

META-COGNITION-BASED SIMPLE AND EFFECTIVE APPROACH TO OBJECT DETECTION

Sannidhi P Kumar*

Chandan Gautam*

Suresh Sundaram*

* WIRIN Autonomous Systems and Robotics Lab
Department of Aerospace Engineering
Indian Institute of Science, Bangalore, Karnataka

ABSTRACT

Recently, many researchers have attempted to improve deep learning-based object detection models, both in terms of accuracy and operational speeds. However, frequently, there is a trade-off between speed and accuracy of such models, which encumbers their use in practical applications such as autonomous navigation. In this paper, we explore a meta-cognitive learning strategy for object detection to improve generalization ability while at the same time maintaining detection speed. The meta-cognitive method selectively samples the object instances in the training dataset to reduce overfitting. We use YOLO v3 Tiny as a base model for the work and evaluate the performance using the MS COCO dataset. The experimental results indicate an improvement in absolute precision of 2.6% (minimum), and 4.4% (maximum), with no overhead to inference time.

Index Terms— Object Detection, Meta-cognition, YOLO v3 Tiny, Deep Learning

1. INTRODUCTION

Object detection is one of the fundamental computer vision tasks used in a plethora of applications such as autonomous navigation, surveillance and security, and facial detection. Such widespread use necessitates the need to build accurate real-time object detection networks. Object detection combines the tasks of image classification and object localization. Thus, an object detector's task is to return the bounding box coordinates of the objects of interest in an image and assign them class labels. Modern deep learning-based object detectors comprise two parts, a backbone pre-trained on ImageNet and a detection head. Several backbones, such as [1, 2, 3, 4, 5] have been proposed. Detection heads are classified into two types, two-stage object detectors, and one-stage object detectors. Two-stage detectors, such as Region-based Convolutional Neural Networks (R-CNN) [6], Fast R-CNN [7], Faster R-CNN [8], and Region-based Fully Convolutional Networks [9], generate region proposals in the first stage. In the second

stage, feature extraction from these region proposals is followed by object classification and bounding box regression. These detectors have achieved high accuracy rates but low operational speeds. To address this issue, single-stage detectors have been proposed. These detectors skip the region proposal stage and make bounding box predictions directly from the input image. As a result, these models are faster but not as accurate as two-stage detectors. Among the single-stage detectors, Single Shot MultiBox Detector (SSD) [10], the You Only Look Once (YOLO) series of networks, and RetinaNet [11] are the most representative. SSD directly produces class predictions and adjustments to default bounding boxes of different aspect ratios and scales for each location in a feature map. YOLO [12] frames object detection as a regression problem wherein network evaluation is performed just once to predict the bounding boxes and class probabilities. YOLO v2 [13] improves the Average Precision (AP) of YOLO significantly by introducing batch normalization for convolutional layers, a high-resolution classifier, and dimension clusters. YOLO v3 [14] proposes a Darknet-53 architecture that further improves upon the performance of YOLO v2. YOLO v4 [15], being the latest addition to the family, introduces new features such as weighted residual connections, mosaic data augmentation, cross mini-batch normalization, cross-stage partial connections, self adversarial training, and Complete Intersection over Union (CIoU) loss.

Variants of the above described YOLO series have been released for constrained environments that meet real-time requirements. One such model is YOLO v3 Tiny, which clocks at 220 FPS, but gives an AP of 33.1%. The objective is to improve accuracy while preserving the detection speed. It is recently shown in cognitive psychology that human meta-cognition learning principles help achieve better generalization ability. Over the last decade, these principles have been used across various disciplines such as disease classification, science education, human emotion recognition, and search algorithms for optimization problems. Various researchers have amalgamated the concept of meta-cognition with neural networks in the past [16, 17, 18, 19]. Such neural networks have two components, cognitive and meta-cognitive.

We would like to thank WIRIN for funding this work.

The latter component scrutinizes the former’s knowledge and devises efficient learning strategies, namely the sample delete, neuron growth, parameter update, and sample reserve strategies to improve its learning ability. [16] implements the sample delete strategy for a classification problem as follows - when the predicted class label of the new training sample is the same as the actual class label, and the confidence level (estimated posterior probability) is greater than the expected value, the new training sample does not provide additional information to the classifier. It is deleted from the training pipeline without being used in the learning process. For a regression problem as in [17], if the predicted error for the current sample is less than the delete threshold, then the knowledge content of the sample is similar to the knowledge present in the network. Thus, the sample is removed from the training sequence.

This paper has explored a sample selection strategy, which focuses on what object instances in each training image can be learned for object detection. The proposed method is based on the self-regulative learning process, which selects appropriate learning samples based on current knowledge, and the knowledge present in the training batch. The self-regulative learning process is evaluated using YOLO v3 Tiny architecture for object detection. The performance of the proposed meta-cognitive YOLO v3 Tiny is evaluated using the MS COCO dataset [20] and compared against the performance of the baseline model. The results clearly highlight that meta-cognition helps achieve better performance without increasing the computational complexity of the model.

The rest of the paper is organized as follows. Section 2 discusses our baseline model. Section 3 describes the proposed method. Section 4 contains the experimental setup and performance evaluations. Finally, we conclude our paper in section 5.

2. BACKGROUND

This section discusses the details of the baseline model. The architecture of YOLO v3 Tiny broadly consists of 13 convolutional layers: feature extractor (seven layers) and feature detector (six layers). The backbone of YOLO v3 Tiny performs feature extraction, where each convolutional layer is followed by a batch normalization layer and Leaky ReLU activation. The feature extractor uses max-pooling layers to reduce the dimensionality of these convolutional layers. The network’s backbone is followed by the detection head, which makes detections at two scales, i.e., output images of sizes 13×13 and 26×26 for an input image of resolution 416. Each grid cell of an output image predicts three bounding boxes. The first output image of size 13×13 , which detects large objects, is obtained at the tenth convolutional layer of the network. The feature map from the sixth pooling layer is upsampled and then concatenated with the feature map obtained at the fifth convolutional layer of the network. The concatenated feature

map is then passed through the next two convolutional layers to give rise to the 26×26 image to detect medium sized objects. YOLO v3 Tiny performs detection based on these attributes, namely the class labels, object confidence score, and bounding box coordinates.

3. META-COGNITIVE YOLO V3 TINY

In this section, the proposed method Meta-Cognitive **YOLO v3 Tiny** (MC-YOLOv3T) is discussed in detail. In MC-YOLOv3T, we explore a sample selection strategy of meta-cognitive learning in the training phase. Initially, we pre-process the training images by performing data augmentation and resizing. The pre-processed images are then passed to the network to obtain the predicted bounding boxes, following which we apply meta-cognitive thresholding to the class predictions of each bounding box. Meta-cognitive thresholding selectively samples the object instances in the training image. The learning threshold is an exponentially decaying function and is defined as follows:

$$N(t) = N_0 e^{-\lambda t}, \quad (1)$$

where $N(t)$ is the threshold at epoch t , $N_0 = N(0)$ is the threshold at epoch 0, and λ is the decay constant. Since loss is higher in the initial phase of training, we set the meta-cognitive threshold to a higher value. As training proceeds and the loss decreases, we need to decrease the threshold value as well. Therefore, an exponentially decaying threshold function is used instead of a constant threshold.

The meta-cognitive threshold selectively samples the instances based on an error term. For every predicted bounding box assigned to a ground truth object, we calculate the error term which penalizes MC-YOLOv3T for assigning low scores to the correct classes and attaching high scores to the incorrect classes. The error term is defined as follows:

$$E = \max(|p(0) - \hat{p}(0)|, |p(1) - \hat{p}(1)|, \dots, |p(c) - \hat{p}(c)|), \quad (2)$$

where the error term E is the maximum of the absolute difference of the ground truth conditional class probability $p(0), p(1), \dots, p(c)$ and predicted conditional class probability $\hat{p}(0), \hat{p}(1), \dots, \hat{p}(c)$, and c is the maximum number of classes. If the error term is less than the meta-cognitive threshold, we set the classification loss contributed by the bounding box prediction to zero. This ensures that MC-YOLOv3T does not learn to classify an object instance if its prediction is within the error limit prescribed, thereby reducing overfitting. If the error term is more than the meta-cognitive threshold, we retain the classification loss contributed by such object instances. The loss function in MC-YOLOv3T comprises of three components - classification loss [14] which penalizes the class predictions made by bounding boxes that predict objects, bounding box regression

loss [21] which penalizes the model for the predicted bounding box coordinates, and objectness loss [14] which penalizes the model for the confidence with which it makes object predictions. These three losses are weighted as follows:

$$\text{Total Loss} = \alpha \text{ Classification loss} + \beta \text{ Bounding box regression loss} + \gamma \text{ Objectness loss} \quad (3)$$

Finally, we train the network by using the back-propagation algorithm. This training procedure is continued until the predefined number of epochs are reached, following which we use the trained network for detection. As the proposed method is based on YOLO v3 Tiny, its detection phase is identical to [14]. The complete training procedure is mentioned in Algorithm 1.

Algorithm 1 MC-YOLOv3T Training Procedure

Input: Number of epochs T , Initial meta-cognitive threshold $N(0)$, Final meta-cognitive threshold $N(T)$, Hyper-parameters

Output: Trained network

- 1: $t \leftarrow 0$
 - 2: **while** $t \neq T$ **do**
 - 3: Pre-process training images and pass through MC-YOLOv3T network
 - 4: Obtain bounding box predictions
 - 5: Compute meta-cognitive threshold, $N(t)$, by using Equation 1 and error term, E , by using Equation 2
 - 6: Compute classification loss for every bounding box assigned to a ground truth object
 - 7: Compute regression loss for every bounding box assigned to a ground truth object
 - 8: Compute objectness loss for every bounding box
 - 9: **for** Bounding box predictions which are assigned to ground truth objects **do**
 - 10: **if** $E < N(t)$ **then**
 - 11: Classification loss $\leftarrow 0$
 - 12: **end if**
 - 13: **end for**
 - 14: $t \leftarrow t + 1$
 - 15: Compute total loss by using Equation 3
 - 16: **end while**
-

4. EXPERIMENTAL RESULTS

We evaluate MC-YOLOv3T on MS COCO dataset [20]. For evaluation, we use the metrics [14] of $AP_{50:95}$, AP_{50} , and AP_{75} , which denote the average precision at IoU thresholds of 0.5 : 0.95, 0.5 and 0.75, respectively. We also use AP_S , AP_M and AP_L , which denote the average precision for small, medium, and large objects, respectively. For pre-processing the images, we use the technique of mosaic data augmentation [15] and resize images for multi-scale training, where

Hyper-parameter	Values
Weight for regression loss (β)	3.54
Weight for classification loss (α)	37.4
Weight for objectness loss (γ)	64.3
Batch size	32
Momentum	0.937
Weight decay	0.0005
Initial Learning Rate	0.01
Final Learning Rate	0.0005
Epochs	600

Table 1: Hyper-parameter values used for training

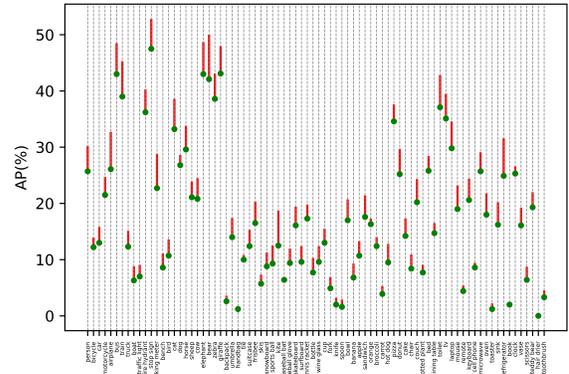


Fig. 1: AP Improvement per class of COCO dataset. Green dots indicate performance of baseline model, red lines indicate performance gain using MC-YOLOv3T

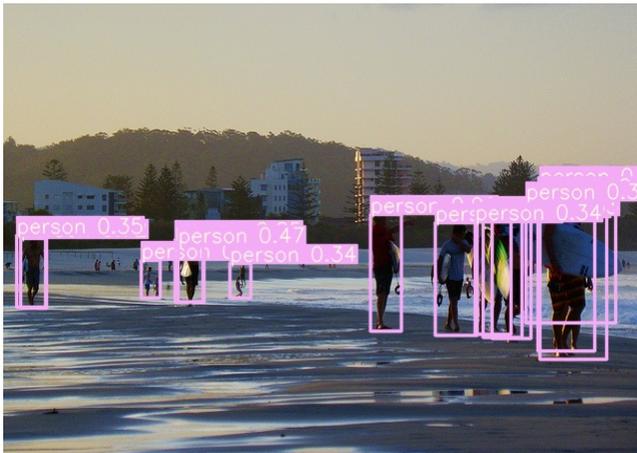
sizes are chosen from $\{320, 352, \dots, 608\}$ in the multiples of 32. Hyper-parameters of the network are listed in Table 1. Moreover, we use cosine scheduling for the learning rate and stochastic gradient descent optimizer. We also set the initial meta-cognitive threshold to 0.5 and final threshold to 0.05. First, MC-YOLOv3T is trained on the complete training set and 88% of the validation set of MS COCO 2014. We then perform validation on the remaining 12% of the MS COCO 2014 validation set. Further, the performance is evaluated using MS COCO 2017 test set. The hardware acquired for training is an Intel Xeon 1.8 GHz NVIDIA RTX Titan GPU with 128 GB RAM and 500 GB SSD for training the network.

4.1. Results

In Table 2, we compare the performance of MC-YOLOv3T against the baseline model at input image resolutions of 320, 416, 512, and 608 to juxtapose at different scales. As MC-YOLOv3T has used a meta-cognition approach during training, it has been exposed to fewer object instances than the baseline model. We find that the loss of the baseline model saturates at 300 epochs, whereas MC-YOLOv3T saturates at 600 epochs. Although MC-YOLOv3T saturates at 600 epochs, it has outperformed the baseline model at 300 epochs only. The results evaluated on MS COCO 2017 test

	AP _{50:95}				AP ₅₀				AP ₇₅			
	Baseline model	MC-YOLOv3T 300 epochs	MC-YOLOv3T 600 epochs	Absolute delta	Baseline model	MC-YOLOv3T 300 epochs	MC-YOLOv3T 600 epochs	Absolute delta	Baseline model	MC-YOLOv3T 300 epochs	MC-YOLOv3T 600 epochs	Absolute delta
YOLO v3 Tiny@320	15.5	17.6	17.8	+2.3	29.5	30.5	31	+1.5	14.6	17.6	18.1	+3.5
YOLO v3 Tiny@416	17.5	20.2	20.4	+2.9	33	35	35.3	+2.3	17	20.2	20.8	+3.8
YOLO v3 Tiny@512	18.7	21.2	21.5	+2.8	35.3	37.3	37.9	+2.6	18	21.2	21.6	+3.6
YOLO v3 Tiny@608	19.4	21.5	21.7	+2.3	37	38.5	38.7	+1.7	18.6	21.2	21.5	+2.9

Table 2: Overview of improvement in AP achieved by meta-cognitive training, evaluated on MS COCO 2014 validation set at multiple image resolutions



(a) Detection results using the baseline model



(b) Detection results using MC-YOLOv3T

Fig. 2: Comparison of Detection Results as seen on a COCO 2017 Test set image

	Baseline model	MC-YOLOv3T 600 epochs	Absolute delta
AP _{50:95}	17.2	20.3	+3.1
AP ₅₀	32.5	35.1	+2.6
AP ₇₅	16.6	20.9	+4.4
AP _S	4.2	4.6	+0.4
AP _M	16.6	20.3	+3.7
AP _L	29.9	34.4	+4.5

Table 3: Overview of improvement in AP achieved by meta-cognitive training, evaluated on MS COCO 2017 test set at image resolution of 416

set are presented in Table 3. Here, it has been observed that MC-YOLOv3T significantly outperforms the baseline model. On the primary detection metric of AP_{50:95}, we see an improvement of 3.1%. Moreover, there is an improvement of 2.6% (minimum) and 4.4% (maximum) on the metric of AP₅₀ and AP₇₅, respectively. These results suggest that MC-YOLOv3T yields better improvements at higher IoU thresholds, i.e., AP₇₅. Further, we compute AP_S, AP_M, and AP_L. In Table 3, there is a marginal increment of 0.4% for small objects detection (AP_S) due to the absence of a detection layer for small objects in MC-YOLOv3T. Moreover, significant improvements of 3.7% and 4.5% have been observed for medium (AP_M), and large (AP_L) objects detection, respectively.

The results in Table 3 are the average precision over 80 classes of the dataset and do not fully reflect per class AP changes. Therefore, we plot per class AP improvements in Figure 1. In this figure, it can be observed that MC-YOLOv3T has exhibited an improvement for all classes. The highest increment of 7.9% is observed for the class 'bear'. However, there is a marginal improvement for a few classes. The marginal improvement for such classes is due to data imbalance, as these classes have fewer annotated instances in the training set, as compared to other classes.

We have also illustrated the detection results in Figure 2a and 2b for the baseline model and MC-YOLOv3T, respectively. The latter makes detections on the 'surfboard' category, which are absent in the former model. Furthermore, confidence scores for the 'person' category are higher in MC-YOLOv3T. These results indicate the superiority of our proposed method over baseline due to meta-cognitive principles.

5. CONCLUSION

In this paper, we have employed a meta-cognitive learning-based sample selection strategy with YOLO v3 Tiny for object detection. The proposed MC-YOLOv3T, shows significant performance improvement across all input image resolutions. In particular, it yields greater improvements as the size of the object increases. Therefore, it has exhibited the best performance in case of large object detection. Finally, it is noteworthy that the MC-YOLOv3T achieves maximum 4.4% improvement without any computational burden in the inference.

6. REFERENCES

- [1] S. Liu and W. Deng, "Very deep convolutional neural network based image classification using small training sample size," in *Proceedings of the Asian Conference on Pattern Recognition*. IAPR, 2015, pp. 730–734.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 770–778.
- [3] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 5987–5995.
- [4] G. Huang, Z. Liu, L. Van Der Maaten, and K.Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 2261–2269.
- [5] A. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *ArXiv*, vol. abs/1704.04861, August 2017.
- [6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, pp. 142–158, January 2016.
- [7] R. Girshick, "Fast r-cnn," in *Proceedings of the International Conference on Computer Vision*. IEEE, 2015, pp. 1440–1448.
- [8] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, pp. 1137–1149, June 2017.
- [9] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," in *Proceedings of the Advances in neural information processing systems*, 2016, pp. 379–387.
- [10] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C-Y. Fu, and A.C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [11] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, pp. 318–327, February 2020.
- [12] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 779–788.
- [13] J. Redmon and A. Farhadi, "Yolo9000: Better, faster, stronger," in *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 6517–6525.
- [14] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *ArXiv*, vol. abs/1804.02767, 2018.
- [15] A. Bochkovskiy, C.Y. Wang, and H. Liao, "Yolov4: Optimal speed and accuracy of object detection," *ArXiv*, vol. abs/2004.10934, 2020.
- [16] G.S. Babu and S. Suresh, "Meta-cognitive neural network for classification problems in a sequential learning framework," *Neurocomputing*, vol. 81, pp. 86–96, April 2012.
- [17] G.S. Babu, X. Li, and S. Suresh, "Meta-cognitive regression neural network for function approximation: Application to remaining useful life estimation," in *Proceedings of the International Joint Conference on Neural Networks*, 2016, pp. 4803–4810.
- [18] A. Das, Nguyen Anh, S. Suresh, and N. Srikanth, "An interval type-2 fuzzy inference system and its meta-cognitive learning algorithm," *Evolving Systems*, vol. 7, pp. 95–105, March 2016.
- [19] N. Anh, M. Prasad, N. Srikanth, and S. Sundaram, "Wind speed intervals prediction using meta-cognitive approach," *Procedia computer science*, vol. 144, pp. 23–32, 2018.
- [20] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C.L. Zitnick, "Microsoft coco: Common objects in context," in *Proceedings of the European conference on computer vision*. Springer, 2014, pp. 740–755.
- [21] H. Rezaatofghi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.