

---

# Vehicle Classification on Low-resolution and Occluded images: A low-cost labeled dataset for augmentation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Video image processing of traffic camera feeds is useful for counting and classifying  
2 ing vehicles, estimating queue length, traffic speed and also for tracking individual  
3 vehicles. Even after over three decades of research, challenges remain. Vehicle  
4 detection is especially challenging when vehicles are occluded which is common  
5 in heterogeneous traffic. Recently *Deep Learning* has shown remarkable promise  
6 in solving many computer vision tasks such as object recognition, detection, and  
7 tracking. We explore the promise of deep learning for vehicle detection and classification.  
8 However, training deep learning architectures require huge labeled datasets  
9 which are time-consuming and expensive to acquire. We circumvent this problem  
10 by data augmentation. In particular, we show that by properly augmenting an existing  
11 large general (non-traffic) dataset with a small low-resolution heterogeneous  
12 traffic dataset (that we collected) we can obtain state-of-the-art vehicle detection  
13 performance. This result is expected to further encourage the wide-spread use of  
14 deep learning for traffic video image processing.

## 15 1 Introduction

16 Traffic cameras play a crucial role in Intelligent Transport Systems. They can be used for counting  
17 vehicles, estimating queue length, traffic speed, and also for classifying and tracking individual  
18 vehicles. Here, we focus on the task of detecting and classifying vehicles from frames acquired from  
19 a traffic video stream.

20 Even after over three decades of research in the field, challenges remain. Vehicle detection is  
21 especially challenging when vehicles are occluded which is commonly observed in heterogeneous  
22 traffic. In heterogeneous traffic, size and type of vehicles vary significantly and vehicular traffic  
23 density is high which leads to frequent occlusion. Another issue that adds to the challenge is the low  
24 quality of the traffic camera feeds and lack of standardization of cameras and camera positions.

25 Traditionally, in the computer vision community, object detection is done in three steps: a sliding  
26 window phase where we search for the object at various scale and positions, followed by feature  
27 extraction at each window and finally classifying each window as either containing or not containing  
28 the desired object [3]. Commonly used features for object detection are histogram of oriented  
29 gradients (HoG) [3], scale-invariant feature transform (SIFT) [12], and speeded up robust features  
30 (SURF) [1]. This is usually followed by Support Vector Machine (SVM) based classification.

31 Recently, deep learning based approaches have shown extraordinary performance in many computer  
32 vision tasks such as object recognition [4], detection [16] [15], tracking [18], and image segmentation  
33 [10]. For certain tasks such as object recognition [4] and face recognition [11] deep learning has out-  
34 performed humans. The main reason behind its superior performance is, unlike traditional methods

35 which use hand-engineered features such as HoG, SIFT, and SURF, deep networks automatically  
36 learn discriminative features from the training data directly.

37 In this paper, we explore the promise of deep learning for doing vehicle detection in the challenging  
38 context of heterogeneous traffic that contains significant fraction of occluded and truncated images of  
39 vehicles. Though deep learning approaches have shown state-of-the-art results for object detection,  
40 they need to be trained on huge datasets such as Imagenet [4] which has millions of images. This is  
41 because the network itself has millions of parameters to learn. However, it is very time-consuming  
42 and expensive to collect such large labeled dataset of heterogeneous traffic. The main bottleneck is  
43 the task of labeling which is required for training the deep networks. For labeling, bounding boxes  
44 need to be manually drawn around all the vehicles present in any given frame and the vehicles need  
45 to be labeled into different classes. Thus, instead of collecting a large labeled dataset for our task,  
46 we propose to use clever data augmentation techniques. We show that by augmenting a large but  
47 general (non-traffic) dataset with a small labeled traffic dataset and by training a deep network on this  
48 augmented dataset, we easily out-perform traditional approaches for vehicle detection and vehicle  
49 classification.

50 We collected a dataset of 1417 images from traffic cameras installed in the city of Chennai, India.  
51 This is a very small dataset to train a deep network. Thus, we have augmented the PASCAL VOC  
52 dataset [5] with our heterogeneous traffic dataset. The PASCAL VOC dataset has around 10000  
53 images of 20 different classes including cats, dogs, trains, bottles, person along with few relevant  
54 classes such as car, truck, and bus. It is interesting to note that though PASCAL VOC has only  
55 a few relevant classes, still by augmenting it with our traffic dataset, we outperform a traditional  
56 approach of applying SIFT/SURF features followed by SVM classification. Though the proposed  
57 data augmentation can work with any deep network architecture for object detection, we have shown  
58 our results on Faster RCNN [16] which is a popular deep learning architecture.

59 Our specific contributions are as follows: (i) We are providing a labeled dataset for vehicle detection  
60 in heterogeneous traffic with significant occurrence of occlusion; (ii) We implement an extended  
61 deep learning architecture for the task of vehicle detection and classification in heterogeneous traffic  
62 scenario; (iii) We achieve high accuracy levels with limited data; and (iv) We demonstrate the superior  
63 performance of developed algorithm compared to a traditional object classification technique.

## 64 **2 Related Work**

65 Computer vision based methods for analyzing traffic systems are gaining in popularity. Vehicle  
66 detection and vehicle tracking have tremendously benefited from the advancements in computer vision  
67 techniques. Earlier work in vehicle detection are based on motion based algorithms (background  
68 subtraction [17], optical flow [7]) to detect vehicles and then use support vector machines [2] on the  
69 detection to classify them. [13] is one such method where authors have proposed to define a grid  
70 structure over the road in order to detect vehicles in heterogeneous traffic. These approaches are not  
71 robust with respect to illumination, occlusions, and scale changes [7] [17]. Also, the SVM classifier  
72 is heavily dependent on hand crafted features such as SURF [1] and SIFT [12].

73 Recently proposed deep learning models are free from these disadvantages. The most important  
74 feature of a deep learning model is: they identify useful features automatically which are quite  
75 robust to illumination and scale changes given enough training data. Authors proposed region based  
76 networks [16] [10] [9] [8] in which a network identifies possible object proposals and then a classifier  
77 classifies them. There are few studies which proposed object detection as an end to end regression  
78 problem [14] [15] [11]. All the deep learning models have been trained on huge datasets [4] which  
79 allows them to generalize well for a given task. Our method is based on one such deep learning  
80 model: Faster R-CNN (Region-based Convolutional Neural Networks) [16].

## 81 **3 Methodology**

82 Deep learning models have a large number of parameters to be tuned, which require a large number of  
83 labeled data samples. For example, in the Faster RCNN architecture, the feature extraction network  
84 (VGG-16) needs millions of high-quality images for tuning the parameters. VGG-16 is trained on  
85 Imagenet dataset [4]. The other components of Faster RCNN, region proposal network and fully



(a)

Classes	Total Samples	Occluded Samples
Light Motor Vehicles	2746	848
Heavy Motor Vehicles	279	157
Two Wheelers	3294	568

Dataset-1

Classes	Total Samples	Occluded Samples
Auto Rickshaw	598	219
Car	2148	629
Heavy Motor Vehicles	279	157
Two Wheelers	3294	568

Dataset-2

(b)

Figure 1: (a) Proposed data augmentation approach. This figure shows addition of a new class (Auto Rickshaw) from low resolution small dataset in high resolution standard dataset, (b) Statistics of our collected dataset. We have created two datasets; first one has one less class because of merging *auto-rickshaw* and *car* classes.

86 connected layers also require carefully annotated large datasets of images for their training. Getting  
 87 such a huge labeled dataset representing every object class is very expensive and time-consuming.

88 Fine-tuning a pre-trained deep neural network is a standard practice in computer vision community.  
 89 We have shown that doing finetuning with such a small task specific dataset performs poorly.

90 We test four approaches sequentially. First, the pre-trained Faster RCNN model is directly applied to  
 91 our dataset. Second, the pre-trained model is fine-tuned with data from our dataset. Third, the model  
 92 is trained from scratch using the collected data only. Finally, the model is trained with the existing  
 93 large dataset and our collected dataset.

94 Our dataset is quite different from the Pascal VOC dataset on which the Faster RCNN model has been  
 95 trained. In Figure 2 we have shown few of the sampled images. Pascal VOC images are high-quality  
 96 images captured using high-resolution cameras whereas images in our dataset are collected from  
 97 traffic surveillance camera feeds. Pascal VOC images contain fewer object instances per image  
 98 compared to our dataset. One more major difference is that PASCAL VOC dataset has 20 different  
 99 categories which are largely diverse. However, our dataset contains only vehicles and has different  
 100 sub-categories of vehicles and vehicle classes. Due to all these differences, we can not directly deploy  
 101 the existing, or even a fine-tuned, Faster RCNN model to our dataset.

102 Finally, we augmented the Pascal VOC dataset with our dataset. There are few vehicle classes in our  
 103 dataset that are not present in Pascal VOC dataset such as auto rickshaw. We can also perform data  
 104 augmentation in such a case as shown in Figure 2(b). Resultant augmented data will contain all the  
 105 20 classes of Pascal VOC and one additional class, Auto Rickshaw. While applying the model, we  
 106 can ignore the irrelevant classes. This is because the model is benefiting from the high-quality images  
 107 of Pascal VOC and also optimizing the loss according to our dataset. This way of augmenting the  
 108 dataset with our specific dataset is leading to improved learning of parameters in the model as shown  
 109 in the results section.

## 110 4 Dataset Collection

111 We generated our own dataset from cameras monitoring road traffic in Chennai, India. To ensure that  
 112 data are temporally uncorrelated, we have sampled frames at 0.5 fps from multiple video streams.  
 113 We extracted 2400 frames in total.



Figure 2: First row shows few images from Pascal VOC dataset [5]. Second row shows few images from our dataset. From these set of images it is clear that Paccal VOC images are of higher quality compared to the images of our dataset.

Table 1: Object detection results on Faster RCNN architecture using different ways of training (AP @ 0.5).

	Model	TW	HMV	LMV
(i)	Pre-trained Model	0.256	0.273	0.600
(ii)	Pre-trained Model + Fine-tuning	0.114	0.043	0.163
(iii)	Training Only on Our Dataset	0.082	0.004	0.055
(iv)	Augmented Data Training	<b>0.887</b>	<b>0.968</b>	<b>0.905</b>

114 We manually labeled 2400 frames under different vehicle categories. The number of available frames  
 115 reduced to 1417 after careful scrutiny and elimination of unclear images. We initially defined eight  
 116 different vehicle classes commonly seen in Indian traffic. Few of these classes were similar while  
 117 two classes had less number of labeled instances; these were merged into similar looking classes. For  
 118 example, in our dataset, we had different categories for small car, SUV, and sedan which were merged  
 119 under the light motor vehicle (LMV) category. Figure 2(b) shows brief statistics of our dataset.

120 A total of 6319 labeled vehicles are available in the collected dataset (see figure 2(b)). This includes  
 121 3294 two-wheelers, 279 heavy motor vehicles (HMV), 2148 cars, and 598 auto-rickshaws. A second  
 122 dataset was created by merging cars and auto-rickshaws together into light motor vehicle (LMV)  
 123 class. Approximately 25.2% of vehicles were occluded.

124 We have released the heterogeneous traffic dataset that we collected<sup>1</sup> for public use.

## 125 5 Experimental Results

126 In this section, we show the results of proposed data augmentation approach and performance obtained  
 127 by extending faster RCNN model for new classes. The results of data augmentation are compared  
 128 with the performance of four different ways of training Faster RCNN on our dataset: (i) training from  
 129 scratch using collected dataset alone, (ii) fine-tuning the pre-trained model with collected dataset,  
 130 and (iii) using pre-trained model directly, and (iv) model trained from scratch using augmented  
 131 dataset. Performance of extended Faster RCNN model is compared with three different ways of  
 132 training: (i) using pre-trained Faster RCNN model for object proposals alone and then using SVM for

<sup>1</sup>[https://www.dropbox.com/s/j1gr0d4w8u57jfv/dataset\\_vehicle\\_detection\\_ilds\\_iitm.tar.gz?dl=0](https://www.dropbox.com/s/j1gr0d4w8u57jfv/dataset_vehicle_detection_ilds_iitm.tar.gz?dl=0)

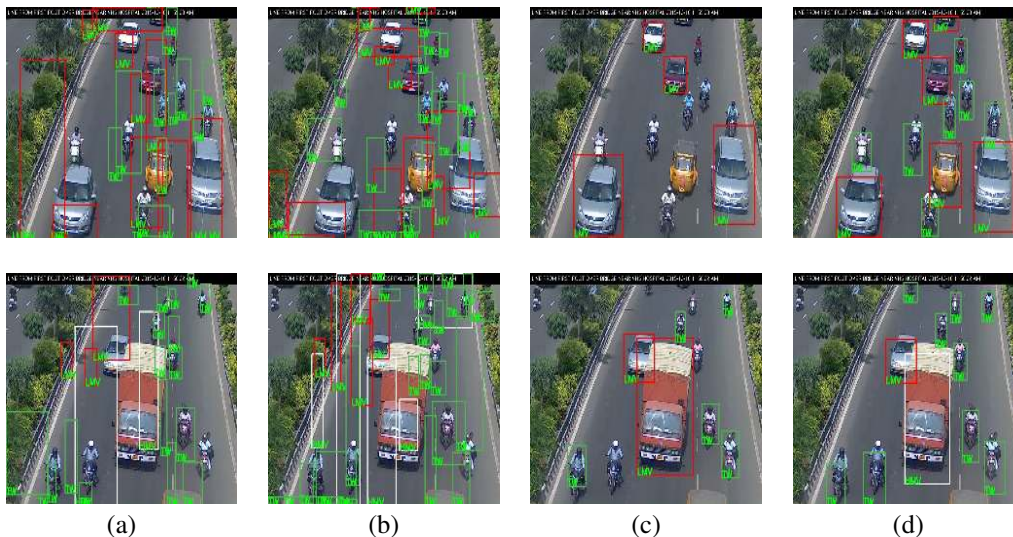


Figure 3: Vehicle detection results on Dataset-1. (a) Faster RCNN model fine-tuned on our data, (b) Faster RCNN model trained on our data from scratch, (c) Faster RCNN pre-trained on PASCAL VOC data, (d) Extended Faster RCNN model trained on 3-class augmented data.

Table 2: Results on adding a new class to the model (AP @ 0.5).

Model	AR	TW	HMV	LMV
(i) Pre-Trained Model + SVM	0.195	0.132	0.417	0.58
(ii) Data augmentation (Dataset-1) + SVM	0.609	0.783	0.653	0.87
(iii) Data Augmentation (Dataset-2)	<b>0.983</b>	<b>0.883</b>	<b>0.987</b>	<b>0.905</b>

133 classification, (ii) training Faster RCNN on augmented dataset and then using SVM for classification  
 134 and (iii) extending Faster RCNN with a new class: auto-rickshaw.

135 All the experiments have been performed on a machine with dual core Intel Xeon processor (2.20  
 136 GHz) having 256 GB of DDR4 RAM with one TitanX graphics processing unit (GPU). Using Faster  
 137 RCNN model we achieved processing speed at 5 frames per second.

## 138 5.1 Data augmentation

139 Table 1 shows results of Faster RCNN architecture using different types of training on Dataset-1  
 140 that has three classes: 1) Two wheelers (TW), 2) Light Motor Vehicle (LMV), and 3) Heavy Motor  
 141 Vehicle (HMV). From this table, we can infer that pre-trained model gives poor results. The poor  
 142 performance of the fine-tuned model can be attributed to the difference in quality and content of the  
 143 collected data compared to Pascal VOC. The model trained only on the collected dataset is performing  
 144 poorly because of limited data. The model trained from scratch on augmented data is performing  
 145 best because it is learning from both datasets; it is benefiting from the good features present in Pascal  
 146 VOC dataset and also optimizing parameter values according to our dataset. Image outputs from each  
 147 approach are shown in Figure 3.

## 148 5.2 Extending Faster RCNN Model for new classes

149 To compare deep learning approaches with traditional approaches we have trained different SVM  
 150 models for vehicle classification. In order to generate feature vectors for SVMs' training, we extracted  
 151 SIFT features from the image patches, where each patch contains only one vehicle. Once we have  
 152 cropped a patch from an image, we change its color space from RGB to gray-scale. Then, SIFT and  
 153 SURF features are extracted for all the patches. K-means clustering is done separately on the SIFT  
 154 and SURF features extracted from all the patches. The final feature vector for each patch is then

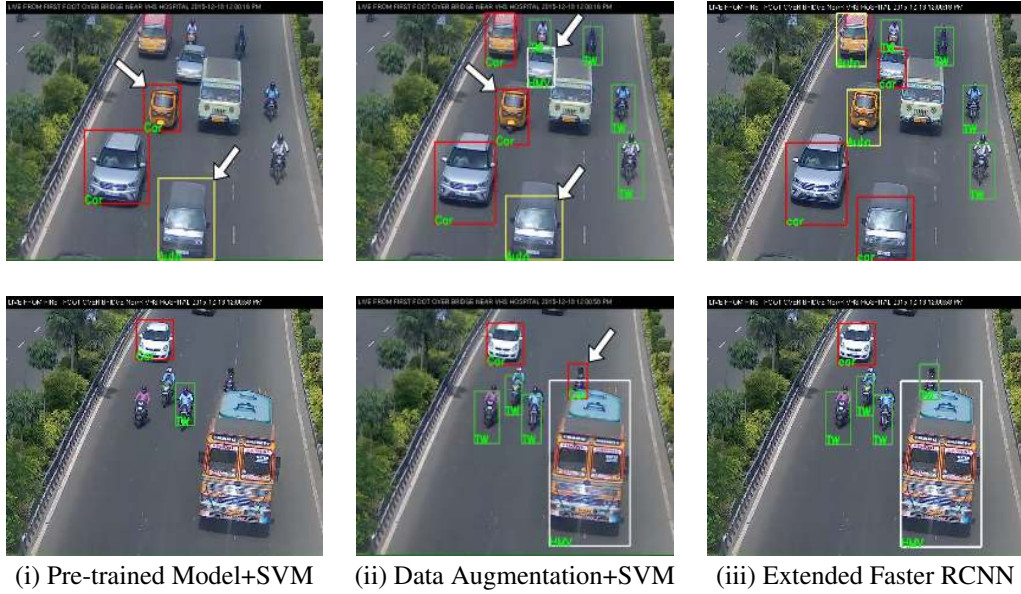


Figure 4: Image outputs of extended Faster RCNN model on Dataset-2.

155 generated following bag-of-words [6] approach, *i.e.*, for a given ‘ $k$ ’ we compute a number of the  
 156 SIFT features getting assigned to a particular cluster center. Experiments were done with different  
 157 values of ‘ $k$ ’. Similarly, another feature vector corresponding to SURF features was generated. After  
 158 getting these feature vectors, SVMs were trained on top of it.

159 As explained in the dataset section, we have merged the eight vehicle categories into three vehicle  
 160 categories. This allowed us to make the best use of Faster RCNN architecture with the existing  
 161 pre-trained model with minimum modifications. The results are shown in Table 2. Faster RCNN  
 162 architecture is able to detect all vehicles well; however, it is unable to classify auto-rickshaw since  
 163 it is not trained on our data. One solution is to train an SVM model to do the classification instead.  
 164 Therefore, in this setting, we get the object proposals from the Faster RCNN model to detect vehicles  
 165 and then employ SVM to classify the detected vehicles into different classes. Finally, we extended  
 166 the Faster RCNN model to incorporate a new class. Adding a new class in Faster RCNN model and  
 167 then training with augmented data gives the best results. Image outputs from each model are shown  
 168 in Figure 4.

## 169 6 Conclusion

170 Deep learning has emerged as a significant new paradigm in object identification and classification.  
 171 However, training deep learning networks requires large datasets. In this paper, we demonstrate the  
 172 use of a limited traffic dataset that augments existing large scale datasets and uses an existing deep  
 173 learning network (Faster RCNN) for detecting and classifying vehicles several of which are truncated  
 174 or occluded. The extended faster RCNN model is also able to deduct a new class of vehicles with  
 175 high degree of accuracy. The results obtained are promising for heterogeneous traffic scenario where  
 176 occlusion is common. This result is expected to encourage the wide-spread use of deep learning for  
 177 traffic video image processing since it is economical in terms of cost and time.

178 The results open up significant avenues for further research. For example, the present model works at  
 179 5 fps on TitanX GPU because of the high computation time of Faster RCNN. To make this model run  
 180 in real-time is one future work direction. A larger dataset with more instances of each class can be  
 181 used to train an eight- or ten-vehicle class model. Given the dissimilarities particularly among vehicle  
 182 types grouped under heavy vehicles, such a finer classification may result in significant improvements  
 183 to overall accuracy. Testing the robustness of developed models with multiple video inputs with  
 184 varying environmental parameters is on-going.

## References

- 185
- 186 [1] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. *Computer*  
187 *vision–ECCV 2006*, pages 404–417, 2006.
- 188 [2] Olivier Chapelle, Patrick Haffner, and Vladimir N Vapnik. Support vector machines for  
189 histogram-based image classification. *IEEE transactions on Neural Networks*, 10(5):1055–  
190 1064, 1999.
- 191 [3] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *Com-*  
192 *puter Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference*  
193 *on*, volume 1, pages 886–893. IEEE, 2005.
- 194 [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
195 hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009.*  
196 *IEEE Conference on*, pages 248–255. IEEE, 2009.
- 197 [5] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman.  
198 The pascal visual object classes (voc) challenge. *International journal of computer vision*,  
199 88(2):303–338, 2010.
- 200 [6] AG Faheema and Subrata Rakshit. Feature selection using bag-of-visual-words representation.  
201 In *Advance Computing Conference (IACC), 2010 IEEE 2nd International*, pages 151–156.  
202 IEEE, 2010.
- 203 [7] David Fleet and Yair Weiss. Optical flow estimation. In *Handbook of mathematical models in*  
204 *computer vision*, pages 237–257. Springer, 2006.
- 205 [8] Ross Girshick. Fast r-cnn. In *International Conference on Computer Vision (ICCV)*, 2015.
- 206 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep  
207 convolutional networks for visual recognition. In *European Conference on Computer Vision*,  
208 pages 346–361. Springer, 2014.
- 209 [10] Yi Li, Kaiming He, Jian Sun, et al. R-fcn: Object detection via region-based fully convolutional  
210 networks. In *Advances in Neural Information Processing Systems*, pages 379–387, 2016.
- 211 [11] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang  
212 Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European Conference on*  
213 *Computer Vision*, pages 21–37. Springer, 2016.
- 214 [12] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal*  
215 *of computer vision*, 60(2):91–110, 2004.
- 216 [13] Gitakrishnan Ramadurai Manipriya, S. and VV Bhavesh Reddy. Grid-based real-time image  
217 processing (grip) algorithm for heterogeneous traffic. In *In Communication Systems and*  
218 *Networks (COMSNETS), 2015 7th International Conference on Intelligent Transportation*  
219 *Systems*, pages 1–6. IEEE, 2015.
- 220 [14] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified,  
221 real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and*  
222 *Pattern Recognition*, pages 779–788, 2016.
- 223 [15] Joseph Redmon and Ali Farhadi. Yolo9000: Better, faster, stronger. *arXiv preprint*  
224 *arXiv:1612.08242*, 2016.
- 225 [16] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time  
226 object detection with region proposal networks. In *Advances in Neural Information Processing*  
227 *Systems (NIPS)*, 2015.
- 228 [17] S Cheung Sen-Ching and Chandrika Kamath. Robust techniques for background subtraction  
229 in urban traffic video. In *Electronic Imaging 2004*, pages 881–892. International Society for  
230 Optics and Photonics, 2004.
- 231 [18] Lijun Wang, Wanli Ouyang, Xiaogang Wang, and Huchuan Lu. Stct: Sequentially training  
232 convolutional networks for visual tracking. In *Proceedings of the IEEE Conference on Computer*  
233 *Vision and Pattern Recognition*, pages 1373–1381, 2016.