

# Joint Event Recognition and Image Selection

Anonymous ECCV submission

Paper ID \*\*\*

**Abstract.** Automatic image organization of personal photos is a problem with many real world applications, and can be divided into two main tasks: recognizing event types of the photo collections, and selecting interesting images from the collections. The two tasks are both challenging, in that the photos from the same event type are highly varied, and there is a great deal of overlap in photos from different event types. In this paper, we looked into the possibility of simultaneously solving both tasks: album-wise event recognition and ~~single image event-specific~~ importance score prediction. We collected an event album dataset with both event type labels and image importance labels, refined from the existing CUFED dataset. We propose a hybrid system consisting of three parts: Siamese network based event-specific image importance prediction, Convolutional Neural Network(CNN) based event recognition, and Long Short-Term Memory(LSTM) based sequence level event recognition, and we propose an iterative updating procedure for event type and image importance score prediction. We show with experiments that the proposed method outperforms the classical approach based on static image classification, and more importantly, we verified that image importance score prediction and event type recognition can in turn help the performance of each other.

**Keywords:** Event Recognition, Image Importance, Convolutional Neural Network, Long Short-Term Memory

## 1 Introduction

With the advent of cheap cameras in nearly all of our devices, and automated uploading to the cloud, and practically unlimited storage, it becomes painless to take photos frequently in daily life, resulting in the explosion of personal photo collections. However, the oversized image collections make it difficult for us to organize the photos, and thus automatic organization algorithms are highly desirable. The organization of personal photo collections can be decomposed into two stages: recognition of the event types of a photo collection, and suggesting the most interesting/important images in the photo collection to represent the album, e.g. making an album cover, or suggesting to the user a set for a further use, e.g. making a photo book. The two stages assist users in keeping the photo collections organized and free of irrelevant images.

Both image importance prediction and event recognition have been studied independently.

Studies of event recognition can be separated into several types. The most popular branches of study use videos as input [1–3]. Spatial and temporal features of videos are usually used for this task. The success of Long Short-Term Memory (LSTM) networks for sequence tasks have also extended to video based event recognition. In [4], the LSTM network and visual feature extraction network is stacked, and the network can deal with both event recognition and description. Reiter *et al.* [5] also combine LSTM and HMMs for video meeting analysis. At the other end of the spectrum, event recognition for single static images has also been attempted [6–8]. In contrast to video based event recognition, there is no temporal information to exploit, and there is no need to consider relevant frame importance or contribution to the event, since there is only one “frame.” This problem can be viewed as a special case of scene recognition, and both object level features and scene level features are utilized [6]. The recent success of Deep Neural Networks, especially Convolutional Neural Networks(CNN) provides us with an outstanding visual representation, that can be used for the single static image event recognition task [7, 8].

Album-wise event recognition lies in the middle between single-image-based and video-based event recognition. Images in an album can be thought of a very sparse samples from the event video. Photo albums differ from videos in that consecutive images from the photo album are no longer continuous, and can have very different visual and semantic information. However, there is still sequential information in time-stamped albums, and the images in an album are of varied importance. In [9], an HMM-based model is proposed to use the sequential information in albums for the recognition task. It was shown that the temporally-sensitive HMM outperforms the simple aggregation of predictions from all the images in an album. This indicates that the sequential information in an album is helpful for album-wise event recognition.

Image importance is a complex image property which is related to various factors, such as aesthetics [10], interestingness [11] and image memorability [12]. It has been shown [13] that importance of an image is modulated by the context of the image, making image importance album-dependent, or event-specific. Event-specific image importance is a highly subjective judgment, and learning it is a very challenging task, due to the very high intra-class variability, and the underlying uncertainty caused by the subjectiveness of the property. Nevertheless, it is still possible to predict it.

Returning to the task of photo collection organization, we ask the question: can we simultaneously recognize the event type for an album, and discover important images in the event album? To answer this question, this paper makes the following contributions: 1. We refine the existing event curation dataset CUFED by collecting more human annotations for album event types for more reliable ground-truth, and allow for multilabel annotation for an album. The multilabel and ground-truth ambiguity between event types provides us with more training information, and allows for a fairer evaluation at the testing stage; 2. We propose a joint event recognition/image importance prediction algorithm. We use a CNN for image level event recognition, and a Siamese Network for event-specific im-

age importance prediction. An iterative update scheme is conducted during the test stage, and we find that event recognition and image importance prediction can improve each other’s performance; 3. We further boost the performance of event recognition with an LSTM network that leverages sequence information in labeling the album.

## 2 Related Works

Our work is partly inspired by [13], who proposed a novel image property: event-specific image importance. In this work, it is claimed that image importance or interestingness is contextual and is related to the album it is in. For example, a photo of a beautiful work of architecture is important in an album of an urban trip, yet not so important in a wedding event. A Siamese network is used to predict the relative score difference between an input image pair, and is jointly trained on all event types. However, in this work, the event-specific image importance score is predicted given the ground-truth event type of the album. In our work, we extend this idea to training a system simultaneously for event recognition and ~~importance~~ image curation, so that additional user input of event type for testing is not required.

Our work is closely related to the study of event recognition for a personal album. The model in [14] classifies a personal album into 8 social events and 10 sports events simply by aggregating the SVM classification result from single images in the album. In [15], Tsai *et al.* exploit object level patterns for event type recognition. Object patterns are learnt from single images, and then an album-wise SVM is trained on the frequency distribution of different object patterns appearing in an album. Similarly, Imran *et al.* [16] use the Pagerank technique to mine the most useful features for an event, and an album-wise SVM classifier is used for recognition. The above works treat albums as an unordered collection of images. On the other hand, in [9], Bossard *et al.* exploit the sequential nature of personal albums and use an HMM based sub-event approach for event recognition. They use temporal sequence of the images, and model an album with successive latent sub-events to boost the recognition performance. They collected a 14 class dataset consisting of 807 albums for the task.

Event recognition for single photos has also been studied. Li *et al.* [6] use a generative graphical model to recognize event types of a database with 8 sports events. Their model integrates cues from scene and object categorization to classify the sports events. Salvador *et al.* [8] focus on cultural event recognition. They integrate cues from visual features extracted by a CNN and from the time-stamp of a photo, inspired by the fact that photos of a cultural event are mostly taken in the same period of time. However, in personal photo collections, the relevance of an image within an event album varies a great deal. These approaches for single images are useful, but not sufficient for album-wise event recognition.

Convolutional Neural Network(CNN) methods have greatly boosted performance in image understanding tasks, such as image classification, object detec-

tion and scene recognition [17–20]. Many studies have switched their focus to higher-level image properties, such as event recognition [21], semantic segmentation [22], multilabel image annotation [23], and image captioning [24]. Long Short Term Memory(LSTM) networks [4] have been proposed for sequence prediction and sequence labeling. LSTM networks have advantages over traditional Recurrent Neural Networks (RNN) in that they can maintain contextual information across an extended sequence of data. LSTMs have achieved success for tasks such as handwritten text recognition [25] and speech recognition [26]. Relevant to our work, the Long-term Recurrent Convolutional Network (LRCN) model [24] has been proposed to stack CNN feature extractor and LSTM networks for sequential learning of videos or images.

### 3 The ML-CUFED Dataset

In order to train and evaluate the joint curation-recognition model, we use the Curation of Flickr Events Dataset (CUFED), and refine it by collecting more human opinions on the event types in the dataset. We ~~can~~ can the new dataset MultiLabel-CUFED (ML-CUFED). In this section, we describe the dataset, and provide a consistency analysis of the labels collected from Amazon Mechanical Turk (AMT). The dataset will be made available to the public.

#### 3.1 The CUFED Dataset

The CUFED Dataset is an image curation dataset extracted from the Yahoo Flickr Creative Commons 100M Dataset (YFCC100M). 20,000 albums were segmented by user tags and timestamp, and then their event types were collected from workers on AMT. Each album received 3 workers’ labels. The dataset contains 23 most common event types from our daily life, ranging from nature trips to weddings. For each event type, 50-200 albums are further randomly sampled from the 20,000 albums, forming the CUFED dataset of 1883 albums. The importance score for each image in the albums was then decided by workers from AMT. The final ground-truth event-specific importance scores of the images were obtained from the average of 5 workers’ votes.

One problem the CUFED Dataset has is that the event type of an album is decided by only three workers, and ~~the~~ the workers were constrained to give a single label to each album. However, for an album with ambiguous event types or with multiple event types, it is overly restrictive to give the album a single event label. For example, the two albums in Fig 1 are both birthday events, but they can also fall into the category of casual family/friends gathering. Those two event types are not mutually exclusive. Moreover, intuitively, we would say the album on the right is a more typical birthday event, with most images focusing on the little boy celebrating his birthday, while the album on the left is more of a casual family/friends gathering rather than an obvious birthday event. Therefore, collecting the event types and their proportion in one album from multiple people’s view is necessary. This results in a multi-label event recognition dataset with richer information.



Fig. 1: Example of two birthday albums (both have the photo uploader’s tag “birthday”).

### 3.2 Data collection

On top of the three votes the dataset already has, we collected 9 more workers’ opinions, and allowed for multiple choices for each album. One worker could select up to three event types for an album. There were totally 299 distinct workers who participated in the task.

The quality of different AMT workers’ submissions vary. Therefore, we need to do quality control in order to collect high quality annotations. Before the real task, only workers who passed an album event recognition test (which is very similar to the actual task) were allowed proceed to work on the actual task. During the tasks, there was another round of quality control. After workers submitted the tasks, the results they turned in were compared with other workers’ submissions, and submissions that highly diverged from others were further manually inspected. If the divergence was unreasonable, the submission was rejected. After all the annotations from workers were collected, we further cleaned the annotations by eliminating the labels with minor votes: for each album, all the event types with only one vote were discarded.

For the final ground-truth event types and their proportion of one album, we use the the proportion of votes among all the votes. For workers who gives votes to more than one label, each of the votes he or she gives is normalized so that all the votes sum to one.

### 3.3 Dataset Analysis

To check the validity of the dataset we collected, we analyzed the annotations in several ways. Each album can between 9 and 27 votes (because we allow for multiple choices from one worker). 76% of the albums received votes for two or fewer event types. This suggests the high coherence of those albums. 95% of the albums received votes for three or fewer event types. To check the consistency among workers, we randomly split the 299 workers into two halves, and for each album we checked whether the annotations from one half was consistent with the other half. For each album, we examined whether the top event types suggested by these two independent groups were the same. We repeated the random split 100 times, and on average, for 89.6% of the albums, the event type receiving most votes were the same for both groups. This suggests that despite the ambiguity of some album types, we got consistent opinions from different AMT workers.

## 4 Approach

In this section, we describe our approach to jointly attain image importance prediction and album event-type recognition. The system is shown in Figure 2.

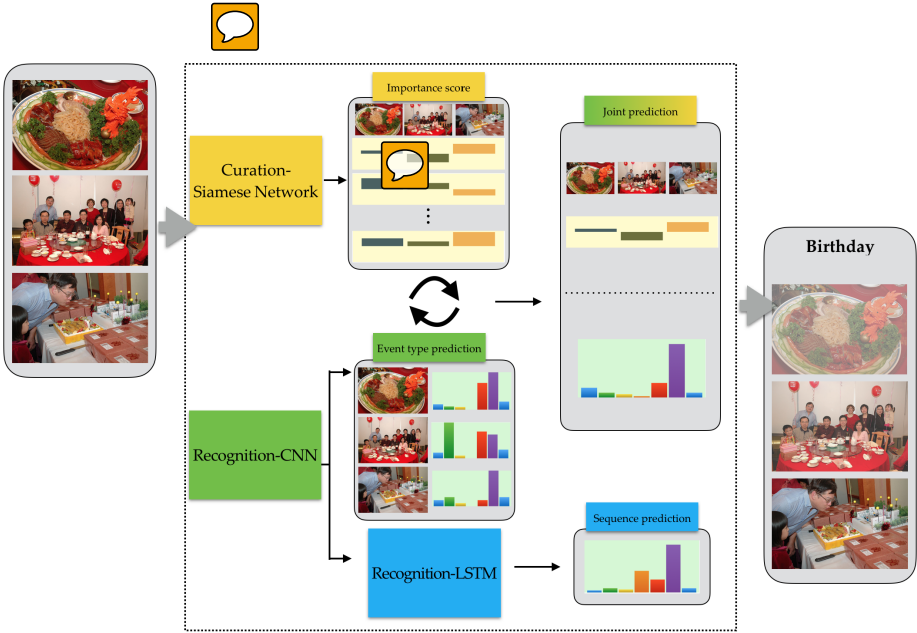


Fig. 2: The joint album recognition-curation system. The system consists of three parts: a Siamese network for image importance score prediction, a CNN for single image event-type recognition, and a LSTM network for album event-type recognition. During the test stage, the three components interact and jointly produce the prediction of album event type and image importance score.

### 4.1 Event curation network


For event curation purposes, we followed the approach in [13], using Piecewise Ranking (PR) loss to jointly train a Siamese network to predict the importance score difference between an image pair given the ground-truth event type of the input image pair. The  architecture can predict the relative event-specific importance score for a set of images, and is found to perform better than a traditional CNN that directly predicts the absolute important score of an image. One difference between our implementation and [13] is that the ground-truth event type for the input image pair is not a one-hot representation; instead, the ground-truth event type label is a soft distribution. The event type label is used to gate the output and gradient of the Siamese network: for an input image pair of a certain event type  $c \in C$ , only the part of the network which corresponds to event type  $c$  is back-propagated. Here, to train the Siamese network on the multi-event type label, we change the 0/1 gating to a soft gating, thus for an input



image pair with multiple event type labels, the error signals from all the possible labels are back-propagated, but with a weight: the ground-truth probability of that event type. Another scheme is to still use 0/1 gating, but for an input image pair with multiple event type labels, there is equal probability for the network to view the input image pair as from one of the possible labels. This gives more relaxation on the ground-truth event type label.

## 4.2 Event recognition network

One of the properties of an “event album” that makes it different from a simple collection of images is that it is a sequence, and this provides us with the temporal relationship between the images. LSTMs have been successfully applied to sequential tasks, and its ability to do long-range context memorization is suitable for our task of album-wise event recognition. Therefore, we use the LSTM network to capture the sequential information, in addition to a classical CNN that captures the visual features of a single image.

We start with a CNN pre-trained on ImageNet [27] [17], and fine-tune it on the CUFED Dataset to recognize a single image’s event type. We then extract the high level CNN features for each image from the adapted network, and use them as the input features to train the LSTM network for album-wise event recognition. The LSTM network consists of a single LSTM layer, a mean pooling layer, and a softmax prediction layer.

The target for both fine-tuned the CNN and the LSTM network is the soft distribution of the ground-truth event labels. Again, another scheme is to use the one-hot target, but treat each training example as from one of the possible event types with equal probabilities. FROM GARY. WE WILL PROVIDE BOTH APPROACH IN THE RESULTS PART. AND ACTUALLY THE SECOND APPROACH HAS A LITTLE BETTER RESULT. Another approach would be to train each album multiple times, with different event labels each time. We deemed this approach to be less elegant than using the distribution as the target.

During the testing stage, the album-wise prediction from the LSTM network and the image-wise prediction from the CNN are combined to produce the final prediction for the album type.

## 4.3 The iterative curation-recognition procedure

For an “event album”, more interesting images or important images give us more information about the event type of this album. For example, although a candle blowing image may only appear in an album once, it is very helpful for deciding the event type of the album. However, as shown in [13], the importance of an image is related to the context it is in, and is event-type dependent. Therefore, we propose that the image importance score can help with event recognition of an album, while event recognition of an album will in turn improves the image importance score prediction.

We denote an  $N$ -image album as  $\mathbf{S} = \{I^1, ..., I^N\}$ . From the event curation network in Section 4.1, we obtain the importance score of an image  $I^n \in \mathbf{S}$  given

its event type:  $W^n = [w_1^n, \dots, w_C^n]^T$  where  $C$  is the number of event types. From the event recognition CNN, we can also get the prediction of the event type for a single image  $P^n = [p_1^n, \dots, p_C^n]^T$ . We then conduct the iterative curation-recognition procedure:

$$\begin{cases} Q(k+1) = [P^1, P^2, \dots, P^N] \cdot V(k)^\alpha \\ V(k+1) = \left\{ [W^1, W^2, \dots, W^N]^T \circ \mathbf{I} \{p_c^n \geq m \cdot \max_{c'}(p_{c'}^n)\}_{n,c} \right\} \cdot Q(k+1) \end{cases} \quad (1)$$

Or equivalently, it can be written in vector form:

$$\begin{cases} q_c(k+1) = V^T(k)^\alpha \cdot P^n \\ v^n(k+1) = \left\{ [W_n^T \circ \mathbf{I} \{p_c^n \geq m \cdot \max_{c'}(p_{c'}^n)\}_{(c,1)}] \right\} \cdot Q(k+1) \end{cases} \quad (2)$$

where the  $N$ -dimensional column vector  $V(k) = [v^1, \dots, v^N]^T$  is the  $k$ -th step prediction for all images' importance score in album  $S$ , and the  $C$ -dimensional column vector  $Q(k) = [q_1, q_2, \dots, q_C]^T$  is the  $k$ -th step prediction for the album's event type.

$\mathbf{I} \{p_c^n \geq m \cdot \max_{c'}(p_{c'}^n)\}_{n,c}$  denotes the binary mask that eliminates the event type predictions with low confidence, and makes sure only event types with high probability contribute to the image importance prediction.  $\alpha$  is the weight factor that reflects our emphasis on the image importance score.

By iteratively conducting the procedures in Equation 1, we obtain the album-wise event prediction  $Q$  and image importance score prediction  $V$ .

Note that Equation 1 is not guaranteed to converge. In the case of oscillation between states, we set a maximum number of iterations, and when the iteration number hits the threshold, predictions for  $Q$  and  $V$  are obtained by averaging over the three previous steps.

#### 4.4 Joint prediction

The iterative curation-recognition procedure takes image importance into account for event type prediction, while the LSTM network is trained to do the same task, taking the sequence of images into account. These two processes stress two distinct properties of event albums, and are complementary to each other. Therefore, we average the predicted probability density from the LSTM network  $Q_{LSTM}$  and from the iterative curation-recognition procedure  $Q_{iter}$ , and get the final event type prediction for an album. The entire system is illustrated in Figure 2.



## 5 Experiments

### 5.1 Baselines

### 5.2 Details

### 5.3 Results

### Importance

### Recognition

## 6 Conclusion

## References

1. Over, P., Awad, G., Michel, M., Fiscus, J., Kraaij, W., Smeaton, A.F., Quénot, G., Ordelman, R.: Trecvid 2015 – an overview of the goals, tasks, data, evaluation mechanisms and metrics. In: Proceedings of TRECVID 2015, NIST, USA (2015)
2. Tang, K., Fei-Fei, L., Koller, D.: Learning latent temporal structure for complex event detection. In: CVPR. (2012)
3. Xu, Z., Yang, Y., Hauptmann, A.G.: A discriminative cnn video representation for event detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2015) 1798–1807
4. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Comput. (1997)
5. Reiter, S., Schuller, B.W., Rigoll, G.: A combined LSTM-RNN - HMM - approach for meeting event segmentation and recognition. In: 2006 IEEE International Conference on Acoustics Speech and Signal Processing, ICASSP 2006, Toulouse, France, May 14-19, 2006. (2006) 393–396
6. Li, L.J., Fei-Fei, L.: What, where and who? classifying events by scene and object recognition. In: Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on. (2007)
7. Park, S., Kwak, N.: Cultural event recognition by subregion classification with convolutional neural network. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops. (2015)
8. Salvador, A., Zeppelzauer, M., Manchon-Vizuet, D., Calafell-Orós, A., Giró-i Nieto, X.: Cultural event recognition with visual convnets and temporal models. In: CVPR ChaLearn Looking at People Workshop 2015. (06/2015 2015)
9. Bossard, L., Guillaumin, M., Van, L.: Event recognition in photo collections with a stopwatch hmm. In: Computer Vision (ICCV), 2013 IEEE International Conference on. (2013)
10. Lu, X., Lin, Z., Jin, H., Yang, J., Wang, J.Z.: Rapid: Rating pictorial aesthetics using deep learning. In: Proceedings of the ACM International Conference on Multimedia. (2014)
11. Dhar, S., Ordonez, V., Berg, T.L.: High level describable attributes for predicting aesthetics and interestingness. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2011)
12. Isola, P., Xiao, J., Torralba, A., Oliva, A.: What makes an image memorable? In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2011)
13. Wang, Y., Lin, Z., Shen, X., Mech, R., Miller, G., Cottrell, G, W.: Event-specific image importance. In: In Proc. Computer Vision and Pattern Recognition Conference (CVPR). (2016)
14. R.MattiviandJ.R.R. HüllingsandF.deNataleandN.Sebe: Exploitation of time constraints for (sub-) event recognition. In: ACM Workshop on Modeling and Representing Events (J-MRE 11). (2011)
15. Tsai, S., Cao, L., Tang, F., Huang, T.S.: Compositional object pattern: a new model for album event recognition. In: Proceedings of the 19th International Conference on Multimedia 2011. (2011)
16. Imran, N., Liu, J., Luo, J., Shah, M.: Event recognition from photo collections via pagerank. In: Proceedings of the 17th ACM International Conference on Multimedia. MM '09, New York, NY, USA, ACM (2009) 621–624
17. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25. Curran Associates, Inc. (2012)

18. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2015)
19. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2014)
20. Zhou, B., Lapedriza, A., Xiao, J., Torralba, A., Oliva, A.: Learning deep features for scene recognition using places database. In: Advances in Neural Information Processing Systems 27 (NIPS). (2014)
21. Xiong, Y., Zhu, K., Lin, D., Tang, X.: Recognize complex events from static images by fusing deep channels. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2015)
22. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. CVPR (to appear) (November 2015)
23. Gong, Y., Jia, Y., Leung, T., Toshev, A., Ioffe, S.: Deep convolutional ranking for multilabel image annotation. CoRR **abs/1312.4894** (2013)
24. Donahue, J., Hendricks, L.A., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., Darrell, T.: Long-term recurrent convolutional networks for visual recognition and description. In: CVPR. (2015)
25. Graves, A., Liwicki, M., Fernández, S., Bertolami, R., Bunke, H., Schmidhuber, J.: A novel connectionist system for unconstrained handwriting recognition. IEEE Trans. Pattern Anal. Mach. Intell. (2009)
26. Sak, H., Senior, A.W., Beaufays, F.: Long short-term memory recurrent neural network architectures for large scale acoustic modeling. In: INTERSPEECH 2014, 15th Annual Conference of the International Speech Communication Association, Singapore, September 14-18, 2014
27. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding. arXiv preprint arXiv:1408.5093 (2014)