# Hard Hat Detection Masked With Color Detection

Ashish Kumar Singh(2019BCS-010), Gurjot Singh(2019BCS-021), Vibhor Sharma(2019BCS-070)

*Dept. of Computer Science and Engineering*

*Atal Bihari Vajpayee Indian Institute of Information Technology and Management*

https://www.overleaf.com/project/6259acb3b77aeacb034b205f Gwalior, India

aks15o8o2001@gmail.com, gurjotsingh21003@gmail.com, vibhors352@gmail.com

*Abstract*—In the United States in 2019, there were 1061 fatalities and many more injuries at construction sites. This was a 5.3 % increase from 2018. The same number in India was around 13000 fatalities. A considerable number of people died due to incidents such as falls, slips, and trips, as well as getting struck by objects that fell on them. This was frequently the outcome of not wearing safety equipment, which happened either because of a lack of knowledge or ignoring the significance of protective gear. In order to decrease fatalities resulting from such accidents, the Occupational Safety and Health Administration (OSHA) provides Fall Prevention and OSHA-10 training to construction workers. This training is intended to help prevent such incidents from occurring at construction sites. Furthermore, safety professionals keep an eye on whether workers are appropriately using personal protective equipment (PPE). Construction deaths have reduced by 2% annually since 1994, according to data; yet, the owners are still dissatisfied with this result. According to many researches, falls are the leading cause of construction fatalities. According to one study, half of the fatalities from falls were caused by workers who either did not use PPEs or did not utilize them effectively. Studies have demonstrated that wearing hard helmets appropriately can prevent deaths resulting from incidents such as falls, slips, trips, and objects falling on workers. To address this issue, researchers developed and assessed a hard-hat detection program that utilizes deep learning and image processing methods to identify whether or not employees are wearing protective helmets. This research should help authorities to monitor the safety of workers at construction sites by automating the process that will reduce the number of accidents

*Index Terms*—PPE,hard-hat detection,image processing,deep learning.

## I. INTRODUCTION

Construction is one of the most dangerous industries where day to day work requires workers to perform hazardous tasks which pose injuries and fatalities to the workers. A considerable proportion of these deaths are the result of accidents like falls, slips, and trips, as well as being hit by falling objects. According to OSHA, around 34% of all construction-related fatalities in the United States in the years around 2012-13 were caused by falls. This figure was as high as 50% in the 1980s and 1990s. Workers fall from great heights and smash their heads on hard floors in the majority of fall accidents. According to an investigation report, half of the fall accidents occurred at a height of less than 3 m; additionally, 57 percent of the fall incidents occurred on ladders, rooftops, construction sites, platforms, or scaffolding. Construction accidents vary in severity, and even a "minor" mishap can have long-term consequences. The impact of traumatic brain injuries on the skull can be categorized. Even if there is no fracture, a TBI can result in internal bleeding, degenerative brain disorders, and death. Another consideration, especially in the case of falls, is the sort of impact that produced the injury. Straight-on collisions with the head are linked to conditions including skull fractures.

Hard helmets can face up to and offer protection towards incidents including shock, object penetration, and contact with electrical hazards. therefore a few of the fatalities from falls and a widespread number of fatalities from slips, journeys, and being struck by using falling items can be decreased if personnel wore hard hats as it should be. according to the findings of one investigation into the frequency of production fatalities and the usage of private shielding system (PPEs), 47.3 percentage of fatally injured sufferers both did now not use PPEs or did not use them appropriately.

consequently a want arises to offer innovation and automation into the safety system. through utilizing real-time video streaming from the worksite to display employees, protection engineers and relevant authorities can be able to put in force difficult hat safety guidelines greater efficiently. This approach has the ability to reduce the frequency of accidents. This studies created a technique for identifying employees who do not put on safety helmets on construction sites.

## II. LITERATURE REVIEW

### A. Convolutional Neural Network

A CNN [1], or convolutional neural network, or convnets, is a deep learning neural network designed to analyze structured arrays of data. CNNs are good at detecting features such as lines, gradients, circles, and even eyes and faces in input image. This makes them pretty useful in computer vision tasks. The main advantage of CNNs over Regular Neural Networks is that they detect important features in image without any human intervention. Also the number of parameters that are learned during the training process in CNNs are comparatively less compared to regular NN. Regular NN in spite of being computationally more expensive than CNNs are not efficient for computer vision tasks.

A CNN consists of mainly 4 types of layers: **Input layer**, **Convolution layer**, **Activation layer** and **Pooling layer**.

*1) Input layer:* This layer contains the image as raw input.

*2) Convolution layer:* This layer gives the output by calculating the dot product between the convolution filter and image patches. The dimensions of the output layer are reliant upon

the height and width of the input, height, width and depth of the convolution layer, padding and stride length.

*3) Activation layer:* This layer applies an activation function to every element of the output of the Convolution layer. Some examples of activation functions include Sigmoid Activation function, ReLu Activation function and tanh.

*4) Pooling layer:* This layer is used to reduce the volume in a CNN which helps reduce memory usage and makes the computation fast. There are usually two types of pooling layers: Max Pool and Average Pool.
The dimensions of the output after applying pooling depends upon the height and width of the input, height, width and depth of the pooling layer, padding and stride length.
**Max Pool** takes the maximum of all the elements over an image patch for a given pool layer size.
**Average Pool** takes the average of all the elements over an image patch for a given pool layer size.

These layers can be made to pile up on top of each other in different combinations to form different Convolution Neural Network Architectures(for example, ResNets [2], Inception Network [3]).

### B. EfficientDet

EfficientDet architecture [4] is a Google Brain-created object detection algorithm. EfficientDet is built on EfficientNet, a classification convolutional neural network that was pre-trained on the ImageNet image database.
EfficientDet pools and blends distinct picture granularities to build features, which are subsequently sent through a NAS-FPN feature fusion layer. The NAS-FPN will integrate number of characteristics at different levels of granularity and feeds them to the detection head, which predicts bounding boxes and class labels.The study [4] closely checks the scaling and object detection model tradeoffs. A comparison is shown in 3.

### C. SSD (Single Shot Detector)

SSD is another popular algorithm used in object detection. It provides high speed as well as accuracy using relatively low resolution images. High detection accuracy is accomplished by utilizing multiple boxes or filters of various sizes and aspect ratios.
Wei Liu [5] presents this approach; to achieve high detection accuracy, prediction for the bounding boxes for different objects in the image is done by multiple feature maps of different sizes that represent multiple scales rather than a single feature map. For feature extraction VGG-16 base network is used as it is standard CNN architecture for image classification. To this base network additional convolutional layers are added for detection. The size of the layers decreases on moving towards the end of the network.
As shown in 1 in the input image a bounding box is formed for each object present in that image. These boxes of different aspect ratio are evaluated by convolutional layers. This is done to find the default box that overlaps with the bounding box containing objects. Two default boxes are shown below 2.
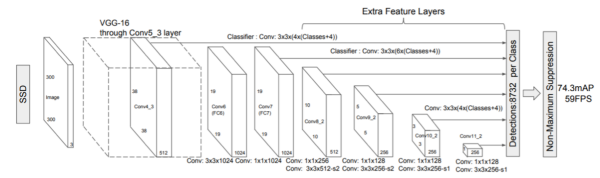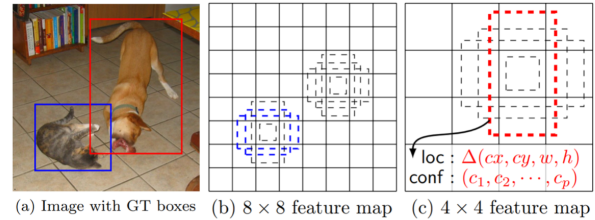


Fig. 1. SSD Architecture



(a) Image with GT boxes  (b) $8 \times 8$ feature map  (c) $4 \times 4$ feature map

Fig. 2. SSD Bounding Box

These boxes that are containing objects are considered as positive boxes and the rest are ignored. $\Delta cx$, $\Delta cy$, h and w, represent the offsets from the center of the default box and its height and width.

### D. YOLO

YOLO [6] stands for **You Only Look Once**. YOLO models are well known for giving high performance at the same time as additionally requiring much less area for its operation, making them perfect applicants for real time object detection eventualities. This technique uses a unmarried neural community to evaluate the overall image, then separates it into sections and predicts bounding boxes and possibilities for each thing. those bounding containers are weighed the usage of the predicted chance. The predictions are made after best
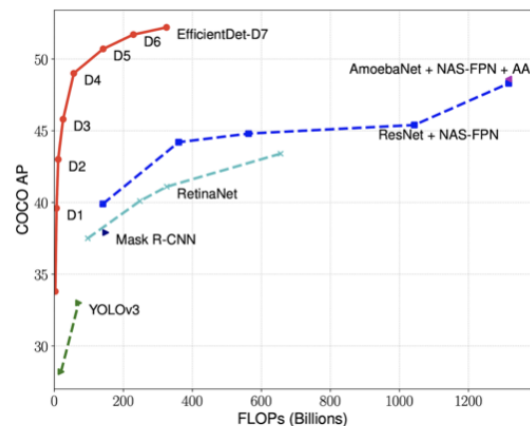


Fig. 3. Comparison between different object detection algorithms

one forward run thru the neural community. After non-max suppression, it presents found objects.

## III. Summary

YOLO, EfficientDet and Single Shot Detection are pioneer steps in the field of Object Detection. They all provide with real time object detection with newer versions of these techniques providing better Speed vs Accuracy tradeoff than the older versions.

However, these techniques are incomplete in addressing the camouflage problem. They also fail to separately segregate objects which are visually similar and are closer to each other. Also, these techniques do not yield the same results for datasets containing images taken from irregular angles and having low lighting.

## IV. Methodology

We have used the YOLOv5(small) model architecture for training our safety helmet detector. The **Model Architecture** has been described in A, B and C subsections.

### A. Backbone

In convolutional neural networks, specifically with admire to the item detection duties, backbones refers back to the a part of the network that is accountable for characteristic extraction. it is part of the community that sees the input. function extraction is used to reduce a variable sized photograph to a fixed set of visual features.

*1) CSPDarknet53:* 5 shows CSPDarknet53 architecture. This has been used for feature extraction. In this section the working of CSPDarknet53 is explored.

**Conv n\*f\*f, stride s**, is nothing but a Convolution layer with filter size as f, number of filters n and stride length s.

**CSPBlock X num** in CSPDarknet53 represents that num CSPBlocks are stacked in the CSPDarknet53 one after each other.

In the CSPBlock where the Base Layer(height\*width\*channels) gets separated into Part 1 and Part 2, where the data in the base layer is divided into two equal parts, one part going to Part 1(height\*width\*(channels/2)) and the other part going to Part 2(height\*width\*(channels/2)).

**Add** in CSPBlock is used to add two different inputs of the same dimensions element by element.

**Mish** is an activation function used in Convolution Neural Networks.

$$f(x) = x.tanh(softplus(x) \tag{1}$$

$$tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \tag{2}$$

$$softplus(x) = ln(1 + e^x) \tag{3}$$

$$tanh(softplus(x)) = \frac{e^{softplus(x)} - e^{-softplus(x)}}{e^{softplus(x)} + e^{-softplus(x)}}$$
$$= \frac{e^{ln(1+e^x)} - e^{-ln(1+e^x)}}{e^{ln(1+e^x)} + e^{-ln(1+e^x)}} \tag{4}$$

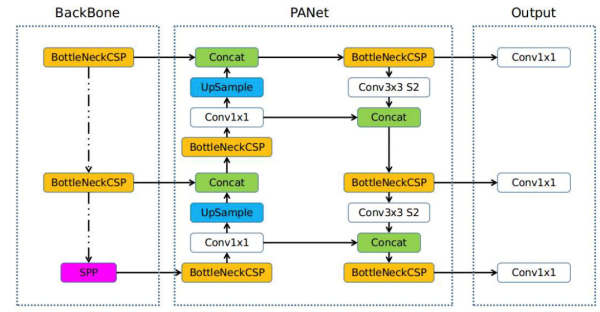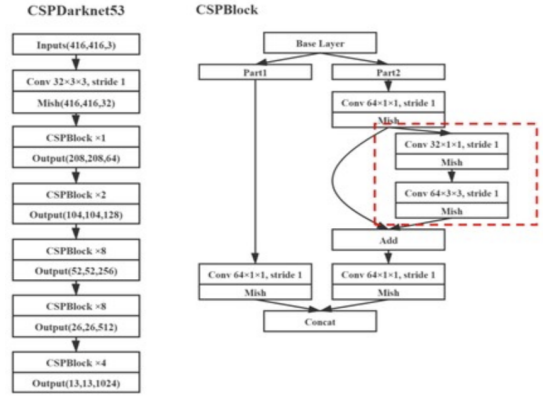

Fig. 4. Overview of YOLOv5 Model Architecture



Fig. 5. Overview of CSPDarknet53 Architecture

$$tanh(softplus(x)) = \frac{1 + e^x - \frac{1}{1+e^x}}{1 + e^x + \frac{1}{1+e^x}}$$
$$= \frac{(1 + e^x)^2 - 1}{(1 + e^x)^2 + 1} \tag{5}$$

$$f(x) = x.\frac{(1 + e^x)^2 - 1}{(1 + e^x)^2 + 1} \tag{6}$$

**Concat** is used to stack two inputs of the same dimensions one over another.

### B. Neck

Neck is used for creating a pyramid feature in the input image. It helps in determining the scaling factor of observed items of similar nature but different scales.

*1) Spatial Pyramid Pooling(SPP):* we've got the capabilities map as the output of the convolution neural networks, which is generated by our various filters acting convolution operation. to place it every other manner, we will have a clear out that can discover circular geometric styles and construct a feature map that highlights these shapes even as preserving the shape's role in the photo. whatever the length of our feature maps, the Spatial Pyramid Pooling Layer lets in us to generate constant length features, this allows in reducing the scale of the map. that is executed via the use of pooling layers, such as Max Pooling and common Pooling, to provide distinct
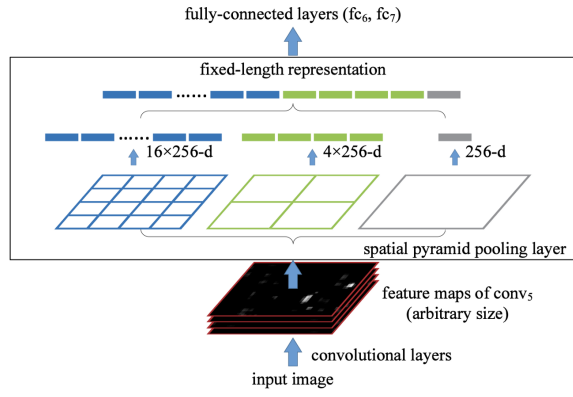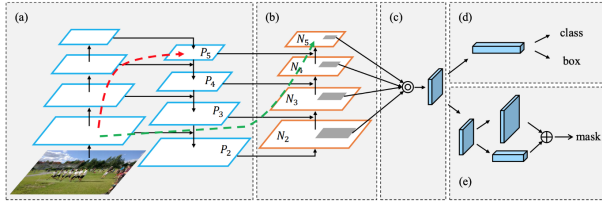
Fig. 6. Overview of SPP



Fig. 8. YOLOv3 head



Fig. 7. Overview of PaNet



Fig. 9. IoU

representations of our function maps in order to generate a hard and fast length. The pooling process is applied to each feature map to reduce it to a single value, resulting in a vector with a size of 1256. Subsequently, each feature map undergoes pooling to produce vectors of sizes 4256 and 16256 by reducing the number of values to 4 and 16, respectively. These vectors, along with another vector produced in the same manner, are then combined or concatenated to generate a vector of a fixed size. This vector serves as the input for the fully connected network.

*2) PANet:* PANet stands for Path Aggregation Network. This architecture allows the extra efficient propagation of layer information both from the bottom to the pinnacle or vice versa. The neck additives normally glide between tiers in an up and down path, linking best the final few layers of the convolutional community. In figure 7, the first layer's facts is mixed with layer P5 (indicated by way of the purple arrow) and propagated to layer N5 (indicated by using the inexperienced arrow). This serves as a shortcut for conveying low-level statistics to the upper levels of the community. The authentic PaNet technique combines the contemporary layer with records from a previous layer to generate a new vector. In our Yolo implementation, we have changed PaNet by concatenating the input vector with the vector from a previous layer to produce a new vector. This allows us to include statistics into the top layers of the network. subsectionHeadHead is the top of a neural network. it's miles the layer in the neural network that is accountable for generating the final output.

*3) YOLOv3 Head:* YOLOv3 has a multi-head detection system. At 3 different particular instances in our neural net-
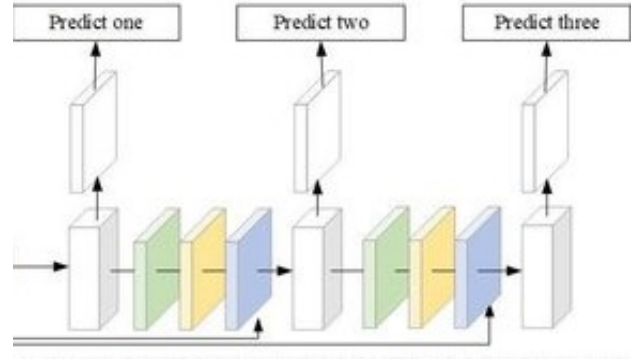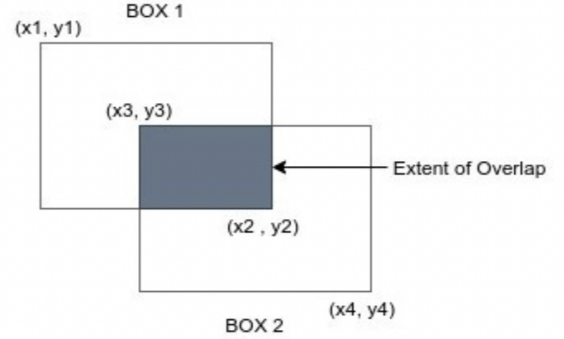
work a prediction is made..

*C. Loss Function*

$$Loss = \lambda_1 L_{loc} + \lambda_2 L_{cls} + \lambda_3 L_{obj}$$

The YOLOv5 loss consists of 3 parts:

*1) Location Loss ($L_{loc}$):* **IOU** (Intersection over union) as shown in 10 is defined as:

$$IoU = \frac{Area of Intersection}{Area of Union}$$

For this particular scenario, BOX 1 and BOX 2 can be considered as predicted Bounding Box and ground truth Bounding Box.

$$L_{IoU} = IoU$$

**GIoU**(Generalised IoU Loss) maximizes the ground truth and anticipated bounding box overlap area. For non-overlapping scenarios, it increases the predicted box's size to overlap with the target box by moving gradually towards the target box.

$$L_{GIoU} = 1 - IoU + \frac{|C - Area of Union|}{|C|}$$

where C is the area of the smallest box covering the predicted and ground truth bounding boxes.

Fig. 10. DIoU

**DIoU**(Distance IoU Loss)

$$L_{DIoU} = 1 - IoU + \frac{d^2}{C^2}$$

here d is the distance between the center points of ground truth bounding box and predicted bounding box and C is the diagonal length of the smallest enclosing box covering both the boxes.

**CIoU** Complete IoU takes three crucial geometric factors, namely overlap area, center point distance, and aspect ratio, should be considered when calculating a good loss for bounding box regression. CIoU Loss is calculated by adding all three:

$$L_{CIoU} = 1 - IoU + \frac{d^2}{c^2} + \alpha\nu$$

here d is the distance between the center points of ground truth bounding box and predicted bounding box and C is the diagonal length of the smallest enclosing box covering both the boxes as mentioned in DIoU above.

$$\nu = \frac{4}{\pi^2}(arctan(\frac{w^{gt}}{h^{gt}}) - arctan(\frac{w}{h}))^2$$

where $w^{gt}$ and $h^{gt}$ is the width and height of ground truth bounding box and w and h is the width and height of the predicted bounding box.

$$\alpha = \frac{\nu}{1 - IoU + \nu}$$

$L_{loc} = L_{CIoU}$

*2) Classes Loss ($L_{cls}$)* :

$$L_{cls} = \sum_{i=0}^{S^2} \chi_i \sum_{c \epsilon classes} (BCE(p_i(c), p_ic))$$

$$BCE(x,y) = -xlog(y) - (1-x)log(1-y)$$

here $S^2$ represents the total number of grid cells our input image is divided into. $\chi_i$ is 1 if an object appears in cell i otherwise 0.
$p_i(c)$ denotes the conditional class probability for class c in cell i.

$$ccp \equiv P_r(Class_i|object)$$

$P_r(Class_i|object)$ is the probability an object belongs to class i given an object is present and $ccp$ is the conditional class probability.

*3) Objectness Loss ($L_{obj}$)* :

$$L_{obj} = \sum_{i=0}^{S^2} \sum_{j=0}^{B} (\chi_{ij}^{obj} * BCE(C_i, C_i')) +$$

$$\lambda_{noobj} \sum_{i=0}^{S^2} \sum_{j=0}^{B} (\chi_{ij}^{noobj} * BCE(C_i, C_i') * (1 - \chi_{ij}^{obj}))$$

where

$$BCE(x,y) = -xlog(y) - (1-x)log(1-y)$$

here $S^2$ represents the total number of grid cells our input image is divided into and B denotes the total number of bounding boxes.
$\chi_{ij}^{obj}$ is 1 if jth boundary box in cell i has an object else 0.
$C_i$ denotes the box confidence score of box j in cell i.

$$bc \equiv P_r(object) \cdot IoU$$

$P_r(object)$ is the probability the box contains an object and $bc$ is the box confidence score.

*D. Training*

We performed training for 100 epochs on a GPU Hardware Accelerator and took batch size as 100.For every epoch we keep track of precision, recall, mAP@0.5 and mAP@0.5:0.95 on our training dataset.
**True Positive** is when the model predicts the positive class correctly.
**False Positive** is when the model predicts the positive class incorrectly.
**False Negative** is when the model predicts the negative class incorrectly.
**True Negative** is when the model predicts the negative class correctly.
**Precision** is a measure of how accurate our predictions are and **Recall** is a measure of how well we find out all the positives.

$$Precision = \frac{TP}{TP + FP} \qquad (7)$$

$$Recall = \frac{TP}{TP + FN} \qquad (8)$$

True Positives, False Positives, True Negatives and False Negatives are calculated using IoU for a given IoU threshold value. This means that for every IoU threshold value we will have different values for Precision and Recall.
AP(Average Precision) for a class is calculated by finding out the area under the precision recall curve for a given class label.
**mAP(mean Average Precision)** is calculated by finding out the average of AP over all the classes.

$$mAP = \sum_{c \epsilon classes} \frac{(AP_c)}{total number of classes} \qquad (9)$$
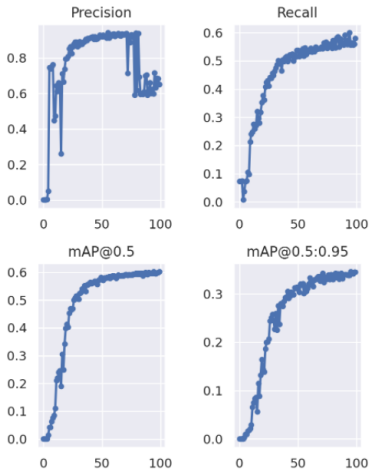
Fig. 11. Precision, Recall, map@0.5, map@0.5:0.95 values on our training dataset from 0 to 100 epochs

**mAP@0.5** means mAP calculated for an IoU threshold value of 0.5.

**mAP@0.5:0.95** means average mAP calculated over different IoU threshold values(0.5, 0.55, 0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95).

### E. Helmet Color Detection

In this section we have explained how we have tried to detect the color of our helmet in case helmet is detected. Here, in an attempt to produce the bounding box with the same color as the helmet, we tried to capture ROI(Region of Interest) in that bounding box and tried to get the most dominant color in that bounding box. As of now we have taken the top 25% of the bounding box as our region of interest. We can toy around with this metric.

Our most dominant color is returned in the form of RGB values. We use the color indicated by those RGB values and plot the box(see Figure 12).

We have produced black color bounding boxes for cases when our model detects a head(not helmet)(see Figure 13).

## V. EXPERIMENTAL SECTION

### A. Dataset

We took 4750 images of workers in a construction site along with their annotation in Pascal VOC Format. 70% of those images were used for training, 20% for validation and 10% for testing. We also performed a preprocessing step by stretching our image to a specific dimensions(416*416 in our case) so as to make sure that all the input images in our input data are of the same dimensions for our model.

### B. Results

Before we started training our model mAP@0.5 value on our training dataset was 0.000417. After training for mAP@0.5 score on training dataset improved to 0.667. Our validation and training dataset gave a mAP@0.5 score of 0.6



Fig. 12. Blue coloured bounding box plotted in our image for a detected helmet of the same color
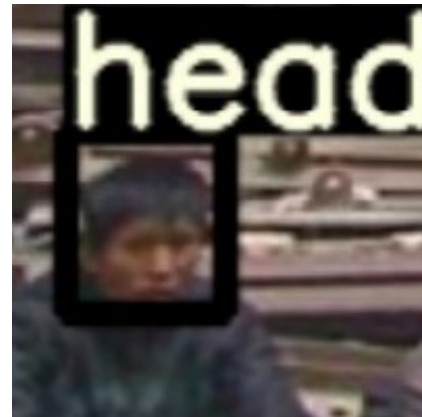


Fig. 13. Black coloured bounding box plotted in our image when our model detects head

and 0.602 respectively. Training our entire model took 2.319 hours.

| Sets | Precision | Recall | mAP@0.5 | mAP@0.5:0.95 |
|------|-----------|--------|---------|--------------|
| **train** | 0.661 | 0.658 | 0.667 | 0.382 |
| **val** | 0.587 | 0.568 | 0.6 | 0.345 |
| **test** | 0.932 | 0.557 | 0.602 | 0.353 |

TABLE I
PRECISION, RECALL, MAP@0.5, MAP@0.5:0.95 VALUES ON OUR TRAINING, VALIDATION AND TEST DATASET FOR OUR TRAINED MODEL.

## VI. CONCLUSION

YOLOv5 even though hasn't been officially published like previous versions of YOLO can be used for our object detection tasks and can be fine tuned to our specific application by training its weights with images of our use case instances. The

model could be trained further on a larger dataset and for more epochs to give better results in real world scenarios. There is also a possibility for exploration on whether doing so(training our model for more epochs and on larger and diverse dataset) would address the problems which have been left unanswered by previous Object Detection Techniques with respect to our task.

Also for Hard Hat color detection, if we have been given possible colors of safety helmets in our dataset, the performance of our model with respect to detecting the color of detected hardhat could be improved. We could compare the most dominant color RGB values in our Region of Interest with RGB values for different colors and use the color which has the least distant RGB values.

## REFERENCES

[1] S. Albawi, T. A. Mohammed, and S. Al-Zawi, "Understanding of a convolutional neural network," in *2017 International Conference on Engineering and Technology (ICET)*, 2017, pp. 1–6.

[2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1–9.

[4] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," 2019. [Online]. Available: https://arxiv.org/abs/1911.09070

[5] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot MultiBox detector," in *Computer Vision – ECCV 2016*. Springer International Publishing, 2016, pp. 21–37.

[6] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 779–788, 2016.