

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 videos<sup>1</sup> 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

<sup>1</sup>YouTube 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  
119 a crucial problem due to its immense multimedia applications.  
120 Early methods [1] in the area focus on learning  
121 a common multi-modal space in order to jointly represent  
122 language and vision. They are also extended to learning  
123 higher level relations between object segments and words  
124 [24]. Zitnick et al.[30, 29] used abstracted clip-arts to further  
125 understand spatial relations of objects and their language  
126 correspondences. Kong et al. [14] and Fidler et  
127 al. [5] both accomplished the same task by using the image  
128 captions only. Relations extracted from image-caption  
129 pairs, are further used to help semantic parsing [28] and activity  
130 recognition [20]. Recent works also focused on generating  
131 image captions automatically using the input image.  
132 These methods range from finding similar images and using  
133 their captions [21] to learning language modal conditioned  
134 on the image [13, 25, 4]. And the methods to learn language  
135 models vary from graphical models [4] to neural networks  
136 [25, 13, 9].  
137

138 All aforementioned methods are using supervised labels  
139 either as strong image-word pairs or weak image-caption  
140 pairs. On the other hand, our method is fully unsupervised.  
141  
142

#### 143 Activity Detection/Recognition:

144 **Recipe Understanding** Following the interest in community  
145 generated recipes in the web, there have been many attempts  
146 to automatically process recipes. Recent methods on natural  
147 language processing [18, 26] focus on semantic parsing of  
148 language recipes in order to extract actions and the objects in  
149 the form of predicates. Tenorth et al.[26] further process the  
150 predicates in order to form a complete logic plan. Mori et al.[19]  
151 also learns the relations of the actions in terms of a flow graph  
152 with the help of a supervision. The aforementioned approaches  
153 focus only on the language modality and they are not applicable  
154 to the videos. We have also seen recent advances [2, 3] in robotics  
155 which uses the parsed recipe in order to perform cooking tasks.  
156 They use supervised object detectors and report a successful  
157 autonomous cooking experiment. In addition to the lan-  
158  
159  
160  
161

162 guage based approaches, Malmaud et al.[17] consider both  
163 language and vision modalities and propose a method to align  
164 an input video to a recipe. However, it can not extract  
165 the steps/actions automatically and requires a ground truth  
166 recipe to align. On the contrary, our method uses both visual  
167 and language modalities and extracts the actions while  
168 autonomously constructing the recipe. There is also an approach  
169 which generates multi-modal recipes from expert  
170 demonstrations [6]. However, it is developed only for the  
171 domain of *teaching user interfaces* and are not applicable  
172 to the videos.  
173

### 174 3. Method

175 In this section, we explain the highlevel components of  
176 the method we develop to jointly represent multi-modal  
177 instructions. We explain the details of the each sub-system  
178 in the following sections. As shown in the Figure ??, our  
179 proposed method consists of three major components; on-  
180 line query and filtering, frame-wise multi-modal represen-  
181 tation and joint clustering to extract activities. **Query:** Our  
182 system starts with querying the YouTube with an *How to*  
183 question for top 100 videos. The text descriptions of the  
184 returned videos are represented as bag-of-words and clus-  
185 tered to eliminate outliers. **Framewise Representation:** In  
186 order to represent the frames of the returned videos, we pro-  
187 cess both the visual and language content of the videos. We  
188 extract object proposals and jointly cluster them in order to  
189 detect the salient objects of the recipe. For the language  
190 descriptions, we use the top salient words of the corpus  
191 generated as concatenation of the all subtitles. We repre-  
192 sent each frame in terms of the resulting salient objects and  
193 words. **Unsupervised Activity Detection:** After describing  
194 each frame by using the salient objects and words, we apply  
195 a non-parametric Bayesian method in order to find tempo-  
196 rally consistent clusters occurring over multiple videos.  
197 The resulting clusters are used to label activities of the test-  
198 videos and summarize them as a flow graph.  
199  
200

### 201 4. Semantic Multi-Modal Frame Representa- 202 tion

#### 203 4.1. Learning Atoms to Represent Multi-Modal 204 Data

205 We represent each frame in our video set as a distribution  
206 over set of language and visual atoms. Language atoms are  
207 the salient words detected by using a tf-idf like measure,  
208 and the visual atoms are found via clustering over object  
209 proposals which we extract from frames. We explain the  
210 details of finding the atoms and representing the frames in  
211 the subsequent sections.  
212  
213  
214  
215

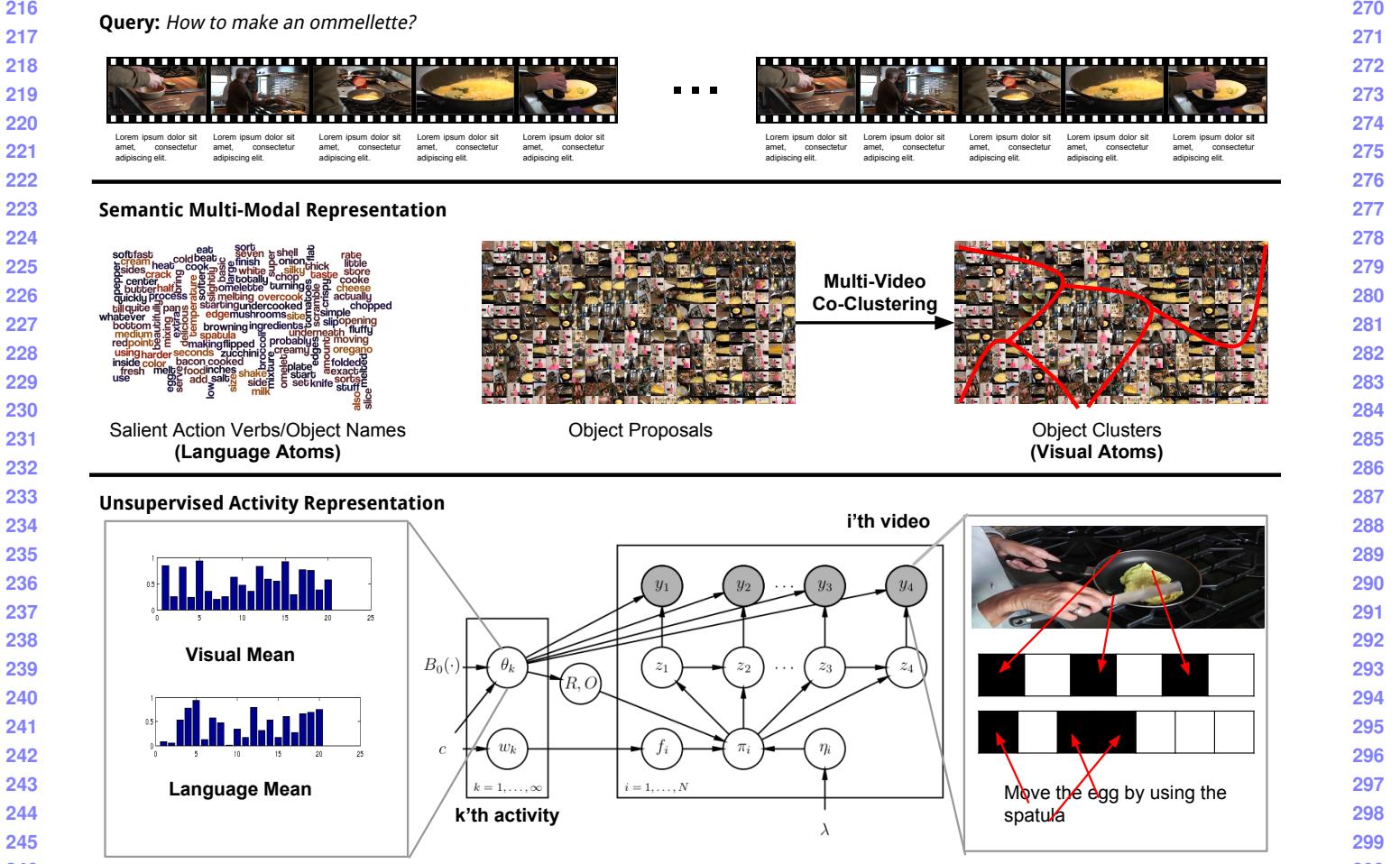


Figure 1: Components of our recipe understanding method. **Query:** We query the YouTube for top 100 *How To* videos and filter the outliers; **Framewise Representation:** We automatically extract object clusters and salient word in order to find multi-modal representation of each frame. **Unsupervised Activity Detection:** We jointly cluster videos in order to learn activities/steps related to the recipe.

#### 4.1.1 Learning Language Atoms

In order to detect salient words, we use tf-idf like measure. For each recipe, we concatinate all subtitles into a single term corpus. As a document corpus, we use all the words extracted from all recipes. Moreover, we compute the tf-idf as  $tfidf(w, d, D) = f_{w,d} \times \log \frac{N}{n_w}$  where  $w$  is the word,  $d$  is the corpus of the recipe and  $f_{w,d}$  is the frequency of the word in the recipe. Moreover,  $N$  represents the total number of videos in all recipes and  $n_w$  is the number of videos whose subtitle include the word  $w$ . After computing the tf-idf, we sort all words with their tf-idf values and choose top  $K$  words as set of salient words.

#### 4.1.2 Learning Visual Atoms

In order to learn visual atoms, we initially generate object proposals from each frame of each videos.

#### Generating Framewise Object Proposals

#### Joint Proposal Clustering to Detect Visual Objects

#### 4.2. Multi-Modal Representation via Learned Atoms

#### 5. Unsupervised Activity Representation

Here we explain the generative model which we use in order to jointly learn the activities from videos. We start

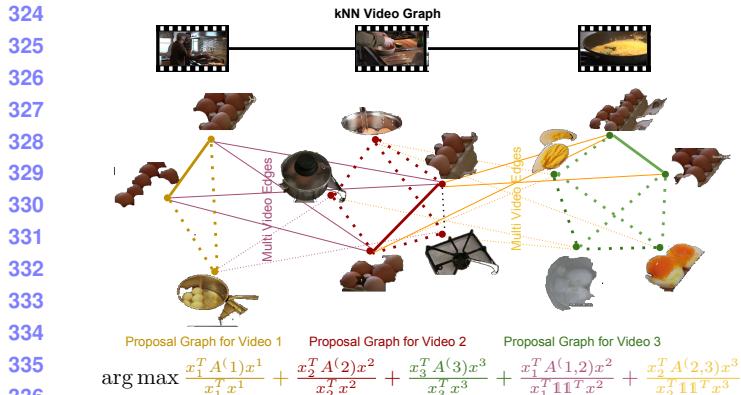
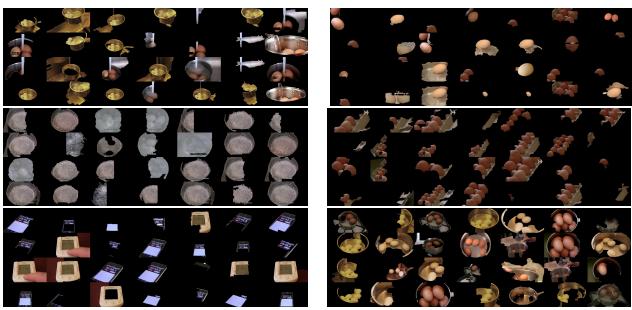


Figure 2: Visualization of the joint proposal clustering.

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

with explaining the basic notation we use in the model. We denote the  $t^{th}$  frame of the  $i^{th}$  video as  $I_t^i$  and its subtitle as  $L_t^i$ . 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 occurrence 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  $\Theta$  from its conjugate distribution *Beta distribution*.

In the following sections, we explain how these models

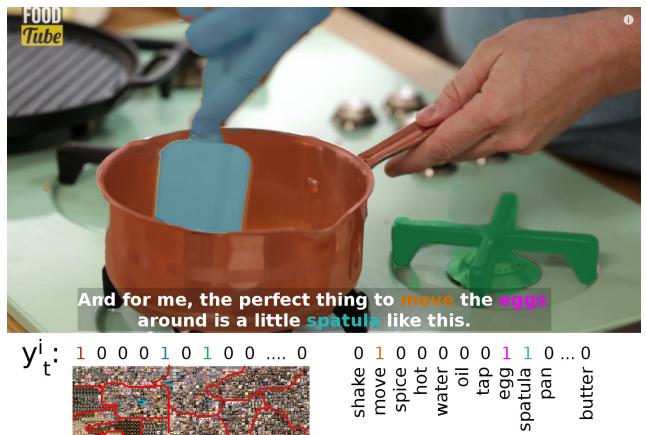


Figure 4: Visualziation of the representation of the Frame.

can be jointly learned and inferred by using the Beta Process Hidden Markov Models. s

## 5.1. 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  $f^i$  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  $F$  such that  $i^{th}$  row of the  $F$  is the activity vector  $f^i$ . Moreover, each feature  $k$  also has an activity frequency  $b_k$  and a distribution parameter  $\Theta_k$ . Distribution parameter  $\Theta_k$  is the Bernoulli distribution as explain in the Section 5. Moreover, its base distribution ( $B_0$ ) is *Beta random variable* since it is the conjugate of Bernoulli. Moreover, in this setting, the activity paremeters  $\Theta_k$  and  $b_k$  follow the *beta process* as;

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

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  $\Theta_k$ , this distribution is shown to be equivalent to Indian Buffet Process [?]. Where videos are customers and activities are dishes in the buffet. The

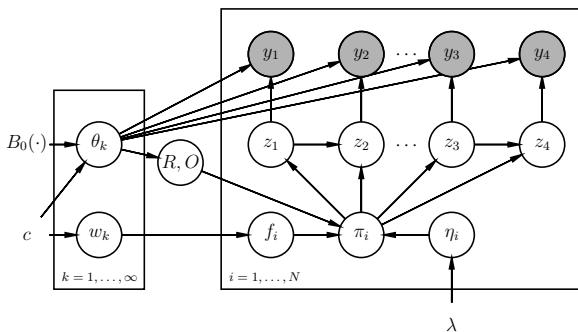
432 first video chooses a Poisson( $\gamma$ ) unique dishes. The fol-  
 433 lowing video  $i$  chooses previously sampled activity  $k$  with  
 434 probability  $\frac{m_k}{i}$  proportional to number of videos ( $m_k$ ) cho-  
 435 sen the activity  $k$ , and it also chooses Poisson( $\frac{\gamma}{i}$ ) new activi-  
 436 ties. Here,  $\gamma$  controls the number of active features in each  
 437 video and  $\beta$  controls the likelihood of the features getting  
 438 shared by multiple videos.  
 439

440 After each video chooses a subset of activities, we model  
 441 the videos as an Hidden Markov Model (HMM) over the  
 442 selected videos. Each frame has the hidden state activity  
 443 id( $z_t^i$ ) and we observe the binary representation  $y_t^i$ . Since  
 444 we model each activity as a Bernoulli distribution, the  
 445 emmission probabilities follow the Bernoulli distribution as  
 446  $p(y_t^i|z_t^i) = Ber(\Theta_{z_t^i})$ . Following the construction of the  
 447 Fox et al.[?], we sample the transition probabilities from a  
 448 normalized Gamma distribution. For each video  $i$ , we sam-  
 449 ple a Gamma random variable for the transition between ac-  
 450 tivity  $j$  and activity  $k$  if both of the activities are included by  
 451 the video (if  $f_k^i$  and  $f_j^i$  are both 1). After sampling these ran-  
 452 dom variables, we normalize them to have proper transition  
 453 probabilities. This procedure can be represented formally  
 454 as  
 455

$$\eta_{j,k}^i \sim Gam(\alpha + \kappa\delta_{j,k}, 1), \quad \pi_j^i = \frac{\eta_j^i \circ f^i}{\sum_k \eta_{j,k}^i f_k^i} \quad (2)$$

458 Where  $\kappa$  is the persistence parameter promoting self state  
 459 transitions to have more coherent temporal boundries,  $\circ$  is  
 460 the element-wise product and  $\pi_j^i$  is the transition probabili-  
 461 ties in video  $i$  from state  $j$  to all states in the form of a  
 462 vector.  
 463

464 This model is also presented as a graphical model in Figure 5  
 465



478 Figure 5: **Graphical model for BP-HMM:**Some explana-  
 479 tion.  
 480

## 482 5.2. Gibbs sampling for BP-HMM

484 We employ Markov Chain Monte Carlo (MCMC)  
 485 method for learning and inference of the BP-HMM. We

486 base our algorithms on the MCMC procedure proposed by  
 487 Fox et al.[?]. It marginilizes over blah and blah and sam-  
 488 ple blah and blah. For faster convergence, we also utilize  
 489 a series of data driven samplers. Here we only discuss the  
 490 proposed data driven samplers and move the details of the  
 491 remainin samplers to the Supplementary Material.  
 492

### 493 Data-Driven Sampler1

### 494 Data-Driven Sampler2

## 495 6. Experiments

### 496 6.1. Dataset

497 We collected blah blah videos from YouTube and did  
 498 blah blah... All numbers etc. These are recipes, 5 of them  
 499 are evaluation set and labelled temporally and labels are  
 500 matched.  
 501

### 502 6.2. Baselines

503 We compare against BP-HMM-O with local feauters and  
 504 HMM with Semantic Features, HMM with local features,  
 505 HMM with CNN feautes...  
 506

### 507 6.3. Qualitative Results

### 508 6.4. Accuracy over Activity Detection

### 509 6.5. Accuracy over Activity Learning

## 510 7. Discussions and Conclusions

511 Discuss which recipes worked and why. Discuss the im-  
 512 portance of semantic representation, scaling features and  
 513 multi-modality.  
 514

## 515 References

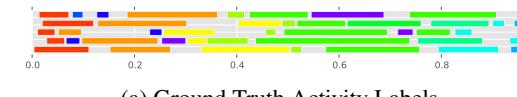
- [1] K. Barnard, P. Duygulu, D. Forsyth, N. De Freitas, D. M. Blei, and M. I. Jordan. Matching words and pictures. *The Journal of Machine Learning Research*, 3:1107–1135, 2003. 2
- [2] M. Beetz, U. Klank, I. Kresse, A. Maldonado, L. Mosenlechner, D. Pangercic, T. Ruhr, and M. Tenorth. Robotic roommates 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

540  
541

Table 1: Notation of the Paper

$y_t = [y_t^v, y_t^l]$	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	$\pi_i^{k,k'}$	$\eta_i^{k,k'} \times f_i^k \times f_i^{k'}$		

546



(b) Activity Labels extracted by Our Method



Crack the eggs one at a time into a bowl.



You can either use a fork or wire whisk to beat the eggs into a bowl.



Remove the omelette onto a plate.



Eggs cook quickly, so make sure the pan gets very hot first; the butter melt completely.

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

- [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 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

- 594
- 595
- 596
- 597
- 598
- 599
- 600
- 601
- 602
- 603
- 604
- 605
- 606
- 607
- 608
- 609
- 610
- 611
- 612
- 613
- 614
- 615
- 616
- 617
- 618
- 619
- 620
- 621
- 622
- 623
- 624
- 625
- 626
- 627
- 628
- 629
- 630
- 631
- 632
- 633
- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647
- [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. Rabivovich, 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 International 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*

- 648        eighth ACM international conference on Multimedia, pages                      702  
649        105–115. ACM, 2000. 1    703  
650        [24] R. Socher and L. Fei-Fei. Connecting modalities: Semi-                          704  
651        supervised segmentation and annotation of images using un-                              705  
652        aligned text corpora. In *Computer Vision and Pattern Recog-                              706  
653        nition (CVPR), 2010 IEEE Conference on*, pages 966–973.                                      707  
654        IEEE, 2010. 2    708  
655        [25] R. Socher, A. Karpathy, Q. V. Le, C. D. Manning, and A. Y.                      709  
656        Ng. Grounded compositional semantics for finding and de-                                      710  
657        scribing images with sentences. *Transactions of the Associa-                              711  
658        tion for Computational Linguistics*, 2:207–218, 2014. 2                              712  
659        [26] M. Tenorth, D. Nyga, and M. Beetz. Understanding and ex-                          713  
660        ecuting instructions for everyday manipulation tasks from the                              714  
661        world wide web. In *Robotics and Automation (ICRA), 2010                      715  
662        IEEE International Conference on*, pages 1486–1491. IEEE,                              716  
663        2010. 2    717  
664        [27] B. T. Truong and S. Venkatesh. Video abstraction: A sys-                          718  
665        tematic review and classification. *ACM Transactions on                              719  
666        Multimedia Computing, Communications, and Applications*,                      720  
667        3(1):3, 2007. 1    721  
668        [28] H. Yu and J. M. Siskind. Grounded language learning from                          722  
669        video described with sentences. In *ACL (1)*, pages 53–63,                              723  
670        2013. 2    724  
671        [29] C. L. Zitnick and D. Parikh. Bringing semantics into fo-                          725  
672        cus using visual abstraction. In *Computer Vision and Pat-                              726  
673        tern Recognition (CVPR), 2013 IEEE Conference on*, pages                              727  
674        3009–3016. IEEE, 2013. 2    728  
675        [30] C. L. Zitnick, D. Parikh, and L. Vanderwende. Learning                          729  
676        the visual interpretation of sentences. In *Computer Vi-                              730  
677        sion (ICCV), 2013 IEEE International Conference on*, pages                              731  
678        1681–1688. IEEE, 2013. 2    732  
679    733  
680    734  
681    735  
682    736  
683    737  
684    738  
685    739  
686    740  
687    741  
688    742  
689    743  
690    744  
691    745  
692    746  
693    747  
694    748  
695    749  
696    750  
697    751  
698    752  
699    753  
700    754  
701    755