# FLICKR GROUPS: MULTIMEDIA COMMUNITIES FOR MULTIMEDIA ANALYSIS

Radu-Andrei Negoescu      Daniel Gatica-Perez

Idiap-RR-18-2010

JULY 2010

# Internet Multimedia Search and Mining

Xian-Sheng Hua, Marcel Worring, and Tat-Seng Chua

June 24, 2010

# Contents

# List of Figures

# List of Tables

# Chapter 1.  Flickr Groups: Multimedia Communities for Multimedia Analysis

**Radu-Andrei Negoescu**
**Idiap Research Institute and Ecole Polytechnique Fédérale de Lausanne**
**Daniel Gatica-Perez**
**Idiap Research Institute and Ecole Polytechnique Fédérale de Lausanne**

**Abstract:**    We present in this chapter a review of current work that leverages on large online social networks' meta-information, in particular Flickr Groups. We briefly present this hugely successful feature in Flickr and discuss the various ways in which metadata stemming from users' interactions with and within groups has been exploited by researchers to improve on state-of-the-art search and browsing algorithms. We then review recent works that have already made use of Flickr Groups, either as a data source, as a way of filtering content, or as a way to reach users for automatic analysis or user studies, and conclude by pointing out to potential directions of future research.

## *Introduction*

It is common nowadays to hear people talking about how "on the internet you can find anything". While this may not hold true for really anything, the amount of information available electronically is staggering. As electronic devices become more ubiquitous, from mobile phones to internet-enabled refrigerators, more and more consumers become producers of content. As of the writing of this piece, YouTube, arguably the world's most popular video sharing website, was boasting 20 hours worth of video uploads each minute [22]. Flickr on the other hand, one of the most popular photo sharing websites, showed an average of 2000 photos uploaded per minute. This average was computed by counting the time it took to upload 1 billion photos starting from the $3^{rd}$ billionth one in November 2008: less than a year.

One of the more interesting aspects of social media today, apart from the sheer amount of data available to the research community, is the even richer metadata. People create, upload, and then annotate content, be it through tags, regions of interest, or comments. Annotation is just one of the explicit ways through which metadata is created - and the most visible, but many other actions also lead to the generation of metadata, such as rating, adding photos or videos to favorite lists, organizing content in sets or collections, and even simply viewing content. This type of implicit metadata is also useful for the understanding of the dynamics of social media. Another way in which social media collections can be enriched is through

the social sphere itself, like user-to-user relationships, or online communities, such as Flickr Groups. These groups are user-managed communities that form around a specific interest, such as a photographic technique, geographical location, event, or simply as a result of a desire for social interaction. There is a great number of groups on Flickr, currently in the range of 300,000. By nature, Flickr Groups become content filtering mechanisms, and this has great potential for multimedia mining.

In this book, the reader has already seen Flickr being used as a data source in Chapters 5 and 7, and some of the following chapters will present further research based on Flickr data. In this chapter we take a comprehensive look at Flickr Groups as one of the most importantn instances of social media communities. While other online photo sharing communities also exist (PicasaWeb, Photobucket, Zooomr), none is as popular as Flickr and indeed none of them offers the same ease of access to their data as Flickr does, through its API. This is the main reason why a lot of multimedia mining research uses Flickr as a data source in the first place.

After briefly describing Flickr Groups, we first analyze how users make use of groups, and then review how the research community has recently begun to take advantage of them in order to address fundamental multimedia problems like image retrieval or search results ranking and re-ranking. We also discuss ways in which the research community, by combining visual content, textual descriptions, and social data, are making progress towards providing users with better content and entity discovery options or personalized content filtering options. The main advantage of metadata over visual content is, still, the computational cost and relatively lower complexity, and this explains why it has been so far the most common source of data for analysis. On the other hand, one of the drawbacks is the relatively higher noise of such metadata, like poor or incomplete annotations. Fortunately, by taking advantage of the massive amount of data available, noise can be somewhat reduced through aggregation techniques. All of these issues are discussed in this chapter. To our knowledge, this represents the first review of the use of Flickr Groups for multimedia analysis.

This chapter is structured as follows: first we present Flickr Groups, from both a functional and a data source perspective. We then review research that has taken advantage of Flickr Groups so far. We end this hapter with a brief discussion in the last section about future directions.

### *What are Flickr Groups?*

Flickr was founded in 2004 and in the following five years became one of the major players in the online photo management market. Part of its success seems to be related to the way photos are central to the website experience: rather than becoming an online photo album storage site like other existing systems, Flickr from the very

beginning encouraged users to share their photos with the rest of the world. Why exactly users share photos online is a question that received tentative answers in several recent studies like those of Van House [19] and Ames and Naaman [1]. Apart from simply needing a place to store their photos online, most users seem driven by social motives like self-expression and self-promotion, and social interaction is indeed one of the key aspects of the Flickr experience. Users can explore random photos automatically ranked by "interestingness", leave comments, add tags, mark-up regions within a photo with notes, list a specific photo as a favorite, or add other users as contacts. All these public displays of interaction have a strong impact on community building.

In Flickr, apart from the previously described interaction mechanisms, users can organize themselves in self-managed communities, called Flickr Groups. As the name suggests, Flickr Groups are sets of users who are brought together by a common feature, and who share photos in the so called *group pool*. There is quite a wide range of interests that bring people together. As a few examples, it may be an interest in a specific kind of photographic technique, leading to groups like *Black and White*, *Closer and Closer Macro Photography*, or *Digital Infrared*. It can also be an interest in a specific geographical location, with groups like *New York City*, *Paris*, or even smaller geographic communities, like the *2010 Olympic Athletes' Village*, shown here in Fig. 1.1. Yet again, it may be an interest in promoting one's images, which leads to groups like *Views 7-25* or *Views 1250-1500*, groups in which people share their photos that have reached the respective number of views on Flickr. Within the same category there are groups like *Nature's Finest (Invited Images Only)*, or *Your Best Shot 2009*, shown in Fig. 1.2. These are communities of sometimes tens of thousands of people, whose common goal is the gathering and exposure of high quality photography. The list of interests that brings people together in groups does not end here: there are the charity oriented groups, the photojournalism ones, the groups for political activism, and even the groups for corporate marketing.

In previous research on a Flickr dataset[12], we analyzed the way users of Flickr use the system, and its groups. First, there are two types of users, ones who have free accounts, and ones who pay in order to have unlimited access to the site's features. The two types of users (paying and non-paying) were almost equally represented (51.4% and 48.6% respectively) in our data set. We show in Fig. 1.3 the relation between the size of the users' photo collections (in number of photos) and the fraction of photos shared in groups. As a first observation, the sizes of the photo collections for users who shared no photos at all were evenly spread over the entire range of sizes (the thick line overlapping the $x$ axis). Furthermore, the sharing fractions for the users who had the maximum number of photos allowed in our dataset (we collected at most 500 photos per user), were also evenly spread over the entire interval $[0, 1]$ (the thick line at $x = 500$). The correlation coefficient

Figure 1.1: Home page of the group *2010 Olympic Athletes' Village*, a very small group with 21 members and just over 150 images. This group is clearly a special interest group, focused on the 2010 Olympic Village.

Figure 1.2: Home page of the group *Your Best Shot 2009*, a group with over 20,000 members and more than 17,000 images. The focus in this group is on social interaction via exposure of high quality photography rather than on a specific subject.

between the two measures was 0.1417, indicating a weak correlation. While the restrictions on free accounts do seem to influence the number of photos users had in their accounts (with an average of around 220 photos for non-paying members as opposed to 450 for paying members) and also the number of groups they shared photos with (on average 60 for paying members, with a median of 23 and an average of 24.7 with a median of 7 for non-paying members), we found that the ratio of photos shared in groups was similar for both categories of users: paying members in our data shared on average 29.4% of their photos (median 17.2%) and non-paying members shared on average 30% (median 17.1%). Our results on the same data showed that, on average, a user shared a small number of photos with each group (mean 9.6, median 5.1) and shared the same photo in multiple groups in even smaller numbers (mean 3.1, median 1.5), with small differences between paying and non-

Relation between the number of photos and the photo sharing fractions



Figure 1.3: The fraction of shared photos ($y$-axis) vs. the number of photos of each user (the $x$-axis): the size of the collection of photos for users who do not share any photos at all is evenly spread over the entire range of sizes $[1, 500]$; the sharing fractions for users who have the maximum number of photos (500) is evenly spread over the full interval $[0, 1]$.

paying members, despite the large differences in the average number of groups noted above. This was an interesting result, showing that users' group-sharing behavior was not influenced by their paying or non-paying status, or by the amount of photos they uploaded. Overall, the analysis showed that through relatively modest photo repurposing, small but persistent group loyalty and active participation in groups, Flickr users contribute a significant proportion of their content to communities, which emerge as rich Flickr entities through the aggregation of their members' contributions.

Thus, from a research perspective, Flickr Groups are interesting for several reasons. First and foremost, many of the groups act like content funnels, gathering in a single place - the group pool - references to photos that match a specific criterion, be it photographic technique, photographic subject, semantic subject, or aesthetic quality. Group administrators, and sometimes even regular group mem-

bers, act as content filters. This is a great resource to tap into for the research community, and several studies have already used groups as a starting point for data collection [10, 9, 16].

Secondly, groups are natural paradigms of the "content+relations" model of social media. That means that not only there is content that is in one way or another consistent, but there is also information about relations between users. Several research groups have started to take advantage of this information in recent studies [8, 5, 18, 12, 13, 11, 17].

Finally, membership in a specific group also brings additional metadata if information about the group itself is included in the dataset, such as the group name, group size, or group type. Group names can be seen as commonly accepted tags by their respective members, or they can be used as groundtruth during evaluation of automatic analysis methods. This kind of metadata has also been exploited in several works [9, 2, 4, 3].

We will have a more in-detail look at these three different ways researchers have been taking advantage of group-related information in the following section.

## *Multimedia Analysis through Flickr Groups*

We start off this section with a table summarizing the existing works which use, in one way or another, Flickr Groups. This list has been compiled by searching online and by manual inspection of the proceedings of the main conferences on multimedia, web and social media, and human-computer interaction. We list in Table 1 the reference, the approximate size of the dataset used, the task dealt with, the approach used for solving the task, and the way in which Flickr Groups were used. Upon inspection, three major uses emerge:

- *groups as metadata*: membership to specific groups becomes metadata used by researchers to model user relations or image similarity, to quantify a user's degree of social participation, or simply to treat group names as additional textual data;

- *groups as funnels for data collection*: starting from certain Flickr groups, researchers gather images for model training with more accurate metadata, select users for in-depth studies of their behavior, or recruit users for ethnographic studies;

- *groups as a research domain on its own*: this is the case when groups themselves constitute a part of (or the totality of) the dataset. Researchers use information such as the number of members, the number of photos, the tags present in the group pool, or the name of the group in order to analyze, understand, or model Flickr Groups as entities of interest themselves.

| Paper | Dataset | Task | Approach | Group Usage |
|---|---|---|---|---|
| Choudhury [3] | 200 groups | predicting group activity over time | interaction metadata | data source |
| Choudhury [4] | 925 groups, 15K images | recommending groups for new images | visual, textual, interaction | data source |
| Egger [5] | 300 groups | extracting automatic hierarchy of groups | membership and group names | data source |
| Lerman [8] | 13K images | personalizing image search | textual and social metadata | additional metadata |
| Negoescu [11] | 10K groups, 8K users | automatic clustering of groups | textual and social metadata | data source |
| Negoescu [13] | 8K groups, 6K users | modeling of groups and users | textual and social metadata | data source |
| Negoescu [12] | 22K users, 51K groups, 6.9M images | analysis of sharing behavior | textual and social metadata | data source |
| Negoescu [14] | 10K groups, 8K users | modeling of groups and users | textual and social metadata | data source |
| Prieur [17] | 72K groups, 4.7M users | analysis of social interaction | metadata | data source |
| Lin [23] | 52 groups, 50K images | modeling group theme evolution over time | visual, social, textual, temporal | data source |
| Chen [2] | 600 groups | recommending groups and tags | visual content | content gathering source |
| Lerman [9] | 55K users | analysis of content discovery | metadata | user selection |
| Miller [10] | 10 users | analysis of photo sharing practices | user studies | user recruitment |
| Plangprasopchok [16] | 21K users | constructing of folksonomies | sets and collection names | user gathering source |
| Nov [15] | 237 users | analysis of tagging motivations | user studies | as measure of social presence |
| Singla [18] | 2.1M users | detecting camera brand congruence | camera and temporal metadata | as prior for friendship strength |
| Van Zwol [20] | 1.83M photos | analysis of social interaction | metadata | as measure of social presence |
| Wang [21] | 103 groups | learning image similarity | visual content | as prior for image similarity, image gathering source |

Table 1.1: Existing works in the literature that use, in one way or another, Flickr Groups. For display reasons, "M" indicates millions, and "K" thousands of items. Table is ordered by group usage and name of first author.

Let us have a more detailed look at each particular usage of Flickr Groups.

### Groups as Metadata

As we mentioned earlier, the simple fact of being a member in a group can be used as metadata. In a study on tagging motivations, Nov et al. [15] examined the tagging behavior of users on Flickr. Starting from findings of previous studies which suggest that social presence influences tagging behavior, the authors used multiple data sources in order to examine which motivations are associated with different tagging levels. One of their data sources were user studies, while for a more computational approach they used the number of Flickr groups a user belongs to as one of the social presence indicators, along with the number of contacts that the user has. Their findings show that, amongst other variables, the variance in tagging activity is best explained by the number of groups a user belongs to. Here, a relatively simple piece of information - the number of groups - helped establish statistically significant evidence to support qualitative research. In a closely related study, Van Zwol [20] tracked the evolution in the Flickr ecosystem of roughly 1.8 million photos uploaded within a 10 day period. The main question in this work is also aimed at understanding user behavior, but instead of targeting the producers of the content, it targeted the consumers (that is, the viewers), trying to identify the factors that explain photo popularity. One of the factors that correlates with the numbers of views a photo receives is the number of groups a photo is shared in. In this case too, metadata stemming from group membership helped verify the research hypothesis, namely that photo popularity is closely related to the social networking behavior of the users. The lesson from these two studies is that the number of groups a user belongs to is a good predictor of their behavior in other areas of the website, such as tagging activity and photo discovery.

In a different study, Singla and Weber [18] looked at Flickr's social network in order to determine whether social relations influence the brand of the camera a user owns. Their dataset included roughly 2.1 million users, and almost 44 milion user-user relationships. Their basic findings are that two users who have a link in Flickr are more likely to own the same brand of camera, with a probability of 0.27 as opposed to 0.19 for two random users. At the same time, the authors made an interesting assumption with respect to the strength of a relationship between two users, taking into account the number of common groups they belong to. Thus, the authors considered users who share a greater number of groups to be *closer* to each other than to other users they share less groups with. In their study however, there appears to be no correlation between the number of common groups and the chances of owning a same-brand camera. The use of group membership information as a measure of friendship strength remains nevertheless an interesting option, to be further explored in subsequent studies. Just as Singla and Weber [18]

used overlapping groups as a measure of similarity for users, Wang et al. [21] used overlapping groups as a measure of similarity for images. The authors used the group membership information in order to solve the problem of image retrieval and organization in user collections that are not necessarily tagged. Their idea was to use images belonging to the same Flickr groups as training data for a classification algorithm that uses visual features, using images that share no groups at all with the training set as negative examples. For their experiments they used as starting point 130 groups that have specific interests, such as objects, like *aquarium* and *car*, specific scene types like *urban* and *sunset*, or abstract concepts like *Christmas* and *smiles*. For these 130 categories, the authors showed an improvement of the classification when the "Flickr similarity" is used as opposed to just visual features without any context.

We have seen in these examples a number of different ways in which the simple fact of belonging to a group can be used as additional information, either as a measure of sociability for the users themselves, or as measures of similarity for users or images sharing the same groups. Whatever the main task, taking into account such information has the advantage of being relatively simple from a computational point of view, and whenever possible it should probably be used.

### Groups as Funnels

A more straightforward way of using Flickr Groups is data gathering. Indeed, as certain groups tend to have very specific themes, they end up acting like content funnels. For example, as shown in Fig. 1.4, the *Everything GREEN!* group requires that any photo shared in the group by its contributing members (920 of them at the time of writing) contain something green. Similar groups can be found for nearly every color in the spectrum, as well as animal species, cars, trains, buildings, airplanes, airports, chairs, trees, leaves ... the list goes on and on. This plethora of *concept groups* has encouraged researchers to start using their photo pools as data sources.

This is for example the case for the work of Chen et al. [2] in which the authors proposed a group and tag recommender. The authors first defined a list of concepts, and then retrieved training data from Flickr by using photo-level and group-level searches. Their experiments showed that group-based training sets lead to better results when compared to photo-based ones, and this is most likely explained by the curating function of group moderators - users who make sure photos submitted to the group abide by the group rules. A photo tagged with concept "X" may or may not contain the specific concept (for instance Kennedy et al. [7] talked about a 50/50 chance of a concept being present in Flickr photos tagged with that concept), but when a photo is added to a group about concept "X" it is more likely to actually contain the concept due to group rules. It might therefore be more pertinent to

Figure 1.4: A partial view of the photo pool for the group *Everything GREEN!*. Members are required to only post photos that contain at least one green element.

gather content starting from groups matching the search term rather than from individual photo search results.

In a study on the ways people discover new content on Flickr, Lerman and Jones [9] have used three different sets of data: one coming from the most popular photos on Flickr, one from a random sample of photos, and a third one from the *Apex* group, a group dedicated, as the name implies it, to high quality photos. They concluded that most of a photo's popularity can be explained by what they called "social browsing", that is users who discover new content mostly through their social network of contacts. In this particular case, a very large group focused on exposing high quality photos was used as a starting point for content gathering.

Other authors have recently used groups as a starting point in user selection. Plangprasopchok and Lerman [16] attempted to construct a common folksonomy by aggregating shallow hierarchies created by many distinct Flickr users. To somehow limit the spread of concepts, they turned to Flickr Groups for user selection, gathering all members of 17 groups interested in a given category (insect photography). For each user the authors then collected information on their photo sets and collections - two photo management structures available in Flickr. Photos can be grouped in sets, and sets can be grouped in collections, therefore forming the shallow hierarchies the authors used for folksonomy extraction. Results seem to be promising, as the authors extract meaningful non-trivial hierarchies from this aggregation effort, with roughly 3200 concepts from data collected from over 21,000 users. A somewhat similar approach for user selection was also used by Miller and Edwards [10] in their ethnographic study on photo-sharing culture and practices of users of Flickr. Although at a much smaller scale - their work involved a user study with a group of 10 people, the authors made use of Flickr's groups in order to recruit participants in the study. The recruiting methods included word-of-mouth, email campaigns, Craiglist ads, and ads within a local Flickr group. It turned out that by far the most successful method was through the Flickr group - another indication that the social aspect is one of the driving forces of Flickr users. They describe a class of users they call *Snaprs* whose photo-sharing practices seem to be deeply tied to the way Flickr Groups work, for online sharing as well as offline activities.

In conclusion, whether looking at collecting a higher quality training set for image retrieval problems, or finding users interested in a certain theme or closely located geographically, Flickr Groups clearly provide an excellent starting point for research.

### Groups as a Subject of Research

Last, but by no means least, we take a look at the research that makes use of groups as primary data source, and treats them as research entities in their own right. We left this discussion to the end as it is the widest spread use of Flickr Groups.

Prieur et al. [17] analyzed Flickr from data collected in 2006. They looked at basic statistics of the Flickr population and its communication patterns, and also looked at groups as a coordination tool. Their analysis, based on a sample of roughly 70,000 groups, seems to support the fact that groups are important community-building entities, and are used by users both as content pools and as social circles. In our own previous work [12], using a dataset of roughly 50,000 groups, we also analyzed Flickr users and groups from a participatory point of view. Our findings showed that users who shared photos in groups did so with roughly 30% of their photo collections, and they participated on average in 50 groups, showing the importance groups have for Flickr users. In the second part of our study we used a probabilistic topic model, namely Probabilistic Latent Semantic Analysis (PLSA), that enabled us to create a topic-based representation for groups. On a dataset of roughly 8,000 groups, each group was regarded as a text document defined by its composing words - the totality of the tags present in the group pool. In Tables 1.2 and 1.3 we show example latent topics extracted from the model, along with groups that are highly probable within that topic, and also some of the photos found in the group pools for topics 1 and 12.　Most topics in the model isolate a certain concept, illustrated here by three instances: topic 1 is about *Plants*, topic 3 is mainly about *New York City*, and topic 12 is about *Beach Landscapes*. The main advantage is that they are learnt automatically from the aggregated tags for all groups, with a global vocabulary of roughly 10,000 words. This topic-based representation can be used as a browsing tool for groups, or as a search improvement, particularly in a query expansion scenario, by going through the topic model and retrieving related tags based on the most likely topic for the search term.

Also dealing with the rather large amount of similar, competing groups, is the work of De Choudhury [3]. The author approaches group recommendation from a slightly different perspective, based not so much on the content itself, as similar groups will host similar content, but from a more social point of view: which is the group that will best reward a user's participation in it? To achieve this, the author proposed a method based on Hidden Markov Models to characterize groups' activity on two planes: content-related activity, such as uploads from users or the process of adding photos as favorites, and user-related activity, such as comments on each other's photos. A dataset of 200 groups and roughly 52,000 users was used for experiments which showed good correlation between the proposed activity prediction model and real group attributes tracked over a period of 30 days, such as the number of new members, the number of comments, the number of favorite images, and the number of new uploads. In a related study, De Choudhury et al. [4] also tried to recommend groups to users for a specific image, however in this case the recommendation took into account visual features of the image, textual tags, and user interaction through comments. Using roughly 15,000 images that belonged to 925 groups, the authors evaluated experimentally the predictive power of their

| Topic 1 | |
|---|---|
| $P(t \mid z)$ | **Tag** |
| 0.0766 | flower |
| 0.0555 | flowers |
| 0.0550 | nature |
| 0.0431 | ilovenature |
| 0.0323 | spring |
| 0.0295 | garden |
| 0.0243 | green |
| 0.0221 | yellow |
| 0.0212 | macro |
| 0.0204 | pink |
| 0.0168 | white |
| 0.0136 | plant |
| 0.0126 | blue |
| 0.0122 | purple |
| 0.0112 | red |
| 0.0110 | flora |
| 0.0109 | canon |
| 0.0095 | rose |

| Topic 1 | |
|---|---|
| $P(z \mid G)$ | **Group Name** |
| 0.9715 | 1-Plants World |
| 0.9456 | Flickr Gardens |
| 0.8783 | In my garden |
| 0.8718 | My Garden |
| 0.8347 | Daffodil World |
| 0.8337 | What plant is that? |
| 0.8214 | Gardening for Fun |
| 0.8102 | Garden Flowers |
| 0.7993 | grow |
| 0.7377 | Backyard Nature |

| Topic 3 | |
|---|---|
| $P(t \mid z)$ | **Tag** |
| 0.1094 | nyc |
| 0.0767 | newyork |
| 0.0441 | newyorkcity |
| 0.0393 | brooklyn |
| 0.0384 | gothamist |
| 0.0363 | manhattan |
| 0.0354 | montreal |
| 0.0344 | ny |
| 0.0206 | quebec |
| 0.0119 | canada |
| 0.0101 | urban |
| 0.0099 | york |
| 0.0096 | new |
| 0.0091 | coneyisland |
| 0.0089 | street |
| 0.0083 | city |
| 0.0079 | usa |
| 0.0067 | subway |

| Topic 3 | |
|---|---|
| $P(z \mid G)$ | **Group Name** |
| 0.9943 | Mermaid Parade, 2007 |
| 0.9866 | Coney Island Mermaid Parade |
| 0.9849 | Coney Island |
| 0.9771 | (718) Brooklyn |
| 0.9224 | N.Y.C. |
| 0.9067 | Curbed |
| 0.8989 | Gothamist |
| 0.8988 | NYC Social |
| 0.8642 | NYC from A to Zed |
| 0.8623 | The NYC Subway |

| Topic 12 | |
|---|---|
| $P(t \mid z)$ | **Tag** |
| 0.0556 | sky |
| 0.0542 | sunset |
| 0.0460 | clouds |
| 0.0411 | beach |
| 0.0349 | sea |
| 0.0267 | water |
| 0.0248 | ocean |
| 0.0233 | blue |
| 0.0177 | sun |
| 0.0124 | sand |
| 0.0108 | sunrise |
| 0.0101 | landscape |
| 0.0089 | cloud |
| 0.0082 | silhouette |
| 0.0076 | boat |
| 0.0076 | coast |
| 0.0073 | desert |
| 0.0070 | nikon |

| Topic 12 | |
|---|---|
| $P(z \mid G)$ | **Group Name** |
| 0.8866 | Beaches & Sunset |
| 0.8480 | wave porn (pls nominate for best of) |
| 0.8430 | Fotos Caribe |
| 0.8338 | *I Love the Ocean/Sea* (Ocean Only, No People Shots) |
| 0.8326 | the sea and its spectrum |
| 0.8126 | Boating |
| 0.8107 | Sky and Sea |
| 0.8070 | Sea |
| 0.8033 | Caribbean perspective |
| 0.8023 | Atlantic Ocean |

Table 1.2: Some of the topics in the PLSA model, characterized by their most probable tags (ranked by the probabilities of the tags given the topic, $P(t \mid z)$), and by their most probable groups (ranked by the probabilities of the topics given the group, $P(z \mid G)$) (taken from [12]).

photos from group *grow*, by *docman*(1), *Ben McLeod* (2,3), *gailf548* (4)

photos from group *Flickr Gardens*, by *Lorika13*, *egg.*, *annethelibrarian*, *Somerslea*

photos from group *Beaches & Sunset*, by *Mallmus*, *marj k*, *The Life of Bryan*, *cakecosas*

photos from group *Sea*, by *Martin Burns*, *mnadi*, *carf*, *Ennor*

Table 1.3: Example photos from group pools, that are highly probable for topics 1 (top row) and 12 (bottom row), shown in Table 1.2 (taken from [12]).

model and found that, when social features were included in addition to the visual and textual ones, predictive performance improved.

Another study that considers Flickr Groups as a central research problem is our work [13] on modeling Flickr users and groups in a joint fashion. In this work we made the simplifying assumption that users and groups can eventually be seen as equivalent entities, as both have collections of photos and inherently tag vocabularies belonging to those photos. We learned a probabilistic topic model on the joint corpus of Flickr entities, much like we had already done with just groups in our previous work [12]. Each group and each user were seen as textual documents composed of their photos' aggregated tags, and the result of the model was a topic-based representation for users and groups alike. The joint model brings users and groups onto a common ground, and as such direct comparison is straightforward. We computed distances between all entities in our dataset of roughly 8,000 users and almost 11,000 groups and could thus simply rank all entities with respect to any other, obtaining recommendation lists of both groups and users. An example is shown in Fig. 1.5, where the "query user" and his topic-based distribution are shown on the left, the top recommended users in the middle, and the top recommended groups on the right. For a comprehensive study on joint modeling of Flickr groups and users we refer the reader to [14].

In a different direction, Lerman et al. [8] used a dataset collected from Flickr of roughly 13,000 images tagged with prototypical terms like *tiger*, *newborn* and *beetle*, on which they attempted to filter the search results based on the searching user's

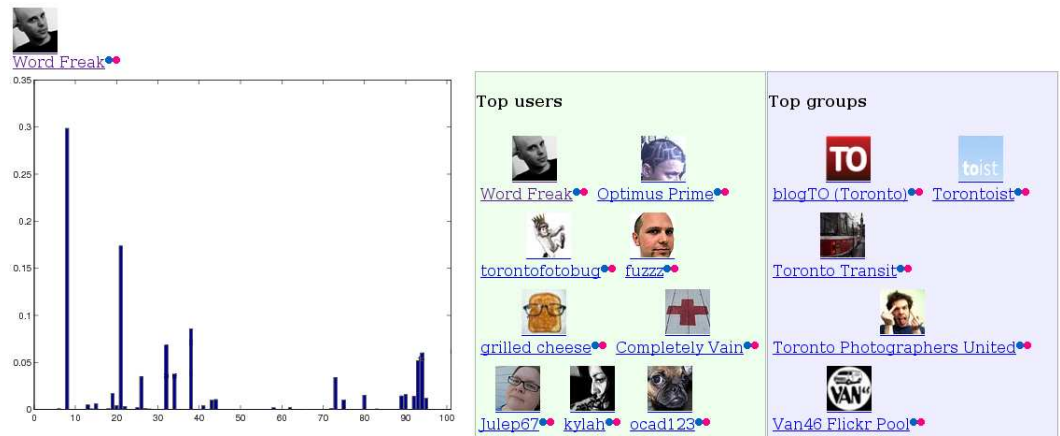Figure 1.5: A user's topic-based representation is shown on the left, with the topics' probability distribution. Based on this representation and a similarity measure, the most similar users and most similar groups are then displayed.

personal information, such as their list of contacts - other Flickr users - and their list of previously used tags. In addition, they explored integrating information from the groups the photos are submitted to in their tag-based model, considering the group name as a commonly agreed-upon tag. Personalization by contacts yielded significant improvement, although as search results are strictly filtered based on whether the images come from the user's social network, it is limited to content retrieval rather than content discovery. Personalization by tags used a topic model for each search set with 10 topics, and showed relative improvements over plain search as well. In this case, including group names as additional tags did not improve filtering results, on the contrary, it sometimes hurt. The authors believe this is mainly an effect of highly social groups, which do not necessarily have a focused photographic interest, and as such are rather noisy. This raises a question for the research community: how to know which groups can be useful in enriching a dataset, and when? This question was at least partially answered by Lin et al.[23] in a paper that tried to extract temporal patterns of themes of interest in Flickr Groups. The authors used visual features as well as social features in order to describe each group as a mixture of themes - each theme is a pattern of image content and context. Context in this case was information related to the owner of the photo, to the tags associated with the photo, and the timestamp of the photo upload. The authors used a dataset of 52 groups and roughly 50,000 images and evaluated their theme extraction method through a tag-prediction task, in which their method outperformed baselines that only used textual or visual features. From a user perspective, theme extraction can represent a novel way to browse group

content. From a research perspective, knowing the rates at which a group's interest changes may become valuable metadata that describes the suitability of a given group for different tasks, like content gathering: the more stable a group's themes, the more chances of the content being homogeneous and suitable for training sets, and reversely, the more dynamic a group's themes, the higher the chances of mixed content.

We conclude with two recent works that deal with the lack of macro-structure related to Flickr Groups. At the time of writing, discovering new groups in Flickr was still a matter of searching by keywords or of serendipitous discovery while browsing someone else's photos. There is no hierarchy per se, nor any other kind of classification. In a study using a dataset of 300 groups, Egger et al.[5] used a membership-based measure they termed *GroupConnectivity* in order to perform community segmentation. This measure is simple to compute, as the fraction between the number of shared members of two groups and the total number of members of the smallest of the two groups. As such, this ratio is bounded by 0 if two groups have no members in common, and by 1 if all members of the smaller group are also members of a larger group. Taking this a step further the method builds, using the same measure of connectivity, a tree of groups. Their assumption was that larger groups are semantic parents of smaller groups with which they are highly connected. Through their experiments their assumption seemed to be generally confirmed, although some counterexamples were also found. The authors obtained semantically meaningful taxonomies, partially shown here in Fig. 1.6. Every node in the tree is a group, and edges imply dependence. With respect to the computational effort involved, this method of automatically extracting taxonomies of Flickr Groups seems very appealing. In a similar study [11] we also looked at the problem of organizing Flickr Groups, by finding what we called *hypergroups*. Our approach was to use Latent Dirichlet Allocation, a probabilistic topic model, to describe each group as a mixture of topics of interest. We then computed a similarity measure based on the Jensen-Shannon divergence, and fed the similarity matrix into an Affinity Propagation (AP) [6] algorithm. The advantages of the AP clustering are manifold: (i) the number of clusters is determined automatically at run-time from the data points by message passing between all exemplars; (ii) at the end of the algorithm, each cluster "center" is an actual data point; (iii) the algorithm is fast to converge. In Fig.1.7 we show clustering results on a dataset of over 10,000 groups. The hypergroups are generally homogeneous, and bring together groups that are not similar by name, but rather at a more abstract level, like for example *RUSTY and CRUSTY* and *Things that Moved*. For a human observer, making the link between things that moved and things that are rusty may be straightforward, but in the absence of keywords, this link is hard to find through a computer search. An even better example of the advantage brought by topic-based clustering is hypergroup 578, which corresponds to groups dedicated to art *except* photography.
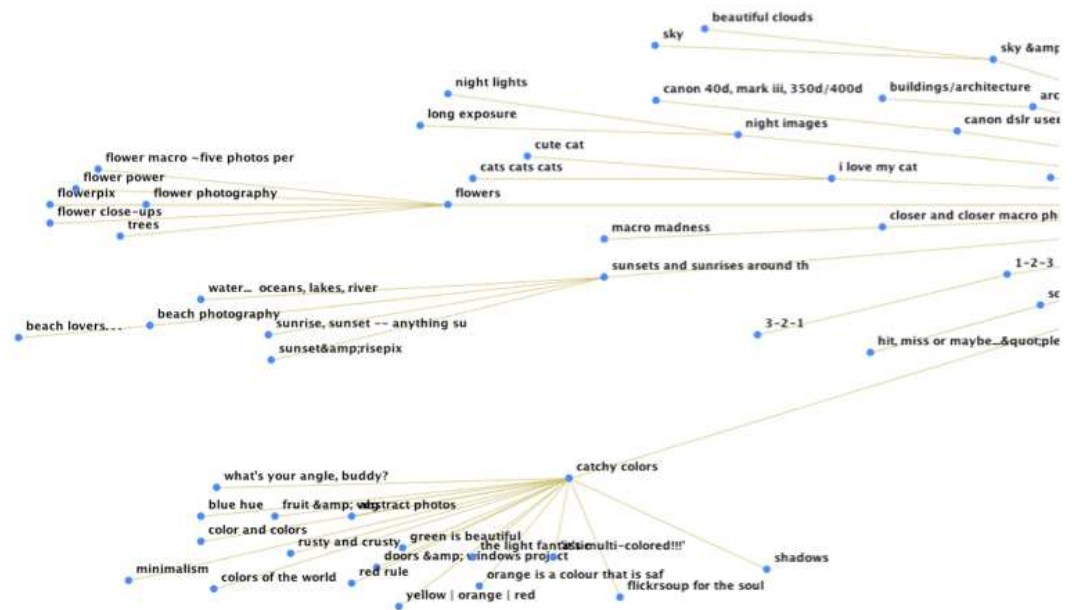
Figure 1.6: Partial view of the automatically discovered taxonomy on 300 Flickr groups. Every node is a group, and edges between groups imply dependence; image courtesy of Egger et al. [5]

The group names are much less homogeneous, but quite clearly the hypergroup is homogeneous from an interest point of view.

## Concluding remarks

In this chapter we have reviewed the existing ways in which Flickr Groups have been used in research in order to either better understand the dynamics of online communities, to improve traditional multimedia problems such as image search, or to help users navigate Flickr Groups themselves. We have seen that, more often than not, including metadata stemming from the social activities of users is beneficial, be it in the form of user interaction through comments, views, and votes, or through membership in specific interest groups. In addition, from a computational point of view this kind of metadata is also almost always cheaper than the content itself (that is, visual features).

We also believe there are several new and interesting avenues to explore, as metadata gets richer and more social media websites will provide access to it. Understanding how this data can be used in order to enrich the user experience or facilitate exploration of ever growing content is of high relevance. Assuming that

Figure 1.7: Hypergroups obtained after clustering through Affinity Propagation. In hypergroup 65, *Strobist.com* is the Flickr group started by a hugely popular US based photographer, David Hobby, who happens to be a Nikon user. Hypergroups 16 and 889 leave no doubt about their member groups' interests. Hypergroup 578 brings together groups dedicated to art drawings and paintings, much less homogeneous from a group names perspective, but quite clearly homogeneous from an interest point of view (taken from [11]).

ways of federating user identities across multiple social media platforms become feasible, future work may investigate whether user behavior and metadata are consistent across social media platforms, or whether specific factors (such as system design and affordances) may impact the way users create, use, and annotate content.

We hope this short review of current research will help practitioners, researchers, and graduate students in social (multi)media processing to better understand the potential of rich communities like Flickr Groups.

## *Acknowledgements*

## *References*

[1]   Morgan Ames and Mor Naaman. Why We Tag: Motivations for Annotation in Mobile and Online Media. In *CHI '07: Proc. of the SIGCHI Conf. on Human Factors in Computing Systems*, San Jose, CA, USA, 2007.

[2]   Hong M. Chen, Ming H. Chang, Ping C. Chang, Ming C. Tien, Winston H. Hsu, and Ja L. Wu. Sheepdog: Group and Tag Recommendation for Flickr Photos by Automatic Search-based Learning. In *MM '08: Proc. of the 16th ACM Int'l Conf. on Multimedia*, Vancouver, Canada, 2008.

[3]   Munmun De Choudhury. Modeling and Predicting Group Activity over Time in Online Social Media. In *HT '09: Proc of the 20th ACM Conf on Hypertext and Hypermedia*, Torino, Italy, 2009.

[4]   Munmun De Choudhury, Hary Sundaram, Yu-Ru Lin, Ajita John, and Doree Duncan Seligmann. Connecting Content to Community in Social Media via Image Content, User Tags and User Communication. In *Intl. Conf. on Multimedia and Expo (ICME)*, New York, NY, USA, 2009.

[5]   Marc Egger, Kai Fischbach, Peter Gloor, Andre Lang, and Mark Sprenger. Deriving Taxonomies from Automatic Analysis of Group Membership Structure in Large Social Networks. In *Lecture Notes in Informatics, vol 154, Proc. of Informatik 2009*, Lubeck, 2009.

[6]   Brendan J. Frey and Delbert Dueck. Clustering by passing messages between data points. *Science*, 315:972–976, 2007.

[7]  Lyndon S. Kennedy, Shih-Fu Chang, and Igor V. Kozintsev. To Search or to Label?: Predicting the Performance of Search-based Automatic Image Classifiers. In *Mir '06: Proc. of the 8th ACM Int'l. Workshop on Multimedia Information Retrieval*, Santa Barbara, CA, USA, 2006.

[8]  K. Lerman, A. Plangrasopchok, and C. Wong. Personalizing Results of Image Search on Flickr. In *AAAI workshop on Intelligent Techniques for Web Personlization*, Vancouver, Canada, 2007.

[9]  Kristina Lerman and Laurie Jones. Social Browsing on Flickr. In *Proc. of Intl. Conf. on Weblogs and Social Media (ICWSM)*, Boulder, CO, U.S.A., March 2007.

[10]  Andrew D. Miller and W. Keith Edwards. Give and Take: a Study of Consumer Photo-sharing Culture and Practice. In *CHI'07: Proc. of the SIGCHI Conf. on Human Factors in Computing Systems*, San Jose, CA, USA, 2007.

[11]  Radu Andrei Negoescu, Brett Adams, Dinh Phung, Svetha Venkatesh, and Daniel Gatica-Perez. Flickr Hypergroups. In *MM '09: Proc. of the 17th ACM Intl. Conf. on Multimedia*, Beijing, China, October 2009.

[12]  Radu Andrei Negoescu and Daniel Gatica-Perez. Analyzing Flickr Groups. In *CIVR '08: Proc. of the Intl. Conf. on Image and Video Retrieval*, Niagara Falls, Canada, July 2008.

[13]  Radu Andrei Negoescu and Daniel Gatica-Perez. Topickr: Flickr Groups and Users Reloaded. In *MM '08: Proc. of the 16th ACM Intl. Conf. on Multimedia*, Vancouver, Canada, October 2008.

[14]  Radu Andrei Negoescu and Daniel Gatica-Perez. Modeling Flickr Communities through Probabilistic Topic-based Analysis. *IEEE Transactions on Multimedia*, 2010.

[15]  Oded Nov, Mor Naaman, and Chen Ye. What drives content tagging: the case of photos on flickr. In *CHI '08: Proc of the 26th SIGCHI Conf. on Human Factors in Computing Systems*, Florence, Italy, 2008.

[16]  Anon Plangprasopchok and Kristina Lerman. Constructing Folksonomies from User-specified Relations on Flickr. In *WWW '09: Proc. of the 18th Intl. Conf. on World Wide Web*, Madrid, Spain, 2009.

[17]  Christophe Prieur, Dominique Cardon, Jean-Samuel Beuscart, Nicolas Pissard, and Pascal Pons. The Strength of Weak Cooperation: a Case Study on Flickr. Retrieved on Jan 21, 2010 from http://arxiv.org/abs/0802.2317, Feb 2008.

[18] Adish Singla and Ingmar Weber. Camera Brand Congruence in the Flickr Social Graph. In *WSDM '09: Proc. of the 2nd ACM Intl. Conf. on Web Search and Data Mining*, Barcelona, Spain, 2009.

[19] Nancy A. Van House. Flickr and Public Image-sharing: Distant Closeness and Photo Exhibition. In *CHI'07: Extended Abstracts on Human Factors in Computing Systems*, San Jose, CA, USA, 2007.

[20] Roelof van Zwol. Flickr: Who is Looking. In *WI '07: Proc. of the Intl. Conf. on Web Intelligence*, San Jose, CA, USA, 2007.

[21] Gang Wang, Derek Hoiem, and David Forsyth. Learning Image Similarity from Flickr Groups Using Stochastic Intersection Kernel Machines. In *Proc. of the 12th Int'l Conf. on Computer Vision*, Kyoto, Japan, 2009.

[22] YouTube. Youtube Factsheet. Retrieved on Jan 21, 2010, http://www.youtube.com/t/fact_sheet, January 2010.

[23] Yu-Ru Lin, Hari Sundaram, Munmun De Choudhury, and Aisling Kelliher. Temporal Patterns in Social Media Streams: Theme Discovery and Evolution using Joint Analysis of Content and Context. In *ICME'09: Proc. of IEEE Int'l. Conf. on Multimedia & Expo*, New York, NY, USA, 2009.