

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Youtube2Story: Unsupervised Joint-Representation of Instructional Videos

Anonymous ICCV submission

Paper ID ****

Abstract

The ABSTRACT

1. Introduction

Leaning the instructions of a novel non-trivial task is both a challenge and a necessity for both humans and autonomous systems. This necessity resulted in many community generated instruction collections [?, ?] and expert curated recipe books[?, ?]. However, these instructions are generally based on a language modality and explains a single way of performing the task although there are variety of ways. On the other hand, online video storage services are full of unstructured instructional videos¹ covering variety of ways, environment conditions and view angles. Although there have been many successful attempts in detecting activities from videos [?, ?], structural representation of such a large and useful video collection is not possible. In this paper, we focus on joint semantic representation of YouTube videos as a response to a single query. We specifically study the unsupervised joint-detection of the activities from a collection of YouTube videos.

Understanding of the instructional videos, requires the careful processing of two complementary modalities namely language and the vision. Luckily the target domain, YouTube videos, has unstructured subtitles as well. They are either generated by the content developer (5% of the time) or automatically generated by using the Automatic Speech Recognition (ASR) softwares. The main limitations of the existing activity detection literature for this problem is scalability and representation level. Existing approaches are mainly supervised and requires extensive training set which is not tractable in the scale of YouTube videos. Moreover, current activity detection research focuses on the low-level visual features. However, such videos in the wild have objects with completely different texture and shape characteristics from wide range of views. Instead, we focus on extracting high-level visual semantic representations and us-

ing salient words occurring among the videos.

We rely on the assumption that the videos collected as the response of a same instructional query, share similar activities performed by the similar objects. We start with the independent processing of the videos in order to create a large collection of visual object proposals and words. After the proposal generation, we jointly process the proposal collections and words to detect the visual objects and words which can be used to represent the unstructured information. Since we rely on high-level information instead of the low-level features, the resulting objects represent the semantic information instead of visual characteristic. By using the extracted objects, we compute the holistic representation of the multi-modal information in each frame.

Moving from frame-wise visual understanding to activity understanding, requires the joint processing of all the videos with the temporal information. In order to exploit the temporal information, we model each video as a Hidden Markov Model using state space of activities. Since we assume that the videos share some of the activities and we have no supervision, we use a model based on *beta process mixture model*. Our model jointly learn the activities and detect them in the videos. Moreover, it does not require prior knowledge over the number of activities.

2. Related Work

Video Summarization Summarizing an input video as a sequence of keyframes (static) or as a sequence of video clips (dynamic) is useful for both multi-media search interfaces and retrieval purposes. Early works in the area are summarized in [27] and mostly focus on choosing keyframes for visualization. Idea of choosing key-frames is also extended by using the video tags by Hong et al[8] and using the spatio-temporal information by Gygli et al.[7].

Summarizing videos is crucial for ego-centric videos since the ego-centric videos are generally long in duration. There are many works which successfully segment such videos into a sequence of important shots [15, 16]; however, they mostly rely on specific features of ego-centric videos. Rui et al. [23] proposed another dynamic summarization method based on the excitement in the speech of the

¹YouTube has 281.000 videos for "How to tie a bow tie"

108 reporter. Due to their domain specific designs, these algorithms are not applicable to the general instructional videos.
109

110 Same idea is also applied to the large image collections by recovering the temporal ordering and visual similarity
111 of images [12]. This image collections are further used to choose important view points for video key-frame selection
112 by Khosla et al.[10]. And further extended to video clip selection by Kim et al[11] and Potapov et al.[22]. Although
113 they are different from our approach since they do not use any high-level semantic information or the language information,
114 we experimentally compare our method with them.
115
116

117 **Understanding Multi-Modal Information:** Learning
118 the relationship between the visual and language data is a crucial problem due to its immense multimedia applications.
119 Early methods [1] in the area focus on learning a common multi-modal space in order to jointly represent
120 language and vision. They are also extended to learning higher level relations between object segments and words
121 [24]. Zitnick et al.[30, 29] used abstracted clip-arts to further understand spatial relations of objects and their language correspondences. Kong et al. [14] and Fidler et
122 al. [5] both accomplished the same task by using the image captions only. Relations extracted from image-caption
123 pairs, are further used to help semantic parsing [28] and activity recognition [20]. Recent works also focused on generating
124 image captions automatically using the input image. These methods range from finding similar images and using
125 their captions [21] to learning language modal conditioned on the image [13, 25, 4]. And the methods to learn language
126 models vary from graphical models [4] to neural networks [25, 13, 9].
127

128 All aforementioned methods are using supervised labels either as strong image-word pairs or weak image-caption
129 pairs. On the other hand, our method is fully unsupervised.
130

131 Activity Detection/Recognition:

132 **Recipe Understanding** Following the interest in community generated recipes in the web, there have been many attempts to automatically process recipes. Recent methods on natural language processing [18, 26] focus on semantic parsing of language recipes in order to extract actions and the objects in the form of predicates. Tenorth et al.[26] further process the predicates in order to form a complete logic plan. Mori et al.[19] also learns the relations of the actions in terms of a flow graph with the help of a supervision. The aforementioned approaches focus only on the language modality and they are not applicable to the videos. We have also seen recent advances [2, 3] in robotics which uses the parsed recipe in order to perform cooking tasks. They use supervised object detectors and report a successful autonomous cooking experiment. In addition to the lan-

133 guage based approaches, Malmaud et al.[17] consider both language and vision modalities and propose a method to align an input video to a recipe. However, it can not extract the steps/actions automatically and requires a ground truth recipe to align. On the contrary, our method uses both visual and language modalities and extracts the actions while autonomously constructing the recipe. There is also an approach which generates multi-modal recipes from expert demonstrations [6]. However, it is developed only for the domain of *teaching user interfaces* and are not applicable to the videos.
134
135

136 3. Method

137 In this section, we explain the high-level components of the method we develop to jointly represent multi-modal instructions. As shown in the Figure ??, our proposed method consists of three major components; online query and filtering, frame-wise multi-modal representation and joint clustering to extract activities. **(1) Query subsystem:** Our system starts with querying the YouTube with an *How to* question for top 100 videos. The text descriptions of the returned videos are represented as bag-of-words and clustered to eliminate outliers. **(2) Framewise Representation:** In order to represent the frames of the returned videos, we process both the visual and language content of the videos. We extract the object proposals and jointly cluster them in order to detect the salient objects of the recipe. For the language descriptions, we use the top salient words of the corpus generated by concatenation of the all subtitles. We represent each frame in terms of the resulting salient objects and words. **(3) Unsupervised Activity Detection:** After describing each frame by using the salient objects and words, we apply a non-parametric Bayesian method in order to find the temporally consistent video clip clusters occurring over multiple videos. Our empirical results suggest that the resulting clusters mostly correspond to the activities. We now explain the details of the each sub-system in the following sections.
138
139

140 3.1. Semantic Multi-Modal Frame Representation

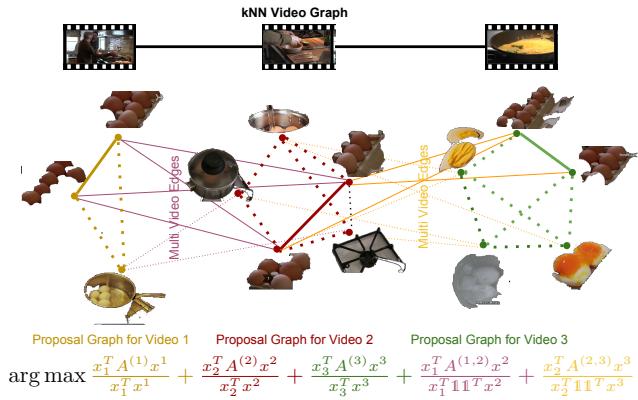
141 We represent the each frame of input videos as in terms of set of language and visual atoms. Language atoms are the salient words detected via ranking tf-idf like measure, and the visual atoms are found via clustering over object proposals which we extract from frames. We explain the details of finding the atoms and representing the frames in the subsequent sections.
142
143

144 3.1.1 Learning Language Atoms

145 In order to detect the salient words, we use tf-idf like measure. For each recipe, we concatenate all subtitles into a single term corpus. As a document corpus, we use all the
146
147

148 162
149 163
150 164
151 165
152 166
153 167
154 168
155 169
156 170
157 171
158 172
159 173
160 174
161 175
162 176
163 177
164 178
165 179
166 180
167 181
168 182
169 183
170 184
171 185
172 186
173 187
174 188
175 189
176 190
177 191
178 192
179 193
180 194
181 195
182 196
183 197
184 198
185 199
186 200
187 201
188 202
189 203
190 204
191 205
192 206
193 207
194 208
195 209
196 210
197 211
198 212
199 213
200 214
201 215

324
 325 **Joint Proposal Clustering to Detect Visual Objects**
 326 Since our proposals are generated from multiple videos,
 327 combining them into a single region collection and cluster-
 328 ing it is not desired for two reasons; (1) the objects have
 329 large visual differences among different videos and accu-
 330 rately clustering them into a single object is hard, (2) objects
 331 clusters are desired to have region proposals from multiple
 332 videos in order to relate videos easily. Hence, we propose a
 333 joint version of the spectral region clustering algorithm.



347 Figure 2: Visualization of the joint proposal clustering.
 348

349 Here, we first explain the original spectral graph cluster-
 350 ing algorithm and then extend it to our joint version. Con-
 351 sider set of region proposals extracted from a single video
 352 r_t^k , and a similarity metric $d(\cdot, \cdot)$ which assigns similarity
 353 between each region pair. We follow the Single Cluster
 354 Graph Partitioning (SCGP)[?] to find the dominant cluster
 355 which maximizes the inter-cluster similarity. In other
 356 words, we solve
 357

$$\arg \max \frac{\sum_{(k_1, t_1), (k_2, t_2) \in K \times T} x_{t_1}^{k_1} x_{t_2}^{k_2} d(r_{t_1}^{k_1}, r_{t_2}^{k_2})}{\sum_{(k, t) \in K \times T} x_t^k} \quad (1)$$

361 Where, $x_t^{i,k}$ is a binary variable which is 1 if $r_t^{(i),k}$ is
 362 included in the cluster. Moreover, we drop the video index
 363 for clarity. When we use the vector form of the indicator
 364 variables as $\mathbf{x}_{tK+k} = x_t^k$ and the pairwise distance ma-
 365 trix as $\mathbf{A}_{t_1 K + k_1, t_2 K + k_2} = d(r_{t_1}^{k_1}, r_{t_2}^{k_2})$, this equation can
 366 be compactly written as $\arg \max_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$. Moreover, it can
 367 be solved by computing the dominant eigenvector of \mathbf{x} af-
 368 ter relaxing it to $[0, 1]$ from binary [?, ?]. After finding the
 369 maximum, the remaining clusters can be found by remov-
 370 ing the selected region proposals from the collection, and
 371 re-applying the same algorithm for the second dominant clus-
 372 ter.
 373

374 Our extension of the SCGP into multiple videos is based
 375 on the assumption that the dominant objects occur in most
 376 of the videos. Hence, we re-formulate the problem by
 377 relating videos to each other. We start with creating the kNN

378 graph of videos based on their language description simi-
 379 larities. Moreover, we also create the kNN graph of region
 380 proposals in each video. This hierarchical graph structure is
 381 visualized in Figure 2. After creating this graph, we impose
 382 the similarity of regions coming from each video as well
 383 as the similarity of regions coming from neighbour videos.
 384 Hence, given the pairwise distance matrix $\mathbf{A}^{(i)}$, binary in-
 385 dicator vectors $\mathbf{x}^{(i)}$ for each video and pairwise distance
 386 matrices for video pairs as $\mathbf{A}^{(i,j)}$, we define our optimiza-
 387 tion problem as;

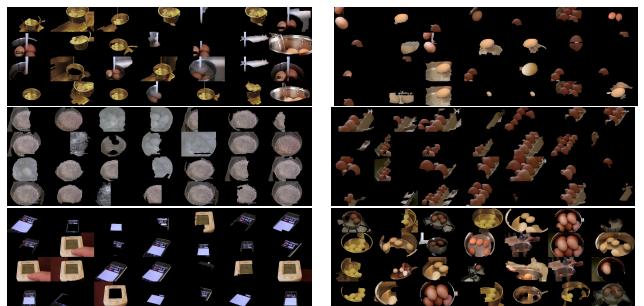
$$\arg \max \sum_{i \in N} \frac{\mathbf{x}^{(i)^T} \mathbf{A}^{(i)} \mathbf{x}^{(i)}}{\mathbf{x}^{(i)^T} \mathbf{x}^{(i)}} + \sum_{i \in N} \sum_{j \in \mathcal{N}(i)} \frac{\mathbf{x}^{(i)^T} \mathbf{A}^{(i,j)} \mathbf{x}^{(j)}}{\mathbf{x}^{(i)^T} \mathbf{1} \mathbf{1}^T \mathbf{x}^{(j)}} \quad (2)$$

392 where $\mathcal{N}(i)$ is the neighbours of the video i in the kNN
 393 graph and $\mathbf{1}$ is vector of ones. We visualize this optimiza-
 394 tion objective in Figure 2 for the case of 3 videos and 1-NN
 395 video graph, 2-NN region proposal graph.

396 After changing the optimization function, we can not use
 397 the efficient eigen-decomposition based approach from [?,
 398 ?]; however, we can use stochastic gradient descent (SGD)
 399 since the cost function is quasi-convex when it is relaxed.
 400 Hence, we use SGD with the following gradient function;
 401

402 Some vector matrix multiplication (3)

403 In Figure 3, we visualize some of the clusters which our
 404 algorithm generated after applied on the videos returned by
 405 the query *How to Hard Boil an Egg*. As shown the figure,
 406 the resulting clusters are highly correlated and correspond
 407 to semantic objects&concepts.



420 Figure 3: Randomly selected images of randomly selected
 421 clusters learned for *How to boil an egg*?
 422

4.2. Multi-Modal Representation via Learned Atoms

4. Unsupervised Activity Representation

428 Here we explain the generative model which we use in
 429 order to jointly learn the activities from videos. We start
 430 with explaining the basic notation we use in the model. We
 431 denote the t^{th} frame of the i^{th} video as I_t^i and its subtitle as

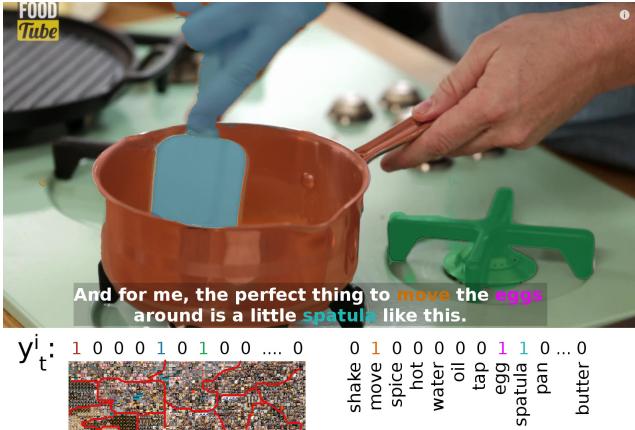


Figure 4: Visualziation of the representation of the Frame.

L_t^i . Moreover, we note the extracted frame representation as y_t^i . We model our algorithm based on activities and the note the activity of the t^{th} frame of the i^{th} video as z_t^i . Since our model is non-parametric, number of activities are not fixed and $z_t^i \in \mathcal{N}$.

Frame: Representation of each frame is computed as the occurence vector over the detected language and visual atoms. Formally, frame y_t^i is represented as $y_t^i = [y_t^{i,l}, y_t^{i,v}]$ such that k^{th} entry of the $y_t^{i,l}$ is 1 if the frame as language atom k and 0 otherwise. $y_t^{i,l}$ is also a binary vector similarly defined over visual atoms. We consider y_t^i as the observed variable and consider the underlying activity label as the hidden state z_t^i . A sample state is visualized in the Figure 4.

Activity: We represent each activity as a Bernoulli distribution over the visual and language atoms. In other words, each frame's representation y_t^i is sampled from its activity distribution as $y_t^i|z_t^i = k \sim Ber(\Theta_k)$. For the sake of clarity, we sample the Θ from its conjugate distribution *Beta distribution*.

In the following sections, we exlpain how these models can be jointly learned and inferred by using the Beta Process Hidden Markov Models. s

4.1. Beta Process Hidden Markov Model

For joint understanding of the time-series information, Fox et al.[?] proposed the Beta Process Hidden Markov Models (BP-HMM) using the Indian Buffet Process[?] of time-series sequences over features. It assumes that there exist a set of features(activities in our case) which can explain the behaviour of all time-series data (all videos in our case). Each time-series data exhibits a subset of available features. This setup is similar to Hughes et al.[?]. However, we differ in the choices of the underlying distributions since we based our model on semantic multi-modal information.

In our model, each video i chooses a set of activities through a activity vector \mathbf{f}^i such that f_k^i is 1 if i^{th} video has activity k , 0 otherwise. When the feature vectors of all videos in the corpus is concatenated, it becomes an activity matrix \mathbf{F} such that i^{th} row of the \mathbf{F} is the activity vector \mathbf{f}^i . Moreover, each feature k also has an activity frequency b_k and a distribution parameter Θ_k . Distribution parameter Θ_k is the Bernoulli distribution as explain in the Section 4. Moreover, its base distribution (B_0) is *Beta random variable* since it is the conjugate of Bernoulli. Moreover, in this setting, the activity paremeters Θ_k and b_k follow the *beta process* as;

$$B|B_0, \gamma, \beta \sim BP(\beta, \gamma B_o), B = \sum_{k=1}^{\infty} b_k \delta_{\Theta_k} \quad (4)$$

where B_0 and the b_k are determined by the underlying poisson process [?] and the feature vector is determined as independent Bernoulli draws as $f_k^i \sim Ber(b_k)$. After marginilizing over the b_k and Θ_k , this distribution is shown to be equivalent to Indian Buffet Process [?]. Where videos are customers and activities are dishes in the buffet. The first video chooses a Poisson(γ) unique dishes. The following video i chooses previously sampled activity k with probability $\frac{m_k}{i}$ proportional to number of videos (m_k) chosen the activity k , and it also chooses Poisson($\frac{\gamma}{i}$) new activities. Here, γ controls the number of active features in each video and β controls the likelihood of the features getting shared by multiple videos.

After each video chooses a subset of activities, we model the videos as an Hidden Markov Model (HMM) over the selected videos. Each frame has the hidden state activity $id(z_t^i)$ and we observe the binary representation y_t^i . Since we model each activity as a Bernoulli distribution, the emmittion probabilities follow the Bernoulli distribution as $p(y_t^i|z_t^i) = Ber(\Theta_{z_t^i})$. Following the construction of the Fox et al.[?], we sample the transition probabilities from a normalized Gamma distribution. For each video i , we sample a Gamma random variable for the transition between activity j and activity k if both of the activities are included by the video (if f_k^i and f_j^i are both 1). After sampling these random variables, we normalize them to have proper transition probabilities. This procedure can be represented formally as

$$\eta_{j,k}^i \sim Gam(\alpha + \kappa \delta_{j,k}, 1), \quad \pi_j^i = \frac{\eta_j^i \circ f^i}{\sum_k \eta_{j,k}^i f_k^i} \quad (5)$$

Where κ is the persistence parameter promoting self state transitions to have more coherent temporal bounderies, \circ is the element-wise product and π_j^i is the transition probabilities in video i from state j to all states in the form of a vector.

This model is also presented as a graphical model in Figure 5

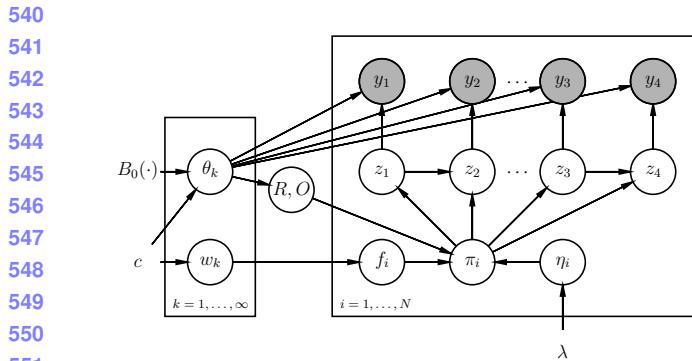


Figure 5: **Graphical model for BP-HMM:** Some explanation.

4.2. Gibbs sampling for BP-HMM

We employ Markov Chain Monte Carlo (MCMC) method for learning and inference of the BP-HMM. We base our algorithms on the MCMC procedure proposed by Fox et al.[?]. It marginilizes over blah and blah and sample blah and blah. For faster convergence, we also utilize a series of data driven samplers. Here we only discuss the proposed data driven samplers and move the details of the remainin samplers to the Supplementary Material.

Data-Driven Sampler1

Data-Driven Sampler2

5. Experiments

5.1. Dataset

We collected blah blah videos from YouTube and did blah blah... All numbers etc. These are recipes, 5 of them are evaluation set and labelled temporally and labels are matched.

5.2. Baselines

We compare against BP-HMM-O with local feauters and HMM with Semantic Features, HMM with local features, HMM with CNN feautes...

5.3. Qualitative Results

5.4. Accuracy over Activity Detection

5.5. Accuracy over Activity Learning

6. Discussions and Conclusions

Discuss which recipes worked and why. Discuss the importance of semantic representation, scaling features and multi-modality.



Crack the eggs one at a time Remove the omelette onto a into a bowl.



You can either use a fork or Eggs cook quickly, so make wire whisk to beat the eggs into sure the pan gets very hot first; a bowl. the butter melt completely.

(c) Sample images and the automatically generated captions for some of the clusters.

References

- [1] K. Barnard, P. Duygulu, D. Forsyth, N. De Freitas, D. M. Blei, and M. I. Jordan. Matching words and pictures. *The Journal of Machine Learning Research*, 3:1107–1135, 2003. 2
- [2] M. Beetz, U. Klank, I. Kresse, A. Maldonado, L. Mosenlechner, D. Pangercic, T. Ruhr, and M. Tenorth. Robotic roommates making pancakes. In *Humanoid Robots (Humanoids), 2011 11th IEEE-RAS International Conference on*, pages 529–536. IEEE, 2011. 2
- [3] M. Bollini, J. Barry, and D. Rus. Bakebot: Baking cookies with the pr2. In *The PR2 Workshop: Results, Challenges and Lessons Learned in Advancing Robots with a Common Platform, IROS*, 2011. 2
- [4] A. Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every picture tells a story: Generating sentences from images. In *Computer Vision–ECCV 2010*, pages 15–29. Springer, 2010. 2
- [5] S. Fidler, A. Sharma, and R. Urtasun. A sentence is worth a thousand pixels. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 1995–2002. IEEE, 2013. 2
- [6] F. Grabler, M. Agrawala, W. Li, M. Dontcheva, and T. Igarashi. Generating photo manipulation tutorials by

Table 1: Notation of the Paper

$y_t = [y_t^v, y_t^f]$	feature representation of t^{th} frame	I_t	t^{th} frame of the video	$x_{i,r}^p$	1 if p^{th} cluster has r^{th} proposal of i^{th} video
x^p	binary vector for p^{th} cluster	L_t	subtitle for t^{th} frame	$O^{k,k'}$	1 if $\#(z_t = k, z_{t'} = k') = 0 \forall t \leq t'$
$\Theta_k = [\Theta_k^v, \Theta_k^l]$	emmition prob. of k^{th} activity	z_t	activity ID of frame t	f_i^k	1 if i^{th} video has k^{th} activity 0.o.w.
$\eta_i^{k,k'}$	$P(z_{t+1} = k' z_t = k)$ for i^{th} vid	$\pi_i^{k,k'}$	$\eta_i^{k,k'} \times f_i^k \times f_i^{k'}$		

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

- demonstration. *ACM Transactions on Graphics (TOG)*, 28(3):66, 2009. 2
- [7] M. Gygli, H. Grabner, H. Riemenschneider, and L. Van Gool. Creating summaries from user videos. In *Computer Vision–ECCV 2014*, pages 505–520. Springer, 2014. 1
- [8] R. Hong, J. Tang, H.-K. Tan, C.-W. Ngo, S. Yan, and T.-S. Chua. Beyond search: event-driven summarization for web videos. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 7(4):35, 2011. 1
- [9] A. Karpathy and L. Fei-Fei. Deep Visual-Semantic Alignments for Generating Image Descriptions. *ArXiv e-prints*, Dec. 2014. 2
- [10] A. Khosla, R. Hamid, C.-J. Lin, and N. Sundaresan. Large-scale video summarization using web-image priors. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2698–2705. IEEE, 2013. 2
- [11] G. Kim, L. Sigal, and E. P. Xing. Joint summarization of large-scale collections of web images and videos for storyline reconstruction. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 4225–4232. IEEE, 2014. 2
- [12] G. Kim and E. P. Xing. Reconstructing storyline graphs for image recommendation from web community photos. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 3882–3889. IEEE, 2014. 2
- [13] R. Kiros, R. Salakhutdinov, and R. Zemel. Multimodal neural language models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 595–603, 2014. 2
- [14] C. Kong, D. Lin, M. Bansal, R. Urtasun, and S. Fidler. What are you talking about? text-to-image coreference. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 3558–3565. IEEE, 2014. 2
- [15] Y. J. Lee, J. Ghosh, and K. Grauman. Discovering important people and objects for egocentric video summarization. In *CVPR*, volume 2, page 6, 2012. 1
- [16] Z. Lu and K. Grauman. Story-driven summarization for egocentric video. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 2714–2721. IEEE, 2013. 1
- [17] J. Malmaud, J. Huang, V. Rathod, N. Johnston, A. Rabinovich, and K. Murphy. What’s Cookin’? Interpreting Cooking Videos using Text, Speech and Vision. *ArXiv e-prints*, Mar. 2015. 2
- [18] J. Malmaud, E. J. Wagner, N. Chang, and K. Murphy. Cooking with semantics. *ACL 2014*, page 33, 2014. 2
- [19] S. Mori, H. Maeta, Y. Yamakata, and T. Sasada. Flow graph corpus from recipe texts. In *Proceedings of the Ninth In-*

ternational Conference on Language Resources and Evaluation, pages 2370–2377, 2014. 2

- [20] T. S. Motwani and R. J. Mooney. Improving video activity recognition using object recognition and text mining. In *ECAI*, pages 600–605, 2012. 2
- [21] V. Ordonez, G. Kulkarni, and T. L. Berg. Im2text: Describing images using 1 million captioned photographs. In *Advances in Neural Information Processing Systems*, pages 1143–1151, 2011. 2
- [22] D. Potapov, M. Douze, Z. Harchaoui, and C. Schmid. Category-specific video summarization. In *Computer Vision–ECCV 2014*, pages 540–555. Springer, 2014. 2
- [23] Y. Rui, A. Gupta, and A. Acero. Automatically extracting highlights for tv baseball programs. In *Proceedings of the eighth ACM international conference on Multimedia*, pages 105–115. ACM, 2000. 1
- [24] R. Socher and L. Fei-Fei. Connecting modalities: Semi-supervised segmentation and annotation of images using unaligned text corpora. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 966–973. IEEE, 2010. 2
- [25] R. Socher, A. Karpathy, Q. V. Le, C. D. Manning, and A. Y. Ng. Grounded compositional semantics for finding and describing images with sentences. *Transactions of the Association for Computational Linguistics*, 2:207–218, 2014. 2
- [26] M. Tenorth, D. Nyga, and M. Beetz. Understanding and executing instructions for everyday manipulation tasks from the world wide web. In *Robotics and Automation (ICRA), 2010 IEEE International Conference on*, pages 1486–1491. IEEE, 2010. 2
- [27] B. T. Truong and S. Venkatesh. Video abstraction: A systematic review and classification. *ACM Transactions on Multimedia Computing, Communications, and Applications*, 3(1):3, 2007. 1
- [28] H. Yu and J. M. Siskind. Grounded language learning from video described with sentences. In *ACL (1)*, pages 53–63, 2013. 2
- [29] C. L. Zitnick and D. Parikh. Bringing semantics into focus using visual abstraction. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3009–3016. IEEE, 2013. 2
- [30] C. L. Zitnick, D. Parikh, and L. Vanderwende. Learning the visual interpretation of sentences. In *Computer Vision (ICCV), 2013 IEEE International Conference on*, pages 1681–1688. IEEE, 2013. 2