

Class-based tag recommendation and user-based evaluation in online audio clip sharing

Frederic Font^a, Joan Serrà^b, Xavier Serra^a

^a*Music Technology Group, Universitat Pompeu Fabra, Barcelona, Spain*

^b*Artificial Intelligence Research Institute (IIIA-CSIC), Spanish National Research Council, Bellaterra, Spain*

Abstract

Online sharing platforms often rely on collaborative tagging systems for annotating content. In this way, users themselves annotate and describe the shared contents using textual labels, commonly called tags. These annotations typically suffer from a number of issues such as tag scarcity or ambiguous labelling. Hence, to minimise some of these issues, tag recommendation systems can be employed to suggest potentially relevant tags during the annotation process. In this work, we present a tag recommendation system and evaluate it in the context of an online platform for audio clip sharing. By exploiting domain-specific knowledge, the system we present is able to classify an audio clip among a number of predefined audio classes and to produce specific tag recommendations for the different classes. We perform an in-depth user-based evaluation of the recommendation method along with two baselines and a former version that we described in previous work. This user-based evaluation is further complemented with a prediction-based evaluation following standard information retrieval methodologies. Results show that the proposed tag recommendation method brings a statistically significant improvement over the previous method and the baselines. In addition, we report a number of findings based on the detailed analysis of user feedback provided during the evaluation process. The considered methods, when applied to real-world collaborative tagging systems, should serve the purpose of consolidating the tagging vocabulary and improving the quality of content annotations.

Keywords: Tag recommendation, User study, Folksonomy, Freesound

1. Introduction

Free-form semantically-meaningful textual labels, called tags, are extensively used in online sharing platforms for describing and annotating contents. Systems that provide the functionality for making these annotations are normally referred to as collaborative tagging systems. Several problems arise when users annotate shared and/or online resources [9]. The most typical ones are tag scarcity, the use of different tags to refer to a single concept (synonymy), the

8 ambiguity in the meaning of certain tags (polysemy), the commonness of ty-
9 pographical errors, the use of user-specific naming conventions, or the use of
10 different languages. To minimise some of these problems, tag recommendation
11 systems can be employed to suggest potentially relevant tags during the an-
12 notation process [14]. As users are exposed to the suggestions of the system,
13 the annotation process partially shifts from the creation of textual labels to the
14 recognition of tags in a list [23], and thus all users receive a certain common
15 influence from the system. Hence, tag recommendation serves the purpose of
16 consolidating the vocabulary of collaborative tagging systems [13].

17 In general, tag recommendations are either based on content analysis of on-
18 line resources or in the other tags that users introduce during the annotation
19 process. In the case of content-based recommendations, a typical approach con-
20 sists in, given a resource to be described, defining a neighbourhood of other
21 resources (based on some similarity measure) and then recommending tags that
22 are used to annotate resources in this neighbourhood [12, 24]. Another approach
23 is the use of machine learning techniques to learn mappings between tags and
24 content features [15, 25, 26]. On the other side, there are tag recommendation
25 strategies which are based on the tags that users introduce during the annota-
26 tion process itself, prior to the moment of the recommendation. Disadvantages
27 of these strategies compared to content-based recommendation methods are
28 that they require the existence of at least one tag to provide recommendations,
29 whereas content-based recommendation systems can provide recommendations
30 to resources with no associated tags or other metadata. Nevertheless, tag re-
31 commendation methods based on the tags that users introduce during the an-
32 notation process have the advantage of not requiring any specific processing of
33 the content of the resources being annotated, thus being typically less expensive
34 in terms of computation resources and being more easily generalisable to other
35 multimedia domains. These methods usually consider the *folksonomy* (i.e., the
36 set of associations between tags, users and content resources) of a collaborative
37 tagging system to estimate tag similarity from their resource co-occurrence. In
38 this way, candidate tags can be selected according to their similarity to the
39 introduced tags, and a sorting algorithm can rank them in terms of estimated
40 relevance [4, 8, 14, 22]. In previous work, we described and evaluated a gen-
41 eral scheme for folksonomy-based tag recommendation in collaborative tagging
42 systems [7]. Out of that scheme, eight particular methods were proposed which
43 form the basis of the method presented in this work.

44 Besides content-based and folksonomy-based tag recommendation systems,
45 other approaches have been described in the literature. Anderson et al. [1]
46 describe a tag recommendation system for Flickr¹, a well known photo shar-
47 ing site, which combines both content-based recommendations (by training a
48 predictive model that learns the mapping between tags and extracted content
49 image features) with folksonomy-based recommendations (following an strategy
50 very similar to [22]). Naaman and Nair[19] describe another tag recommen-

¹www.flickr.com

51 dation system for Flickr, which takes advantage of the geolocation metadata
52 attached to images and recommends tags that other users employed in close
53 areas. Chen et al. [3] describe a tag recommendation system for video resources
54 which crawls the web for information about these videos and identifies keywords
55 to recommend as tags.

56 Although it is quite common to personalise tag recommendation systems
57 to the tagging behaviour of particular users by promoting, for example, tags
58 that users introduced in past annotations [2, 8, 14, 16, 18, 20], most of the
59 current systems do not introduce direct user feedback in the evaluation loop.
60 Thus recommendations are generally evaluated using traditional information
61 retrieval approaches based on the comparison of tag rankings produced by dif-
62 ferent methods, or using precision and recall metrics computed after a tag pre-
63 diction task [2, 7, 8, 16, 18, 20]. To the best of our knowledge, only three stud-
64 ies perform some kind of user-based evaluation. Sigurbjörnsson and Zwol [22]
65 automatically generate tag recommendations for several images from a Flickr
66 dataset and then ask users to rate, in a four-point scale, whether the recommen-
67 dations are appropriate to a given image. Similarly, De Meo et al. [4] extend
68 the annotations of Delicious’ bookmarks² and then ask users to evaluate the
69 relevance of every tag/resource association. Jäschke et al. [13] perform a small
70 evaluation based on a real-world scenario where users have to tag bookmarks in
71 BibSonomy³. Specifically, precision and recall metrics are computed by compar-
72 ing tag recommendations performed to every bookmark and the final taglines
73 that users introduced. Due to its subjectiveness and many different ways to be
74 accomplished, tag recommendation is not an easy task to evaluate, and some
75 advantages and disadvantages can be found in both user-based and information
76 retrieval evaluation approaches [8]. However, there is a clear lack of user-based
77 evaluation in previous work, and we believe that every recommendation system
78 should be validated at some point using both evaluation strategies. Proper user
79 feedback should be helpful not only to compare tag recommendation methods
80 but also to better understand the nature of the task and learn how can systems
81 be improved.

82 The contribution of the present work is twofold. First, we propose an ex-
83 tended version of the best performing tag recommendation method found in
84 our previous work [7]. The main idea behind this extended method is to exploit
85 the automatic classification of the resources to be annotated into a number of
86 predefined classes to further adapt the tag suggestions to the context of these
87 classes. This classification is based on the tags that users start introducing
88 during the annotation process. In this way, instead of personalising recom-
89 mendations for particular users, we “personalise” them to particular classes of
90 resources. Next, as a second contribution, we perform a comprehensive user-
91 based evaluation through an online experiment where participants are presented
92 with some resources which have to be annotated with the help of a tag recom-

²www.delicious.com

³www.bibsonomy.org

93 mendment system. These kind of user-based evaluations are very costly and we
94 have seen that they are not very common in the tag recommendation literature.
95 For that reason, we believe our contribution is of great valuable to the commu-
96 nity. In our evaluation, we compare the recommendation method we proposed
97 in previous work and the extended version we describe here along with two ran-
98 dom baselines. Moreover, we perform a complementary evaluation based on a
99 tag-prediction task following common information retrieval methodologies. In
100 our previous work [7], the tag recommendation methods were evaluated using
101 a tag-prediction task and compared favourably against four baselines and two
102 state of the art methods [8, 22]. For this comparison, we used data from the
103 folksonomies of Freesound⁴, an online audio clip sharing site with more than
104 3,5 million registered users and 180,000 uploaded sounds [5], and Flickr. There-
105 fore, the recommendation methods were tested in the audio and image domains.
106 Similar results were obtained in both scenarios. In this work, evaluations are
107 carried out in the context of Freesound. Results show that the newly proposed
108 recommendation method brings a statistically significant improvement over the
109 previous method, according to both user-based and prediction-based evalua-
110 tions. Analysing user-based evaluation results we find that participants which
111 are experienced in working with sound libraries tend to better appreciate the
112 improvements of the new tag recommendation method we describe here. More-
113 over, we see that the more familiarised the users are with Freesound, the more
114 the number of tag suggestions they accept as valid annotations. User feed-
115 back reveals that tag recommendation methods tend to be more useful when
116 recommending broad tags (i.e., referring to generic concepts). Participants also
117 recognise tag annotation as a particularly difficult task, specially if the resources
118 being annotated are not authored by themselves.

119 The rest of the paper is organised as follows. First, we summarise the steps
120 of the tag recommendation method we proposed in previous work and describe
121 the new approach based on the classification of input tags (Sec. 2). Then, we
122 describe the online experiment we designed for user-based evaluation (Sec. 3).
123 Results of the online experiment are reported in Sec. 4, and the complementary
124 prediction-based evaluation is described and reported in Sec. 5. We conclude
125 the paper with a discussion about our findings and future work (Sec. 6).

126 2. Tag recommendation methods

127 The two tag recommendation methods we describe in this work are based
128 on tag-tag similarities derived from the folksonomy of Freesound. Given a set
129 of input tags Γ_I , the methods output a set of recommended tags Γ_R .

130 2.1. General tag recommendation

131 The general tag recommendation method presented in [7], which we denote
132 by GEN, consists of three steps (Fig. 1):

⁴www.freesound.org

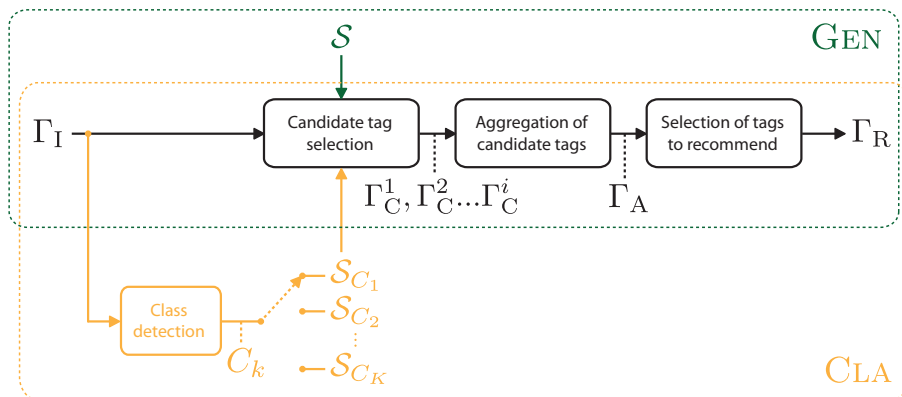


Figure 1: Schematic block diagram of the general (GEN) and class-based (CLA) tag recommendation methods.

- 133 1. Candidate tag selection: Given a set of input tags Γ_I , this step uses a
134 tag-tag similarity matrix \mathcal{S} derived from the Freesound folksonomy to
135 select a set of N candidate tags Γ_C^i for each input tag Γ_{I_i} . The tag-tag
136 similarity matrix \mathcal{S} is constructed by computing the association matrix
137 $\mathcal{D} = \{d_{i,j}\}$, which represents the associations between tags and audio
138 clips in the Freesound folksonomy ($d_{i,j} = 1$ if audio clip a_i is labeled with
139 tag t_j , and $d_{i,j} = 0$ otherwise). Hence, \mathcal{D} is a sparse matrix that has as
140 many columns as audio clips in Freesound and as many rows as the set of
141 distinct tags being used to label these audio clips⁵. Given \mathcal{D} , the tag-tag
142 similarity matrix is obtained as $\mathcal{S} = \mathcal{D}\mathcal{D}'$ ($'$ indicates transposition), and
143 we apply a simple normalisation to the elements $\{s_{t_i,t_j}\}$ of \mathcal{S} so that s_{t_i,t_j}
144 corresponds to the cosine similarity between tags t_i and t_j on the basis of
145 their co-occurrence in audio clips. Tags in Γ_C^i are selected as the N
146 most similar tags to a given input tag Γ_{I_i} .
- 147 2. Aggregation of candidate tags: Given the sets Γ_C^i from the first step,
148 candidates are assigned a score ϕ and aggregated into a single list of tags
149 with scores Γ_A . Such score is determined by the candidate similarity-
150 based ranking so that $\phi = 1$ for the most dissimilar candidate to a given
151 input tag and $\phi = N$ for the most similar one. The scores of tags that
152 are present in different sets of candidates Γ_C^i are added when aggregated
153 in the final set Γ_A .
- 154 3. Selection of tags to recommend: Considering the scores in Γ_A , this step

⁵In order to reduce the computational cost of the operations performed in this step and to get rid of potentially noisy tags, when building the association matrix we only consider tags whose frequency of occurrence is higher than a threshold $\omega = 10$ (i.e. we only consider tags that are used at least 10 times in the Freesound folksonomy). In this way the number of rows of the association matrix is reduced by $\approx 80\%$, with only around $\approx 10\%$ of the associations between tags and audio clips being actually ignored [7].

155 determines a threshold ϵ to select the tags that are finally recommended.
 156 Here we use the strategy of determining the threshold ϵ as a percentage
 157 of the maximum score in Γ_A [7]. Tags in Γ_A are sorted by their score
 158 and those that satisfy $\phi \geq \epsilon$ are outputted as Γ_R , the final set of recom-
 159 mended tags.

160 In this way, the method GEN can generate a sorted list of recommended
 161 tags Γ_R given a set of input tags Γ_I and a tag-tag similarity matrix \mathcal{S} which
 162 is derived from previous tag associations. Given that this method does not
 163 take into account any audio-specific information such as content features, it is
 164 general enough to be applied to other kinds of multimedia domains. Example
 165 applications for audio and images, as well as more detailed explanations, are
 166 provided in [7].

167 2.2. Class-based tag recommendation

168 The proposed class-based tag recommendation method, which we refer to as
 169 CLA, is a variation of GEN based on the classification of the input tags Γ_I into
 170 a set of K predefined audio classes. For every class C_k , $k \in [1, K]$, a tag-tag
 171 similarity matrix \mathcal{S}_{C_k} is built in the same way as in the GEN method, except
 172 that in this case only the tag assignment information corresponding to the audio
 173 clips of the current class is considered (see below). As a result, a different tag-
 174 tag similarity matrix can be computed for every audio class, and the matrix \mathcal{S}_{C_k}
 175 that is used in the candidate tag selection step of the recommendation process
 176 depends on the classification of the input tags Γ_I (Fig. 1). Once the candidates
 177 are selected, the other two steps (aggregation of candidate tags and selection of
 178 tags to recommend) are computed exactly in the same way as in GEN.

179 2.2.1. Classification of input tags

180 The classification of input tags is performed using a supervised learning
 181 model trained with the original tag annotations of audio clips in Freesound.
 182 We defined $K = 5$ audio classes (Table 1) and manually built a ground truth
 183 of 1,200 audio clip examples of each class. Then, we trained a multivariate
 184 Bernoulli Naive Bayes classifier feeding it with the taglines of the audio clips in
 185 the ground truth. Given a set of input tags Γ_I , the classifier can predict which
 186 category C_k better fits the input. Details on the class detection step and the
 187 process we followed for defining the audio classes, building the ground truth and
 188 evaluating the classifier can be found in [6]. The resulting classification system is
 189 able to classify a set of input tags Γ_I within the five defined classes with different
 190 accuracies depending on the length of Γ_I . The lowest accuracy, obtained when
 191 $|\Gamma_I| = 1$ (i.e., only one tag is given to the classifier), is approximately 75%. For
 192 $|\Gamma_I| \geq 4$ the classification accuracy reaches a plateau between 90 and 95%.

193 2.2.2. Computation of tag-tag similarity matrices

194 As mentioned, the process of building the tag-tag similarity matrices \mathcal{S}_{C_k} is
 195 the same as the one for building \mathcal{S} , except that for every matrix \mathcal{S}_{C_k} we only
 196 consider tag assignment information from audio clips belonging to C_k . For that

Class name	Description and examples
SOUNDFX	Sound effects (including <i>foley</i>), footsteps, opening and closing doors, alarm sounds, cars passing by, animals, and all kinds of noises or artificially created glitches.
SOUNDSCAPE	Environmental recordings, street ambiances or artificially constructed complex soundscapes.
SAMPLE	Instrument samples including single notes, chords and percussive hits (e.g. single notes of a piano recorded one by one and uploaded as different audio clips, or samples from a complete drum set).
MUSIC	Musical fragments such as melodies, chord progressions, and drum loops. This class is to SAMPLE what SOUNDSCAPE is to SOUNDFX.
SPEECH	All sorts of speech-related audio clips such as text reading, single words or recordings of text-to-speech processors.

Table 1: Audio classes.

197 we reused the classification system described in Sec. 2.2.1 to classify all audio
198 clips in Freesound in one of the five audio classes, with input tags corresponding
199 to the original taglines of audio clips in Freesound. Then, matrices \mathcal{S}_{C_k} can be
200 built by only considering the columns of \mathcal{D} corresponding to the audio clips of
201 C_k . Hence, $\mathcal{S}_{C_k} = \mathcal{D}_{C_k} \mathcal{D}'_{C_k}$, where \mathcal{D}_{C_k} is a subset of \mathcal{D} where the columns cor-
202 responding to audio clips not in C_k are removed. Each matrix \mathcal{S}_{C_k} is normalised
203 using the same process we use for \mathcal{S} (Sec. 2.1).

204 Notice that the similarity value between two tags t_i and t_j will be different
205 in every matrix \mathcal{S}_{C_k} and in \mathcal{S} , with \mathcal{S}_{C_k} being tailored to the particular context
206 of the k -th class. Also notice that the number of distinct tags resulting from
207 considering all audio clips belonging to C_k will be smaller than the total number
208 of distinct tags resulting from considering all audio clips from all classes (the size
209 of the *class vocabulary* will be smaller than the size of the *general vocabulary*).
210 Therefore, there will be some “all-zeros” rows in \mathcal{S}_{C_k} , corresponding to the tags
211 that are not used in the context of the particular class C_k . Hence, these tags
212 are never recommended when using \mathcal{S}_{C_k} .

213 3. User-based evaluation

214 We designed an online experiment where participants have to tag a set of
215 audio clips from Freesound with the help of the tag recommendation systems
216 of Sec. 2. The experiment was online for 15 days during June 2013, and was
217 publicised in the Freesound front page. The goal of this experiment is twofold.
218 First, we want to assess which of the recommendation methods is more useful
219 for users when tagging audio clips. Second, we want to get qualitative user
220 feedback to better understand the strengths and weaknesses of the considered
221 tag recommendation systems and, in a further stage, to understand the poten-
222 tial strengths and weaknesses of tag recommendation processes in general. As

FREESOUND DATASET

Number of audio clips	140,622
Number of unique tags†	43,696
Number of contributor users‡	6,948
Number of tag assignments	990,574
Average tags per audio clip (tagline length)	7.044

TAG-TAG SIMILARITY MATRICES

	Num. audio clips	Vocabulary size
General matrix (\mathcal{S})	140,622	7,710
Matrix for class SOUNDFX	29,725	4,584
Matrix for class SOUNDSCAPE	38,001	5,768
Matrix for class SAMPLE	26,452	3,280
Matrix for class MUSIC	34,139	4,303
Matrix for class SPEECH	15,305	3,557

Table 2: General statistics of the Freesound dataset and the resulting tag-tag similarity matrices. †Some of these tags are not semantically unique, and may include synonyms and typographic errors. ‡Users that have contributed by uploading at least one audio clip.

223 mentioned in Sec. 1, this is yet an under-explored area.

224 Along with GEN and CLA, we also evaluate two random variants of them,
 225 named RGEN and RCLA, respectively. These differ from the original variants in
 226 that, in the final step of the recommendation process, the set of recommended
 227 tags Γ_R is replaced with an alternative set of the same length containing ran-
 228 domly selected tags either from the general vocabulary (RGEN) or from the
 229 corresponding particular class vocabulary (RCLA). Notice that the general vo-
 230 cabulary is always bigger than any of the individual classes' vocabulary. Hence,
 231 the random selection in RGEN is performed over a bigger and more diverse pool
 232 of tags. Participants were not aware of the particular recommendation method
 233 underlying tag suggestions nor knew about the five audio classes in which we
 234 classify all annotated audio clips. The dataset we use for the evaluation com-
 235 prises Freesound data⁶ gathered between April 2005 and May 2012 (Table 2).
 236 It includes tag assignment information which relates tags, audio clips and users,
 237 and it is used to build the tag-tag similarity matrices \mathcal{S} and \mathcal{S}_{C_k} , as explained
 238 in Sec. 2.2.2.

239 The online-experiment proceeded as follows:

240 **Instructions page:** First, participants were presented with an introduction
 241 page displaying detailed instructions for the experiment (Fig. 2). Partici-
 242 pants were told they would have to annotate 20 audio clips from Freesound,
 243 using as many tags as they felt appropriate for every clip (we suggested
 244 participants to use five or more tags, but it was not mandatory). Partici-
 245 pants were also told that as soon as they started typing tags, a list of

⁶Freesound data, including audio clips and tag annotations, can be gathered using the public Freesound API (www.freesound.org/help/developers/).

Freesound tagging experiment

Welcome to the Freesound tagging experiment!

Instructions

- In this experiment you will be presented with some sounds from Freesound.org and **you will have to annotate them** with textual labels (tags!). Please use any expressions -even onomatopoeic- that come to your mind. Feel free!
- The number of tags you can use for labeling each sound is up to you, although **we suggest using 5 or more tags**.
- As soon as you start annotating a sound, a tag recommendation system will analyse your input and will **display a list of tags that might be meaningful** for the sound you are describing. You can add tags from this list (if you feel they are appropriate) by clicking on them. You do not necessarily have to add any of these tags if you do not find them relevant.
- Once you have finished annotating a sound, **click on the "Next Sound!"** button and you will be presented with another sound to annotate.
- You will have to annotate a total of **20 sounds**.
- To better appreciate the sounds you will be presented, we recommend **using headphones**.
- We will randomly select two participants in the experiment to receive a [Freesound t-shirt!](#)

Thank you very much for your participation!

Figure 2: Screenshot of the instructions page.

246 tag suggestions would appear and that they could choose tags from this
 247 list if they felt the suggestions were appropriate. We also recommended
 248 participants to use headphones for better listening conditions.

249 **Questionnaire:** After the introduction, a short questionnaire (Fig. 3) was
 250 presented to collect some basic user data and information about their ex-
 251 perience in working with sound libraries, their experience using Freesound
 252 (including the number of uploaded sounds) and their native language (in
 253 particular to be able to differentiate between native and non-native En-
 254 glish speakers).

255 **Audio clip annotation:** Once the questionnaire was completed, participants
 256 started annotating audio clips. From the ground truth we defined when
 257 designing the recommendation system (Sec. 2.2.1), we manually selected
 258 50 audio clips per class⁷. These clips were selected trying to cover a
 259 certain variety of sounds and avoiding those that would presumably be
 260 very hard to annotate. From this pool of 250 clips, every participant
 261 was assigned a random selection of four clips per class. Then, each of
 262 the four clips was assigned a different tag recommendation method that
 263 would be used when the participant annotated the clip. In this way, every
 264 participant was assigned a total 20 audio clips, equally distributed among
 265 audio classes and recommendation methods. Participants were presented
 266 with the first audio clip and had to annotate by typing tags in a text box.

⁷The clips we selected for the annotation phase of the online experiment (a total of 250, 50 per class) were removed from the ground truth and thus were not used to train the classifier described in section 2.2.1.

Before starting, some information about you...

Name: (optional)

Email: (optional | we will use the email to contact you in case you win a t-shirt)

Age: Gender: Male Female (optional)

Check this box if you're a native english speaker.

In case **you're not** a native english speaker, could you please indicate here which is your first language? (optional)

Are you a Freesound user? Yes No

If you're a Freesound user, could you please tell us:

a) How long have you been using Freesound?
"I have been using Freesound for years"

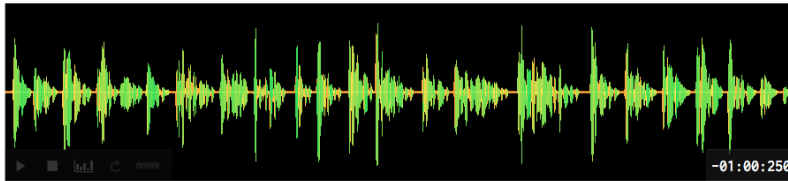
b) How many sounds have you uploaded?
 I have not uploaded any sound
 Between 1 and 10
 Between 10 and 50
 Between 50 and 500
 Between 500 and 1000
 More than 1000

Are you used to working with sound libraries? Yes No

How would you qualify your experience in fields such as sound libraries, sound recording and sound design?
 Accidental
 Amateur
 Advanced
 Professional

Figure 3: Screenshot of the questionnaire page.

Sound number 7



Introduce your tags:
(separate tags with spaces, join multi-word tags with dashes, e.g. first-tag)

Suggestions of other possibly relevant tags given your input: (click on the tags to add them, click here to clear the recommendation)

Figure 4: Screenshot of the sound annotation page.

267 The audio clip could be reproduced using a web player that also showed
 268 a visualisation of the waveform and the spectrogram of the audio clip. As
 269 soon as the participant started typing, a list of suggested tags appeared
 270 below the text box. This list was computed using the tag recommendation
 271 method assigned to the currently annotated audio clip, and was being
 272 updated every time a new tag was written in the text box⁸. Users could
 273 click over the tags shown in the list to automatically append them in the
 274 text box (Fig. 4). Once a participant considered an audio clip was fully
 275 annotated, she could click on the “Next sound” button and be presented
 276 with the following clip. Participants were also provided an URL that they
 277 could save for later resuming the experiment in case they did not want
 278 to annotate all clips in one go. Noticeably, we logged information about
 279 all the keystrokes and mouse clicks that participants performed with the
 280 corresponding timestamps.

281 **Feedback page:** After annotating the 20 audio clips, participants were pre-
 282 sented with a page thanking their participation and offering some space
 283 in a text box to give some feedback about the experiment. Alternatively,
 284 they were also offered to write the feedback in a particular section of the
 285 Freesound forums.

286 Considering the logs resulting of the user experiments we define a simple
 287 measure for evaluating the “usefulness” of every tag recommendation method
 288 in the tagging process. The measure consists in counting, for every set of tags
 289 assigned to an audio clip by a particular participant, the number of these tags
 290 that were recommended by the system during the annotation process (i.e., the
 291 number of recommended tags that were *accepted* by the participant). Let Γ_P be
 292 the set of tags that a participant used to annotate a particular audio clip, and
 293 let Γ_{R_m} be one of the sets of recommended tags that were presented to the user
 294 in the successive M tag recommendations during the tagging process of that
 295 particular audio clip. Then, we can define Λ , the number of accepted tags, as

$$\Lambda = \left| \Gamma_P \cap \left(\bigcup_{m=0}^M \Gamma_{R_m} \right) \right|, \quad (1)$$

296 where $|\cdot|$ measures set cardinality. Notice that Λ is roughly equivalent to
 297 a standard recall measure (without the normalization by $|\Gamma_P|$). We employ
 298 this measure instead of standard precision and recall (e.g., as done in [13])
 299 because the nature of our evaluation has some particularities which make such
 300 metrics less useful. As described above, in our evaluation system several tag
 301 recommendations are performed during the annotation of a single clip (i.e., every
 302 time that a new tag is introduced the recommendation is recomputed). As a

⁸Similarly to the Freesound upload system, tags had to be separated by spaces and multi-words joined with dashes. Hence, the recommendation was updated every time a blank space was introduced.

303 result, the total number of recommended tags for every audio clip is much larger
304 than the final number of assigned tags. If we computed precision and recall by
305 comparing the whole set of recommended tags for every audio clip with the final
306 taglines assigned by users, we would obtain very low precision values which, in
307 our opinion, are not as representative as Λ . In our evaluation (and in a real-
308 world tag recommendation scenario), users are the ones who finally decide which
309 of the recommended tags are relevant for a particular resource. Therefore, the
310 length of the recommendation is not as important as the fact that it contains
311 meaningful suggestions (i.e., recall is much more important than precision).

312 4. Results

313 During the two weeks the experiment was online we gathered a total 201
314 experiment logs from 190 unique participants (some participants decided to re-
315 peat the experiment more than once). Among all these experiment logs, 80
316 correspond to unfinished experiments (i.e., with less than 20 audio clips anno-
317 tated) which we do not consider in the analysis. In addition, we apply a filter
318 to discard logs from experiments that were finished very quickly and with very
319 few calls to the recommendation methods. More specifically, we discard logs
320 from experiments completed in less than 10 minutes (average of 30 seconds per
321 audio clip) and from experiments not reporting a minimum of three calls to the
322 recommendation system for every annotated audio clip. We discard these logs
323 as we consider that participants did not pay enough attention when annotat-
324 ing audio clips and thus contain potentially noisy data. After filtering, we are
325 left with 70 logs that we consider as sufficiently reliable data for analysis. In
326 the following subsections we show the results of different aspects of the online
327 experiment analysis.

328 4.1. Accepted tags per recommendation method

329 First, we report on the basic accuracy of the considered tag recommendation
330 methods (Table 3, leftmost column). We observe that random methods RCLA
331 and RGEN report way lower average Λ than CLA and GEN. Thus, our methods
332 do perform much more meaningful recommendations than the random baselines.
333 Interestingly, we also observe that both class-based methods CLA and RCLA
334 report higher averages than their general counterparts GEN and RGEN. This
335 suggests that tag recommendations improve when using class-based methods.
336 However, the differences are not statistically significant⁹.

337 Next, we repeat the same analysis but considering different groups of ex-
338 periment logs according to the questionnaire that participants had to fill at the

⁹If not stated otherwise, statistical significance is assessed by performing pairwise compar-
isons using the Mann-Whitney U test with $\alpha=0.05$ [17]. When performing multiple compar-
isons, we apply a correction to the rejection criteria in order to reduce the familywise error
rate. In particular, we use the Holm-Bonferroni correction [11]. Notice that these are robust
and stringent criteria for measuring statistical significance (cf. [21]).

	All	Expert	Non-expert	Native	Non-native
CLA	2.414 (2.775)	2.547 (2.988)	2.179 (2.224)	2.950 (3.382)	1.963 (2.027)
GEN	2.154 (2.526)	2.163 (2.663)	2.147 (2.229)	2.656 (3.006)	1.732 (1.938)
RCLA	0.260 (0.671)	0.278 (0.680)	0.211 (0.663)	0.300 (0.705)	0.226 (0.638)
RGEN	0.166 (0.455)	0.139 (0.458)	0.253 (0.458)	0.194 (0.518)	0.142 (0.392)

Table 3: Average number of accepted tags Λ (standard deviation into parenthesis) of the user-based evaluation approach for the following groups of participants. From left to right these correspond to all participants, expert participants, non expert participants, native English speakers and non-native English speakers.

339 beginning of the experiment (Table 3). In particular, we compute Λ for each rec-
340 ommendation method considering groups of logs corresponding to experienced
341 participants (i.e., participants that checked the box marked with the question
342 “Are you used to working with sound libraries?” in the questionnaire; second
343 column in Table 3), non-experienced participants (third column), native English
344 speakers (fourth column), and non-native speakers (fifth column). We again ob-
345 serve that, except for RCLA and RGEN in the non-expert group, all class-based
346 methods report higher averages than the general methods. This further sup-
347 ports the idea that class-based recommendations bring some improvements over
348 the general method. Interestingly, in the case of experienced participants, the
349 difference between CLA and GEN increases with respect to the same comparison
350 when considering all participants. In this case we get a statistically significant
351 increase of 0.38 ($p < 2.91 \cdot 10^{-2}$). Furthermore, the difference between RCLA
352 and RGEN also increases for the experts (with respect to all participants) and
353 becomes statistically significant ($p < 2.47 \cdot 10^{-3}$). This suggests that expert par-
354 ticipants clearly appreciate a difference between CLA and Gen methods (even
355 for the random versions) and find class-based recommenders to be more useful.
356 On the other side, we observe that when analysing the non-experienced par-
357 ticipants group, the differences between class-based and general methods gets
358 blurred, with almost no difference between the two types of recommendation
359 methods. Thus, non-experienced participants are not able to tell the differ-
360 ence between class-based and general recommendations. Overall, these results
361 indicate that the usefulness of class-based tag recommendations compared to
362 general recommendations is slightly higher, and specially in the case of experi-
363 enced participants.

364 Considering the last two groups of participants (native and non-native En-
365 glish speakers), we observe that the differences between class-based and general
366 recommendation systems are quite similar to those obtained when considering
367 all participants. Class-based systems report a higher Λ but the increments are
368 practically the same for both native and non-native groups (there is no statisti-
369 cally significant difference between the increments). Thus, we do not see a direct
370 general implication of language in method preference. Nevertheless, there is a
371 significant difference in the absolute number of accepted tags among the native
372 and non-native participant groups (Table 3). Native English speakers tend to
373 accept an average of 0.96 tags more than non-native ones ($p = 4.61 \cdot 10^{-3}$). Fur-

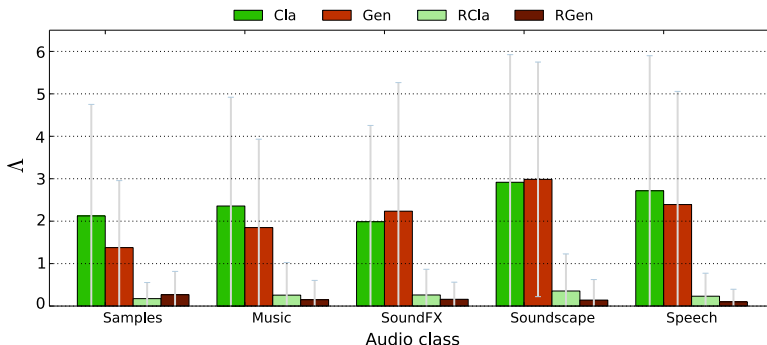


Figure 5: Average accepted tags (Λ) per audio class and recommendation method.

374 furthermore, we observe that native English participants tend to annotate audio
375 clips with an average of 0.32 tags more than non-native ones ($p = 3.24 \cdot 10^{-6}$).
376 Thus, in our experiments, native speakers consistently use more tags for describ-
377 ing audio clips than non-native speakers and tend to accept more recommenda-
378 tions. To the best of our knowledge, this is the first time evidence is reported
379 with regard to the comparison of native’s and non-native’s tagging behaviour.
380 Our results suggest that native speakers use more tags when describing online
381 resources than non-native participants and that, therefore, this aspect should
382 not be overlooked in future studies. Overall, we see that both native and non-
383 native speakers prefer CLA over GEN (and RCLA over RGEN), but that this
384 preference is not stronger than in any of the other user groups.

385 4.2. Accepted tags per audio class

386 To gain insight about how do recommendation methods work for the different
387 audio classes defined above (Table 1), we grouped annotated sounds by class and
388 recommendation method and computed the average number of accepted tags Λ
389 for each group (Fig. 5). In general, clips under SOUNDSCAPE and SPEECH classes
390 reported higher Λ than clips under the other classes. This is probably because
391 there are some tags such as **field-recording**, **nature** or **voice** which are very
392 common in these classes and are very generic (i.e., could be used to annotate
393 almost any clip in SOUNDSCAPE or SPEECH classes).

394 It can be also observed that not all audio classes feature a higher Λ for the
395 CLA method than for the GEN method. SOUNDSCAPE clips report higher Λ for
396 GEN than for CLA, although the difference of 0.07 is not statistically significant
397 ($p = 4.56 \cdot 10^{-1}$). SOUNDFX clips also report higher Λ for the GEN method and,
398 although the difference is still not statistically significant ($p = 3.80 \cdot 10^{-1}$), the
399 increase of 0.25 is this time bigger. SAMPLE, MUSIC and SPEECH classes report
400 higher Λ for CLA recommendations, with larger Λ increases and with improved
401 statistical significance. This suggests that the knowledge-based adaptation that
402 the CLA performs is better exploited in SPEECH, MUSIC and SAMPLE classes
403 than in SOUNDSCAPE or SOUNDFX. We hypothesise that the vocabulary needed

404 to accurately describe clips from the former classes is more reduced than the
 405 vocabulary needed for other audio clips. Therefore, the class-based method can
 406 easily adapt to the class context and produce better recommendations, probably
 407 including less generic tags than the ones that would be recommended using the
 408 general method. On the other side, clips under SOUNDSCAPE and SOUNDFX
 409 classes cover a wider range of sounds and need a larger vocabulary to be well-
 410 described. In this situation, the CLA method does not adapt well and does not
 411 improve the GEN results. Our hypothesis is partially supported by looking at
 412 the actual size of the resulting class vocabularies after computing the tag-tag
 413 similarity matrix per class (\mathcal{S}_{C_k} , Table 2). SPEECH, MUSIC and SAMPLE produce
 414 smaller similarity matrices, with less tags in the vocabulary, than SOUNDSCAPE
 415 and SOUNDFX.

416 4.3. Correlation between number of uploaded sounds and accepted tags

417 All participants in our experiment were Freesound users. However, not all
 418 of them had experience in uploading and tagging audio clips in Freesound. In
 419 order to get some insight as how being used to tagging audio clips affects Λ ,
 420 we computed the correlation¹⁰ between the number of uploaded sounds and the
 421 number of accepted tags, grouping audio clips into the four evaluated recom-
 422 mendation methods (Table 4). We find the strongest correlation for the CLA
 423 method ($\varrho = 0.276$, $p < 3.76 \cdot 10^{-7}$). Thus, in this case, Λ tends to grow along
 424 with the number of uploaded sounds. A less significant correlation is reported
 425 for the GEN method ($\varrho = 0.105$, $p < 5.61 \cdot 10^{-3}$). RCLA and RGEN present no
 426 significant correlations ($\varrho = 0.087$, $p < 1.13 \cdot 10^{-1}$ and $\varrho = 0.063$, $p < 2.55 \cdot 10^{-1}$,
 427 respectively). This finding suggests that the more familiar the participants are
 428 with the Freesound uploading and tagging process, the more recommended tags
 429 they tend to accept, specially when recommendations are generated with the
 430 CLA method. This result is consistent with the previous observation that experi-
 431 enced participants tend to accept more tags than non-experienced ones when
 432 recommendations are generated by CLA (Sec. 4.1). Again, we are not aware of
 433 any study considering user familiarity in the context of resource tagging. There-
 434 fore, our results represent a novel and original contribution with regard to this
 435 aspect.

436 4.4. Timing aspects

437 Timing is also an often unconsidered aspect when evaluating tag recommen-
 438 dation systems. However, it is interesting because it can reveal some insights
 439 about the annotation process. We measured the average time invested for an-
 440 notating an audio clip in our experiments and observed that there exists a
 441 significant correlation between the length of the audio clips and the time in-
 442 vested to annotate them, being shorter clips the fastest to describe ($\varrho = 0.24$,

¹⁰We employ the Spearman’s rank correlation coefficient [10], with ϱ denoting the correlation coefficient and p the p-value associated with it.

Number of uploaded sounds†	COM	GEN	RCOM	RGEM
0	2.105	2.036	0.221	0.126
1 to 10	1.823	2.027	0.293	0.133
11 to 50	2.580	1.820	0.220	0.240
51 to 500	2.289	2.222	0.311	0.133
501 to 1000	4.160	2.035	0.380	0.300

Table 4: Average number of accepted tags Λ per number of uploaded sounds and recommendation method. †The ranges in the number of uploaded sounds are determined in the questionnaire that participants had to fill at the beginning of the experiment (Fig. 3).

443 $p < 5.68 \cdot 10^{-19}$). That could be expected, as shorter clips tend to be less com-
444 plex and need less time for listening to them. Consistently, audio clips belonging
445 to the SOUNDSCAPE class need an average of 15 extra seconds to be described
446 when compared to clips belonging to other classes ($p < 8.12 \cdot 10^{-3}$). On the
447 other side, SAMPLE clips need less time than the rest ($p < 3.15 \cdot 10^{-2}$). This can
448 be explained because SOUNDSCAPE clips are generally longer than clips from
449 other classes, while SAMPLE clips tend to be shorter. We have not observed
450 any statistically significant differences in the average time invested in annotat-
451 ing audio clips when comparing the four different recommendation methods.
452 Therefore, the choice of a recommendation method does not seem to affect the
453 time needed to annotate audio clips.

454 4.5. User feedback

455 In the last phase of the online experiment, participants were provided the
456 opportunity to give some feedback in the form of comments (Sec. 3). We observe
457 some recurring opinions that, if extrapolated, bring also valuable insights into
458 recommendation processes in general. First of all, participants agree in that the
459 process of annotating audio clips (and by extension the process of recommend-
460 ing tags) is a very hard task, and that recommendations are a generally useful
461 tool but not always needed or used. In our case, the 30% of all tag annota-
462 tions performed during the experiment were suggested by the recommendation
463 systems¹¹.

464 A lot of participants point out that annotation is especially hard when the
465 audio clip being described is not recorded/created by the person annotating it
466 (which was always the case in our experiment). In those cases, there is a lot
467 of meaningful information about the sound which most of the times can not be
468 determined without the knowledge of how the clip was created (e.g., software
469 used, recording device, location of a recording, etc.). Some participants also
470 point out that in order to perfectly annotate musical audio clips such as drum
471 loops or instrument notes, a lot of time needs to be invested in determining
472 properties such as beats per minute or the pitch of a note. Those issues are par-
473 ticularly relevant in our context, where participants had to annotate audio clips

¹¹This percentage is computed without taking into account tag recommendations performed with random methods which obviously did not provide meaningful recommendations.

474 not created by themselves. Finally, another repeated comment is that tag sug-
475 gestions are more useful for “nature” and “human-related” audio clips, whereas
476 “abstract” and “synthetic” clips require more tags to be manually introduced
477 before some meaningful suggestions are made. These comments are somehow
478 aligned with the results reported in Fig. 5, where we see that SOUNDSCAPE and
479 SPEECH classes are the ones that report higher Λ .

480 4.6. Tag analysis

481 We perform here a close-look analysis to the experiment logs in order to get
482 some insight in the type of tags that are recommended and in which cases those
483 are accepted by participants. We detect several interesting patterns that we
484 believe also help comprehending in more detail tag recommendation processes
485 in general. First of all, there are some tags which are recommended and ac-
486 cepted a lot of times in the online experiment. These tags correspond to very
487 generic concepts such as **field-recording**, **voice**, **electronic**, **loop**, **nature**
488 or **percussion**. These recommendations are useful to provide some kind of gen-
489 eral categorization to annotated audio clips, but clips only tagged with these
490 kind of tags do clearly lack specificity in the annotations. We observe that an-
491 other recommendation pattern consists in tags that are suggested many times
492 but are rarely accepted. This is the case of tags such as **sound** or **recording**, for
493 which we hypothesise that the meaning is too obvious to be considered as rel-
494 evant information for participants. It is also the case of tags like **soundscape**,
495 **percussion-loop**, **drum-loop** or **natural-reverb** which are normally repre-
496 sented by alternative tags (or combinations of tags) such as **field-recording**
497 (instead of **soundscape**), **loop**, **percussion**, **drum**, **natural** or **reverb**.

498 We also observe that there are some tags with low acceptance because of its
499 subjective meaning (e.g. **groovy**, **threatening**) or because participants can not
500 assess its correctness because they are not the authors of the annotated clips
501 (e.g. **multi-sample**, **improvised**). Obviously there are also some suggested
502 tags which are not accepted because they are simply not appropriate for the
503 clips being described. That could be the case of tags like **piano**, **guitar** or **pad**
504 which are sometimes recommended to audio clips which clearly do not contain
505 piano, guitar or pad-like sounds. Finally, we observe a last group of suggestions
506 which correspond to tags not usually suggested but normally accepted such as
507 **announcement**, **synthesizer**, **footsteps** or **airplane**. We consider these as be-
508 ing very good recommendations as they correspond to not-so-general concepts
509 and are apparently recommended only when they are needed. Overall, recom-
510 mendations provided by our methods tend to be useful when recommending
511 general tags, referring to concepts than can be used as a broad categorisations
512 of the audio clips. However, recommendations are not as useful when they refer
513 to more detailed aspects of the sounds being annotated.

514 5. Complementary evaluation

515 In order to complement the performed user-based evaluation, we also con-
516 sider a more systematic and empirical assessment of the different tag recommen-

517 dation methods (CLA, GEN, RCLA and RGEN) following the methodology we
 518 described in [7]. This complementary assessment follows a typical information
 519 retrieval evaluation setup based on a tag prediction task which we now describe.

520 *5.1. Prediction-based evaluation methodology*

521 For this evaluation we consider audio clips and annotations of the same
 522 Freesound dataset described in Sec. 3. We perform a 10-fold cross-validation
 523 following the methodology described by Salzberg [21] and others. For each fold,
 524 there is a training phase consisting of two steps which preprocesses all the nec-
 525 essary data for performing recommendations with the four evaluated methods.
 526 The first step consists in training a classifier that allows the classification of the
 527 input tags in one of the five defined audio classes, as described in Sec. 2.2.1. To
 528 do that, we feed the classifier only with these audio clips that are present both
 529 in the training set and in the ground truth we built when designing the system
 530 (i.e., we only use audio clips from the training set that we know to which audio
 531 category they belong to).

532 The second step of the training phase consists in building the general tag-tag
 533 similarity matrix \mathcal{S} and the matrices \mathcal{S}_{C_k} for every class C_k . For that we use
 534 information from all the audio clips in the training set. Notice that building \mathcal{S}_{C_k}
 535 requires the classification of all audio clips of the training set in one of the five
 536 defined categories (Sec. 2.2.2). We perform that classification using the same
 537 classifier trained in the first step of the training phase. Hence, this classifier
 538 is not only used in the recommendation process to classify the input tags and
 539 select a similarity matrix \mathcal{S}_{C_k} , but it is also used to build the similarity-matrices
 540 \mathcal{S}_{C_k} by classifying the audio clips of the training set.

541 After the training phase, we pick every audio clip in the evaluation set and
 542 randomly delete a set of tags Γ_D from its originally assigned tags, yielding Γ_I , the
 543 input to our recommendation system. The number of tags we delete is chosen
 544 uniformly at random, with the only constraint of leaving a minimum number of
 545 input tags of $|\Gamma_I| \geq 3$ so that there is presumably enough information for the
 546 recommender systems to provide good recommendations [7]. This constraint
 547 also implies that in order to be able to remove at least one tag for each audio clip
 548 ($|\Gamma_D| \geq 1$), we can only consider for evaluation the audio clips that have at least
 549 four tags¹². After we remove some tags, we run the four tag recommendation
 550 methods using Γ_I as input and the similarity matrices we computed in the
 551 training phase.

As evaluation measures we compute standard precision (P_n), recall (R_n),
 and f-measure (F_n) for each individual audio clip n according to

$$P_n = \frac{|\Gamma_R \cap \Gamma_D|}{|\Gamma_R|}, R_n = \frac{|\Gamma_R \cap \Gamma_D|}{|\Gamma_D|}, \text{ and } F_n = \frac{2P_n R_n}{P_n + R_n},$$

¹²This filtering is done before the whole evaluation process starts, therefore we evaluate the same number of clips in each fold.

552 where Γ_R is the set of recommended tags and Γ_D is the set of deleted tags.
553 Then, global P , R and F measures for each tag recommendation method are
554 calculated by averaging P_n , R_n and F_n across all resources $n \in [1, N]$ evaluated
555 with the chosen recommendation method.

556 The prediction-based evaluation approach is interesting as it allows us to
557 evaluate the different recommendation methods in a systematic way and using
558 a lot of audio clips. In previous work [7], we used this evaluation methodology
559 to exhaustively compare eight variations of the GEN recommendation method
560 (using different sets of parameters for each one of the recommendation steps)
561 against four baselines and two state of the art folksonomy-based tag recom-
562 mendation methods, and using data from the folksonomies of Freesound and
563 Flickr. That number of methods could have not been comprehensively com-
564 pared through a user-based evaluation approach such as the one presented in
565 the above sections. However, prediction-based evaluation has an important lim-
566 itation which is that we need an extensive ground truth to evaluate whether our
567 predictions are correct or not. In our case, this ground truth is composed by
568 the original taglines of sounds in Freesound. This means that the recommenda-
569 tions we evaluate will only be considered as “correct” recommendations if they
570 contain tags that the original author of the sound used to annotate it. As a
571 result, tags that could be subjectively considered as good recommendations for
572 a particular audio clip but are not present in the original annotations do not
573 count as correct predictions. Moreover, prediction-based evaluation does not al-
574 low the collection of qualitative user feedback that, as we have seen before, can
575 shed some light on relevant aspects of the recommendation process. For that
576 reason, we state that the prediction-based evaluation approach may be taken as
577 a complement to the results already described in the previous sections, allowing
578 us to further test our previous findings.

579 5.2. Prediction-based evaluation results

580 Results for the four evaluated tag recommendation methods appear to be
581 very similar to what we observe in the user study (Table 5). We can see that
582 CLA outperforms GEN by a small but statistically significant difference of 0.011
583 ($p < 6.51 \cdot 10^{-8}$). This difference suggests that CLA can successfully take advan-
584 tage of the classification step and the knowledge derived from the ground truth
585 to slightly improve the recommendations of the system. As expected, random
586 methods RCLA and RGEN score much lower F than CLA and GEN. Neverthe-
587 less, it is interesting to note that RCLA also features a statistically significant
588 increase in F with respect to RGEN ($p < 1.57 \cdot 10^{-24}$). This increase can be
589 explained by recalling that the pool of tags from which the random selection is
590 performed in RCLA is different in every audio class and it always contains less
591 tags than the pool in RGEN (Sec. 2.2.2). Hence, these results suggest that at
592 least some tags which are not relevant for a particular audio class are effectively
593 removed when building the similarity matrices \mathcal{S}_{C_k} . We also observe that CLA
594 and GEN feature a very similar number of recommended tags $|\Gamma_R|$, with an
595 average of 3.99 and 3.88 tags, respectively.

	P	R	F
CLA	0.476 (0.428)	0.488 (0.424)	0.440 (0.389)
GEN	0.486 (0.429)	0.467 (0.408)	0.429 (0.372)
RCLA	0.003 (0.031)	0.003 (0.038)	0.002 (0.025)
RGEN	0.002 (0.024)	0.002 (0.031)	0.001 (0.019)

Table 5: Average precision, recall and f-measure (standard deviation in parenthesis) for the prediction-based evaluation approach. Results are sorted by f-measure.

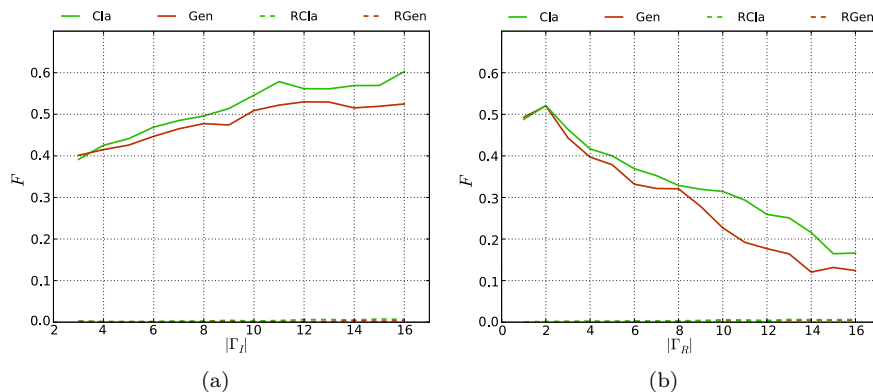


Figure 6: Average f-measure F as a function of the number of input tags $|\Gamma_I|$ (a) and the number of recommended tags $|\Gamma_R|$ (b).

596 If we analyse F as a function of the number of input tags $|\Gamma_I|$ and the number
597 of recommended tags $|\Gamma_R|$ we can get some more insight on the behaviour of
598 the considered recommendation methods (Fig. 6). For instance, we see that
599 both CLA and GEN have a tendency of increasing F as the number of input
600 tags also increases (Fig. 6(a)). This suggests that the recommendation system
601 is able to provide better recommendations when it is feed with more input tags.
602 The opposite happens with the number of recommended tags (Fig. 6(b)). This
603 can be explained as bigger numbers of recommended tags imply lower precision
604 values because more non-relevant tags are recommended. Nevertheless, it is
605 interesting to observe that the increase in F of CLA over GEN is specially
606 notorious for large numbers of recommended tags ($|\Gamma_R| > 8$, Fig. 6(b)). This
607 highlights the superiority of CLA over GEN when larger number of tags are
608 recommended, and suggests that CLA is able to provide more comprehensive
609 and relevant recommendations.

610 6. Conclusion and discussion

611 In this work we describe and evaluate two tag recommendation methods in
612 the audio clip sharing context of Freesound. One general tag recommendation
613 method (GEN) was introduced in previous work by the authors. The other

614 method, which is class-based (CLA), is an original contribution of this article.
615 It extends the former in two main aspects: it automatically determines to which
616 class an audio clip belongs and it produces specific recommendations for differ-
617 ent audio classes. As both tag recommendation methods (GEN and CLA) are
618 folksonomy-based, they are easily generalisable to other multimedia domains.
619 However, the CLA method requires the definition of K classes of resources in
620 the particular domain, and the building of a ground truth to train the classifier
621 needed to perform recommendations. The main bottleneck in terms of scala-
622 bility lies in the computation of the tag-tag similarity matrices that inform the
623 candidate selection step. However, these matrices can be computed offline, and
624 their size can be easily reduced by raising the threshold ω during the construc-
625 tion of the association matrix. This will discard those tags whose frequency of
626 occurrence is below that threshold (Sec. 2.1). That means that our recommen-
627 dation methods can scale well to even bigger amounts of data, as the number
628 of tags above the threshold ω will grow much more slowly than the number of
629 resources.

630 A limitation of the proposed recommendation methods is that they can suffer
631 the *cold-start* problem if deployed to collaborative tagging systems which have
632 not enough data to derive reliable tag-tag similarity matrices. Although our
633 recommendation methods have not been designed for collaborative tagging sys-
634 tems with scarce data, it would be interesting to evaluate how fast the methods
635 could acquire enough data from user annotations in order to provide meaning-
636 ful recommendations. In other words, it would be interesting to investigate
637 how big the folksonomy of a collaborative tagging system should be to enable
638 our tag recommendation methods to provide meaningful recommendations. We
639 hypothesise that, on a first step of the implementation of the system, tag-tag
640 similarity matrices would need to be recomputed very often as relatively small
641 changes in the folksonomy could have a big impact on the resulting similar-
642 ity matrix. In that case, the system would quickly learn from user tagging
643 behaviour and recommendations would quickly start to become more diverse.
644 Besides the similarity matrices, the CLA method also needs annotation data to
645 train the classifier. However, a collaborative tagging system could start using
646 the GEN method until enough data would be collected to build the ground truth
647 and train the classifier.

648 As a second contribution, we perform a user-based evaluation through an
649 online experiment. In it, participants had to annotate several audio clips with
650 the help of the different tag recommendation strategies. We logged the activity
651 of the participants and analysed these logs with the goal of comparing the
652 considered methods and, in addition, getting more insight into the positive and
653 negative aspects of tag recommendation systems in general. To the best of our
654 knowledge, this is one of the very few user-based evaluations carried out for a
655 tag recommendation task. Finally, as a further contribution, we complement
656 the user-based evaluation with a prediction-based evaluation, following a well-
657 established methodology and not considering any user input.

658 As a main result, we have seen that class-based recommendation reports
659 statistically significantly better scores than general recommendation, both in

660 the user-based and prediction-based evaluations. The difference in scoring is,
661 in absolute terms, more prominent for the user-based evaluation. Moreover, it
662 further improves when considering only expert users. This suggests that the
663 class-based method does indeed bring some improvements in the recommenda-
664 tions compared to general recommendation, and that these improvements are
665 more noticeable to expert users.

666 Among all annotations that participants performed during the online exper-
667 iment, approximately one third of them correspond to tags recommended by
668 the system (for both GEN and CLA methods). That by itself brings evidence
669 with regard to the general utility of tag recommendation systems. However,
670 the found results also indicate that tag suggestions referring to generic con-
671 cepts or sound classes tend to be more useful than recommendations of very
672 concrete tags describing specific sound characteristics. Participants found tag
673 suggestions more useful for sounds under SOUNDSCAPE and SPEECH categories.
674 We hypothesise that this happens because these categories are more suited to
675 the use of generic tags. MUSIC and SAMPLE audio classes require of annota-
676 tions describing very specific musical concepts such as pitch, tonality or beats
677 per minute. Participants had difficulties in annotating such concepts, as they
678 are problematic to annotate without having a certain knowledge of the record-
679 ing context (i.e., without being the author of the audio clip) and because tag
680 recommenders tend to produce less meaningful suggestions in these cases. All
681 these often overlooked qualitative evaluation aspects also represent a valuable
682 contribution of the present article.

683 We believe that, in order to build better tag recommendation systems, those
684 should be more aware of the particular contexts of the resources being described
685 and should extensively exploit all available knowledge. To generate tag sugges-
686 tions describing more concrete aspects of sound characteristics we need systems
687 that know about the specifics of the audio domain, such as which are the most
688 relevant properties of audio clips for different audio categories, and how to au-
689 tomatically estimate some of these properties. For that reason, we believe that
690 future tag recommenders should take advantage of knowledge representation
691 mechanisms such as ontologies to be able to include tags describing the audio
692 domain in some structured representation, and to be able to produce informed
693 recommendations based on reasoning and users' input. Such a system should
694 contribute in greatly improving online resource descriptions and thus facilitating
695 and providing new opportunities for content reuse.

696 **Acknowledgements**

697 We would like to thank Perfecto Herrera for his help in designing the online
698 experiment and also all Freesound users that participated. This work has been
699 supported by BES-2010-037309 FPI from the Spanish Ministry of Science and
700 Innovation (TIN2009-14247-C02-01; F.F.), 2009-SGR-1434 from Generalitat de
701 Catalunya (J.S.), JAEDOC069/2010 from CSIC (J.S.), ICT-2011-8-318770 from
702 the European Commission (J.S.), and FP7-2007-2013 / ERC grant agreement
703 267583 (CompMusic; F.F., X.S.).

704 **References**

- 705 [1] Anderson, A., Ranghunathan, K., Vogel, A., 2008. Tagez: Flickr tag recom-
706 mendation, in: Proceedings of the 23rd Conference on Artificial Intelligence
707 (AAAI 2008).
- 708 [2] Cao, H., Xie, M., Xue, L., Liu, C., 2009. Social tag prediction based on
709 supervised ranking model, in: Proceedings of the Conference on Machine
710 Learning and Principles and Practice of Knowledge Discovery in Databases
711 (ECML/PKDD), Discovery Challenge Workshop, pp. 35–48.
- 712 [3] Chen, Z., Cao, J., Song, Y., Guo, J., Zhang, Y., Li, J., 2010. Context-
713 oriented web video tag recommendation, in: Proceedings of the 19th Inter-
714 national Conference on World Wide Web (WWW 2010), p. 1079.
- 715 [4] De Meo, P., Quattrone, G., Ursino, D., 2009. Exploitation of semantic
716 relationships and hierarchical data structures to support a user in his an-
717 notation and browsing activities in folksonomies. *Journal of Information*
718 *Systems* 34, 511–535.
- 719 [5] Font, F., Roma, G., Serra, X., 2013a. Freesound Technical Demo, in:
720 Proceedings of the 21st ACM Conference on Multimedia (ACM MM 13),
721 pp. 411–412.
- 722 [6] Font, F., Serrà, J., Serra, X., 2013b. Audio clip classification using social
723 tags and the effect of tag expansion, in: Proceedings of the 53rd AES
724 Conference on Semantic Audio.
- 725 [7] Font, F., Serrà, J., Serra, X., 2013c. Folksonomy-based tag recommendation
726 for collaborative tagging systems. *International Journal on Semantic Web*
727 *and Information Systems* 9, 1–30.
- 728 [8] Garg, N., Weber, I., 2008. Personalized, interactive tag recommendation for
729 flickr, in: Proceedings of the 2nd ACM Conference Recommender systems
730 (RecSys 08), pp. 67–74.
- 731 [9] Halpin, H., Robu, V., Shepard, H., 2006. The dynamics and semantics of
732 collaborative tagging, in: Proceedings of the 1st Semantic Authoring and
733 Annotation Workshop, pp. 1–21.
- 734 [10] Hogg, R.V., Craig, A.T., 1995. *Introduction to Mathematical Statistics*.
735 5th ed., Prentice Hall.
- 736 [11] Holm, S., 1979. A simple sequentially rejective multiple test procedure.
737 *Scandinavian Journal of Statistics* 6, 65–70.
- 738 [12] Ivanov, I., Vajda, P., Goldmann, L., Lee, J.S., Ebrahimi, T., 2010. Object-
739 based tag propagation for semi-automatic annotation of images, in: Pro-
740 ceedings of the International Conference on Multimedia Information Re-
741 trieval, pp. 497–506.

- 742 [13] Jäschke, R., Eisterlehner, F., Hotho, A., Stumme, G., 2009. Testing and
743 evaluating tag recommenders in a live system, in: Proceedings of the 3rd
744 ACM Conference on Recommender systems (RecSys 09), pp. 369–372.
- 745 [14] Jäschke, R., Marinho, L., Hotho, A., Schmidt-Thieme, L., Stumme, G.,
746 2007. Tag Recommendations in Folksonomies, in: Proceedings of the 11th
747 European Conference on Principles and Practice of Knowledge Discovery
748 in Databases, pp. 506–514.
- 749 [15] Li, J., Wang, J.Z., 2006. Real-time computerized annotation of picture, in:
750 Proceedings of the 14th ACM Multimedia Conference, pp. 911–920.
- 751 [16] Lipczak, M., 2008. Tag recommendation for folksonomies oriented to-
752 wards individual users, in: Proceedings of the Conference on Machine
753 Learning and Principles and Practice of Knowledge Discovery in Databases
754 (ECML/PKDD), Discovery Challenge Workshop, pp. 84–95.
- 755 [17] Mann, H., Whitney, D., 1947. On a test of whether one of two random
756 variables is stochastically larger than the other. The annals of mathematical
757 statistics 18, 50–60.
- 758 [18] Marinho, L.B., Preisach, C., Schmidt-Thieme, L., 2009. Relational classi-
759 fication for personalized tag recommendation, in: Proceedings of the Con-
760 ference on Machine Learning and Principles and Practice of Knowledge
761 Discovery in Databases (ECML/PKDD), Discovery Challenge Workshop,
762 pp. 7–15.
- 763 [19] Naaman, M., Nair, R., 2008. ZoneTags Collaborative Tag Suggestions:
764 What is This Person Doing in My Phone? IEEE MultiMedia 15, 34–40.
- 765 [20] Rendle, S., Schmidt-Thieme, L., 2009. Factor models for tag recom-
766 mendation in bibsonomy, in: Proceedings of the Conference on Machine
767 Learning and Principles and Practice of Knowledge Discovery in Databases
768 (ECML/PKDD), Discovery Challenge Workshop, pp. 235–242.
- 769 [21] Salzberg, S.L., 1997. On Comparing Classifiers: Pitfalls to Avoid and a
770 Recommended Approach. Journal of Data Mining and Knowledge Discov-
771 ery 1, 317–328.
- 772 [22] Sigurbjörnsson, B., Zwol, R., 2008. Flickr tag recommendation based on
773 collective knowledge, in: Proceedings of the 17th International Conference
774 on World Wide Web (WWW 2008), pp. 327–336.
- 775 [23] Sood, S.C., Owsley, S.H., Hammond, K.J., Birnbaum, L., 2007. TagAs-
776 sist: Automatic Tag Suggestion for Blog Posts, in: Proceedings of the 1st
777 International Conference on Weblogs and Social Media (ICWSM 07), pp.
778 1–8.
- 779 [24] Sordo, M., 2012. Semantic Annotation of Music Collections: A Computa-
780 tional Approach. Ph.D. thesis. Universitat Pompeu Fabra, Barcelona.

- 781 [25] Toderici, G., Aradhye, H., Pasca, M., Sbaiz, L., Yagnik, J., 2010. Find-
782 ing meaning on youtube: Tag recommendation and category discovery,
783 in: Proceedings of the IEEE Conference on Computer Vision and Pattern
784 Recognition (CVPR 2010), pp. 3447–3454.
- 785 [26] Turnbull, D., Barrington, L., Torres, D., Lanckriet, G., 2008. Semantic
786 Annotation and Retrieval of Music and Sound Effects. IEEE Transactions
787 On Audio Speech And Language Processing 16, 467–476.