

# Incorporating Object and People Information to Improve Video Activity Recognition

*A thesis submitted in partial fulfillment  
of the requirements for the degree of*

MASTER OF TECHNOLOGY

*in*

Computer Science & Engineering

*by*

Niranjan Viladkar

Entry No. 2012MCS2810

*Under the guidance of*

Dr. Parag Singla



Department of Computer Science and Engineering,  
Indian Institute of Technology Delhi.

June 2014.

# Certificate

This is to certify that the thesis titled **Incorporating Object and People Information to Improve Video Activity Recognition** being submitted by **Niranjan Viladkar** for the award of **Master of Technology in Computer Science & Engineering** is a record of bona fide work carried out by him under my guidance and supervision at the **Department of Computer Science & Engineering**. The work presented in this thesis has not been submitted elsewhere either in part or full, for the award of any other degree or diploma.

---

**Dr. Parag Singla**  
**Department of Computer Science and Engineering**  
**Indian Institute of Technology, Delhi**

# Acknowledgments

My sincere thanks to my project guide, Dr. Parag Singla. I am grateful to have his constant support and motivation during the difficulties that I faced at various stages of the project. He not only guided me towards the solutions to the problems but encouraged me to learn the concepts that gave me insights into the subject. I am thankful to Dr. Subhashis Banerjee, he gave me directions for the project in initial stages.

I thank Ruchin Kukreja and Arunim Samat for their help with acquiring datasets and initial project setup.

I would also like to thank Happy Mittal, a Ph.D scholar at the Department of Computer Science, IIT Delhi, who helped me understand Markov Logic Networks and use of Alchemy.

**Niranjan Viladkar**

# Abstract

Human Activity Recognition is an important vision problem. In recent years, there have been various approaches towards the problem [10, 9, 11, 13]. Most existing approaches perform classification of video clips based on low level features like HoG-HoF. This approach can not exploit the semantic relationship between activities and presence of various kinds of objects (and people) in the underlying domain.

For example, in a video clip showing an action of ‘eating’, even if the low level feature based detector has low confidence about this activity, presence of dining table and other objects such as bottle and chair can boost up the confidence and help make the correct prediction.

This project adds domain knowledge in the form of object and people information to the classification using Markov Logic. Markov Logic captures semantic relationship using weighted first order logic formulae. An experiment where object and people features are added to SVM features, shows significant improvement in the predictions.

This project gives an end to end system for activity recognition. It takes input from a video classifier and an object detector and builds semantic model on top of this information. Activity predictions using this approach shows improvement over the existing frameworks.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Classification based on Video Features . . . . .	3
2.2	Object Detection . . . . .	4
2.3	Markov Logic Networks . . . . .	5
<b>3</b>	<b>Related Work</b>	<b>7</b>
<b>4</b>	<b>Approach</b>	<b>9</b>
4.1	MLN Evidence . . . . .	10
4.2	MLN Rules . . . . .	11
4.3	MLN Learning and Inference . . . . .	12
<b>5</b>	<b>Experiments</b>	<b>14</b>
5.1	Dataset . . . . .	14
5.2	Methodology . . . . .	14
5.3	Results . . . . .	16
<b>6</b>	<b>Conclusion</b>	<b>19</b>
	<b>Bibliography</b>	<b>20</b>

# List of Figures

4.1	Approach of the Project . . . . .	9
5.1	Car detection in a frame . . . . .	15
5.2	Person detection in a frame . . . . .	15

# List of Tables

5.1	Output of object detector with decision values . . . . .	15
5.2	Thresholds on object detector confidence . . . . .	16
5.3	MLN Experiments - Average precisions . . . . .	17
5.4	Average Precision : Object and People Features in SVM . . .	18

# Chapter 1

## Introduction

Human activity recognition is an area of interest in the field of vision. For example, if one wants to extract only a certain type of activity sequence from YouTube videos, like eating or driving etc., a highly accurate and precise video recognition system would drastically reduce the need for human intervention. The work in [11] processes YouTube dataset for activity recognition problem.

Existing activity recognition systems [10, 9, 11] do not take advantage of underlying domain. Domain may have significant influence over the activity types, the objects appearing in the activity etc. For example, presence of a dining table improves the chances of activity eating as compared to driving. Or if, the video is occluded partially, but objects like bottle and chair are visible, then also, chances of activity eating are more as compared to other activities. This relationship between objects and activities can be captured by First Order Logic (FOL).

But merely using FOL is not sufficient. In its basic form, FOL uses hard constraints. To explain this with example, consider a rule

$$\forall clip \text{ } HasObject(clip, bottle) \implies HasActivity(clip, eating) \quad (1.1)$$

In real world scenario, it is perfectly possible that a video clip contains a bottle, but the activity might be different than eating. The rule gets violated in this scenario as it is universally quantified over all the clips. Rule may also get violated in case of noisy training data. Thus, using hard constraints is neither practical nor robust for real world video clips.

Solution to this is to attach weights with rules like (1.1) with weights being proportional to the real world relevance. This idea can be captured using Markov Logic Networks (MLNs). Markov Logic allows us to capture the



relationship between activities and objects using weighted first order logic formulae. This helps in building a better overall model for activity detection.

This project exploits the semantic relationship between activities and objects augmenting the information captured by low level features. Aim of this project is to provide end to end recognition system to improve the existing human activity recognition systems by adding domain knowledge i.e. object and people information to them.

Chapter 2 gives the theoretical background. It explains video classification using only video features. It explains how to extract and use the video features. It further describes object detection in detail and also describes the theory of Markov Logic Networks. Chapter 3 throws light on previous approaches and points out potential improvement areas. Chapter 4 explains design decisions and approach of this project. Chapter 5 shows the result comparisons. Finally thesis concludes with chapter 6 which notes conclusion and future work.

# Chapter 2

## Background

Video clips representing human actions can be classified based on features extracted from the video. Object detection can add the object and people information to the recognition process.

This chapter explains three components in this project. Classifiers based on video features, object detectors and gives a brief introduction to Markov Logic Networks.

Classifiers based on video features use low level features like HoG-HoF features. Using clustering of these features, each clip is represented in bag-of-features representation. Finally, SVM [3] is used to learn a model and classify the testing instances. Object detectors run on frames of videos and detect pre-defined objects along with a confidence value for the detections. Such detections are useful to add object and people information to the existing recognition system. Markov logic allows us to capture the semantic relationship using weighted first order formulae.

### 2.1 Classification based on Video Features

As explained in [8] spatial temporal interest points (STIPs) are 3 dimensional points in space and time where the local video features show significant variations. In [10] the features extracted at STIPs are shown to perform satisfactorily for human activity recognition. At each STIP, histograms of oriented gradient (HoG) and histograms of optical flow (HoF) are evaluated. The number of HoG and HoF features at each STIP are 72 and 90 respectively. Concatenating vectors of HoG and HoF features, we get a single 162 element descriptor for each STIP.

After extraction of STIP features, each video clip roughly has of the order of 1000 such descriptors. Some number of descriptors are sampled randomly

from all the descriptors across all the clips. These random descriptors are clustered into  $k$  clusters using k-means.

Each descriptor in each clip is clustered into one of the  $k$  clusters using least Euclidean distance. Thus each clip is represented by a  $k$  sized vector where  $i^{th}$  element in this vector represents number of descriptors of that clip nearest to the  $i^{th}$  cluster. This vector is called a Bag of Features (BoF) representation of the clip.

The BoF representation of all training clips are used to train a support vector machine (SVM) [3] and BoF representation of a disjoint set of testing clips is classified using learnt model.

## 2.2 Object Detection

Video clips of certain activity class have peculiar set of objects present in it. If one could find objects present in the video clip, the activity prediction can potentially be made more confident. Below subsections introduce the object detection method and its use in this project.

Object detector based on Discriminatively Trained Deformable Part Models [5] was used in the project. It can run on an image at a time. The detectors used in this project are trained on PASCAL 2007 [4] datasets. There are 20 models corresponding to 20 different objects. These objects are : aeroplane, bicycle, bird, boat, bottle, bus, car, cat, chair, cow, dining table, dog, horse, motorbike, person, potted plant, sheep, sofa, train, tv monitor. For this project only 10 relevant models were used : bicycle, bottle, bus, car, chair, diningtable, motorbike, person, sofa, tvmonitor.

Above described object detector models work on a single image at a time. To detect objects in video clips, 1 frame per second was extracted from the video clip and all the object detector models were run on the frame. In the dataset used in this project, shortest clips is about 3-4 seconds long and longest clips are of the order of a minute. Frame rate is 24 frames per second.

The object detector models give a bounding box and a confidence (also called decision value) for each detection on an absolute scale between  $-\infty$  to  $\infty$ .

Usually, a positive decision value represents true positive in all the models. A negative decision value represent lesser confident detection. It might be a true positive or a false positive. One can specify a threshold decision value before running object detector model so that only the detections at least as confident as the threshold are considered. Usually, such thresholds are negative to allow detection of lesser confident objects also.

## 2.3 Markov Logic Networks

As explained in [12], Markov Logic consists of set of First order logic (FOL) formulae associated with weights. The FOL formulae form the structure of Markov Logic Networks (MLNs).

In its basic form, first order knowledge base is a set of hard constraints over a set of possible worlds. Because of these hard constraints, if a world does not satisfy even a single formula, the world is considered as false or impossible.

Markov Logic uses FOL with soft constraints instead of hard. In markov logic, if a world does not satisfy a formula  $F_i$ , it is considered as less probable world in comparison to a world which satisfies  $F_i$  assuming rest of the formulae have same state in both the worlds. Together with set of constants, Markov Logic defines a ground Markov Network distribution.

Definition 2.3.1 from [12] formally defines Markov Logic Network (MLN) as follows

**Definition 2.3.1.** *A Markov Logic Network(MLN)  $L$  is a set of pairs  $(F_i, w_i)$ , where  $F_i$  is a formula in first order logic and  $w_i$  is weight associated with the formula - a real number. Together with a finite set of constants  $C = \{c_1, c_2, \dots, c_{|C|}\}$ , MLN  $L$  defines a markov network,  $M_{L,C}$  as follows :*

1.  $M_{L,C}$  contains one binary node for each possible grounding of each atom appearing in  $L$ . The value of the node is 1 if ground atom is true and 0 otherwise.
2.  $M_{L,C}$  contains one feature for each possible grounding of each formula  $F_i$  in  $L$ . The value of the feature is 1 if the ground formula is true and 0 otherwise. The weight of the feature is  $w_i$  associated with  $F_i$  in  $L$ .

Probability distribution over possible world  $x$  specified by markov network  $M_{L,C}$  is given by

$$P(X = x) = \frac{1}{Z} \exp \left( \sum_{i=1}^F w_i n_i(x) \right) \quad (2.1)$$

where  $F$  is number of formulae in MLN and  $n_i(x)$  is the number of true groundings of  $F_i$  in world  $x$ .

# Chapter 3

## Related Work

In [9], authors have come up with a method to annotate video clips from Hollywood movies using their text scripts. Authors present a method for classification using local space-time features, space-time pyramids and multi-channel non-linear SVM.

Their approach is based on low level space time features such as HoG-HoF features. All video clips are represented in Bag-of-Features representation. These Bag-of-Features vectors are used to train a SVM model. Note that as this method uses low level features, it can not exploit the rich semantic information relating activities with presence of objects and people.

In [10], authors add context in the form of scene information to the classification problem. As quoted by authors, an example is, activity eating is more probable to happen in kitchen while activity running is likely to happen outdoor. Thus scene contexts of these two activities will help in classification.

There approach also consists of forming a Bag-Of-Features representation of video clips using low level features like HoG-HoF. Thus suffering similar problems as explained above. The scene information is retrieved using text scripts. Thus this information is limited by the text available and may not be able to capture the actual relationship between scene and activity as seen in the video.

In [11], authors have used object information present in the video clip. The correlation between activity and object has been found using analysis of large amount of text. Note that this correlation is not found directly via video but via text. Thus it has the limitation that only the correlations appearing in the text can be captured. Also, a naive Bayes assumption is made which assume that video features and object features are independent of each other given

the activity class. Finally an integrated classifier is formed which combines classification based on methods from section 2.1 and probability of activity deduced from object information.

In [13], Markov Logic Networks are used to capture the interaction between humans and vehicles, such as, open door, drive away etc. But their model is limited to human vehicle interaction. For real world videos, a more general system is required.

Thus, it can be seen that most of the previous approaches don't capture semantic relationship between activities and objects. Next chapter explains how this project combines object and people information with existing activity recognition to improve the recognition results.

# Chapter 4

## Approach

This project aims at providing end to end system which takes inputs from video activity recognizer and object detector to enhance the video activity recognition.

Figure 4.1 explains the approach of this project.

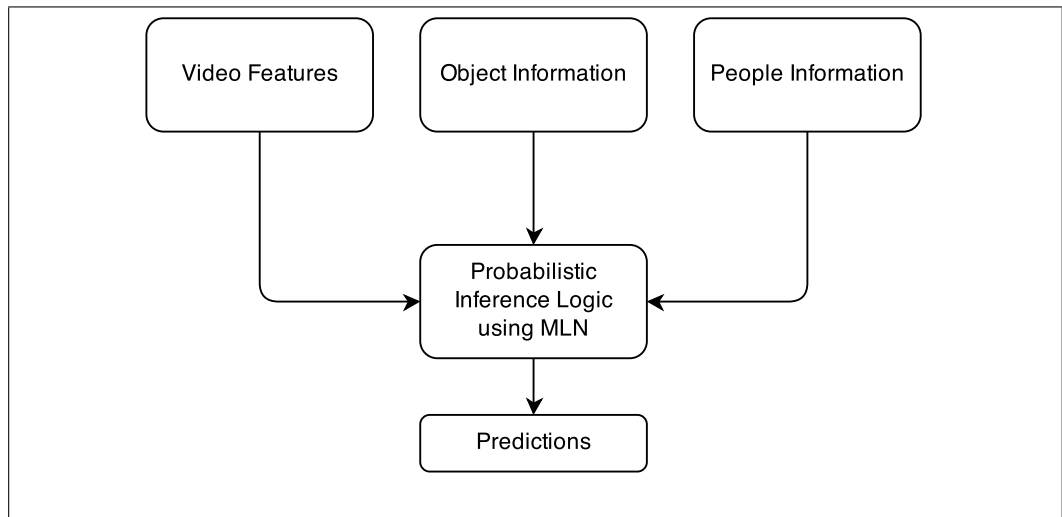


Figure 4.1: Approach of the Project

The video features are extracted as explained in section 2.1. The code used in this phase is publicly made available at [6]. Using one versus all (ovr) matlab interface of libsvm [2], model is learnt to extract the confidence values from video feature classification.

The object features are evaluated using [5]. Each object detection has a confidence value along with it. This value is used to determine whether to consider detection as a viable detection or spurious one. The thresholds for confidence values are found experimentally. Details are given in chapter 5.

Main contribution of this project is in the formulation of markov logic system. Alchemy [1] software is used to form markov logic network. The software



allows convenient representation of rules using first order logic and assigns weights to such rules using the evidence provided. Following sections explain in detail about the evidence creation, rules, etc.

## 4.1 MLN Evidence

The Hollywood2 dataset [7] (described in chapter 5) provides true activity labels for training data. These labels are used to create predicates like

$$HasActivity("actioncliptrain00001", "SitUp") \quad (4.1)$$

### 4.1.1 Evidence from Video Recognizer

The learnt SVM model is used to classify training dataset itself to get confidence values of activities for all clips. These confidence values are partitioned into bins to create predicates like

$$\begin{aligned} ActivityConf\_P1\_TO\_P15("actioncliptrain00001", "SitUp") \\ ActivityConf\_N1\_TO\_N2("actioncliptrain00002", "HandShake") \end{aligned} \quad (4.2)$$

Here, P1\_TO\_P15 represents a bin with confidence value of +1 to +1.5. Similarly other bins are formed and predicates are generated to form the evidence. Each clip has 12 such predicates as there are 12 activity classes.

### 4.1.2 Evidence from Object Detector

Object detector [5] gives object detections in frames along with the corresponding confidence values. Experimentally determined thresholds are used to decide whether to consider object detection and predicates like following are formed

$$\begin{aligned} ObjPresent("actioncliptrain00001", "person") \\ ObjPresent("actioncliptrain00002", "car") \end{aligned} \quad (4.3)$$

Finally, to generate evidence for people information, average number of “person” objects are calculated in each clip. These averages are partitioned into bins and each clip is assigned one bin. Predicates like following are formed using this method

$$\begin{aligned} &NumPersons\_1\_TO\_15(“actioncliptrain00001”) \\ &NumPersons\_0\_TO\_1(“actioncliptrain00002”) \end{aligned} \quad (4.4)$$

Now, in order to find the semantic relationship between activities and object & people information, the MLN model needs to be learnt using MLN evidence and rules.

## 4.2 MLN Rules

For learning MLN model, rules corresponding to bins of confidence from **video classifier** are written. For example,

$$\begin{aligned} &ActivityConf\_N1\_TO\_N05(clip, activity) \implies HasActivity(clip, activity) \\ &ActivityConf\_P1\_TO\_P15(clip, activity) \implies HasActivity(clip, activity) \end{aligned} \quad (4.5)$$

Here, the 1<sup>st</sup> rule captures relation between low confidence (−1 to −0.5) of presence of activity and actual evidence. And 2<sup>nd</sup> rule captures relation between higher confidence (+1 to +1.5) of presence of activity and actual evidence. Thus intuitively, 1<sup>st</sup> rule will get lower weight as compared to 2<sup>nd</sup> rule.

For adding **object information**, both positive and negative rules are written. Positive rules are intuitive rules which enhance the probability of likely inference. Negative rules are counter intuitive rules which suppress the probability of unlikely rules. In both the types, the confidence of classification

improves. Below is example of few positive rules

$$\begin{aligned} \text{ObjPresent}(c, \text{"chair"}) &\implies \text{HasActivity}(c, \text{"Eat"}) \\ \text{ObjPresent}(c, \text{"car"}) &\implies \text{HasActivity}(c, \text{"DriveCar"}) \end{aligned} \quad (4.6)$$

and few negative rules

$$\begin{aligned} \text{ObjPresent}(c, \text{"bus"}) &\implies \text{HasActivity}(c, \text{"StandUp"}) \\ \text{ObjPresent}(c, \text{"car"}) &\implies \text{HasActivity}(c, \text{"HandShake"}) \end{aligned} \quad (4.7)$$

For learning relationship between **people information** and activities, rules consisting of bins of number of people and activities are written. For example,

$$\begin{aligned} \text{NumPersons}_{0\_TO\_1}(\text{clip}) &\implies \text{HasActivity}(\text{clip}, +a) \\ \text{NumPersons}_{1\_TO\_15}(\text{clip}) &\implies \text{HasActivity}(\text{clip}, +a) \end{aligned} \quad (4.8)$$

where, 1<sup>st</sup> rule has a bin that corresponds to average number of people per frame between 0 to 1. “+a” means rule will be expanded to all possible activities. Thus for each bin, rules corresponding to all 12 activities will be used in MLN learning.

## 4.3 MLN Learning and Inference

### 4.3.1 Learning

The evidence generated is used to learn a MLN model according to the rules provided. It learns weights for the rules mentioned in the previous section. The query predicate given is “HasActivity”. The positive rules, which capture intuitive relationship between activity and object & people information, are expected to get higher weights as compared to negative rules, which state counter intuitive relationship.

### 4.3.2 Inference

The learnt MLN model, i.e. a set weighted first order logic formulae, is used with evidence to predict activity classes for test dataset. Inference is made on predicate “HasActivity”. Thus only difference between evidence and inference is absence of “HasActivity” predicates. MLN Inference process provides with the probabilities of all activities for each clip.

**AAP Calculations** - Using the probabilities given by MLN inference, inferred activity for each clip is determined corresponding to maximum probability. The dataset provides with true activity labels for test clips. These true activities and MLN inferred activities are used to calculate average average precision of MLN system. Next chapter presents the experimental results.

# Chapter 5

## Experiments

### 5.1 Dataset

All experiments in this project are performed on standard Hollywood2 (actions) [7] dataset from [10]. It has 823 labeled training video clips and 884 labeled testing video clips. There are 12 activity classes : AnswerPhone, DriveCar, Eat, FightPerson, GetOutCar, HandShake, HugPerson, Kiss, Run, SitDown, SitUp, StandUp. Each clip is 10 to 30 seconds long. Frame rate is 24 fps. Each clip is provided with a true label. Few of the clips have multiple labels, in this project, only single labeled clips are used.

### 5.2 Methodology

**Video recognizer setup** - 100,000 HoG-HoF feature vectors are randomly sampled from set of all feature vectors extracted from video clips using [6]. While doing random sampling, each clip is given equal chance - this avoids biased sampling towards (or against) any particular class of videos if they happen to be longer (or shorter) than other videos on an average. These randomly sampled descriptors are clustered in  $k = 200$  clusters using k-means. All clips are represented in Bag-of-Features representation over these clusters as explained in section 2.1.

**SVM** - Matlab interface of libsvm [2] is used to train 12 one-vs-rest (ovr) models corresponding to 12 activity classes. A RBF kernel is used with parameters,  $\gamma = 0.01$  and  $C = 100$ .

**Object Detector Setup** - Object detector explained in section 2.2 is

used. This run of object detection has a threshold of -0.9. Object specific thresholding is done at a later stage while forming evidence for MLN.

The visual output of object detector looks as shown in figures 5.1 and 5.2. The boxes are called as bounding boxes. For each box, a confidence value is also evaluated. Table 5.1 shows sample confidence output.

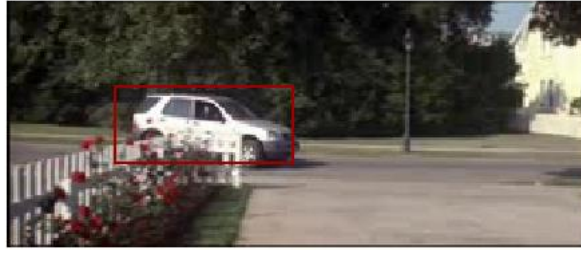


Figure 5.1: Car detection in a frame. *Source:[7]*



Figure 5.2: Person detection in a frame. *Source:[7]*

FRAME1	
Car	-0.181786
⋮	⋮
⋮	⋮
FRAME151	
Person	0.579786
Person	-0.593087

Table 5.1: Output of object detector with decision values

Table 5.2 shows the thresholds used to generate predicates for evidence and inference of MLN.

Object	Threshold
Bicycle	-0.90
Bottle	-0.83
Bus	-0.83
Car	-0.20
Chair	-0.77
Dining Table	-0.80
Motor Bike	-0.90
Person	-0.70
Sofa	-0.87
TV Monitor	-0.80

Table 5.2: Thresholds on object detector confidence

**Alchemy Setup** - MLN model is learnt using Alchemy 2.0 [1] software. All the parameters are kept default. The learning algorithm used is the default Alchemy learning algorithm.

**Classification metric** - Predictions are compared by taking average average precision as a measure. As explained in [10], average precision(AP) approximates area under recall-precision curve. Thus, AP for each activity class is calculated and then finally averaged over all the classes, average AP (AAP) is calculated. AP is taken as a measure of performance in the work of [10].

## 5.3 Results

**MLN** - Using alchemy 2.0 [1] software, various experiments were performed to add the object and people information to the existing video recognizer.

Table 5.3 shows comparison of APs for various strategies using MLNs. As can be seen, Average AP, when only action is considered for evidence and inference, is 28.53%. This number gradually improves as more information is brought into the system. Action information when added with Object information shows Average AP to be 30.30%. Action and People information

Activity Class	MLN			
	Only Action	Action & Object	Action & People	Action Object & People
AnswerPhone	10.64%	11.11%	11.67%	12.73%
DriveCar	66.06%	66.67%	71.57%	68.18%
Eat	32.50%	40.00%	35.00%	40.00%
FightPerson	56.90%	54.84%	61.54%	62.26%
GetOutCar	8.00%	13.79%	17.39%	14.29%
HandShake	21.43%	25.00%	30.77%	41.67%
HugPerson	15.79%	13.79%	14.29%	16.13%
Kiss	18.07%	19.78%	19.79%	20.65%
Run	36.42%	41.48%	40.32%	42.15%
SitDown	38.10%	35.56%	34.78%	39.56%
SitUp	0.00%	5.26%	0.00%	12.50%
StandUp	38.46%	36.29%	38.26%	36.24%
AAP	28.53%	30.30%	31.28%	33.86%

Table 5.3: MLN Experiments - Average precisions

shows Average AP to be 31.28%. And finally, adding both Object and People information to action information, Average AP of 33.86% is achieved.

**Object and People Features in SVM** - The Bag-of-Features vector of each clip can be appended with Object and People features. Term Frequency - Inverse Document Frequency (tf-idf) counts of objects are taken as object features. One feature corresponding to one object. Average number of people per frame are taken as people features. Thus adding 11 more features to each feature vector. These new feature vectors are now used to learn SVM model with similar parameters and methodology explained in section 5.2. It shows significant improvement as described in table 5.4.

Advantage of **MLN approach over SVM** approach is that, MLNs can support feed-backs. A high confidence of an activity can increase probability of object being present in that particular clip. This kind of feedback information cannot be captured using Object and People SVM features.



Activity Class	SVM Basic	SVM Basic + Object	SVM Basic + People	SVM Basic + People + Object
AnswerPhone	11.36%	11.63%	16.67%	16.67%
DriveCar	66.96%	66.69%	70.09%	70.09%
Eat	45.45%	44.12%	35.00%	50.00%
FightPerson	57.63%	57.63%	60.00%	66.04%
GetOutCar	17.86%	17.86%	10.34%	12.12%
HandShake	25.93%	25.93%	36.36%	31.82%
HugPerson	15.15%	15.62%	17.86%	17.86%
Kiss	18.18%	18.18%	17.65%	20.69%
Run	38.78%	38.51%	36.42%	39.35%
SitDown	40.96%	40.96%	42.68%	42.05%
SitUp	5.26%	5.26%	4.55%	8.70%
StandUp	35.20%	35.20%	37.21%	37.88%
Average AP	31.56%	31.49%	31.64%	34.44%

Table 5.4: Average Precision : Object and People  
Features in SVM

# Chapter 6

## Conclusion

Various experiments aimed at improving video activity recognition were performed during this project. Approaches included experiments using MLNs and adding object & people information to SVM learning. We provide a seamless method to combine outputs of multiple noisy detectors to enhance the recognition process. It was found that adding object and people information to the existing recognition framework improves the precision results.

The system developed by this project builds on top of video recognizer and object detector. Higher the precision of both of these systems, higher the precision of overall system.

Future work can target to add more semantic information to the model such as scene information. An important future contribution could be to complete feed-back loop to improve the quality of the object detector. For instance, if dining table is occluded in a video of ‘eating’ activity, the feedback from the activity recognizer can be used to boost up the confidence and predict the presence of dining table.

# Bibliography

- [1] Alchemy-2 : Inference and Learning in Markov Logic. <http://code.google.com/p/alchemy-2/downloads/list>.
- [2] Chih-Chung Chang and Chih-Jen Lin. LIBSVM – A Library for Support Vector Machines. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.
- [3] Corinna Cortes and Vladimir Vapnik. Support-Vector Networks. In *Machine Learning*, pages 273–297, 1995.
- [4] Everingham, M. and Van Gool, L. and Williams, C. K. I. and Winn, J. and Zisserman, A. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. <http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html>.
- [5] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester. Discriminatively Trained Deformable Part Models, Release 4. <http://cs.brown.edu/~pff/latent-release4/>.
- [6] Ivan Laptev. HoG-HoF Feature Extraction at Space Time Interest Points Code. <http://www.di.ens.fr/~laptev/download/stip-2.0-linux.zip>.
- [7] Ivan Laptev. Hollywood 2 Dataset (Action). <http://www.di.ens.fr/~laptev/actions/hollywood2/>.
- [8] Ivan Laptev. On Space-Time Interest Points. *IJCV*, 64(2/3):107–123, 2005.
- [9] Ivan Laptev, Marcin Marszałek, Cordelia Schmid, and Benjamin Rozenfeld. Learning Realistic Human Actions from Movies. In *Conference on Computer Vision & Pattern Recognition*, Jun 2008.
- [10] Marcin Marszałek, Ivan Laptev, and Cordelia Schmid. Actions in Context. In *IEEE Conference on Computer Vision & Pattern Recognition*, 2009.

- 
- [11] Tanvi S. Motwani and Raymond J. Mooney. Improving Video Activity Recognition using Object Recognition and Text Mining . In *Proceedings of the 20th European Conference on Artificial Intelligence (ECAI-2012)*, pages 600–605, August 2012.
  - [12] Matthew Richardson and Pedro Domingos. Markov Logic Networks. *Mach. Learn.*, 62(1-2):107–136, February 2006.
  - [13] Son D. Tran and Larry S. Davis. Event Modeling and Recognition Using Markov Logic Networks. In *Proceedings of the 10th European Conference on Computer Vision: Part II*, ECCV '08, pages 610–623, Berlin, Heidelberg, 2008. Springer-Verlag.