austria-Forum tag dataset contains 13,314 tags, 19,430 resources and 97,908 tag assignments (see also Table 1).

#### 3.2 Measuring the Average and Maximum Number of Clicks and the Drop Rate

In a first experiment we investigated average and maximum tag-resource taxonomy depths for different branching factors  $b$  in order to measure the number of clicks a user would need to reach a designated target resource in the taxonomy. Furthermore we examined



**Figure 3: Example of a branching factor distribution for a tag-resource taxonomy with maximum branching  $b = 5$ .**

the number of missing resources (=drop rate) after the generation of a tag-resource taxonomy from tagging data with different branching factors  $b$ .

Since the resulting tag-resource taxonomies generated by Algorithm 2 are not complete the average nor the maximum depth of the taxonomy can be estimated by formulas. If the tag-resource taxonomy was complete, we could calculate the maximum number of clicks as  $\max\{click(T_{res})\} = \log_{b/2}(n)$ , where  $n$  are the number of nodes in the taxonomy. Hence, these values were conducted empirically through an experiment.

For the experiment three different tag-resource taxonomies named *Res2*, *Res5* and *Res10* with three different branching factors  $b = 2, 5$  and  $10$  were generated. In order to compare the resulting taxonomies against a golden standard taxonomy the DMOZ Open Directory Project (ODP) taxonomy<sup>4</sup> was consulted. This experiment was conducted to determine whether the generated tag-resource taxonomy would be usable or not.

As shown in Table 5 the tag-resource taxonomy with the smallest branching factor  $b = 2$  is the deepest, i.e. a user would need 17 clicks to reach a target resource in the worst case. On the other, and as expected the tag-resource taxonomy with highest branching factor  $b = 10$  is less deepest taxonomy, i.e. in the worst case a user would have to click 8 times to reach a desired resource. For  $b = 5$  the worst case 10 clicks. On average for branching factor  $b = 2$  the mean number of clicks is 12.45. For  $b = 5$  the mean number of clicks is 5.93 and for  $b = 10$  4.44 clicks. The ODP Taxonomy has a mean depth of 6.86 [1]. The maximum depth is 13. Hence, compared to the ODP taxonomy the tag-resource taxonomy with branching factor  $b \geq 5$  will be most usable.

In order to measure the number of missing resources (=drop rate) after the generation process of the taxonomies, we simply calculated the number of resources contained in tag-resource taxonomies

<sup>4</sup><http://www.dmoz.org>

*Res2*, *Res5* and *Res10* and compared it to the number of unique resources contained in the original Austria-Forum tag dataset. As shown in Table 5 and represented as parameter  $n$  none of the resources dropped during the tag-resource taxonomy generation process. The reason for this behavior is the high number of existing links between the resources of the resource-to-resource network of the Austria-Forum tag dataset.

### 3.3 Measuring the Collision Rate

In the second experiment we measured the number of collisions when generating a tag-resource taxonomy with different branching factors  $b$ . As explained the tag-resource generation algorithm is not to 100% collision free, i.e. it could happen that in a tag trail of a given resource the same tags occurs twice or even more often. Hence, the goal of this experiment was to reveal how many collisions occur in general if a tag-resource taxonomy with a given branching  $b$  is created. For this experiment we used the three resource taxonomies from the former experiment. Table 6 shows collision rates for the three generated tag-resource taxonomies. All in all, we observe that the collision rate is relatively small. However to make the approach totally collision free one might use additional metadata as described in Section 2.2.

Name	$b$	$n$	$CR(\%)$
Res2	2	19,430	0.1%
Res5	5	19,430	0.2%
Res10	10	19,430	0.2%

**Table 6: Collision Rates (CR) for different resource taxonomies with different branching factor  $b$ .**

### 3.4 Measuring the Semantic Structure of the Tag-Resource Taxonomy

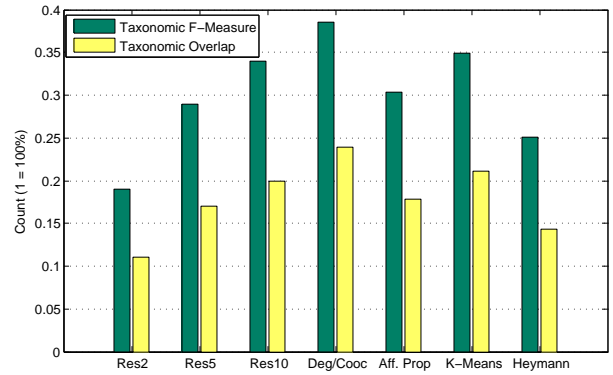
In the third experiment we measured the quality of the semantic structure of three tag-resource taxonomies that were generated for the two former experiments.

For that purpose, we consulted two different semantic measures – the Taxonomic F-Measure (in short  $TF$ ) [10] and the Taxonomic Overlap ( $TO$ ) [22]. Both measures identify the quality of a given taxonomy against a golden standard via common concepts. We used Germanet<sup>5</sup> as golden standard for the experiment since the Austria-Forum tag dataset contains only German tags.

To determine the overall semantic quality of our three generated tag-taxonomies four tag taxonomies on the basis of the following popular tag taxonomy induction algorithms were generated – Hierarchical K-Means [11], Affinity Propagation [12, 25], Heymann [18] and Deg/Cooc [16, 3]. In the experiment,  $TF$  and  $TO$  values for all seven taxonomies were measured and compared against one another. The goal of the experiment was to study how semantic structures generated by the tag-resource induction algorithm (Algorithm 2) compare to semantic structures produced by other popular tag taxonomy induction algorithms such as Hierarchical K-Means, Affinity Propagation, Heymann or Deg/Cooc [16].

Figure 4 shows the results of the semantic evaluation of the experiment. We observe, the higher the branching factor the better the semantic structure of the generated tag-resource taxonomies. The results indicate that tag-resource taxonomies with branching factors between  $b = [5...10]$  perform on average as good as normal

<sup>5</sup><http://www.sfs.uni-tuebingen.de/GermaNet/>



**Figure 4: Results of the semantic evaluation of the three generated tag-resource taxonomies *Res2*, *Res5* and *Res10*.**

tag taxonomies based on a Affinity Propagation tag taxonomy induction algorithm.

### 3.5 Empirical Analysis

In order to conduct whether the approach of a tag-resource taxonomy is also usable for humans we conducted a user study based on the ideas of [27].

First, we took the tag-resource taxonomy with branching factor  $b = 10$  and extracted 100 tag trails uniformly at random from the tag-resource taxonomy. After that, a Deg/Cooc tag taxonomy with a maximum branching factor of  $b = 10$  was generated in order to compare our approach of a tag-taxonomy to an existing method. Again, 100 tag trails were extracted uniformly at random from the generated tag taxonomy. Since shorter concept trails are typically evaluated better, we chose tag trails from both taxonomies that had a minimum tag trail length of 3 concepts (excluding the root node). After that, we presented the trails of both taxonomies in random order and generated an online test containing 200 tag trails, 837 relations and 1,037 concepts. Each of our users had to evaluate the exact same tag trails. The study participants received instructions on how to rate the trails and an exemplary taxonomy. During the test the users were asked to rate the trails according to the classification schema presented in Table 8.

Classification	Description
Correct	Correct hierarchy relation
Related	Correct relation, but not hierarchical or reverse hierarchical
Equivalent	Synonym
Not Related	The relations do not have anything to do with each other
Unknown	The evaluator does not recognize the meaning of the tag(s)

**Table 8: Classification Labels for the User Evaluation.**

Nine test subjects from three different departments at our university participated in the experiment. All participants were experienced computer users and familiar with user studies and the evaluation of concept hierarchies. The study was conducted between April 25<sup>th</sup> and 28<sup>th</sup> of 2011.

Table 7 shows the classification results of the user study. Compared to a tag taxonomy comprising only tags we see that concept relations of a tag-resource taxonomy with a branching factor  $b = 10$  are



Name	b	Correct (%)	Related (%)	Equivalent (%)	Not Related (%)	Unknown(%)
Deg/Cooc10	10	33.2	27.3	13	21.9	5.1
Res10	10	27.3	36.2	12.3	19.8	4.2

**Table 7: Results of the empirical analysis of the tag-resource taxonomy with branching factor  $b = 10$  compared to a Deg/Cooc tag taxonomy with branching factor  $b = 10$ .**

only to 5% less hierarchically arranged than the tag concepts of the in theory best semantically correct tag taxonomy approach the so-called Deg/Cooc tag taxonomy induction algorithm. Regarding the relatedness of the tag concepts we can observe that the tag-resource taxonomy was rated to 9% better than the Deg/Cooc tag taxonomy. Overall, the rating for the not related tags for both taxonomies was relatively small, taking into account that the maximum branching factor in both taxonomies was set to relatively small value of  $b = 10$ .

## 4. RELATED WORK

For the presented work the following research topics on tagging are relevant:

### 4.1 Analysis of Social Tagging Systems

The first analysis of social tagging systems was done by Golder and Huberman [13]. In this work the authors show stable usage patterns within collaborative tagging systems and introduce an initial model of collaborative tagging. Subsequent work by Marlow et al. [23] introduces another model which gives insight into a simple taxonomy of incentives and contribution models within these systems. Hammond et al. [15] give a high level overview of different social tagging tools and examine various aspects such as audience and types of tagged media.

### 4.2 Navigation in Social Tagging Systems

As previously mentioned Chi and Mytkowicz [6] studied Delicious using information theory (entropy) and found that the system becomes harder to navigate over time. The main reason for this is the small increase of tag vocabulary as opposed the vast growth of tagging information in these systems. In previous work [17] we analyzed tag clouds as means to browse tagging systems and showed that tag-resource networks have sufficient navigation properties in theory but also illustrated that user interface restrictions (such as pagination) spoil efficient navigation for all practical purposes.

### 4.3 Tag Semantics

In previous work [4] we compared different methods (such as network centrality, subsumption etc.) to measure the generality of tags in social systems. In [21] we showed that semantics within a social tagging system are heavily influenced by the users' tag usage. Users who are more verbose in the process of social tagging are better candidates for the construction of semantic structures out of folksonomies.

### 4.4 Creating Hierarchies from Social Tagging Data

Heymann et al. [18] converted a large corpus of tags annotating objects into a navigable hierarchical taxonomy of tags by evaluating the centrality of the tags in a similarity graph. In another work Solskinnsbakk et al. [27] constructed tag hierarchies using association rule mining of the corresponding tag set. Kiu and Tsui [19] introduced *TaxoFolk* - an algorithm which integrates a tags and

resources into a taxonomy by applying various data-mining techniques such as formal concept analysis. In another work Plangprasopchok [24] propose a hierarchy generation algorithm based on an examination of user-defined relations within the system. Schmitz [26] gives insight into an algorithm that induces an ontology from tags in the Flickr system using a subsumption-based model. However, contrary to our work, none of these previous approaches examine the implications the resulting structures have on the navigability of the system.

## 5. CONCLUSIONS AND FUTURE WORK

In this paper we introduced a novel approach to enhance the navigability of social tagging system through tag-resource taxonomies. We showed that tag taxonomies are in general well suited for finding related tag concepts, but perform worse in finding resources in an efficient number of clicks. By introducing the notation of the so-called tag-resource taxonomies we presented a method that tackles this issue. We illustrated in theory that with the approach of a tag-resource taxonomy it is possible to navigate to resources efficiently. Additionally to these findings, we evaluated the approach empirically and found that tag-resource taxonomies perform on a semantic level nearly as good or even better than other popular tag taxonomy approaches.

Thus, with the notation of tag-resource taxonomies we have introduced a novel hierarchical method that allows the user to navigate the resources in the tagging system in an efficient and semantically appropriate manner. To the best of our knowledge, this is the first work that describes such an efficient hierarchical navigation tool on the basis of tag-resource hierarchies.

In the future we plan to integrate a prototype visualization of tag-resource sub taxonomies into each article page of the Austria Forum in the form of tag trails in order to support the user in the process of navigating the system.

## 6. ACKNOWLEDGMENTS

We would like to thank Markus Muhr for helping us with the generation of the Affinity and K-Means tag taxonomies.

This work is funded by - BMVIT - the Federal Ministry for Transport, Innovation and Technology, program line Forschung, Innovation und Technologie für Informationstechnologie, project NAV-TAG – Improving the navigability of tagging systems and the European Commission as part of the FP7 Marie Curie IAPP project TEAM (grant no. 251514).

## 7. REFERENCES

- [1] S. Alexaki, V. Christophides, G. Karvounarakis, D. Plexousakis, and K. Tolle. The ics-forth rdfsuite: Managing voluminous rdf description bases. In *SemWeb*, 2001.
- [2] Austria-Forum. Das Österreichische Wissensnetz. <http://www.austria-lexikon.at>, 2011. [Online; accessed 2011-03-01].

- [3] D. Benz, A. Hotho, and G. Stumme. Semantics made by you and me: Self-emerging ontologies can capture the diversity of shared knowledge. In *Proc. of the 2nd Web Science Conference (WebSci10)*, Raleigh, NC, USA, 2010. Web Science Trust.
- [4] D. Benz, C. Körner, A. Hotho, G. Stumme, and M. Strohmaier. One tag to bind them all : Measuring term abstractness in social metadata. In *Proceedings of the 8th Extended Semantic Web Conference (ESWC 2011)*, Heraklion, Crete, May 2011.
- [5] BibSonomy. BibSonomy: A blue social bookmark and publication sharing system. <http://www.bibsonomy.org>, 2011. [Online; accessed 2011-04-21].
- [6] E. H. Chi and T. Mytkowicz. Understanding navigability of social tagging systems. In *proceedings of the SIGCHI conference on Human Factors in Computing Systems (CHI'07)*, 2007.
- [7] E. H. Chi and T. Mytkowicz. Understanding the efficiency of social tagging systems using information theory. In *HT '08: Proceedings of the nineteenth ACM conference on Hypertext and hypermedia*, pages 81–88, New York, NY, USA, 2008. ACM.
- [8] CiteULike. CiteULike: Everyone's library. <http://www.citeulike.org>, 2011. [Online; accessed 2011-04-21].
- [9] A. Clauset, C. R. Shalizi, and M. E. J. Newman. Power-law distributions in empirical data. *SIAM Rev.*, 51:661–703, November 2009.
- [10] K. Dellschaft and S. Staab. On how to perform a gold standard based evaluation of ontology learning. In *Proceedings of ISWC-2006 International Semantic Web Conference*, Athens, GA, USA, November 2006. Springer.
- [11] I. Dhillon, J. Fan, and Y. Guan. Efficient clustering of very large document collections. In *Data Mining for Scientific and Engineering Applications*. Kluwer Academic Publishers, Heidelberg, 2001.
- [12] B. J. J. Frey and D. Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972–976, January 2007.
- [13] S. A. Golder and B. A. Huberman. The structure of collaborative tagging systems. *Journal of Information Science*, 32(2):198–208, 2006.
- [14] H. Halpin, V. Robu, and H. Shepherd. The complex dynamics of collaborative tagging. In *Proceedings of the 16th international conference on World Wide Web*, WWW '07, pages 211–220, New York, NY, USA, 2007. ACM.
- [15] T. Hammond, T. Hannay, B. Lund, and J. Scott. Social bookmarking tools (i): A general review. *D-Lib Magazine*, 11(4), 2005.
- [16] D. Helic, M. Strohmaier, C. Trattner, M. Muhr, and K. Lermann. Pragmatic evaluation of folksonomies. In *Proc. of the 21st International World Wide Web conference*, WWW '11, New York, NY, USA, 2011. ACM.
- [17] D. Helic, C. Trattner, M. Strohmaier, and K. Andrews. Are tag clouds useful for navigation? a network-theoretic analysis. *International Journal of Social Computing and Cyber-Physical Systems*, 2011.
- [18] P. Heymann and H. Garcia-Molina. Collaborative creation of communal hierarchical taxonomies in social tagging systems. Technical Report 2006-10, Stanford InfoLab, April 2006.
- [19] C.-C. Kiu and E. Tsui. Taxofolk: A hybrid taxonomy-folksonomy structure for knowledge classification and navigation. *Expert Systems with Applications*, 38(5):6049 – 6058, 2011.
- [20] J. M. Kleinberg. Small-world phenomena and the dynamics of information. In *Advances in Neural Information Processing Systems (NIPS) 14*, page 2001, Cambridge, MA, USA, 2001. MIT Press.
- [21] C. Körner, D. Benz, M. Strohmaier, A. Hotho, and G. Stumme. Stop thinking, start tagging - tag semantics emerge from collaborative verbosity. In *Proc. of the 19th International World Wide Web Conference (WWW 2010)*, Raleigh, NC, USA, Apr. 2010. ACM.
- [22] A. Mädche and S. Staab. Measuring similarity between ontologies. In *Proc. Of the European Conference on Knowledge Acquisition and Management - EKAW-2002. Madrid, Spain, October 1-4, 2002*, volume 2473 of *LNC3/LNAI*, Heidelberg, 2002. Springer.
- [23] C. Marlow, M. Naaman, D. Boyd, and M. Davis. Ht06, tagging paper, taxonomy, flickr, academic article, to read. In *Proc. Seventeenth Conference on Hypertext and Hypermedia (Hypertext 2006)*, HT'06, pages 31–40, USA, NY, 2006. ACM.
- [24] A. Plangprasopchok and K. Lerman. Constructing folksonomies from user-specified relations on flickr. In *Proc. of 18th International World Wide Web Conference, WWW '09*, May 2009.
- [25] A. Plangprasopchok, K. Lerman, and L. Getoor. From saplings to a tree: Integrating structured metadata via relational affinity propagation. In *Proceedings of the AAAI workshop on Statistical Relational AI*, Menlo Park, CA, USA, July 2010. AAAI.
- [26] P. Schmitz. Inducing ontology from flickr tags. In *Proceedings of the Workshop on Collaborative Tagging at WWW2006*, Edinburgh, Scotland, May 2006.
- [27] G. Solskinnsbakk and J. Gulla. A hybrid approach to constructing tag hierarchies. In *On the Move to Meaningful Internet Systems, OTM 2010*, volume 6427 of *Lecture Notes in Computer Science*, pages 975–982. Springer Berlin / Heidelberg, 2010.
- [28] C. Trattner, I. Hasani-Mavriqi, D. Helic, and H. Leitner. The austrian way of wiki(pedia)!: development of a structured wiki-based encyclopedia within a local austrian context. In *Proceedings of the 6th International Symposium on Wikis and Open Collaboration, WikiSym '10*, pages 1–10, New York, NY, USA, 2010. ACM.
- [29] C. Trattner and D. Helic. Linking related documents: combining tag clouds and search queries. In *Proceedings of the 10th international conference on Web engineering, ICWE'10*, pages 486–489, Berlin, Heidelberg, 2010. Springer-Verlag.
- [30] C. Trattner, M. Strohmaier, D. Helic, and K. Andrews. The benefits and limitations of tag clouds as a tool for social navigation. *Technical Report – IICM, Graz University of Technology*, 2011.