

# Learning to Hash-tag Videos with Tag2Vec

Aditya Singh

Saurabh Saini

Rajvi Shah

P J Narayanan

CVIT, KCIS, IIIT Hyderabad, India<sup>\*</sup>

{(aditya.singh,saurabh.saini,rajvi.shah}@research.,pjn@}iiit.ac.in

<http://cvit.iiit.ac.in/research/projects/tag2vec>



Figure 1. Learning a direct mapping from videos to hash-tags : sample frames from short video clips with user-given hash-tags (left); a sample frame from a query video and hash-tags suggested by our system for this query (right).

## ABSTRACT

User-given tags or labels are valuable resources for semantic understanding of visual media such as images and videos. Recently, a new type of labeling mechanism known as hash-tags have become increasingly popular on social media sites. In this paper, we study the problem of generating relevant and useful hash-tags for short video clips. Traditional data-driven approaches for tag enrichment and recommendation use direct visual similarity for label transfer and propagation. We attempt to learn a direct low-cost mapping from video to hash-tags using a two step training process. We first employ a natural language processing (NLP) technique, skip-gram models with neural network training to learn a low-dimensional vector representation of hash-tags (*Tag2Vec*) using a corpus of  $\sim 10$  million hash-tags. We

then train an embedding function to map video features to the low-dimensional Tag2vec space. We learn this embedding for 29 categories of short video clips with hash-tags. A query video without any tag-information can then be directly mapped to the vector space of tags using the learned embedding and relevant tags can be found by performing a simple nearest-neighbor retrieval in the Tag2Vec space. We validate the relevance of the tags suggested by our system qualitatively and quantitatively with a user study.

## CCS Concepts

•Computing methodologies → Computer vision; Visual content-based indexing and retrieval;

## Keywords

Tag2Vec; Video Tagging; Hash-tag recommendation

## 1. INTRODUCTION

Over the last decade, social media websites such as Twitter, Instagram, Vine, YouTube have become increasingly popular. These media sites allow users to upload, tag, and share their content with a wide audience across the world. In case of visual media such as images and videos, the user-given tags often provide rich semantic information about the visual context as well as affective appeal of the media, otherwise hard to recognize and categorize. Most popular image

<sup>\*</sup>Center for Visual Information Technology (CVIT), Kohli Center on Intelligent Systems (KCIS), International Institute of Information Technology (IIIT) Hyderabad

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

ICVGIP, December 18-22, 2016, Guwahati, India

© 2016 ACM. ISBN 978-1-4503-4753-2/16/12...\$15.00

DOI: <http://dx.doi.org/10.1145/3009977.3010035>

and video search engines still heavily rely upon user-given tags for relevant retrieval. Apart from search and retrieval, hash-tags also facilitate browsing and content management. Varied web content can be organized easily using hash-tags which defines new trending concepts. From user perspective this makes it easier to share content and follow trends.

With increasing success of object category recognition algorithms on large data such as ImageNet, automatic image tagging and captioning is showing remarkable progress. However, for videos with dynamic events, fewer attempts have been made at unsupervised video tagging.

The methods for video annotation and tagging can be classified into two categories, (i) model based methods, and (ii) data-driven methods. Model based methods mainly apply several concept classifiers, pre-trained with low-level video features, and use the resulting concept labels for effective tagging. The shortcoming of this approach is revealed by the fact that it is impossible to enumerate all possible concept categories and their inter-relationships as perceptually understood. Data-driven approaches on the other hand do not explicitly discover or recognize visual concepts. They, instead directly propagate or transfer them from tagged videos to query videos using some measure of visual similarity.

In this paper, we present a hybrid approach for tag suggestion. We do not explicitly use pre-trained concept classifiers, nor do we use video-to-video similarity based measures. We propose a method to directly embed video features into a low-dimensional vector space of tag distribution. Given a query video, the relevant tags can then be retrieved by a simple nearest neighbour strategy. The *video-to-tag* training is carried out in two stages. First, we learn a 100-dimensional vector space representation of popular tag words using a corpus of  $\sim 10$  million hash-tags using the algorithm of Mikolov et al. [11]. This algorithm trains a two-layer neural network with skip-gram representation to learn word embeddings in vector space. Extending the terminology of [11], we call this vector space *Tag2Vec* space in this paper. Second, we learn a nonlinear embedding of high-dimensional video features to the low-dimensional Tag2Vec space using a separate neural network Socher et al. [18]. For this task, we use  $\sim 2740$  short video clips from 29 categories and their associated hash-tags. Once trained, the final video-to-tag embedding can be leveraged to suggest tag words for query videos. Our approach is pictorially summarized in Figure 3. We evaluate the performance of our system qualitatively and quantitatively with a user study and show the method is promising. We also discuss limitations and future directions to improve effectiveness of this simple approach.

The main contributions of this paper are as follows, (i) We study the problem of hash-tag suggestion for videos from a novel perspective and present a mechanism for direct embedding of videos to a vector space of hash-tags; (ii) We present a new dataset consisting of 3000 random wild short social video clips spanning 29 categories with associated user-given hash-tags.

## 2. RELATED WORK

In this section we present a brief discussion of the relevant literature related to our problem. First, we discuss the methods related to the two core sub-systems of our method, tagging and word embedding, separately. Later, we discuss the recent works on visual semantic joint understanding that leverage similar methodology.

### 2.1 Tagging

Image and video tagging approaches can be coarsely categorized as model-based or data-driven. Qi et al. [15] combine the strengths of the two separate paradigms using separate binary classifiers learning and concept fusion. They focus on multi-label results and use gibbs random field based mathematical model. Lavrenko et al. [7] construct a joint probability of visual region-based words with text annotations, incorporating co-occurrent visual features, and co-occurrent annotations to demonstrate that statistical methods can be used to retrieve videos by content. However, the main drawback of model-based approaches is the limit on detectors that can be trained. As there are thousands of concepts for which it is difficult to gather large training data for reliable learning. Due to this, there has been a shift from model-based methods to data-driven similarity based methods.

Data-driven methods utilize the abundance of videos shared by users and transfer tags based on similarity measures. Ballan et al. [2] proposed a video retagging approach based on visual and semantic consistency. This approach however only acknowledges tags which are nouns in the WordNet lexicon. Often the hash-tags are informal internet slang-words which rarely occur in proper language documents and hence do not have semantic consistency. Tang et al. [20] use a graph based semi-supervised learning approach for manifold ranking. They use partial differential based anisotropic diffusion for label propagation. Zhao et al. [28] focuses on fast near duplicate video retrieval for automatic video annotation. They find near duplicates by indexing local features, fast pruning of false matches at frame levels, and localization of near duplicate segments at video levels. Then a weighted majority approach is used for tag recommendation. Moxley et al. [12] devices a graph reinforcement framework to propagate tags developed by crawling tags of similar videos for annotation by using text and visual features. Wang et al. [25] computes similarity between two samples along with the difference in their surrounding neighbourhood sample & label distribution. The neighbourhood sample similarity is computed using KL divergence and label similarity is based on difference of label histograms of the two samples. Yao et al. [26] utilize the user click-through data along with the similarity based measures to tackle the problem of video tagging. Our approach is also based on data-driven similarity but instead of directly measuring video similarity, we learn a direct mapping to embed the videos in a lower-dimensional Tag2Vec space then use a nearest-neighbour classifier for tag suggestion.

### 2.2 Word Embedding

Common approaches for word embedding is through neural networks [17], dimensionality reduction of co-occurrence matrix. [8, 9, 10], explicitly constructed probabilistic models [5] etc. In our method we use neural network based embedding as will be focusing only on such approaches here. Amongst the early approaches, Bengio et al. [4] reduces the high dimensionality of words representations in contexts by learning a distributed representation for words. They use a feed forward neural network with a linear projection layer and a non-linear hidden layer which jointly learns a word vector representation and a statistical language model. Currently, widely popular technique of Mikolov et al. [11] proposes method for learning word vectors from a large amount

of unstructured data. They also show that the learned space is a metric space and meaningful algebraic operations can be performed on the word vectors. Barkan [3] proposes a bayesian skip-gram method which maps words to densities in a latent space rather than word vectors which results in less effort in hyperparameter tuning.

### 2.3 Visual-semantic Joint Embedding

Vector representation of words is used by many recent approaches for joint visual semantic learning [17, 16, 23, 27, 6]. Socher et al. [17] learns an embedding function which performs a mapping of image features to semantic word space. Then they utilize this for categorization of seen and unseen classes. Zhang et al. [27] utilize linear mappings and non-linear neural networks to tag an image. They define the problem of assigning tags as identification of a principal direction for an image in word space. This principal direction ranks relevant tags ahead of irrelevant tags. Similar to all these approaches we also used a neural network for learning embedding function but we focus only on hash tags. [16] maps the video representations to semantic space for improving action classification. Recent image captioning methods [23, 6] have shown remarkable progress in generating rich descriptive sentences for natural images. These methods use recurrent neural network architectures with large data for training. Visual question answering (VQA) systems [1] also employ deep learning to train an answering system for natural language queries. Though deep features have shown promise in image based captioning and VQA systems, one either needs a huge amount of data to train deep networks or needs to fine-tune a pre-trained network. For videos such pre-trained networks are not readily available yet and though we have collected a dataset of nearly 3000 short videos, it is not sufficient to train a deep network. Hence, we use the state-of-the-art hand-crafted features for our application. On a related note, recently Tapaswi et al. [21] released an interesting video based question answering dataset by aligning book descriptions to movie scenes. Our work however has a different focus in its application to hash-tags and wild social video clips.

The rest of the paper is organized as follows: [section 3](#) explains the methodology and generation of tag space. [section 4](#) provides the experimental details, results and analysis. Finally in [section 5](#) we conclude our work and present future research directions.

## 3. METHOD

The proposed system is trained using a large number of videos and associated hash-tags scraped from social media platform [vine.co](#). Videos shared on this platform (commonly known as vines) are six seconds long, often captured by hand-held or wearable devices, with cuts and edits, and present a significantly wilder and more challenging distribution than traditional videos. For each uploaded video, the original poster also provides hash-tags. Unlike tag words typically used as meta-data, hash-tags serve more of a social purpose to improve content visibility and to associate content with social trends. Many hash-tags do not adhere to the commonly understood semantics of the natural language. Due to this reason, we learn a new tag space representation *Tag2Vec* instead of directly using semantically structured Word2Vec space of [11]. [Figure 1](#) shows examples of a few vine videos and associated hash-tags. It can

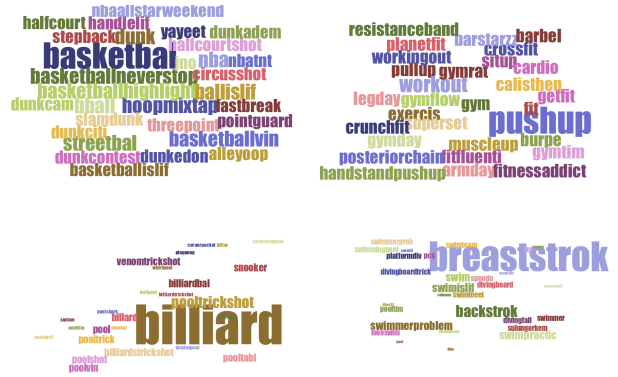


Figure 2: Illustration of effectiveness of Tag2Vec representation. The word clouds represent 29 nearest neighbours in Tag2Vec space for the following queries : ‘basketball’, ‘billiards’, ‘pushup’, ‘baseball’, ‘kayaking’, and ‘breast stroke’. The word sizes are proportional to similarity (inversely proportional of  $L_2$  distances), query word being the largest due to 100% self-similarity. Note that tag words are stemmed, not spelled incorrectly.

be observed that #FitFluential and #Mr315 are non-word hash-tags but understandable social media jargons given the content.

As mentioned previously, our end-to-end training for *video-to-tag* mapping consists of two stages, of learning the Tag2Vec representation, and of learning the video features to tag vector space embedding. Finally, given a query video, the learned embedding projects it to the tag space and nearby hash-tags are retrieved as suggestion. This process is outlined in [Figure 3](#). In the following subsections, we explain (i) the hash-tag data and learning of Tag2Vec space, (ii) the video data and learning of visual to Tag2Vec space embedding, and finally (iii) retrieval for tag recommendation.

### 3.1 Hash-tag Data and Pre-processing

To gather hash-tags data, we first use the 17,000 most common English words (as determined by n-gram frequency analysis of the Google’s Trillion Word Corpus <sup>1</sup>) as queries and retrieve a total of about 2.7 million videos ( $\sim 150$  videos per query). We scrape the hash-tags corresponding to each retrieved video, remove all special characters, and perform *stemming* on the hash-tags.

Stemming is a popular technique in natural language processing (NLP) community for reducing words to a root form such that multiple inflections of a word reduce to the same root e.g. ‘fish’, ‘fished’, ‘fishing’, and ‘fish-like’ reduce to the stem ‘fish’. The stemmed hash-tag words for each video form a *hash-tag sentence* leading to a total of 2.7 million sentences. These 2.7 million hash-tag sentences together form a text-corpus that we use to learn the Tag2Vec representation.

### 3.2 Learning Tag2Vec Representation

Mikolov et al. [11] proposed efficient unsupervised neural network based methods to learn embeddings of semantic words in vector space using a large corpus of text data (web-based) consisting of 1.6 billion words. These methods either use continuous Bag of Words representation or skip-gram

<sup>1</sup><https://github.com/first20hours/google-10000-english>

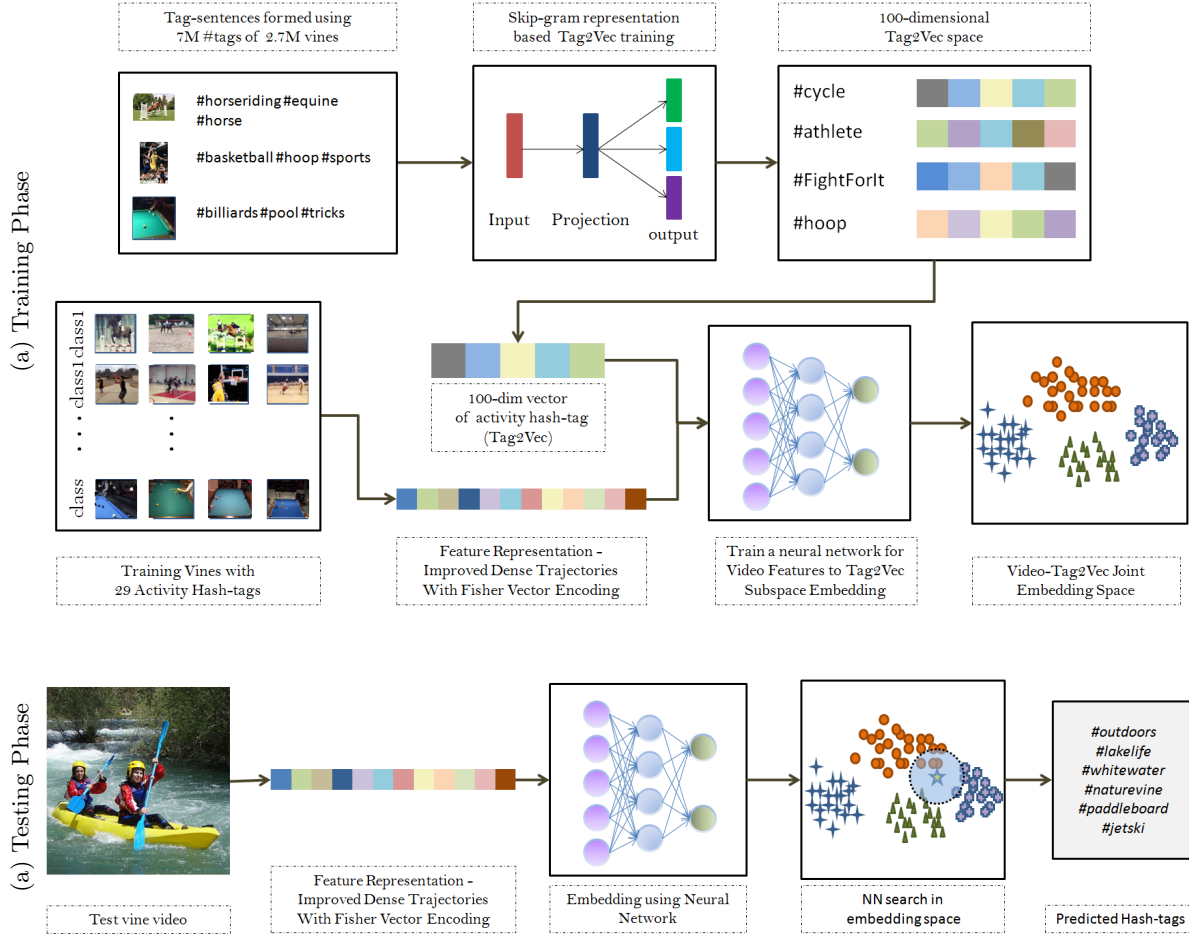


Figure 3: Pictorial representation of our tag recommendation system

representation of text sequences to train a two-layer neural network. The resulting word embedding assigns every word in the corpus to a vector in a 300-dimensional space. This word embedding is known as Word2Vec and the final vector space is popularly referred to as Word2Vec space. The main advantage of this representation is that vector operations can be performed on words. This property is extremely useful, e.g. to compute similarity between words using vector space distances.

Similar to this, we also train a neural network to learn hash-tag embeddings in a vector space and call it *Tag2Vec*. Since, we work with much smaller data ( $\sim 7$  million unique tags and 2.7 million sentences), we use skip-gram representation which is more effective on small data and we also restrict the resulting space to be 100-dimensional. We use publicly available code<sup>2</sup> for learning the Tag2Vec embeddings. The training converges quickly ( $\sim 10$  minutes) as we have a relatively small corpus of hash-tags. The resulting vector space enables us to perform vector operations on hash-tags. Figure 2 shows a word cloud representation of 30 nearest neighbours in Tag2Vec space for query vectors corresponding to stemmed tag words, ‘basketball’, ‘billiards’,

‘pushup’, ‘baseball’, ‘kayaking’, and ‘breast stroke’. The word sizes are inversely related to the distance from the respectively query words, the closer the words in vector space, the larger the font size. It can be observed that ‘basketball’ tag has many tags with high similarity whereas ‘billiards’ has fewer. It is also worth noting how contextual similarity is also well captured, for example ‘kayaking’ tag has strong neighbours such as ‘state park’ and ‘summer adventure’.

To demonstrate that the learned Tag2Vec space is different from Word2Vec space, we list the top-10 near-neighbours for four action word queries in both spaces (see Table 1). It can be clearly seen that the tags retrieved using Tag2Vec space are more diverse and socially relevant. See particularly the results for ‘Polevault’ and ‘Basketball’ queries. We also measure similarity between pairs of tag vectors and see that the Tag2Vec space models social jargons well. For example, we noticed that tags like #lol and #laugh have high similarity, #lol is also has high similarity with #fail owing to users tagging funny videos showing people failing at doing something, #fight is closer to both #win and #fail. This shows that our tag2vec space is able to capture meaningful tag relationships and similarity.

<sup>2</sup><https://code.google.com/p/word2vec/>



Query Words	Vector Space	top-10 retrieval results									
Pushups	Word2Vec	jumping jacks	pushup	situps	calisthenics	abdominal crunches	pushups situps	burpees	pullups	ab crunches	squat thrusts
	Tag2Vec	workout	gym	burpees	gymflow	gymday	exercise	muscleup	calisthenics	superset	gymrat
Polevault	Word2Vec	Ivan Ukhov	Gulfiya Khanafeyeva	Yaroslav Rybakov	Tatyana Lysenko	Anna Chicherova	Andrey Silnov	champion Tatyana	Croatia Blanka Vlasic	Svetlana Feofanova	Olga Kuzenkova
	Tag2Vec	USATF	athlete	tracknation	trackandfield	discusthrow	tooathlete	highjump	maxvelocity	blockstart	longjump
Kayaking	Word2Vec	canoeing	Kayaking	kayak	paddling	sea kayaking	rafting	whitewater kayaking	kayaking canoeing	rafting kayaking	canoing
	Tag2Vec	statepark	lake	whitewater	paddleboard	outdoorsfinland	lagoon	lakelife	outdooraction	emeraldbay	boat
Basketball	Word2Vec	baskeball	volleyball	basketbal	basketball	hoops	soccer	softball	football	bas ketball	roundball
	Tag2Vec	dunk	ballislife	hoopmixtape	basketballvine	NBA	basketballneverstops	streetball	basketballhighlight	bball	dunkcity

Table 1: Comparison of Tag2Vec and Word2Vec spaces. Top-10 nearest neighbour results shown for four query words. It can be seen that the tag words retrieved from Tag2Vec space are more diverse and socially relevant.

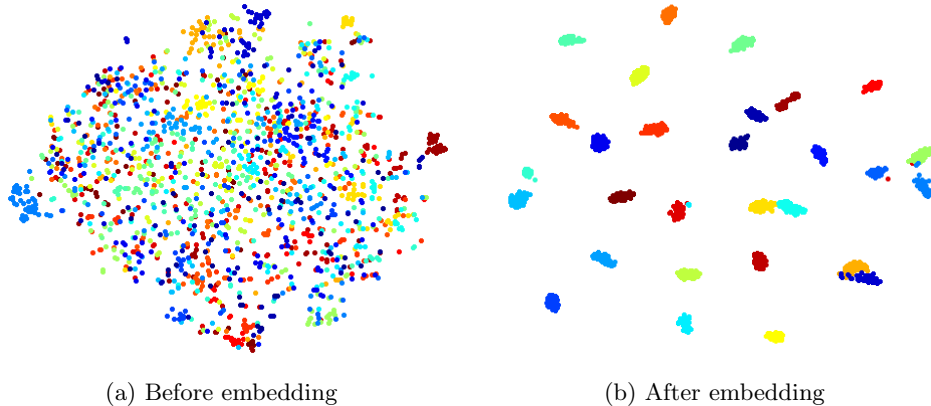


Figure 4: t-SNE visualization of the video features in fisher vector space and Tag2Vec space (after embedding)

### 3.3 Video Data and Feature Representation

Our video data consists of vine video clips (vines) spanning 29 categories. These categories are listed in Table 2. As mentioned before, vines are short but wild and complex amateur video clips with heavy camera motion, cuts, and edits. Hence, they have high intra-class variance and direct similarity based approaches don’t work well. We obtain 3000 useful videos with hash-tags after removing duplicates. We split these into training and testing sets of 2740 and 260 videos respectively. For both sets of videos, we compute visual and motion features.

In particular, we compute Improved Dense Trajectory (IDT) [24] features which consist of HoG (histogram of oriented gradients), HoF (histogram of optical flow), and MBH (motion boundary histograms) features. These low level features capture global scene, motion and rate of motion information respectively. We encode the low-level IDT features using 1003676 dimensional fisher vectors for better generalization [13, 14]. Fisher encoding relies on gaussian mixture model (GMM) computed over a large vocabulary of low-level features. We use UCF51 action recognition dataset [19] to compute generalized vocabulary for GMM estimation. For IDT extraction, we use the code<sup>3</sup> made available by the authors. For computing GMM parameters and fisher vectors, we use the VLFeat Computer Vision library [22].

<sup>3</sup><https://lear.inrialpes.fr/people/wang/improved-trajectories>

### 3.4 Learning Video to Hash-tag Embedding

In the first step, a 100-dimensional vector space, representative of the hash-tag distribution has been learned. Next we need to learn a mapping function that can project the 1003676 dimensional video features (fisher encoded IDTs) to the 100 dimensional tag vector space.

Socher et al. [18] proposed a method to learn a mapping from visual features (images) to word vectors (word2vec space) for detecting objects in a cross-modal zero-shot framework. We adopt this cross-modal learning approach to learn a mapping from fisher vectors to tag vectors. Similar to [18], we train a neural network with (fisher vector, tag word) pairs for each of the 2770 training videos to learn a non-linear embedding function from video features too Tag2Vec space. The tag word is the same as the category/class label of the training video. For learning this embedding function, we use the publicly available code for zero-shot learning<sup>4</sup>. The neural network is set up to have 600 hidden nodes and maximum iterations are set to 1000 as we have more categories. Training this network with our data took approximately 2 hours.

Figure 4 shows the t-SNE (t-Stochastic Neighborhood Embedding) visualization of training features in fisher vector space and Tag2Vec space. It can be clearly seen that after embedding the the training vectors form distinct clusters around their category words. Once the embedding function

<sup>4</sup><https://code.google.com/p/word2vec/>

Baseball	Basketball	Benchpress	Biking	Billiards	Boxing	Breaststroke	Diving	Drumming	Fencing
Golf	HighJump	Horseriding	HulaHoop	Juggling	Kayaking	Lunges	Nunchucks	Piano	PoleVault
Pushups	Yoyo	Salsa	Skateboarding	Skiing	Soccer	Swing	Tennis	Volleyball	

Table 2: 29 video categories used for training

is learned, a query video (belonging to these 29 categories) can be directly mapped to the tag space and relevant hash-tags can be recommended. In the next section, we explain the hash-tag recommendation mechanism.

### 3.5 Tag Suggestion Metrics

Given a query video, we first compute its fisher vector representation. We then use the learned embedding function to project the query fisher vector in the learned Tag2Vec space. We utilize a simple nearest neighbour approach based on  $L_2$  distance to retrieve potentially relevant hash-tags for a given query video. It can be seen that we do not directly compare the test/query vector to any of the training video vectors, neither in fisher vector space, not after the embedding. This is advantageous in terms of retrieval time and memory because, (i) there is no need to store the training set but only the Tag2Vec model and the fisher vectors to Tag2Vec space embedding function; and (ii) redundant comparisons are avoided. The tag words retrieved from the Tag2Vec space are stem words. Since we cannot suggest stems as hash-tags we need to convert a stemmed tag to its proper form. However, each stem corresponds to multiple tags, e.g. ‘beauty’, ‘beautiful’, ‘beautifully’ would all map to stem word ‘beauti’. In our system, for a particular stemmed tag, we pick the most commonly used word in our hash-tag corpus from among all corresponding inflections of that stem. This de-stemming approach is a bit limiting as many variations of the same stem words would always be rejected. A better approach based on parts of speech (PoS) tagging and edit distance can replace this. The time complexity of tag suggestion depends on the dimensionality of the tag space ( $d$ ) and number of tags ( $N$ ),  $O(N * d)$ . Our MATLAB implementation takes around 1 second for video-to-tag space embedding and less than a second for near-neighbour tag retrieval. For large-scale application, the simple nearest neighbour approach can be replaced by better retrieval mechanisms for efficiency and robustness.

## 4. EXPERIMENT & RESULTS

We conduct experiments to both qualitatively and quantitatively validate the tags suggested by our approach. Users who are familiar with social networking jargons are asked to take a survey. In user survey, each user is shown a set of vines with 15 recommended hashtags per vine. Users can pause/replay the vine and select the hashtags they consider can be used with the shown vine. Figure 5 shows a snapshot of the user study session. We perform these experiment with 14 users where every user marks each vine once. The testing set contains 270 vines with approximately 9 vines per class. Vines are wild and intra-class variations are quite drastic. Hence, we compute the average relevance scores per class for better understanding the cases where the system performs better or worse. As our dataset consists of wild, unconstrained, and unfiltered vines, there are cases where users don’t find any suggested hashtags as relevant. To see

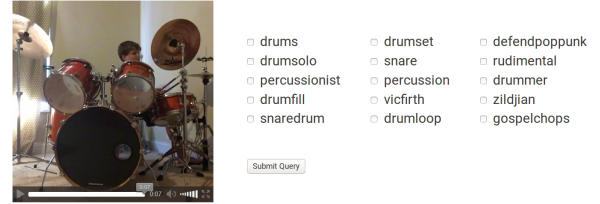


Figure 5: A screen shot of our user annotation (survey) platform.

whether this happens for specific classes or it is uniform across classes, we also collect the data for the number of vines per class for which users didn’t find any recommended tag as relevant.

The plot in Figure 7 (left) shows the average number of relevant tags suggested for each class. The plot in Figure 7 (right) shows the class-wise distribution of vines with no relevant hashtags. Benchpress contains the highest number of tags, 7, on an average followed by Volleyball at 6.88. Yoyo & Salsa are the worst performers with 1.25 & 1.45 tags respectively on an average. In total 52 vines out of the 270 didn’t contain a single relevant hashtag. Upon viewing these vines, we noticed that these vines in majority were the ones which visually and textually contained no information pertaining to the action tag. For example, a person talking about how good boxing is might contain relevant hashtags as assigned by the uploader but as we don’t process auditory modality and training of the embedding function relies only on visual features, the system is unavailable to correctly map such vines to the relevant concepts. Figure 6 shows some examples of success and failure cases for qualitative evaluation.

The overall number of relevant tags suggested for a vine is 4.03 out of 15 which is 27% of the tags suggested for each vine. Based on the hash-tag statistics collected by scraping the vine platform, we observed that 4.79 is the average number of hashtags associated with a typical vine (based on 2.5 million entries). One thing to note is that not all the tags in the ground-truth data are relevant hence this number is likely to come down. By suggesting 15 tags we are able to reproduce a similar number where an uploader finds 4 tags relevant to the vine which suggests that our system performs well even for such unconstrained videos.

## 5. CONCLUSION AND FUTURE WORKS

In summary, we present a method to automatically suggest hash-tags for short social video clips. Hash-tags are noisy and have ambiguous semantics. We learn a vector space which we show is able to capture these semantics. We call this vector space Tag2Vec. Also for automatically annotating hash-tags for a given video we learn a neural network based embedding function. We work on a self gathered dataset of wild short video clips of 29 categories. The em-

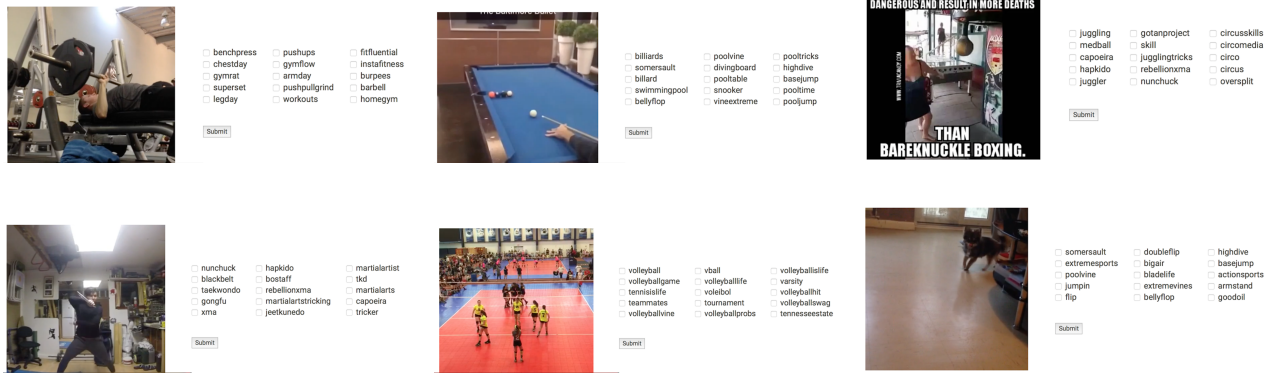


Figure 6: Hash tags suggested by our framework for the given video clip. First two columns show relevant hash-tag suggestions as predicted by our proposed model (like armday, benchpress, instafitness etc. for top left image). Last column shows two failure cases.

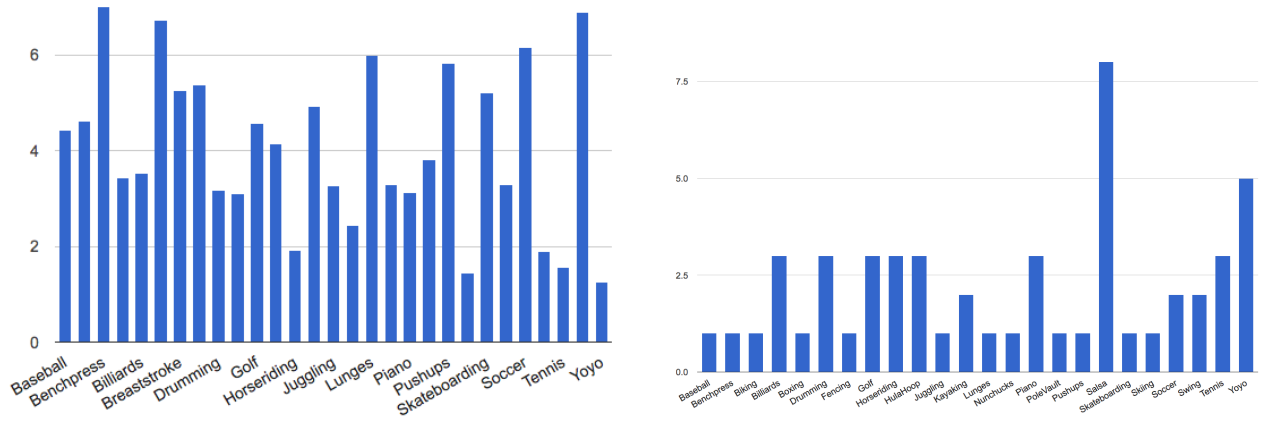


Figure 7: Left image shows the average number of relevant tags marked by the users for each class out of 15 suggest tags. Right image shows average number of videos across all users per class for which there were no relevant tags found out of entire test dataset of 50 videos per class

bedding function embeds any query video to our Tag2Vec space from which we propose hash-tags using simple nearest neighbour retrieval. We show that our Tag2Vec space has desired semantic structure and we are able to suggest relevant hash-tags for the query videos. In future, we would like to explore sentimental vs. contextual relevance in the Tag2Vec space and would also like to incorporate relevance metric that depends on the evolving popularity and trends of the hash-tags.

## 6. ACKNOWLEDGEMENT

We thank Shivam Kakkar and Siddhartha Gairola for their crucial contributions in developing the tool for user study. Rajvi Shah and Saurabh Saini are supported respectively by Google and TCS PhD fellowships. We thank these organizations for their valuable support.

## References

- [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *In-*

*ternational Conference on Computer Vision (ICCV)*, 2015.

- [2] Lamberto Ballan, Marco Bertini, Alberto Del Bimbo, Marco Meoni, and Giuseppe Serra. Tag suggestion and localization in user-generated videos based on social knowledge. In *Proceedings of Second ACM SIGMM Workshop on Social Media, WSM '10*, pages 3–8. ACM, 2010.
- [3] Oren Barkan. Bayesian neural word embedding. *CoRR*, abs/1603.06571, 2016. URL <http://arxiv.org/abs/1603.06571>.
- [4] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003. ISSN 1532-4435.
- [5] Amir Globerson, Gal Chechik, Fernando Pereira, and Naftali Tishby. Euclidean embedding of co-occurrence data. In *Advances in Neural Information Processing Systems 17*, pages 497–504, 2005.

- [6] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3128–3137, 2015.
- [7] V. Lavrenko, S. L. Feng, and R. Manmatha. Statistical models for automatic video annotation and retrieval. In *Acoustics, Speech, and Signal Processing, 2004. Proceedings. (ICASSP '04). IEEE International Conference on*, volume 3, pages iii–1044–7 vol.3, 2004.
- [8] Rémi Lebrete and Ronan Collobert. Word embeddings through hellinger pca. In *14th Conference of the European Chapter of the Association for Computational Linguistics*, 2014.
- [9] Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2177–2185, 2014.
- [10] Yitan Li, Linli Xu, Fei Tian, Liang Jiang, Xiaowei Zhong, and Enhong Chen. Word embedding revisited: A new representation learning and explicit matrix factorization perspective. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJ-CAI'15*, pages 3650–3656, 2015.
- [11] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- [12] E. Moxley, T. Mei, and B. S. Manjunath. Video annotation through search and graph reinforcement mining. *IEEE Transactions on Multimedia*, 12(3):184–193, 2010.
- [13] F. Perronnin and C. Dance. Fisher kernels on visual vocabularies for image categorization. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2007.
- [14] Florent Perronnin, Jorge Sánchez, and Thomas Mensink. Improving the fisher kernel for large-scale image classification. In *Proceedings of the 11th European Conference on Computer Vision: Part IV, ECCV'10*, pages 143–156, 2010.
- [15] Guo-Jun Qi, Xian-Sheng Hua, Yong Rui, Jinhui Tang, Tao Mei, and Hong-Jiang Zhang. Correlative multi-label video annotation. In *Proceedings of the 15th ACM International Conference on Multimedia*, MM '07, pages 17–26. ACM, 2007.
- [16] Aditya Singh, Saurabh Saini, Rajvi Shah, and P J Narayanan. From traditional to modern: Domain adaptation for action classification in short social video clips. In *Proceedings of German Conference on Pattern Recognition (GCPR)*, pages 245–257, 2016.
- [17] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. Zero-shot learning through cross-modal transfer. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 935–943, 2013.
- [18] Richard Socher, Milind Ganjoo, Hamsa Sridhar, Osbert Bastani, Christopher D. Manning, and Andrew Y. Ng. Zero-shot learning through cross-modal transfer. *CoRR*, abs/1301.3666, 2013.
- [19] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *CoRR*, abs/1212.0402, 2012.
- [20] Jinhui Tang, Xian-Sheng Hua, Guo-Jun Qi, Meng Wang, Tao Mei, and Xiuqing Wu. Structure-sensitive manifold ranking for video concept detection. In *Proceedings of the 15th ACM International Conference on Multimedia*, MM '07, pages 852–861, 2007.
- [21] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Movieqa: Understanding stories in movies through question-answering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [22] Andrea Vedaldi and Brian Fulkerson. Vlfeat: An open and portable library of computer vision algorithms. In *Proceedings of the 18th ACM International Conference on Multimedia*, MM '10, pages 1469–1472, 2010.
- [23] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. *CoRR*, abs/1411.4555, 2014.
- [24] Heng Wang and C. Schmid. Action recognition with improved trajectories. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 3551–3558, 2013.
- [25] M. Wang, X. S. Hua, J. Tang, and R. Hong. Beyond distance measurement: Constructing neighborhood similarity for video annotation. *IEEE Transactions on Multimedia*, 11(3):465–476, 2009.
- [26] Ting Yao, Tao Mei, Chong-Wah Ngo, and Shipeng Li. Annotation for free: Video tagging by mining user search behavior. In *Proceedings of the 21st ACM International Conference on Multimedia*, MM '13, pages 977–986. ACM, 2013.
- [27] Yang Zhang, Boqing Gong, and Mubarak Shah. Fast zero-shot image tagging. *CoRR*, abs/1605.09759, 2016. URL <http://arxiv.org/abs/1605.09759>.
- [28] W. L. Zhao, X. Wu, and C. W. Ngo. On the annotation of web videos by efficient near-duplicate search. *IEEE Transactions on Multimedia*, 12(5):448–461, 2010.