

# Case Study: Performance Of Yolov 4 Is Better Than Yolov3

**Prof. Dr.Kamal Alaskar<sup>1</sup>, Dr.Firoj A. Tamboli<sup>2</sup>, Dr.Rajendra Jadhav<sup>3</sup>, Mrs.Manjushri A Kadam<sup>4</sup>**

<sup>1</sup>Professor, Department of Computer Applications Bharati Vidyapeeth (Deemed to be University) Institute of Management Kadamwadi,Kolhapur,Maharashtra,India.Pin-416003.

<sup>2</sup>HOD- Pharmacognosy Bharati Vidyapeeth college of Pharmacy, Kolhapur, Maharashtra, India

<sup>3</sup>Associate Professor, Department of Management Bharati Vidyapeeth (Deemed to be University) Institute of Management Kadamwadi,Kolhapur,Maharashtra,India.Pin-416003.

<sup>4</sup>Assistant Professor, Department of Management Bharati Vidyapeeth (Deemed to be University) Institute of Management Kadamwadi,Kolhapur,Maharashtra,India.Pin-416003.

---

**Abstract-**In computer vision, there are so many applications and uses, one of which is object detection. Object detection is a subset of computer vision that is used to detect the presence, location, and type of objects in images. Object detection is also a combination of three functions; Object recognition, to find objects in an image, Object localization, to find where exactly in the image the objects are located, and Object classification, to detect what particular objects are in that image.

There are lots of ways used to pick or grasp object through robotic hand but there are some hardly worked done with help the of Deep Learning approaches. To solve this issue, a solution is proposed which involves human strategies of picking up an object using Neural Network classifier. Classifier uses help of object detection model to detect object in environment and classifier classifies into picking strategy as per objects shape and orientation. Strategy detected by classifier can be used by soft hand as anticipatory action and reactive grasp. To increase accuracy number of primitive measures taken into consideration, our bounded and some of limitation are taken into mind while proposing architecture.

**Keywords-**Object Detection, Deep learning, Grasping Strategies, Neural Network, Soft Hand.

### I. INTRODUCTION

Soft hand as proven to be more efficient when used in supervision of human[1][2]. But such approach is still lagging in performance hence Data Driven approach can be used to improve the performance (see [3]).For detection of object on which grasping is done YOLO is embedded (You Only Look Once) technology YOLOv4, unlike used in [5] YOLOv3. YOLOv4 has shown good accuracy with respect to Single Shot MultiBox Detector (SSD). The code is online at <https://pjreddie.com/yolo/>. [4]. YOLOv4 replaces soft max function with binary classification which helps in reducing complexity and output is same as multiclass classification only[6].

Machine Learning approach had proved positive results in detection of the object or say grasping details of object [7][8][9]. In [10]Neural network helps in predicting unseen object and action to be perform on that object with the help of learning from pre labeled data. In [11] convolution neural network had played good role in detecting objects strategy output of whose network can be used to determine controlling robotic soft hand. Using such network objects are trained on 45 objects and tested with 10 objects which gives approximately accuracy of 84% [3]. Especially focus on updating object detection technique used in [3] from YOLO9000 to YOLOv4 which are faster than as used in [3]Below Block Diagram shows overall system of project except change/updating in object detection technology to YOLOv4 [3]

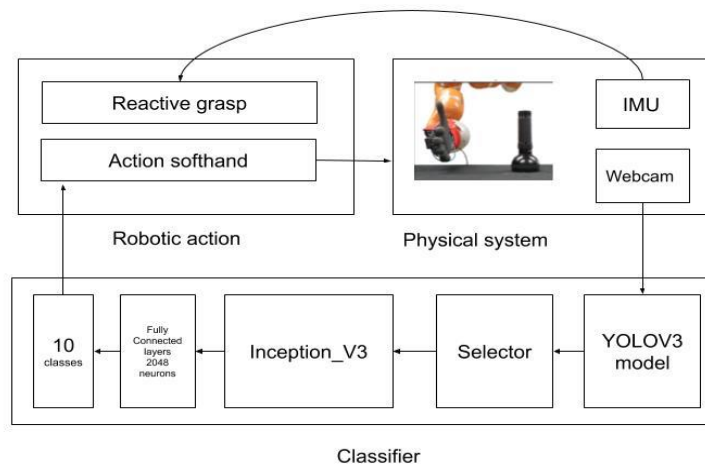


Fig I architechure for primitve selection

### II. LITERATURE SURVEY

As a main reference to project [3] whose system is used for further increment in approach is used in Data Driven approach to control Anthropomorphic soft hand Which uses Deep learning technology significance for increasing performance of control strategy of soft hand. It uses InceptionV3 module for grasping object and concluding strategy to be used by [11] reactive strategies have become significant in human robot exchanging of objects. It is tested over 10 objects by training on 45 objects which corresponds to accuracy of 84% on test objects. Where it is decided to update object detection method input given to model of Inceptionv3.

Paper achieves these goals by:

- i) Using Deep Neural Network model Inceptionv3 model predicts or decides which action human would take to pick certain type of objects.
- ii) Understanding through which action would be performed by human using robotic hand containing soft hand capability.
- iii) Testing over on 10 objects which are not used in model training model would be able to predict action to perform on object by soft hand or robotic arm with accuracy of prediction 84%.

#### **A. DIFFERENT MODELS**

For Deep Learning InceptionV3 module will be used which is pre trained model for object grasping and third version of inception module (see [12]). Module is trained over ImageNet Dataset .InceptionV3 model with 144 crops gained top-5 error rate is 4.2%, which pullback PReLU-Net and Inception-v2 which were used in 2015. With 42 layers deep, the parameter complexity increases by only 2.5% Google Net [12].It consists of 313 layers of neuron where some of layers will be retrained to get our performance output. PyTorch version Inception-v3:[14] <https://github.com/pytorch/vision/blob/master/torchvision/models/inception.py>

#### **B. OBJECT DETECTION MODELS**

YOLOv4 are proved to be faster than YOLOv3 used in [3] as said in[4]It's very fine on the old detection metric of 5 IOU. YOLOv4 uses multiclass classification method for object detection. It uses boxes which are predefined, and models are trained on it with lots of images label with boxes. While Detecting object it uses 10,000 approximately boxes to predict one of which have more significant level and use that box as detected object box and label that box depending on pre trained images.

	backbone	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP <sub>S</sub>	AP <sub>M</sub>	AP <sub>L</sub>
<i>Two-stage methods</i>							
Faster R-CNN+++ [5]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [8]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [6]	Inception-ResNet-v2 [21]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [20]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	<b>52.1</b>
<i>One-stage methods</i>							
YOLOv2 [15]	DarkNet-19 [15]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [11, 3]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [3]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
RetinaNet [9]	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
RetinaNet [9]	ResNeXt-101-FPN	<b>40.8</b>	<b>61.1</b>	<b>44.1</b>	<b>24.1</b>	<b>44.2</b>	51.2
YOLOv3 608 × 608	Darknet-53	33.0	57.9	34.4	18.3	35.4	41.9

Table I Object Detection score working of block in detail for figure shown above (Figure) .

### III. METHODOLOGY

- Training Phase
- Test Phase

#### A. TRAINING PHASE

To reduce training time, pre trained model is used to get more accuracy in less time YOLOv4 pre trained model weights and model to detect object which is been trained on COCO datasets which detects 80 classes if required to train can train using link to code (<https://github.com/qqwweee/keras-yolo3>). Training InceptionV3 module which is best for classification which can perform better than human vision also.

- **MODEL ARCHITECTURE**

Instead of creating whole model from scratch pre trained model InceptionV3 is used and use its weights and retrain some of its layers. Figure II shows detailed description of v3 module. To train model Dataset is required. Pop last layer of Inceptionv3 module and add two layers each of 2048 neuron and SoftMax layer of 6 classes, the indication probability of each object which strategy to be used more as shown in figure

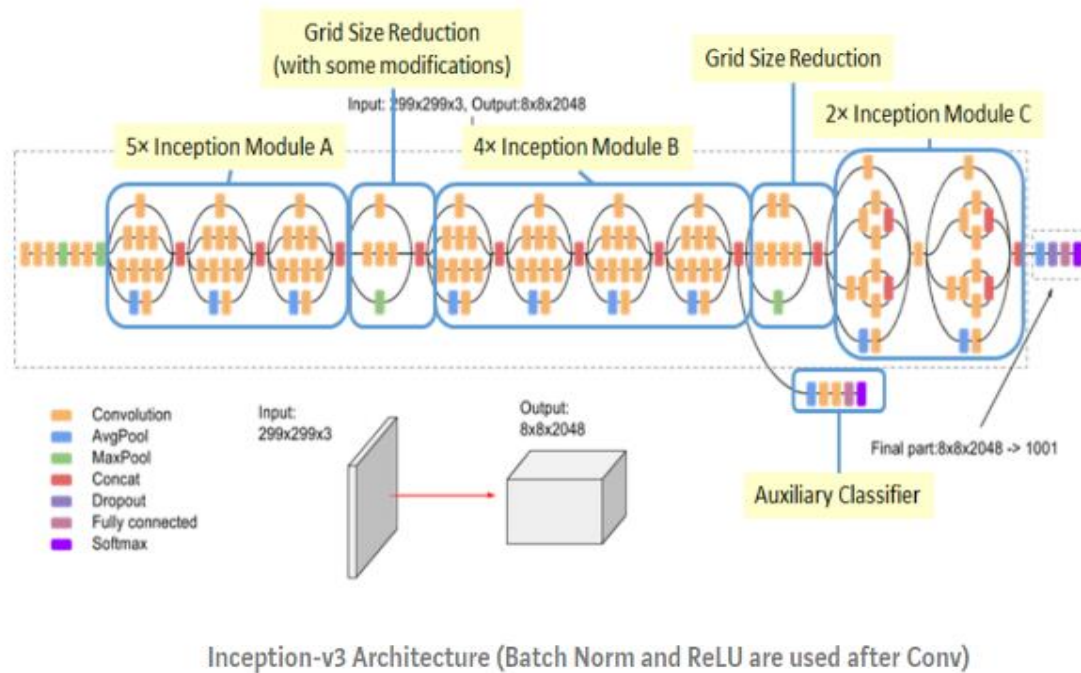


Figure II Inceptionv3 model

- **DATASET**

For creation of dataset as suggested in paper [3] objects are trained on 45 objects that are mostly different than used in [3]. Objects are kept on table at center position, different people are asked to pick up object and keep it on side and record this video. Each video is labeled as per the primitive taken by human in that video and identified six actions namely unlike of [4] to increase accuracy of model as top, left and right pinches can be replaced

**Top:** In top grasp approach, hand is approached from top with palm parallel to object on top it. This approach is used to pick object which are big in size like box used in following example.

**Bottom:** In this primitive hand are approach from right side and picked object with 4 fingers and thumb in opposite of 4 fingers providing negative force to hold object e.g. bowl.

**Pinch:** This type of approach is like top approach only, but have some difference which is commonly used for small objects like match stick.

**Slide:** In this approach object which are thin in size like CD are slide to corner of table and flipped and picked up

**Flip:** This is primitive use for small object with thin in size like coin where slide primitive cannot be used here and just flip object perpendicular to table.

**Lateral:** In this primitive hand approaches towards objects from right side with palm perpendicular to table and parallel to object like picking up bottle

These are 45 objects that are used here

**Note:** Objects which are used are also orientation dependent.

- **TRANSFER LEARNING**



Figure III Training Data images

Transfer Learning approach is where used weights of pre trained model to predict our aim can add some of more layers at the end of model that are fully connected or expand model and train concatenate part can also retrain some of the layers of models.

Adam as our optimizer.

- **STEPS IN TRAINING**

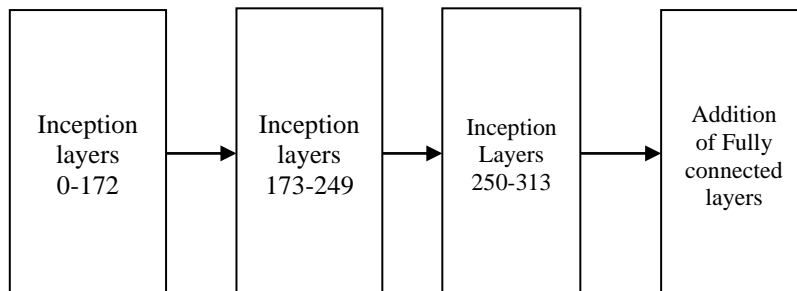
Initially it is set are some of the parameters of our build model taking from [3] paper outcomes

BATCH\_SIZE =10

EPOCHS = 10

Regularizer = 0.01

Pdrop = not used in our case



#### Fig IV Inception block diagram

Learning rate = Used Adam doesn't require learning rate.

Learning rate for fine tuning = Used Adam doesn't require learning rate.

Optimizer = Adam.

Activation function for last FFC Layers = relu.

Prediction layer activation function = SoftMax.

Next freeze all layers in model except the last three layers in model as shown in figure and train last three layers with learning rate initialized in step1.

After training of last layers freeze all layers in model except unfreezing middle layers from 173-249 This layers are used for fine tuning to get inner attributes.

Using less amount of time using Kaggle platform to run our network on GPU which is freely provided to us by Kaggle of 12 GB ram of NVIDIA graphics.

#### **TEST PHASE**

Once our model is trained, it can now integrate in our system to test new objects.

#### **FLOWCHART**

For implementation two flow charts are suggested one for Training and other for Real Time implementation. Dataset is created using Videos generated by human strategy for picking up objects Instead of generating videos we have used objects image directly images. Model is created using pre trained model of InceptionV3 popping last layer and adding 2 fully connected layers and SoftMax layer for multi class classification. Dataset is resized to input require size of model i.e. 416\*416 and then split into 80% as training data and 20% as validation data.

Freeze all layers except added layers and compile model with batch size of 10 and learning rate of 0.001 Adam as optimizer and train model with 30 epochs on gpu. Freeze all layers except layers from 72-249 layers and compile model with the same parameters as mentioned above except learning rate of 0.00001 and train model.

Model is trained and ready for prediction.

#### **IMPLEMENTATION FLOWCHART**

Input is taken from Webcam and use OpenCV library to process on frames got from webcam and capture video using webcam is setup the size of frames as per YOLOv4 model requirement and send frames to Object Detection Model. YOLOv4 method is selected for detecting of object frames of webcam are given to YOLOv4 whereby it randomly generates boxes in frames and each frame is passed through Darknet 53. Model which has convolution layers which gives 3d tensors giving parameter for detected object box. After completion of object detection object is selected which is bounded by boxes present near to center of frames Image is selected and

resized to required size of trained model i.e.  $416 \times 416$  Model which is trained previously is used here to predict one of 6 primitives defined earlier. The Class selected which gives more probability below graph shows confusion matrix for trained model.

Anticipatory action is divided into two phase

- Approach Phase
- Grasp Phase

Approach phase can be studied from [3] and further robotic action is not performed

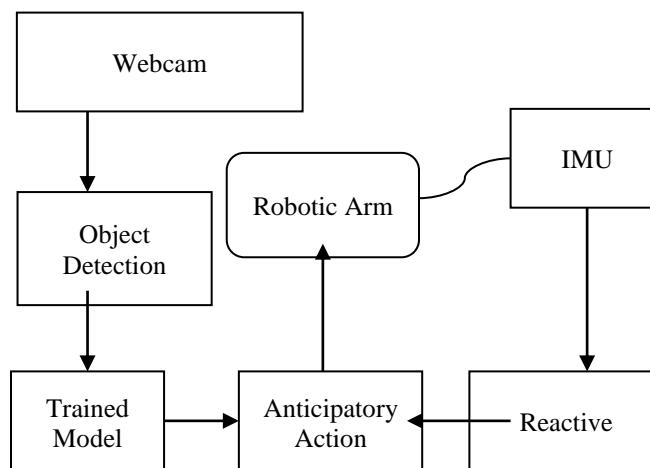


Figure VI Implementation Flowchart

#### IV) RESULTS

- Model is trained on 45 objects which gives high accuracy to 97% and validation accuracy to 93% as shown above
- Model is Tested our trained model on 10 objects which gives results of overall average of 84% which works very good with lateral top and top pinch but slight less accuracy as shown in confusion matrix with slide and flip, it works worst with bottom primitive giving accuracy of 46%
- YOLOv4 which gave 51ms of response which is more than YOLOv3 and show good performance on gpu



```
Train for 584 steps, validate for 156 steps
Epoch 1/5
584/584 [=====] - 144s 247ms/step - loss: 2.8217
- accuracy: 0.5205 - val_loss: 2.6481 - val_accuracy: 0.3359
Epoch 2/5
584/584 [=====] - 124s 212ms/step - loss: 0.6282
- accuracy: 0.8021 - val_loss: 2.2706 - val_accuracy: 0.6205
Epoch 3/5
584/584 [=====] - 124s 213ms/step - loss: 0.3554
- accuracy: 0.9332 - val_loss: 0.1948 - val_accuracy: 1.0000
Epoch 4/5
584/584 [=====] - 124s 212ms/step - loss: 0.1832
- accuracy: 0.9800 - val_loss: 0.2815 - val_accuracy: 0.9109
Epoch 5/5
584/584 [=====] - 124s 213ms/step - loss: 0.1783
- accuracy: 0.9783 - val_loss: 0.1976 - val_accuracy: 0.9333
```

Figure VII Training Output

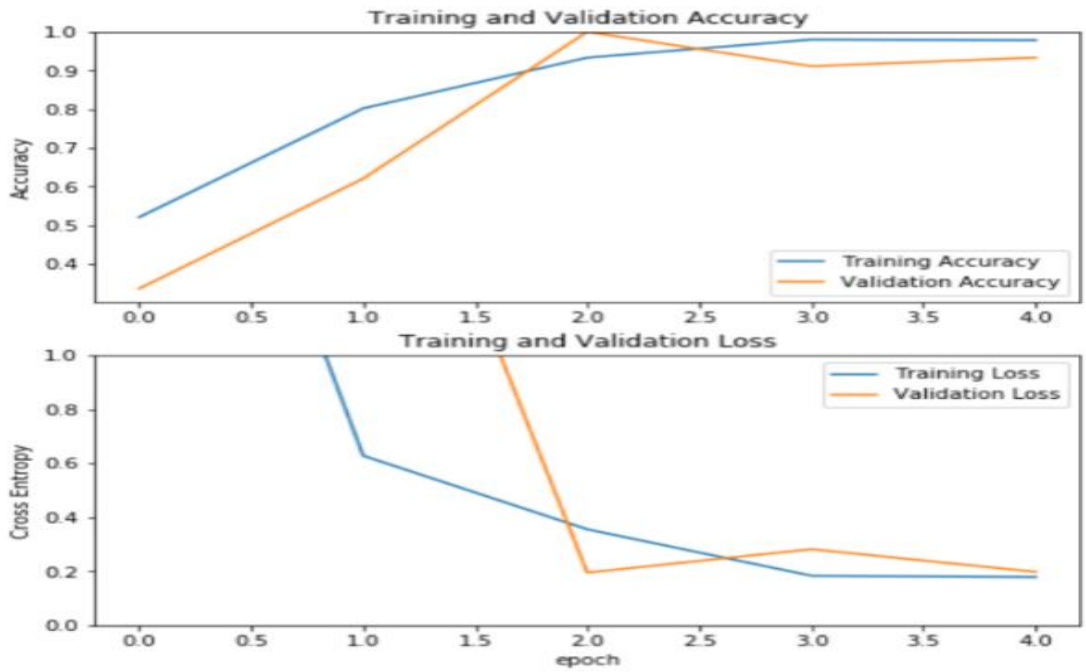


Figure VIII Training Graph

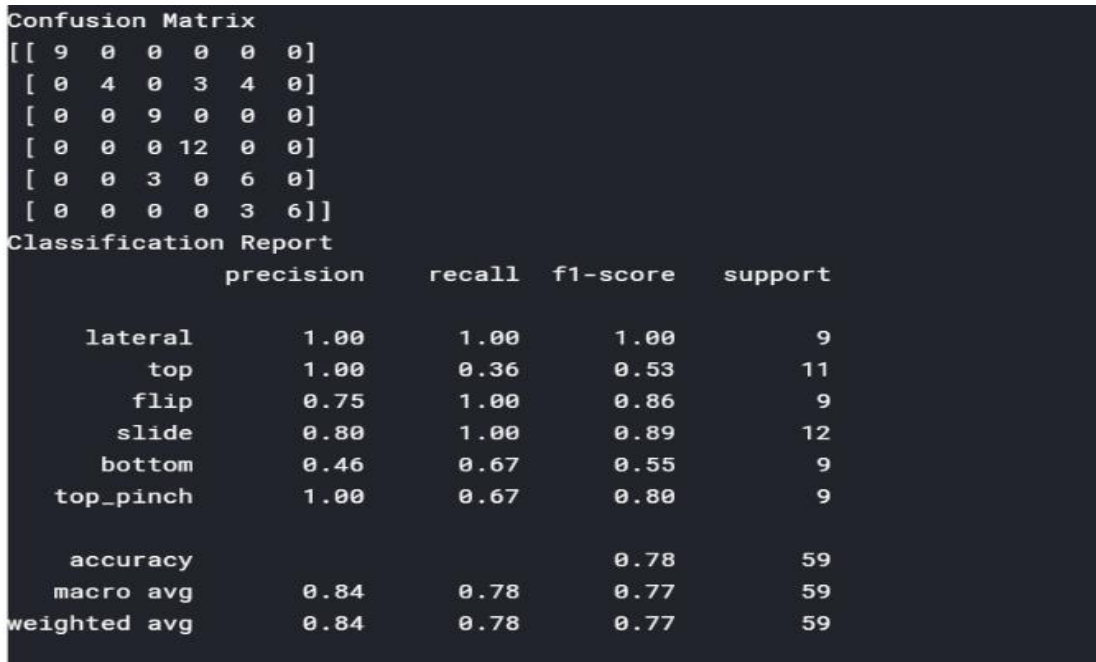
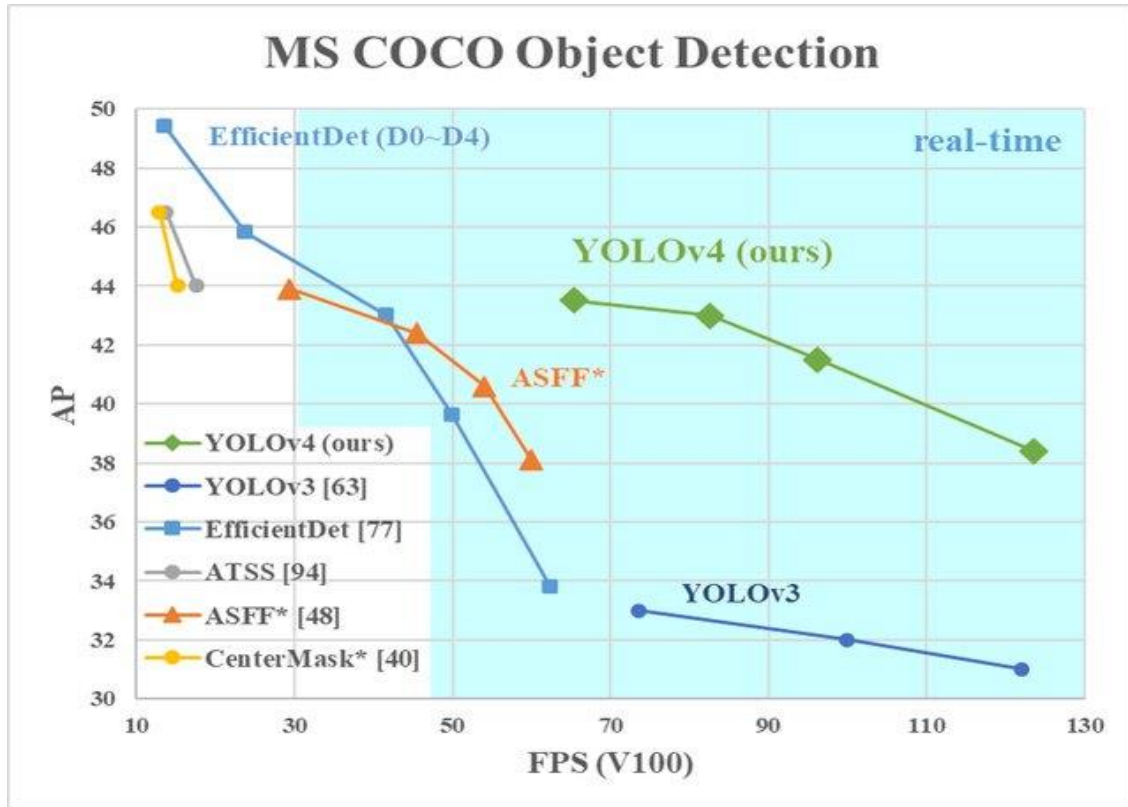


Figure IX Predicted Data



Source:- <https://ai-pool.com/a/s/yolov3-and-yolov4-in-object-detection>

## V) CONCLUSION

Our work includes use of YOLOv4 instead of YOLOv3 which detects objects more accurately. Proposed and validated a Primitive detection using deep learning and performing using soft hand Goal can be gained by:

- i) Creating new model by using transfer learning methodology and using new Adam optimizer.
- ii) Creating new primitive to increase accuracy of model to predict actions and reducing primitive used in previous work,
- iii) Using Neural Network methodology in robotic system to make it more intelligent.
- iv) Testing on wide range of object that are not used in training data.

Further in the next approach or future work stereo camera will be used to get more accuracy.

## REFERENCES

1. C. Erdogan, A. Schroder, and O. Brock, "Coordination of intrinsic and extrinsic degrees of freedom in soft robotic grasping," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1–6.

2. M. Haas, W. Friedl, G. Stillfried, and H. Hoppner, "Human-robotic variable-stiffness grasps of small-fruit containers are successful even under severely impaired sensory feedback," *Frontiers in neurorobotics*, vol. 12, p. 70, 2018.
3. J. Zhou et al., "A Soft-Robotic Approach to Anthropomorphic Robotic Hand Dexterity," in *IEEE Access*, vol. 7, pp. 101483-101495, 2019.
4. YOLOv3: An Incremental Improvement Joseph Redmon, Ali Farhadi University of Washington
5. J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," arXiv preprint, 2017.
6. [https://medium.com/@jonathan\\_hui/real-time-object-detection-with-yolo-yolov2-28b1b93e2088](https://medium.com/@jonathan_hui/real-time-object-detection-with-yolo-yolov2-28b1b93e2088)
7. X. Yan, J. Hsu, M. Khansari, Y. Bai, A. Pathak, A. Gupta, J. Davidson, and H. Lee, "Learning 6-dof grasping interaction via deep geometryaware 3d representations," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1–9.
8. P. Schmidt, N. Vahrenkamp, M. Wächter, and T. Asfour, "Grasping of unknown objects using deep convolutional neural networks based on depth images," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 6831–6838.
9. A. Zeng, S. Song, K.-T. Yu, E. Donlon, F. R. Hogan, M. Bauza, D. Ma, O. Taylor, M. Liu, E. Romo et al., "Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching," in 2018 IEEE International Conference on Robotics and Automation (ICRA). IEEE, 2018, pp. 1–8.
10. T. Nishimura, K. Mizushima, Y. Suzuki, T. Tsuji, and T. Watanabe, "Thin plate manipulation by an under-actuated robotic soft gripper utilizing the environment," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 1236–1243.
11. M. Bianchi, G. Averta, E. Battaglia, C. Rosales, A. Tondo, M. Poggiani, G. Santaera, S. Ciotti, M. G. Catalano, and A. Bicchi, "Tactilebased grasp primitives for soft hands: Applications to human-to-robothandover tasks and beyond," in *Robotics and Automation (ICRA), 2018 IEEE International Conference on*. IEEE, 2019.
12. [2015 CVPR] [GoogLeNet / Inception-v1] Going Deeper with Convolutions
13. C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
14. <https://medium.com/@sh.tsang/review-inception-v3-1st-runner-up-image-classification-in-ilsvrc-2015-17915421f77c>