140

141

142

143

144

145

146

147

148

149

100

## Microsoft Word Author Guidelines for CVPR Proceedings

# Object Detections by a Super-Resolution Method and Convolution Neural Networks

Paper ID \*\*\*\*< 2072 >

# Abstract

Recently with many of blur-less or slightly blurred images, convolutional neural networks classify objects with around 90 percent regression rates, even if there are variable sized images. However, small object regions or cropping of images make object detection or classification difficult and decreases the detection rates. In many methods related to convolutional neural network (CNN), Bilinear or Bicubic algorithms are popularly used to interpolate region of interests. To overcome the limitations of these algorithms, we introduce a super-resolution method applied to the cropped regions or candidates and this leads to improve recognition rates for object detection and classification. Large object candidates comparable in size of the full image have good results for object detections using many popular conventional methods. However, for smaller candidates, super-resolution region using our preprocessing and region candidates, allows a CNN to outperform conventional methods in the number of detected objects when tested on the VOC2007 and MSO datasets.

## 1. Introduction

Since Krizhevsky et al. [1] introduced specifically designed CNN architectures, there have been many methods to increase the rate of object classification on convolution neural networks (CNN). [2], [3], [4], [5], [6], [7] have shown performances to be increased. Nowadays, with many of blur-less or slightly blurred images, CNNs classify objects with around a 90 percent regression rates.

Recently, there has been research focused on reducing the misdetection and detection failures on convolution neural networks with the help of generative adversarial network (GAN). Furthermore, GANs have been extended by using reinforcement learning, [8].

Frameworks such as Caffe and Tensorflow can help to design CNN models for any specific computer visioning. Also, from a computer infrastructure point of views, graphics processing unit (GPU), such as those from Nvidia substantially increase performance. These advances allow one to design systems for vision detection and classification 161 designed to run in real time. 162

However, even as there have been improvements of 163 convolutional neural networks in multiple ways, there are 164 still misdetections and detection failures for object 165 classifications. For example, randomly 100 images from [9] 166 were selected and preprocessed at three times lower 167 resolution than the original images to build cropped images. 168 Then they were interpolated by Bilinear and Bicubic 169 algorithms, and finally, they were tested by [3]. Figure 1 170 shows the detection failure of the second person from left side of the image. The second person from the left side in 171Figure 1 is missed by the algorithm. We will compare 172 several different algorithms and propose an advanced 173 method in later sections. 174



Figure 1: An example of object classification of "person" by CNN

Alex Krizhevsky et al. in [1] described the number of categories and designed their network with 1,000 object classes from ImageNet. [2], [5] classified 22 object classes from PASCAL VOC 2007 [9]. There is a 200 class data set in ILSVRC2013, ILSVRC2017 contains several datasets and especially, the ImageNet dataset contains 1000 classes and Omniglot contains 1623 classes. But research is often done with VOC2007 20+1 classes even though there are datasets with more than several hundred classes. 196

Recently, the most popular image size on CNN has been 197 around 256x256. Spatial Pyramid Pooling network, Faster 198 R-CNN, and ConvNet are examples testing the relevance of 199 the image size. Advances in camera systems and the popularity of mobile devices cause consumers to demand higher resolution images. Moreover, in medical and

150 151

152

155

156 157

158

159

160

186

187

geometrical satellite imaging systems the demand of object classifications are required in strong. In [11], [12], [13], they shows that recurrent neural network model is capable of extracting information from an image by selecting a sequence of regions and processing by selected region at high resolution.

In practical applications, there are many low quality image processing systems such as surveillance camera systems, car black-box systems, or even mobile phone cameras for taking pictures of long distance. For example, in surveillance systems it may not easy to increase the capacities for better image qualities because of their storage capacities, dark conditions on camera image sensors, and night visioning. In car black box systems, moving vibrations and the requirements of low electric power consumption may cause bad image compressions or lack of fast zooming or focusing devices. Especially, mobile phone pictures are blurred or have small target objects according to the limit of the software zooming algorithms.

Thus, we will introduce our research to improve object classification rates with improvements through super resolution algorithms and convolution neural networks.

# 2. Related works

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

To support fast implementation such as video data, the number of classes in CNN is preferred to be kept small. The number of classes to detect or classify objects is a major factor on the design of systems. The number of neurons, even in a hidden layer, are required to increase to handle many classes and it causes the system to slow down because of several hidden layers associated with heavy computations. Papers such as [2], [14] have shown new methods or new parameters related to lower layers of neural networks. The paper [14] introduced a new feature extraction algorithm for object recognition. The paper [2] also improved CNN in region proposal. Their new methods speed up the detection allowing the image data set to run in almost real time.

Keiming He et al. in [5] had the best results for training and testing with the maximum image side of 392 because they had image dataset from VOC 2007 and ImageNet. Also they showed results indicating scale matters in the classification processing. Thus, they suggested the spatial pyramid pooling model to support various image sizes in convolution layers which work with various image sizes while the standard fully connected layer requires a fixed image size. With Caltech101 image dataset, they found objects that had better performances among several scaled datasets. They noticed that this is mainly why the detected objects usually occupy large regions of the whole images. They evaluated cropped or warped images and got lower accuracy rates than the same model on the undistorted full image.

In [15] Generative Adversarial Networks (GANs) are described as generative models that use supervised learning

to approximate an intractable cost function and can simulate 250 many cost functions, including the one used for maximum 251 likelihood. More specifically, a generative model *g* trained 252 on training data *X* sampled from some true distribution *D* is 253 one which, given some standard random distribution *Z*, 254 produces a distribution *D'* which is close to *D* according to 255 some closeness metric (a sample  $z \in Z$  maps to a sample 256  $g(z) \in D'$ ). 257

There are several points of researches related to the GANs 258 which are training auto-encoder based GANs [4], learning 259 semi-supervised and generating images that humans find 260 visually realistic [16], and training for semi-supervised text 261 classification [17]. Also, there are GAN extensions related 262 to super resolution methods and reinforcement learning. 263

Image super resolution (SR) means a deep learning 264 method which infers a high resolution image from a low 265 resolution image. In [18] SR algorithms are introduced in 266 two groups: single image based and multiple image based 267 ones. That is, in multiple image SR algorithm, to restore an 268 image a couple of low resolution images of the same scene 269 are fed as input and a registration algorithm to find the 270 transformation between them is added while single image SR algorithm employs a training step to learn the 271 relationship between a set of high resolution image and 272 their low resolution counterparts and the relationship is 273 used to predict the missing high resolution details of input 274 images. They can be applied for video data to improve 275 276 action recognition.

As single image SR, [19], [20], [21] present methods 277 using mixture of experts (MoE) which are anchor based 278 local learning approach, sparse coding, and deep 279 convolution neural networks, respectively. Christian Ledig 280 et al. in [22] proposed a GAN method using generator 281 networks and discriminator networks to recover photo 282 realistic natural images with minimizing the mean squared 283 reconstruction error. 284

In SR and CNN algorithms, image interpolations are used to resize images or to crop/ warp image patches. Especially, small region of interest in a big image resolution is processed with a cropping or warping to fit into bounding boxes or region proposal and it results in decreased object detection rate.

Current CNNs, which are using the sliding window 290 method, are able to process variable size images as their 291 input. However, the fully connected layer can work only 292 with a fixed size of feature maps. This is why the CNNs at 293 the first time had applied a fixed size of input images. 294 Kaiming He et al. [5] introduced variable size of images as 295 the input of their CNNs and warped or cropped images to 296 get fixed feature maps. Cropped or warped images might 297 have the poor object recognition results. This caused us to 298 implement super resolution instead of interpolation 299 methods.

## 3. Proposed mage resolution improving algorithm

CNN requires a target image as a fixed size one before processing hidden layers. Thus with a given input image, the network resizes it to the fixed-size of image. In [2], the Bilinear interpolation algorithm is applied to extend the image size to fit as an input for a neural network. As future research, the extended image size directly given by the super-resolution will be used. In this paper, we use the super resolution method to increase the image size to around 492x324.

Thus, we will first describe the super resolution method to improve the images before the input layer of CNN, then, how the preprocessed image can be classified by CNN.

## 3.1. Super resolution image scale-up

Most work on image processing focuses on improving the deep hidden neural networks. In this section, however, we will describe more on pre-processing of image samples which have not enough qualities than on improving of hidden layers. For example, a cropped image whose size of 166x110 is extracted from a normal image is too small to feed into the input layer of our CNN.

In the expert method, component regressors,  $W_i$ , and as an anchor point,  $v_i$  have relations such as:

$$\min_{\{v_1, v_2, \dots, v_k; w_1, w_2, \dots, W_k\}} \sum_{j=1}^{N} \sum_{i=1}^{N} c_{ji} ||h_j - W_i l_j||^2,$$

where  $l_j$  means a low resolution path,  $h_j$  is the corresponding high resolution patch,  $v_j$  is the nearest anchor point for  $l_j$  and  $c_{ji}$  is a continuous scalar value which represents the degree of membership of  $l_j$ .

However, to keep the computational efficiency and the competitive image quality of anchor-based local learning method of multiple regressors, a mixture of experts which is one of conditional combined mixture models [23, 24] is proposed. Like [24] we define the model of mixture of experts for super-resolution images, describing the expectation and maximization (EM) algorithm, and also train and test this model.

In a mixture of expert model, a maximum likelihood estimation should be solved iteratively by the EM algorithm. At every iteration, the posterior probabilities are calculated for patches and then we get the expectation of the log likelihood as a result of the E-step. During the Mstep, anchor points and regressors are updated which is a softmax regression problem. After training, super resolution images can be constructed by collecting all the patches from regressed low resolution patches and averaging the overlapped pixels.

To differentiate the performance of super resolution method from interpolation methods, Bilinear, Bicubic interpolated images will be built in addition to the images processed by super resolution method.

## 3.2. Object detection

351 We implemented the convolution neural network based 352 on the Faster R-CNN [3] because it supports variable image 353 sizes as the input data of CNN and sliding widow proposal 354 scanning for the convolution network. Above all, the 355 detection speed is fast enough to be almost real-time. As 356 mentioned earlier, Faster R-CNN uses region proposal to detect an object. But the size ratio of region proposal to the 357 whole image is critical for the object to be detected. It 358 means we can distinguish our proposed method compared 359 to other interpolated data. 360

The CNN is implemented based on Intel® Xeon CPU 3612.30GHz and 4 NVIDIA Tesla K80 GPU boards. Each 362GPU board has memory of 12GB and two GPUs.363

Our CNN has several convolutional layers to support a 364 region proposal network in addition to the conventional 365 CNN [1], [7]. Also, these convolution layers are shared with 366 object detection networks as in [5]. Thus, our model works 367 as a deep CNN which has more convolution and pooling 368 layers too. 369

As multiple-scale prediction schemes for regression references, there are schemes based on multiple-images, on multiple-filters, and on multiple-anchors in [3]. With multiimage schemes, images are resized at multiple scales and feature maps are computed for each scale. Even though this scheme is time-consuming, our super resolution method can resize images to get better prediction scores. However, for less usage of computational resources, for fast detections, and for feature sharing in fully convolution layer, we choose the multi-anchors scheme.

Additionally, to control the memory usages of CAFFE, 379 we impose constraints on the number of images. In our deep 380 CNN model allowing multiple scaled image sizes, the 381 number of region proposals is w (width of the image) x h 382 (height of the image) x k (the number of anchors of a region 383 proposal). The memory space for the region proposal 384 network needs to be limited. Therefore, we constrain the 385 number of simultaneous images in hidden layers to be less 386 than or equal to 20 images. 387

For model training, we use pre-trained model parameters taken from the Faster R-CNN implementation. Faster R-CNN was trained, validated and tested with the PASCAL VOC 2007 dataset of 9000 images. As the result, we consider that the parameters should be fitted well enough for our model.

In every proposal from an image, an object is roughly considered and the proposal size is chosen by the predefined window size. Thus, this patch of an image may be interpolated. If instead of interpolation methods the super resolution method is adopted, any variable sizes of input images with variable sizes or shapes can be supported with better cropped or warped regions. 393 394 394 395 397 398 398 399

We will not adopt super resolution method to improve only the specific patch image qualities now. We will keep this for future research.

#### 4. Performance evaluations

Before any other processing, we considered several image file formats to get better image quality for preprocessing, training, and testing images. Therefore, we picked two image data formats and with a few images we processed the whole procedures of our proposed model. While most of image datasets are built images based on the JPG file format, this did not seem to offer as good results in super resolution processing compared with the BMP format. As the result of this simple testing, we decided to convert images from a file format of JPG into BMP format as a pre-processing procedure. Then, the image is taken to is further processed with Bilinear, Bicubic interpolations, and the super resolution method.

Also, we scale down image size by 3 times smaller in each width and height to build images with lower quality, instead of collecting images by cropping or warping images.

In our model, we implement the super resolution procedure based on [19]. Unlike training with less than 50 images in which most of super resolution models are published, we have trained our model based on their initial parameters and with 100 PASCAL VOC2007 images. There is not overfitting with this number of images for training. With random images from PASCAL VOC2007 [9], [25] and test images from [26] we have tested our super resolution model. We could not find any differences between them. Therefore we will testify object 450 classification part with 520 images randomly extracted 451 from VOC2007 and 1224 images from MSO [27]. 452

In Figure 2, we found the dataset has a different number 453 of objects compared to our intuition. For example, top-left 454 image is labeled with no object but our proposed model 455 detected an object or objects, like as shown in the image. 456 Thus, we decide not to use the given label from dataset. 457

### 4.1. Comparison of output pictures

Many images from the PASCAL VOC 2007 have 460 multiple objects, as shown in the results given in Table 1 461 with 3 different pre-processing models, which are bilinear, 462 Bicubic interpolation and a super resolution. Our model is 463 set as learning rate of 0.001 and detection scores with 0.8 464 or higher. 465

The first column shows that by bilinear interpolation there 466 are the number of detected objects in each given class. For 467 the second column every rows show that the number of 468 Bicubic interpolated and the number of object-detected 469 images with the regression rate. And so are the same in the 470 third column with super resolution. 471

As shown Table 1, our model detected more objects than the other two models but the average detected scores are not much different. It means that our model made better images as the input image and fed into the CNN then as the results 



Figure 2: Images which have objects detected but labeled without any object or less number of objects.

there are more number of detections. Especially, Figure 3 and Figure 4 show the cumulative number of detected objects and compare them.



Figure 3: Distribution diagram of the number of classified objects for bilinear, Bicubic and super resolution models on VOC 2007 dataset.



Figure 4: Histogram for convolution network models with 3 different image pre-processing on randomly extracted images from VOC 2007 dataset

With image dataset of Micro soft MSO, we present results in Table 2. Figure 5, and Figure 6. Like VOC2007 dataset, our proposed model detected more number of objects than the two other models. Unlike the case of the VOC2007, MSO dataset has zero object images. However, the average detection scores of our model are given as similar with the one of VOC2007.

#### 4.2. Big vs. small ROI pictures and their regression rates

As mentioned previous section, if objects are big enough compared to the size of the image which is containing the object proposal, objects from interpolated images with Bilinear or Bicubic methods have well enough or sometimes may have better performance than objects from



Figure 5: Diagram of the number of classified objects for bilinear, Bicubic and super resolution models on MSO dataset



Figure 6: Histogram for convolution network models with 3 different image pre-processing on randomly extracted images from MSO dataset

### our proposed model.

As mentioned in previous section, if objects are big enough compared to the size of the image which is containing the object proposal, objects from interpolated 587 images with Bilinear or Bicubic methods are good enough 588 or sometimes may have better performance than objects 589 from our proposed model. However, our proposed model 590 has much better results with objects from small bounding 591 boxes or small ratio of objects to the size of the image which 592 contains the object. In Figure 7, our proposed model detects 593 a chair which is a pretty small object in a given image but 594 the other two models did not detect 'chair' object. Even 595 though the other chair is detected in all of three models, our 596 proposed model has a little bit higher score. 

#### 5. Conclusion

We proposed a model which is composed with super resolution processing and a convolution neural network. In this model, several kinds of classes from two different datasets are scored when objects are detected. Our model has benefits on the number of object detections, especially in case of small object detections.



chair detections with p(chair | box) >= 0.8







Figure 7: Comparison of small object detection through Bilinear, Bicubic, and SR models

In this proposal, we did not implement that only bounding box areas are processed with super resolution method. But we can pretty sure that this can save the computational resources and thus we can adopt this in real-time processing for better object detection or classification.

Our proposed model appears powerful in scenarios such as relatively small objects in big pictures, warping on region proposals, and object detection from cropped image.

#### References

 Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton. "ImageNet Classification with Deep Convolution Neural Networks". In Advances in Neural Information Processing Systems (NIPS), (1097-1105). (2012).
 Chi Li Li D. The D. Chilling Market and Construction Statement of the second s

- [2] Girshick Ross. "Fast R-CNN". arXiv: 1504.08083v2 [cs.CV] 27 Sep 2015. 658
- [3] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. "Faster R-CNN: Towards Real-time Object Detection with Region Proposal Networks". arXiv:1506.01497v3[cs.CV] 6 Jan 2016.
- [4] David Berthelot Schumm, Luke Metz Thomas.
   "BEGAN: Boundary Equilibrium Generative Adversarial Networks". arXiv: 1703.10717v2
   [cs.LG]. (2017).
- [5] Keiming He, Xiangyu Zhang, Shaoqing Ren, and Jian 667 Sun. "Spatial Pyramid Pooling in Deep 668 Convolutional Networks for Visual Recogniton". 669 arXiv:1406.4729v4. (2015). 670
- [6] Ross Girshick, Jeff Donahua, Trevor Darrell, Jitendra 671
  Malik, "Region-Based Convolution Networks for Accurate Object Detection and Segmentation", *IEEE 7ransactions on Pattern Analysis and Machine 1ntelligence Vol. 38, NO.1*, 2016.
- [7] Karen Simonyan, Zisserman Andrew. "Very Deep 676 Convolutional Networks for Large-Scale Image 677 Recognition". ICLR. (2015).
- [8] Ian J. Goodfellow Pouget-Abadie, Mehdi Mirza, 679
   Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron 680
   Courvill, Yoshua BengioJean. "Generative 681
   Adversarial Nets". arXiv:1406.2661v1. (2014). 682
- [9] M. Everingham Van Gool, C. K. I. Williams, J. Winn, 683 and A. ZissermanL. "The PASCAL Visual Object 684 Classes Challenge 2007 (VOC2007) Results". 685 (2007). 686
- [10] Achanta R., Estrada F., Wils P., Süsstrunk S. 687
  "Salient Region Detection and Segmentation". In: 688
  Gasteratos A., Vincze M., Tsotsos J.K. (eds) 689
  Computer Vision Systems. ICVS 2008. Lecture 690
  Notes in Computer Science, vol 5008. Springer, 691
  Berlin, Heidelberg. (2008) 692
- [11] Volodymyr Mnih, Nicolas Heess, Alex Graves, and 693 Koray Kavukcuoglu. "Recurrent Models of Visual 694 Attention". arXiv:1406.6247v1[cs.LG] 24 Jun 2014. 695
- [12] Jimmy Lei Ba, Volodymyr Mnih, and Koray 696 Kavukcuoglu. "Multiple Object Recognition with 697 Visual Attention". arXiv:1412.7755v2[cs.LG] 23 698 Apr 2015. 699
- [13] Karol Gregor Danihelka, Alex Graves, Danilo Jimenez Rezende, Daan Wierstra Ivo. "DRAW: A

Recurrent Neural Network For Image Generation". arXiv:1502.04623v2[cs.CV] 20 May 2015.

- [14] J.R.R. Uijlingsvan de Sande, T. Gevers, and A.W.M. Smeulders K.E.A. "Selective Search for Object Recognition". IJCV. (2012).
- [15] Ghosh Nachum and Debiprasad Ofir. https://www.quora.com. "What-are-Generative-Adversarial-Networks-GANs".
- [16] Tim Salimans Goodfellow, Wojciech Zaremba, Vicki Cheung Ian. "Improved Techniques for Training GANs". arXiv:1606.03498v1. (2016).
- [17] Takeru Miyato M Dai, Ian Goodfellow Andrew. "Adversarial Training Methods for Semi-Supervised Text Classification". ICLR. (2017).
- [18] Nasrollahi Kamal, Guerrero Escalera Sergio, Rasti Pejman, Anbarjafari Gholamerza, Baro Xavier, J. Escalante Hugo, Moeslund B. Thomas. "Deep Learning based Super-Resolution for Improved Action Recognition". In International Conference on Image Processing Theory, Tools and Applications (IPTA) IEEE Signal Processing Society. (2015).
- [19] Kai Zhang Wang, Wangmeng Zuo, Hongzhi Zhang, Lei ZhangBaoquan. "Joint Learning of Multiple Regressors for Single Image Super Resolution". IEEE Singnal Processing Letters. Vol. 23 No. 1. (2016).
- [20] Zhaowen Wang Liu, Jianchao Yang, Wei Han, Thomas Huang Ding. "Deep Networks for Image Super Resolution with Sparse Prior". ICCV. (2015).
- [21] Chao Dong Change Loy, Kaiming He, Xiaoou Tang Chen. "Image Super Resolution Using Deep Convolution Networks". arXiv:1501.00092v3.v (2015).
- [22] Christian Ledig Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, Wenzhe ShiLucas. "Photo-Realistic Single Image Super Resolution Using a Generative Adversarial Network". arXiv:1609.04802v5[cs.CV]. (2017).
- [23] Bishop Christopher. "Pattern Recognition and Machine Learning". Springer. (2006).
- [24] Zhang Zhang, Baoquan Wang, Wangmeng Zuo, Hongzhi ZhangKai. "Joint Learning of Multiple Regressors for Single Image Super-Resolution". IEEE Signal processing letters, Vol. 23, No.1. (2016).
- [25] Olga Russakovsky Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg and Li Fei-Fei Jia. "ImageNet Large Scale Visual Recognition Challenge". arXiv:1409.0575. (2014).

- [26] Olga Russakovsky Deng, Hao Su, Jonathan Krause, 750 Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej 751 Karpathy, Aditya Khosla, Michael Bernstein, 752 Alexander C. Berg and Li Fei-Fei Jia, "ImageNet Large Scale Visual Recognition Challenge 2017". (2017). 756
- [27] Zhang Ma, Shuga Sameki, Mehrnoosh Sclaroff, Stan Betke, Margrit Lin, Zhe Shen, Xiaohui Price, Brian Much, Radom Jianming. "Salient Object Subitizing". IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2015).

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

736

737

738

739

740

741

742

743

744

745

746

747

748

749

700

701

702

703

704

705

706

707

708

709

710

	Bilinear		В	icubic	Ours	
Classes	#tp	mean	#tp	mean	#tp	mean
aeroplane	25	0.9569	29	0.9472	27	0.9787
bicycle	16	0.9602	17	0.9697	18	0.9520
bird	18	0.9204	25	0.9178	30	0.9500
boat	16	0.9330	18	0.9260	20	0.9276
bottle	19	0.9222	17	0.9208	15	0.9331
bus	26	0.9626	25	0.9639	26	0.9588
car	93	0.9662	100	0.9671	104	0.9710
cat	11	0.9427	10	0.9665	11	0.9739
chair	25	0.9273	34	0.9368	45	0.9320
cow	11	0.9149	12	0.9216	16	0.9385
diningtable	7	0.9317	7	0.9420	12	0.9193
dog	37	0.9559	36	0.9657	35	0.9565
horse	30	0.9560	34	0.9574	37	0.9694
motorbike	13	0.9467	13	0.9630	16	0.9581
person	405	0.9546	432	0.9575	466	0.9590
pottedplant	10	0.9467	12	0.9194	18	0.8767
sheep	15	0.9275	13	0.9380	14	0.9227
sofa	7	0.9191	8	0.9175	8	0.9475
train	7	0.9193	4	0.9276	4	0.9572
tvmonitor	25	0.9653	25	0.9729	26	0.9479
Total	816	0.9415	871	0.9450	948	0.9465
mis-classified	25		21		39	

 Table 1: Detected objects on 3 different pre-processing and convolution neural networks with PACAL VOC2007. #tp is number of true positives correctly predicted. The column labelled mean is average probability from method.
 850

900	
901	Tal
902	pos
903	
904	
905	
906	а
907	
908	b
909	
910	b
911	h
912	D
913	b
914	_
915	b
916	
917	С
918	
919	С
920	
921	C
922	С
923	_
924	d
925	u
926	d
927	
928	h
929	
930	n
931	
932	р
933	n
934 025	ρ
935	s
930	
937	S
930	tr
939	u u
940 0/1	tv
0/2	
942	Т
944	
945	n
946	
947	
ודע	

Table 2: Detected objects on 3 different pre-processing and convolution neural networks with MSO dataset. #tp is number of true	
positives correctly predicted. The column labelled mean is average probability from method.	

	Bilinear		Bicubic		Ours	
Classes	#tp	mean	#tp	mean	#tp	mean
aeroplane	5	0.9650	6	0.9300	9	0.9263
bicycle	1	0.9992	2	0.9112	1	0.9978
bird	50	0.9512	59	0.9521	75	0.9627
boat	1	0.8301	1	0.9771	1	0.9713
bottle	23	0.9074	24	0.9062	21	0.9117
bus	3	0.9954	3	0.9954	3	0.9871
car	16	0.9641	19	0.9558	18	0.9499
cat	11	0.9639	10	0.9829	11	0.9540
chair	19	0.9299	20	0.9270	22	0.9393
COW	5	0.9452	6	0.9488	7	0.9584
diningtable	3	0.9410	4	0.8927	4	0.9072
dog	42	0.9559	47	0.9636	50	0.9478
horse	11	0.9728	11	0.9726	15	0.9362
motorbike	4	0.9487	4	0.9529	4	0.9588
person	302	0.9715	318	0.9718	340	0.9719
pottedplant	4	0.9023	5	0.8756	9	0.9035
sheep	1	0.8780	1	0.9353	4	0.9064
sofa	3	0.9204	3	0.9099	6	0.9261
train	6	0.9746	7	0.9529	8	0.9522
tvmonitor	6	0.9720	6	0.9804	10	0.9248
Total	516	0.9444	556	0.9447	618	0.9447
mis-classified	86		82		92	