

# Gathering Training Sample Automatically for Social Event Visual Modeling

Xueliang Liu  
EURECOM  
Sophia Antipolis, France  
xueliang.liu@eurecom.fr

Benoit Huet  
EURECOM  
Sophia Antipolis, France  
benoit.huet@eurecom.fr

## ABSTRACT

In recent years, the emergence of social media on the Internet has derived many of interesting research and applications. In this paper, a novel framework is proposed to model the visual appearance of social events using automatically collected training samples on the basis of photo context analysis. While collecting positive samples can be achieved easily thanks to explicitly identifying tags, finding representative negative samples from the vast amount of irrelevant multimedia documents is a more challenging task. Here, we argue and demonstrate that the most common negative sample, originating from the same location as the event to be modeled, are best suited for the task. A novel ranking approach is devised to select a set of negative samples. The visual event models are learned from automatically collected samples using SVM. The results reported here show that the event models are effective to filter out irrelevant photos and perform with a high accuracy on various social events categories.

## Categories and Subject Descriptors

H.3.1 [Content Analysis and Indexing]: miscellaneous

## General Terms

Algorithms, Design, Performance

## Keywords

Events, social media, multimedia semantics

## 1. INTRODUCTION

In the past few years, social media have becomes an integral part of many people's life. Thanks to the rapid increase of websites like Facebook, Flickr, YouTube, people are able to share information in a quick, simple and cheap way. With the increasing popularity of social media, users generated content is being produced and shared online at an unprecedented rate. Real problems are beginning to surface from

this situation and are generating growing interest within the multimedia research community, such as how to analyze the semantic pattern and how to mine valuable information from such big data.

In their daily life, people naturally organize their personal data according to occurring events; holiday, wedding, birthday party, concert, etc... Events are a natural way for referring to any observable and describable occurrence grouping persons, places, times and activities [16]. Events are also observable experiences that are often documented by people through different media (e.g. blog, videos and photos). The intrinsic connection between media and experiences are has been explored in previous work [11, 5], These works aim at associating media data with events by exploring their rich contextual information (metadata). However, it is well known that missing or inaccurate data is a frequent issue in user contributed data, which limits the application of these methods.

Besides metadata, the main content in social media is the visual content, in the format of photos or videos (audio being only present in videos). In the multimedia community, remarkable progress has been made on visual content based analysis. Much work has been done to model concepts using low level visual feature and machine learning techniques [4]. However, the labeling of a large dataset is a compulsory step to any supervised learning process. Furthermore, manual labeling is a particularly expensive and tedious task, which impedes its widespread utilization for social media content. To address this important issue, we present a framework aimed at collecting high quality social event data in an automated way from the Internet. The data is acquired based on the analysis of rich event context metadata and is then used for training, verifying and testing visual models. The positive samples are obtained based on specific tags which identify the events accurately, in the form of machine tag and abbreviation of event title. The negative samples are selected by ranking localized photos candidates based on their associated tag commonness. Finally, both positive and negative samples are employed to train event models, which are then evaluated on a manually labeled ground-truth, demonstrating the effectiveness of the proposed approach compared with 3 baseline experiments. The contributions of this paper are twofold:

- We propose framework to collect the training samples automatically. The collection is done by analysing both social media and events contextual information. It shows a possible solution on how to use the social media data in visual content based analysis.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SAM'12, October 29, 2012, Nara, Japan.

Copyright 2012 ACM 978-1-4503-1586-9/12/10 ...\$15.00.

- Compared with previous work which uses K-NN to classify photos in an event-driven search task where no negative samples are available, we tackle visual modeling of social events as a real classification problem. The visual properties of each event are learned using well known and widely employed SVM.

The remainder of this paper is organized as follows: After briefly reviewing the related work in Section 2, we describe the proposed framework in Section 3. Experimental results obtained on a dataset of several thousand media documents for 14 difference events are reported in Section 4. Finally, we summarize and discuss future work in Section 5.

## 2. RELATED WORK

In the past few years, research on how to mine the relation between the end-user experience and multimedia content has drawn lots of attention. The methods found in the literature addressing this issue cover many multi-modal processing techniques. Therefore, we address the related work from a number of relevant research directions, including: event illustration by media data; event detection from social media data; content based media analysis; as well as automatically training sample collections.

In the past few years, research on how to enrich the end-user experience with multimedia content has drawn lots of attention. A tremendous amount of work has been done in very different areas. Illustrating events with media data studies the problem of how to leverage vivid visual content to share experience. In [5], the authors proposed a framework to generate photos collections of news to enhance user’s experience while reading news articles. They computed the similarity between news text and image tags and obtained the relevant images using text retrieval techniques. In [11], an approach aimed at creating a vivid experience to user browsing public events such as concerts is proposed. They studied the user uploading behaviors on Flickr and matched concert with photos in different modalities, such as text/tags, time, and geo-location. This results in an enriched photos set which better illustrates the event.

The study of event detection from media data has also gained a lot of attention within the multimedia research community in the past years. The objective of event detection is to discover events out of previous record that occur on given location and time. To address the problem, Quack *et al.* [14] presented methods to mine events and object from community photo collections. They clustered the photos with multi-modal feature and then classified the results into events and objects. A similar problems is also studied in [6] where Firan *et al.* focused on building a Naive Bayes event models which classify photos as either relevant or irrelevant to given events. In [1, 2], the authors follow a very similar approach, exploiting the rich “context” associated with social media content and applying clustering algorithms to identify social events.

Much of the previously presented approaches aimed at mining the intrinsic connection between event and media, do so using metadata comparison (i.e. time, location, owner, tags, etc...). Only little work has been done on the analysis the visual content of medias in the context of event, and this is precisely the issue we address with this paper.

The usage of low-level visual features for improving content-based multimedia retrieval systems has made great progress

[4]. TRECVID [13] is a video retrieval benchmark that make its effort on content based video retrieval. To address the problem of web visual data analysis, some large scale datasets have been built using multimedia data crawler from shared portals [3]. Beside those web datasets built very recently, a number of learning techniques performed on these dataset have shown acceptable results[17, 7]. And many works [8, 15] have been done to study how to collect online data for the training purpose. In [8], Li et al proposed their work on how to train visual concept model by data collected from Internet automatically. An improved work is reported in [15], where the authors employed text, meta-data and visual information in order to achieve better performance.

In [9], the authors also present work related to the collection the negative training samples from the semantic analysis of tags and visual features. However, the approach can not be adapted to solve our problem since concepts belonging to negative samples are not known in advance. Although the work presented in this paper also tackles the problems of training sample collection, its objective is radically different. Here, we attempt to semantically associate media with missing or inaccurate metadata to their corresponding social event. In addition, our solution leverages the rich event context to build the collection of online media samples, using an approach inspired from tag based photo ranking techniques. Finally, those automatically learned visual models are employed to effectively classify photos with insufficient and/or erroneous metadata.

## 3. THE VISUAL EVENT MODEL

We propose an original scheme for collecting the training samples for modeling social events visual semantics without human assistance. We define a social event as the specific happening that take place at a given location and time and involve many persons (i.e. concerts, conferences, exhibitions, etc...). Figure 1 depicts the automated steps leading to the creation of the dataset employed for learning event models. The positive samples are collected on social media platforms using identification tag based query. The identification tags are those defining the event accurately (i.e. event machine tag).

Collecting the representative negative sample is a more challenging task due to the vast amount of irrelevant data available. Here, negative samples are retrieved from online social media data using metadata analysis. We have observed while experimenting that when querying for photos originating from an event, based on its date and location, the negative samples (those photos which do not correspond to this particular event) are photos depicting general concepts for this location. Among such photos one typically finds, buildings, objects and portraits, etc... and some of the tags associated with these media are common for this location. For example, the city name is a popular tags in many situation yet it doesn’t allows to accurately define an event. In the work presented here, we consider those photos captured at the same location to events and containing common tags as the most relevant negative samples for this specific event. Common tags, along with their corresponding photos, are identified based on a novel approach inspired from learning to rank [10], which we detail in section 3.2.

Having collected both positive and negative visual examples of a particular event, machine learning approaches can be employed to learn the visual model. The methodology

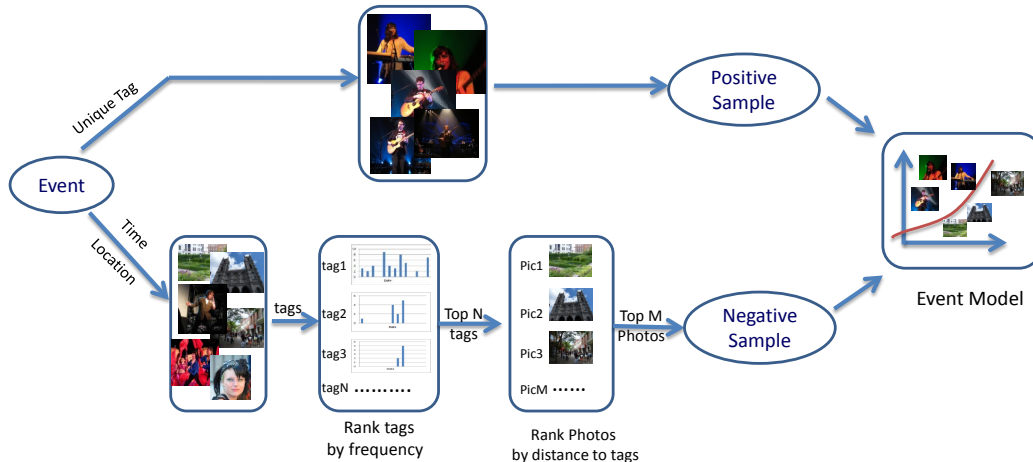


Figure 1: Overview of the framework for modeling events semantic

used to train the Support Vector Machines used in this work is detailed in 3.3.

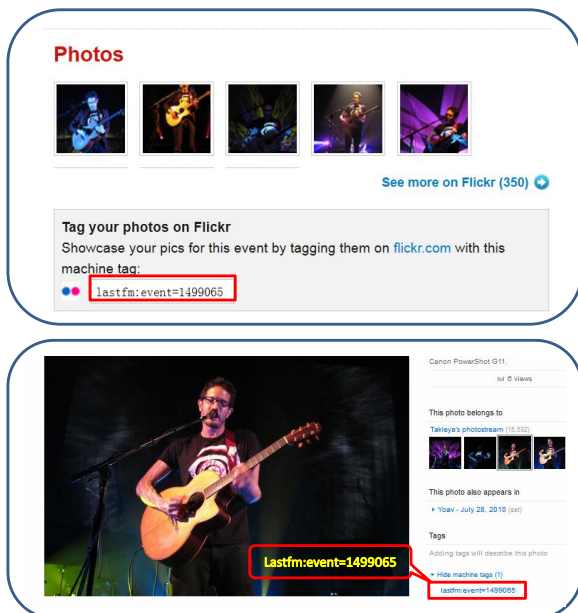


Figure 2: Machine Tags Used in Last.fm(Top) and Flickr(Bottom)

### 3.1 Positive Samples Collection

We collect social event visual positive samples by querying social media platforms with event identification tag. The identification tag of an event is that tag that precisely refers to this unique event. There are different kinds of tags to identify events in social media data. The machine tags is a overlap metadata that is available from some events repositories (such as LastFM<sup>1</sup>, Upcoming<sup>2</sup>, Facebook<sup>3</sup>) and can be used by users to refer the event when they upload media data taken during the event, so it is popularly used to connect events and photo/video in media sharing platforms, such as

<sup>1</sup><http://www.last.fm>

<sup>2</sup><http://www.upcoming.org>

<sup>3</sup><http://www.facebook.com/events/>

Flickr<sup>4</sup>. In these social event website, machine tags are formatted as “\$DOMAIN:event=\$XXX”, where “\$DOMAIN” is the name of website, and “\$XXX” is unique event id provided by the event side, for example, “lastfm:event=1842684” is an event registered in Last.FM whose id is 1842684, and “facebook:event= 108938242471051” is a public event in facebook whose id is 108938242471051. When users take photos during the event, they could upload them to media shared website with such as a tag that provides the background of the photos. The machine tags can be recognized by both kinds of web service and give explicit and accurate links between events and multimedia documents. The media documents containing the appropriate machine tag are taken as positive samples for the corresponding event.

Machine tag is frequently and commonly used nowadays. In the Flickr namespace, there are about 2 million Last.fm event tags, and 400 thousand Upcoming event tags<sup>5</sup>. However, many real world events still do not feature such metadata. To overcome this issue, we use the abbreviated event name to identify such events. Such events abbreviations are well known and popular among the attendees. For example, “ACMMM10” is short for ACM conference on Multimedia 2010, without any ambiguity. All photos with such tag are assumed to be positive samples of this social event.

### 3.2 Negative Samples Collection

Since social events are characterized by a grouping of people at a given time and place, we argue that it is more reasonable to collect the most relevant negative samples from images taken around the same period and location as the event but do not originate from the event. Here is an example to motivate our assumption. Given an event held in a city near a famous landmark, it is likely that among the photos taken by attendees some will show the landmark. As a famous landmark, it is expected to be captured frequently by tourist too. Therefore, it is important that such photos are included in the negative samples in order to differentiate between the event and its surrounding.

In this paper, we collect the most representative negative samples using a ranking approach, that identifies the most representative photos as the ones with tags referring to the common concept in the location. The tags are integrally

<sup>4</sup><http://www.flickr.com>

<sup>5</sup><http://www.husk.org>

considered as carrying such concept. Let  $C$  as target concept, and  $T = T_i$  as the tag list of an image  $I$  containing  $n$  tags. The probability of  $C$  in  $I$  is defined as:

$$P(C|I) = \frac{P(T|C) * P(C)}{P(T)} \quad (1)$$

where the prior probability  $P(C)$  and  $P(T)$  can be viewed as a constant for the purpose of ranking the images. We assume that the concept  $C$  is dominant in the location but different from the event. Our solution to estimate  $C$  and calculate  $P(C|I)$  can be solved in 3 steps, as follows;

The first step consists in gathering the photo candidates. For each event, online services (i.e. LastFM) are used to identify the location and date. These parameters are then employed to query the Flickr API for a photo set ( $P$ ). The location is defined by a circle, whose center is determined by the GPS coordinates of the event venue and radius value ( $R$ ). The time interval is the period of  $D$  days before and after the event’s date. In order to obtain a large set of candidate photos, appropriate values should be set for both  $D$  (days) and  $R$  (kms). The influence of those two parameters will be studied in the experiment section 4.2.

The second step is to estimate the concept  $C$ . We argue that concept  $C$  can be appropriately expressed by the “common tags” associated with photos taken in a given location. The “common tags” are tags which are commonly and frequently associated by users with photos taken at a location, and as such are different from event specific tags. The commonness of a tag can be measured by the number of days it appears within a given period for a given location. More formally, the commonness of tag  $t$  can be calculated as:

$$Score(t) = \sum_{i=1}^D SD(t, i) / D$$

where the value of  $SD(t, i)$  is 1 if tag  $t$  appears on day  $i$ , and 0 if not.

We rank the tags according to their score decreasingly. The top  $N$  tags are kept as a group of common tags  $CTags$ . These tags are prevalently used and highly relevant to the location but do not represent an event due to the fact that they cover a too large time-span. The effect of  $N$ , the number of common tags kept to represent the location, is also studied in the experiment section 4.2.

The last step is to calculate the probability of  $C$  in  $I$  so that negative photo samples could be selected based on commonness ranking. For each photo  $p$  of  $P$ , we extract the title and tags as their text description  $Text(p)$ , and compute the similarity between those terms and the common tags obtained previously. The measure used here is the cosine distance.

$$P(C|I) = \frac{CTags \cdot Text(p)}{\|CTags\| \|Text(p)\|}$$

All of the negative candidate are ranked by their textual similarity to the common tags set ( $CTags$ ) and the top  $M$  photos are kept as negative samples for training the visual model.

Having collected both positive and negative visual examples of a particular event, machine learning approaches can be employed to learn the visual model. The methodology used to train the Support Vector Machines used in this work is detailed in 3.3.

### 3.3 Model Training

Individual event model is obtained as follows; First, 128D Scale Invariant Feature Transform (SIFT) feature is computed over the local region detected by DoG filter, then we cluster all of the visual feature with K-means for each event, and the SIFT description is quantized to generate 300-dimensional Bag of Visual Words. The event model is learned by *Support Vector Machine* with *Radial Basis Function* kernel for learning. Model parameters are optimized using cross-validation methods.

## 4. EXPERIMENTS

### 4.1 Data Set and Experiment Setting

Our proposed algorithm is evaluated on different type of events, including 10 concerts from LastFM, 3 scientific conferences and 1 popular carnival. The photo source used here is Flickr, although other media and source could be easily added to the framework. The details of the dataset for each event could be found in Table 1.

For our experiments, three photo sets are created. The first set contains all the Flickr photos which match the identification tag of the 14 selected events. We randomly split the positive photos originating from each event into two parts according to usage: 50% for training, 50% for verifying.

The second set contains the negative candidates. Photos that are taken within a given spatial distance (less than  $R$  Kms) and within a given temporal interval (less than  $D$  days) of each of the selected events are retrieved from Flickr. The process of common tags generation and photos ranking is performed on this photo set in order to retain only the 200 most common photos for each event as negative samples for training the model. We selected such value in order to balance both positive and negative training samples.

The third set of media is called Real Online data (**RO**) and is used to evaluate our approach in a real life situation. The collection is obtained using Flickr queries combining text, location and time as presented in [11]. The ground truth on this collection is provided by manual human labeling.

The number of photos for each event of the three sets can be found in Table 2.

Table 2: The media collection

EventID	Positive Samples	Negative Candidate	RO	
			Pos	Neg
lastfm:804783	441	1063	466	64
lastfm:1830095	716	748	398	134
lastfm:1858887	408	745	431	266
lastfm:1499065	348	712	16	153
lastfm:1787326	446	913	0	313
lastfm:1351984	307	584	498	19
lastfm:1842684	602	1125	535	78
lastfm:2020655	538	745	750	6
lastfm:1301748	944	541	1157	80
lastfm:1370837	592	1025	592	115
SIGIR2010	100	557	178	23
ACMMM07	30	525	0	201
ACMMM10	118	64	15	44
NICECarnival2011	52	848	60	209
Total	5642	10195	5096	1705

Table 1: the Event DataSet

EventID	Title	Date	Latitude	Longitude
lastfm:804783	Metallica	03/03/2009	54.964053	-1.622136
lastfm:1830095	Hole in the Sky Bergen Metal Festival XII	24/08/2011	60.389585	5.323773
lastfm:1858887	Duran Duran	23/04/2011	41.888098	-87.629431
lastfm:1499065	Osheaga en Ville	28/07/2010	45.509788	-73.563446
lastfm:1787326	The Asylum Tour: The Door	03/03/2011	34.062496	-118.348874
lastfm:1351984	Bospop 2010	10/07/2010	50.788893	5.708738
lastfm:1842684	Buskers Bern	11/08/2011	46.947232	7.452345
lastfm:2020655	Lacuna Coil - Darkness Rising Tour	18/11/2011	50.723090	-1.864967
lastfm:1301748	End Of The Road Festival	10/09/2010	50.951341	-2.082616
lastfm:1370837	Into The Great Wide Open	03/09/2010	52.033333	4.433333
ACMMM10	the ACM conference on Multimedia 2010	25/10/2010	43.777846	11.249613
SIGIR2010	ACM Special Interest Group on Information Retrieval,2010	19/07/2010	46.194713	6.140347
ACMMM07	the ACM conference on Multimedia 2007	24/09/2007	48.334790	10.897200
NICECarnival2011	the Carnival de Nice 2011	05/03/2011	43.701530	7.278240

We use half the positive samples and 200 negative samples to train SVM model for each event, and optimize the parameters  $D$ ,  $R$  and common vocabulary size  $N$  using the verification data.

## 4.2 Location Distance, Time Interval and Tags Size

We investigate the impact of parameter  $R$ , and  $D$ , the location distance and time interval between photo taken and event held, to the final event model. We change the two parameters gradually and test the trained model on the verification dataset. Specifically,  $R$  is chosen from 4 to 20 with step of 4 kms, and  $D$  is set from 5 to 30 with step 5 days. Cross-validation on the two parameters is performed in the process. Figure 3 shows an example of resulting classification accuracy averaged over the different size of common tags vocabulary. Results for all selected events favor the use of rather large values for both time interval and location distance parameters. This finding is supported by the fact that the larger the values of  $D$  and  $R$ , the more photos are retrieved from Flickr and this results in increased diversity within the selected negative samples. As a result, we set  $R$  and  $D$  to 20km and 30 respectively for all further experiments.

We also evaluate the influence of  $N$ , the number of common tags employed, with respect to the resulting event model accuracy. For each combination of parameters  $R$  and  $D$ , we optimize the model with vocabulary size varying from 5 to 50 tags. The results, presented in Figure 4, clearly indicate that the best performance is obtained for a vocabulary of 10 tags.

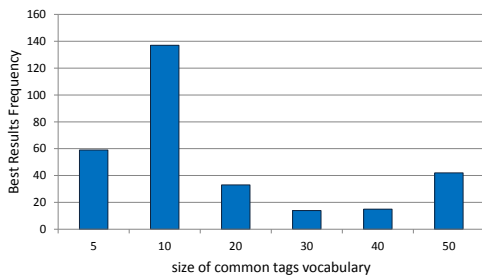


Figure 4: Performance (accuracy) vs size of common tag vocabulary

## 4.3 Performance Evaluation

Having carefully chosen the parameters ( $R$ ,  $D$  and  $N$ ), we evaluate the optimized visual models on manually labeled real online data (**RO**). The results of the evaluation runs are measured in terms of classification accuracy (Acc) [12] and presented in Table 3. Our automatically learned visual event models are compared with four other approaches at the task of mining online media illustrating events and collecting training sample effectively. The first and also the most basic approach, consist in simply running a Flickr query (as the one used to create the real online data) and assuming all returned media are positive. In other words, the accuracy value reported in the column **Query**, indicates the proportion of correct photo event associations in (**RO**). The second approach reported for comparison is one where the SVM model is replaced by a K-NN visual filter proposed in [11]. In addition, we compare different negative sample collection methods. In the third approach (column **Random Sample**), the negative samples are randomly selected from the localized negative candidates. In order to evaluate the influence of “location”, a unique set of 200 negative samples is randomly selected from the entire set of (200 \* 14 events) negative samples and used to train all SVM models (column **Uniform Negative**).

Table 3: Performance Evaluation (Accuracy)

EventID	Query	Our Algorithm	Pruning in [11]	Random Sample	Uniform Negative
lastfm:804783	87.92	88.68	46.98	50.00	75.85
lastfm:1830095	74.81	78.38	80.26	96.62	84.96
lastfm:1858887	61.84	63.41	63.56	76.47	73.89
lastfm:1499065	9.47	90.53	89.94	92.90	89.35
lastfm:1787326	0.00	98.40	92.65	97.12	42.49
lastfm:1351984	96.32	96.32	55.32	86.65	93.81
lastfm:1842684	87.28	87.93	67.86	79.28	87.11
lastfm:2020655	99.21	91.80	71.69	75.00	94.58
lastfm:1301748	93.53	93.53	73.73	64.83	93.21
lastfm:1370837	83.73	85.15	73.83	60.25	80.62
SIGIR2010	0.00	60.19	42.28	16.41	22.38
ACMMM07	25.01	57.62	46.61	28.81	27.18
ACMMM10	85.83	91.04	87.56	86.57	89.05
NICECarnival2011	22.30	76.58	59.10	55.39	56.51
Average	69.41	83.31	68.64	70.07	73.42

From the results presented in table 3, it is interesting to note that the approach proposed in [11] for analyzing visual content achieves, on average, almost the same performance as the Flickr **Query**. When compared with the approach

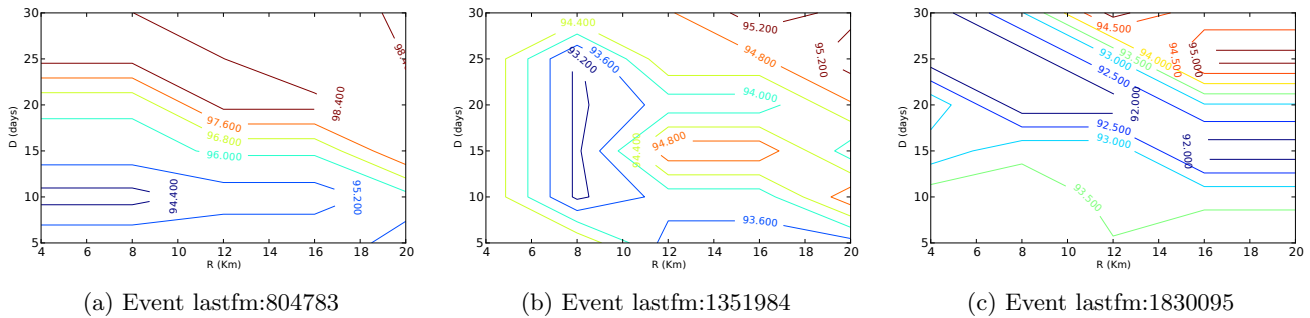


Figure 3: Cross Validation on  $R$  (location distance) and  $D$  (time interval) for 3 events, measured according to accuracy

in [11], the our learned visual model performs significantly and consistently better (83.3% vs 68.6% on average over all 14 events). This result shows the importance of modeling visual content.

In addition, compared with our approach, the models trained using random negative samples expose degraded accuracy (from 83.3% to 70.1%). Moreover, the performance of models trained with the uniform negative dataset is better than when random negative event sample are used, but not as accurate as our approaches. Those results confirm our hypothesis, “location” information plays an important role in negative samples collection and our approach is effective in collecting such negative samples.

Overall, the experiments have clearly shown the value of using visual analysis to model social events content. Furthermore, we have demonstrated that the construction of the event model can be automated without compromising the resulting performance.

## 5. CONCLUSION AND FUTURE WORK

We proposed an novel framework leveraging on the huge number of media documents available on social media website to gather the training data collection necessary to learn social event models. The positive samples are collected using photos with identification tags explicitly referring to the event. The negative samples correspond to those photos taken at the same period and in the vicinity of the event but for which the tags are identified as being common (repeatedly appearing over time). We evaluate the trained visual models on a manually labeled dataset, study the effect of the methodology related parameters and finally report accuracy results of 83% on real world scenarios. As future work, we currently investigate approaches for collecting additional positive samples with extended coverage of the event while preserving accuracy.

## 6. REFERENCES

- [1] H. Becker, M. Naaman, and L. Gravano. Event Identification in Social Media. In *International Workshop on the Web and Databases*, Providence, USA, 2009.
- [2] H. Becker, M. Naaman, and L. Gravano. Learning Similarity Metrics for Event Identification in Social Media. In *ACM Conference on WSDM*, pages 291–300, New York, USA, 2010.
- [3] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y.-T. Zheng. NUS-WIDE: A Real-World Web Image Database from National University of Singapore. In *ACM Conf. on CIVR*, Santorini, Greece.
- [4] R. Datta, D. Joshi, J. Li, James, and Z. Wang. Image retrieval: Ideas, influences, and trends of the new age. *ACM Computing Surveys (CSUR)*, 40, 2008.
- [5] D. Delgado, J. Magalhaes, and N. Correia. Automated Illustration of News Stories. In *IEEE Conference on Semantic Computing*, pages 73–78. IEEE, Sept. 2010.
- [6] C. S. Firan, M. Georgescu, and et. al. Bringing order to your photos: event-driven classification of flickr images based on social. In *ACM conference on CIKM*, New York, USA, 2010.
- [7] R. Hong, G. Li, L. Nie, J. Tang, and T.-S. Chua. Explore Large Scale Data for Multimedia QA. In *ACM conference on CIVR*, Xi’an, China, 2010.
- [8] L.-J. Li and G. Wang. OPTIMOL: automatic Online Picture collecTion via Incremental MOdel Learning. *IEEE Conference on CVPR*, 88(2):1–8, 2007.
- [9] X. Li, C. G. Snoek, M. Worring, and A. W. Smeulders. Social negative bootstrapping for visual categorization. In *ACM Conference on ICMR*, 2011.
- [10] T.-Y. Liu. *Learning to Rank for Information Retrieval*. springer, 2011.
- [11] X. Liu, R. Troncy, and B. Huet. Finding Media Illustrating Events. In *ACM Conference on ICMR*, 2011.
- [12] C. D. Manning, P. Raghavan, and H. Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 1 edition, July 2008.
- [13] P. Over, G. Awad, M. Michel, J. Fiscus, W. Kraaij, and A. F. Smeaton. Trecvid 2011 – an overview of the goals, tasks, data, evaluation mechanisms and metrics. In *Proceedings of TRECVID*. NIST, USA, 2011.
- [14] T. Quack, B. Leibe, and L. Van Gool. World-scale mining of objects and events from community photo collections. In *ACM conference on CIVR*, page 47, New York, USA, July 2008.
- [15] F. Schroff, A. Criminisi, and A. Zisserman. Harvesting Image Databases from the Web. In *IEEE Conference on ICCV*, pages 1–8. IEEE, 2007.
- [16] U. Westermann and R. Jain. Toward a Common Event Model for Multimedia Applications. *IEEE MultiMedia*, 14(1):19–29, 2007.
- [17] Z.-J. Zha, T. Mei, J. Wang, Z. Wang, and X.-S. Hua. GRAPH-BASED SEMI-SUPERVISED LEARNING WITH MULTI-LABEL. *ACM Trans. Program. Lang. Syst.*, 20(5):97–103, 2009.