# SH17: A Dataset for Human Safety and Personal Protective Equipment Detection in Manufacturing Industry

Hafiz Mughees Ahmad[†*] and Afshin Rahimi[§]

Mechanical, Automotive and Materials Engineering Department,
University of Windsor, Windsor, ON, Canada
[†]ahmad54@uwindsor.ca, [§]arahimi@uwindsor.ca

## Abstract

Workplace accidents continue to pose significant risks for human safety, particularly in industries such as construction and manufacturing, and the necessity for effective Personal Protective Equipment (PPE) compliance has become increasingly paramount. Our research focuses on the development of non-invasive techniques based on the Object Detection (OD) and Convolutional Neural Network (CNN) to detect and verify the proper use of various types of PPE such as helmets, safety glasses, masks, and protective clothing. This study proposes the SH17 Dataset, consisting of 8,099 annotated images containing 75,994 instances of 17 classes collected from diverse industrial environments, to train and validate the OD models. We have trained state-of-the-art OD models for benchmarking, and initial results demonstrate promising accuracy levels with You Only Look Once (YOLO)v9-e model variant exceeding 70.9% in PPE detection. The performance of the model validation on cross-domain datasets suggests that integrating these technologies can significantly improve safety management systems, providing a scalable and efficient solution for industries striving to meet human safety regulations and protect their workforce. The dataset is available at https://github.com/ahmadmughees/sh17dataset.

## Index Terms

SH17, Object Detection, Convolutional Neural Network, YOLO, Personal Protective Equipment, Worker, Human Safety, Dataset

## I. INTRODUCTION

Occupational Safety and Health (OSH) is a multidisciplinary field that ensures safety in work environments and avoids risks to a worker's health. This includes preventing work-related injuries, accidents, and illnesses by managing hazards and promoting healthful practices. PPE is crucial in industrial manufacturing for safeguarding workers against various occupational hazards. The importance of PPE lies in its ability to minimize the risks of injuries and illnesses from exposure to dangerous conditions such as toxic chemicals, extreme temperatures, and mechanical hazards. Research emphasizes that PPE is an essential last line of defense in protecting workers when other safety measures are insufficient [1].

Based on the guidelines from the Occupational Safety and Health Administration (OSHA) [2], and Vukicevic et al. [3], this study divides PPE into five categories according to the body parts it protects; 1) Head 2) Upper Body 3) Hands 4) Feet and 5) Whole body. Common examples of these PPEs include helmets, safety glasses, face shields, earmuffs, safety vests, gloves, and safety shoes, each designed to address specific threats. Studies have shown that proper selection, usage, and maintenance of PPE are vital for ensuring its effectiveness and the safety of the employees.

Manufacturing industries are trying their best to ensure the safety of their workforce and use of PPE in their facilities to reduce the injuries causing millions of dollars in damages. Earlier systems use sensor-based approaches for automatic PPE compliance but with the recent popularity of Machine Learning (ML) and Computer Vision (CV) systems, researchers have been proposing non-invasive solutions for PPE compliance for the manufacturing facilities [4], [5]. Vukicevic et al. [6] proposed 2 stage approach for the detection of PPE using person and keypoint detection in the first step and then detecting the PPE type in the second stage using OD model. Yu et al. [5] proposed to use the You Only Look Once (YOLO)v5 OD model for detecting PPE in chemical plants.

In this study, we also propose to use OD-based methods for the PPE compliance. To make this study more accessible and help other researchers pursue advancement in this field, we offer an open-source dataset consisting of 8,099 images and 75,994 instances. We compare this dataset with existing benchmark datasets and relevant studies and provide open-source weights for PPE detection to make it easy for others to recreate the results.

The main contributions of this work are,

- A collection of a high-quality, large-scale dataset from the internet. The dataset contains images from all over the world, removing the location and gender bias and making it inclusive. The dataset is open-sourced for commercial and research purposes.
- Extensive training of the state-of-the-art OD models with publicly available weights for community use.
- Evaluation of trained model with different dataset to test the efficacy across out-of-distribution data.

The remainder of this paper is structured as follows: a literature review and background are provided in Section II. The collected dataset is introduced in Section III. The insights from the experiments and results are provided in Section IV. Finally, Section V offers concluding remarks and future directions.

## II. LITERATURE REVIEW

Human safety inspectors are still the most common and easiest way to ensure PPE compliance at any workplace. However, it is costly and laborious work with a high error margin. Earlier, the wide usage of Internet of Thing (IoT) devices in industrial manufacturing introduced sensors for PPE compliance in the process. However, that is also a very costly and sensitive process [7], [8]. Recently, Computer Vision (CV) methods have emerged as a non-invasive solution, providing cheaper and better alternatives to the sensors [3], [9], [10] and well accepted in construction as well as manufacturing. These methods use the annotated data to train the ML models to detect the required objects. In the following section, we have discussed the existing datasets for PPE detection and OD methods used in the literature.

### A. Existing Datasets

Most existing datasets for PPE detection are tailored towards hardhat (often called helmets) detection in construction engineering.

*1) Hardhat Wearing Detection (GDUT-HWD) Dataset:* The GDUT-HWD dataset [11] contains hardhat images collected from internet sources in 5 different colors labeled individually. It contains 3,174 images and 18,893 object instances within multiple sizes. Due to crowd-sourcing, its publicly available version[1] contains a lot of advertisement images with non-relevant backgrounds.

*2) Safety Helmet Wearing (SHW) Dataset:* The SHW dataset [12] is a comprehensive public dataset[2] used for detecting both safety helmets and heads. It comprises 7,581 images featuring 9,044 instances of humans wearing safety helmets and 111,514 instances of humans without safety helmets. The examples were sourced from search engines and a portion from the SCUT-HEAD dataset [13] and manually annotated afterward.

---

[1]https://github.com/wujixiu/helmet-detection
[2]https://github.com/njvisionpower/Safety-Helmet-Wearing-Dataset

*3) Color Helmet and Vest (CHV) Dataset:* The CHV dataset [14] is publicly available[3] and consists of 1,330 high-quality images labeled into 6 categories of person, vests, and colored helmets in 4 colors. The authors selected images from GDUT-HWD [11] and SHW [12] datasets using the strict criteria of related construction background, gestures of workers, and object angles and distances from the camera.

*4) SHEL5K Dataset:* The SHEL5K dataset [15] is an improved version of Safety Helmet Detection (SHD) dataset[4] [16] consisting of 5,000 images with 75,570 instances in 6 classes, including head, helmet, face, and person with and without a helmet. The dataset is publicly available[5].

*5) Pictor-PPE Dataset:* The Pictor-PPE dataset [4] is one of the earliest datasets collected for the PPE detection. It contains 3 classes, Hat, Vest, and Worker, and contains 774 crowd-sourced and 698 web-mined images. Its publicly available version[6] only contains 784 images available for scientific purposes from construction sites with a total of 2,496 worker instances.

*6) TCRSF Dataset:* The TCRSF dataset [5] is a closed-source dataset[7] and extracted from the video feeds collected from different viewpoints in the chemical plants. It provides more realistic scenarios from complex backgrounds. It contains 50,558 labeled instances from 12,373 images in 7 categories, including helmet, safety clothing, head, etc[8].

All the above datasets, summarized in Table I, are collected to address the safety needs of the construction industry in particular, but human safety in the general manufacturing environment is not fully considered, and in this study, we have explored this holistic aspect by collecting data from diverse environments.

TABLE I: Existing Datasets of helmet and PPE detection.

| Dataset | Classes | Images | Instances | Available | Paper |
|---|---|---|---|---|---|
| Pictor-PPE | 3 | 784 | - | ✓ | [4] |
| SHW | 1 | 7,581 | 120,558 | ✓ | [12] |
| CHV | 6 | 1,330 | - | ✓ | [14] |
| TCRSF | 7 | 12,373 | 50,558 | ✗ | [5] |
| GDUT-HWD | 5 | 3,174 | 18,893 | ✓ | [11] |
| SHD | 3 | 5,000 | - | ✓ | - |
| SHEL5K | 5 | 5,000 | 75,570 | ✓ | [15] |
| SH17 (proposed) | 17 | 8,099 | 75,994 | ✓ | - |

## B. YOLO Models

You Only Look Once (YOLO) is a series of popular real-time object detection models that has been through several iterations from YOLOv1 to YOLOv8 [17]–[23]. The latest version, YOLOv8, builds upon its predecessors by incorporating state-of-the-art techniques and innovations to improve accuracy, speed, and adaptability.

YOLOv8 [21] follows a one-stage OD approach, where the input image undergoes a single pass through the network for bounding box prediction and classification. The architecture consists of three main components: a backbone network, a neck, and a prediction head following redmon et al. [19] and bochkovskiy et al. [20]. While the previous variants used the same head for objectness, classification, and

---

[3]https://github.com/ZijianWang-ZW/PPE_detection

[4]https://www.kaggle.com/datasets/andrewmvd/hard-hat-detection

[5]https://data.mendeley.com/datasets/9rcv8mm682/4

[6]https://github.com/ciber-lab/pictor-ppe

[7]Authors announced in [5] that dataset will be publicly released at https://github.com/sofffty/TCRSF. However, it has not been released as of 2024-04-19.

[8]Information mentioned in their paper [5]

regression tasks. Jocher et al. proposed the decoupled head along with Complete Intersection Over Union (CIoU) loss [24] and Distribution Focal Loss (DFL) [25], especially for the calculation of bounding box loss.

The major drawback of YOLOv8 was reliance on DarkNet-53 [19]. It limited the ability to capture fine-grained features, especially for small and occluded objects. Wang et al. proposed YOLOv9 [26] to address this by introducing Programmable Gradient Information (PGI), a novel concept that ensures the retention of critical information throughout the detection process. PGI integrates a reversible branch that works alongside the main network, preserving essential features and improving training outcomes without additional computational costs. Additionally, they proposed a Generalized Efficient Long-Range Attention Network (GELAN) block that optimizes the balance between parameter count, complexity, accuracy, and inference speed, enabling users to select the best computational blocks for various devices.

Recently, proposed YOLOv10 [27] further builds upon predecessors and improves the inference speed by eliminating the non-max-suppression post-processing step which only picks the most probable bounding box out of 1,000s of boxes produced by OD model. The authors proposed a consistent dual assignment strategy where a one-to-one assignment head is introduced in addition to a one-to-many head [28] with an identical structure. Both heads are optimized together during training utilizing rich features, while only a one-to-one head is used in the inference for better inference speed. They also reduced the number of trainable parameters by carefully analyzing each trainable block and removing non-contributing blocks to the overall efficiency. They achieved comparable performance as of YOLOv9-c variant by reducing the 46% trainable parameters.

The next section discusses the use of these OD methods in the PPE detection for human safety.

### C. *Object Detection (OD) for PPE Detection*

Most previous studies that focus on using the OD for PPE detection have considered only the applications in construction engineering aiming to verify the use of hard hats and safety vests. Isailovic et al. [3] and Vukicevic et al. [6] proposed the PPE compliance using a two-stage approach by using a keypoint detector to detect regions and then passing these regions to the binary classification or OD model for further PPE detection. They used the images from Pictor-PPE [4] and hardhat dataset [29] to achieve this objective.

Wu et al. [30] used the Single Stage Detector (SSD) [31] architecture to identify hardhats of different colors on construction sites. They collected the novel GDUT-HWD dataset to train the model. Delhi et al. [32] and Tran et al. [33] similarly employed YOLOv3 to detect hardhats and safety jackets in the construction environment and further used it to check for PPE compliance on various body parts [33] labeling safe and not-safe to each worker. They gathered the data from web sources. Otgonbold et al. [15] compared the performance of multiple OD models for detecting 6 different classes from person, helmet, head, and face on the novel SHEL5k dataset [15]. Nath et al. [4] used different combinations of OD approaches to address worker safety and PPE compliance. The authors initially detected the hard, vest, and worker and passed these to Machine Learning (ML) classifier, while in another approach, they trained the model to predict the state of the PPE. They passed the cropped RoI for a person to another CNN classifier for PPE classification. The Pictor-PPE dataset [4] is publicly available. Chen and Demachi [34] introduced a method using OpenPose [35] for body landmark detection and YOLOv3 [19] OD model for PPE detection. They analyzed the geometric relationships between the key points and detected PPE to assess compliance. They used the GDUT-HWD dataset for all their experiments. Zhafran et al. [36] explored the Fast R-CNN [37] architecture for checking masks, gloves, hardhats, and vests, noting a decrease in accuracy with greater distance and changes in lighting conditions using the manually collected laboratory data. Additionally, several studies have focused on detecting protective masks due to the Coronavirus Disease (COVID-19) pandemic. Loey et al. [38] used YOLOv2 while [39] combined SSD [31] with MobileNetV2 [40] for medical mask detection during COVID-19.

Ferdous and Ahsan [41] used YOLOX [42] OD model on the CHV dataset. Kim et al. [43] scrapped the internet to collect a novel dataset of 4,844 images and manually annotated them into 3 classes, heavy,

PPE, and worker to train YOLOv5 [21] and YOLOv8 [44]. Some other relevant works using OD for PPE detection includes Lung and Wang [45], Xiang et al. [9], Di et al. [46], Han et al. [47], and Azizi et al. [48].

Some traditional approaches include combining the Histogram of Oriented Gradients with the Circle Hough Transform algorithms [49]. Li et al. [50] implemented a radiomics-based method for helmet detection, while Mneymneh et al. [51] introduced a motion detection-based system that subsequently identifies workers and hardhats. Balakrishnan et al. [52] designed a software system that includes an IoT module and Microsoft Azure's Custom Vision AI and Intelligent AI Services to detect safety glasses in lab settings. Amazon has also introduced the proprietary Amazon Rekognition PPE detection system [53].

It is evident that most studies have been motivated by the needs of the construction industry, which remains one of the least digitized sectors and records a high number of fatal injuries. While these studies typically focus on specific types of PPE such as helmets, vests, and masks, the needs of other industries remain largely unexplored. As per the authors' knowledge, there is no comprehensive study on PPE compliance for manufacturing study. Instead, separate studies have targeted specific PPE types relevant to particular industries. These studies typically used a single architecture on custom or small-scale data, making direct comparisons challenging due to the different datasets used for training.

This is the first study to encompass and directly benchmark recent deep learning object detection architectures, providing an objective comparison using a newly developed dataset. This object detection method aims to ensure more efficient compliance with multiple PPEs on body regions. This study proposes a modular framework for PPE compliance that can be applied to various types of PPE and body parts.

## III. SH17 Dataset

Current datasets for PPE detection often focus on specific scenarios, such as detecting helmets, and may not reflect the variety of situations found in real-world industrial settings. In this study, we propose **S**afe **H**uman dataset consisting of **17** different objects referred to as **SH17 dataset**. We scrapped images from the Pexels[9] website, showcasing a range of human activities across diverse industrial operations. Samples of the dataset are shown in Fig. 1. The following sections discuss data collection, annotation, and further details about the curation process for this dataset.

### A. Collection Process

Many open-source datasets are sourced from platforms such as Flickr[10] or gathered through web crawling with search engines such as Google or Bing. However, such data can be noisy and require substantial cleaning effort, with images often subject to different licensing conflicts, resulting in studies not publicly sharing their data. To streamline this process, we gathered images from Pexels, which offers clear usage rights for all its images. To extract relevant images, we used multiple queries such as *manufacturing worker, industrial worker, human worker, labor, etc*. The tags associated with Pexels images proved reasonably accurate. After removing duplicate samples, we obtained around 11,000 samples, of which around 26% were empty, containing no objects from our target classes. These empty images were excluded during the labeling process, resulting in a final dataset of 8,099 images. The dataset exhibits significant diversity, representing manufacturing environments globally, thus minimizing potential regional or racial biases.

### B. Annotation Process

The annotation process involved four human annotators, three initially completing annotations. The team lead then verified and corrected any mistakes in the annotations. Finally, a graduate student performed a final verification and addressed any remaining mislabeling. All annotators utilized DarkLabel[11] while

---

Fig. 1: Samples from the proposed **SH17 Dataset.** Best viewed online.

graduate student used the LabelImg[12] for data annotation. The selection of the tool is a personal preference and has no impact on the quality of annotation.

### C. Categories

We categorized the objects into 17 classes, each representing a body part and the required PPE items for safety. We chose these classes with the aim that downstream applications could ignore irrelevant classes while still covering a wide range of scenarios in industrial manufacturing. The complete list of classes and their definitions can be found in Table II. Additional tags with each class to match some of the classes in earlier datasets were also mentioned with each annotation and are released as the extended version of the dataset. The tags are also mentioned in the Table II. These can be utilized during training in downstream tasks by advanced users. We also extracted the metadata of each image sample. Its details are explained in Section A.

### D. Data Details

We provide the original images extracted from Pexels in their native resolution with a maximum image size of 8,192×5,462 and a minimum of 1,920×1,002. Data consists of both landscape and portrait-style images. Each image contains an average of 9.38 instances. We have labeled ears and earmuffs, which are

---

[12]https://github.com/HumanSignal/labelImg

TABLE II: List of all the annotated classes and their counts. Tag **off** means the item is present in the scene but not worn by the person, while **on** means the item is present and worn by the person.

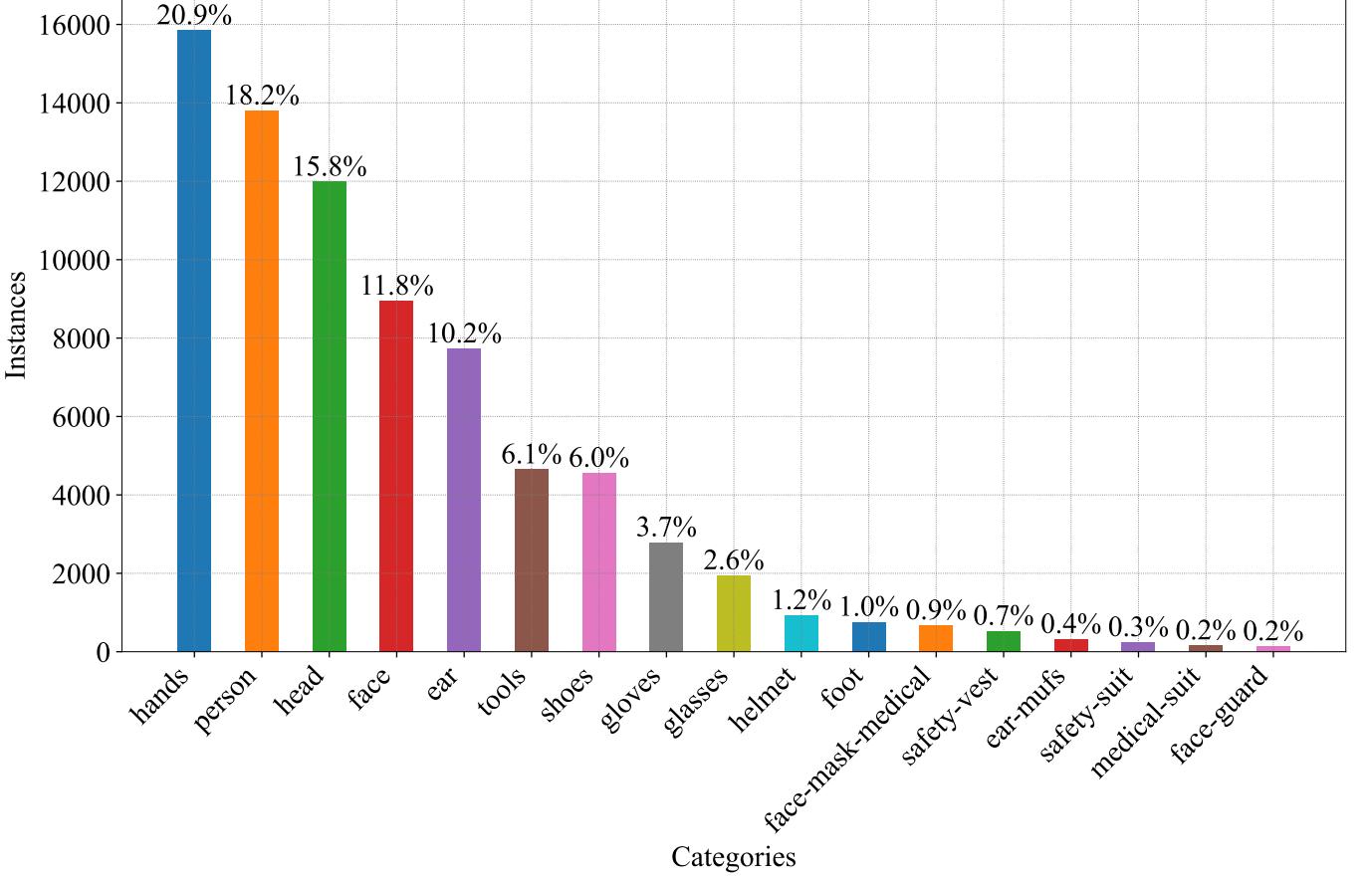| ID | Name | Additional Tags | Instances | Description |
|---|---|---|---|---|
| 1 | Person | male, female, children | 13802 | Uses visible features for classification. |
| 2 | Head | - | 11985 | Includes any view of the head: front, back, top or else. |
| 3 | Face | - | 8950 | Only classified as visible when the nose is visible. |
| 4 | Glasses | on, off, safety, vision | 1945 | Detection of safety glasses. |
| 5 | Face-mask-medical | on, off | 669 | Detection of medical face masks |
| 6 | Face-guard | on, off | 134 | Detection of whether faceguards |
| 7 | Ear | | 7730 | Focused on the ears for safety equipment detection. |
| 8 | Earmuffs | on, off | 318 | Detection of earmuffs. Over-Ear-Headphones are labelled as earmufs. |
| 9 | Hands | | 15850 | Focus on the hands for safety equipment detection. |
| 10 | Gloves | on, off | 2790 | Detection of gloves. |
| 11 | Foot | | 796 | Visible when there are no shoes, and each foot is annotated regardless of with or without socks. |
| 12 | Shoes | on, off, safety, other | 4560 | Shoes detection. Safety: Includes safety shoes and thick joggers. Others: slippers, sneakers, or other types of footwear. |
| 13 | Safety-vest | on, off | 530 | Detection of safety vests. |
| 14 | Tools | on, off | 4647 | Detection of tools being held; on means in hand. If the tool is present in the scene but not in hand, it's off. Pencils and laptops are not considered tools. |
| 15 | Helmet | on, off, white, red, black, yellow, blue | 927 | tags contain the color of a helmet as well. |
| 16 | Medical-suit | on, off | 155 | Detection of medical suits |
| 17 | Safety-suit | on, off | 530 | Detection of safety suits |

Fig. 2: Distribution of all class instances. The percentage of each class is shown at the top of each bar.

very small objects compared to a person. Hence, the dataset contains objects of all sizes. There are 39,764 annotations containing less than 1% of the area, while 59,025 annotations contain less than 5% area of the image. Examples of all object sizes are included in Fig. 1. The dataset contains the maximum instances of hands 15,850, which is 20.9% of the instances. The helmet class contains 927, approximately 1.2% of the data, while there are 134 instances of faceguards, the lowest class instance with 0.2% data only. Table II contains the complete list of instances of each class. Figure 2 shows the distribution of all classes, their instances, and their percentage in the dataset, demonstrating an imbalance in class distribution.

### E. Evaluation Metrics

Object detection models are evaluated using several metrics to assess their performance accurately. Common metrics used by Microsoft Common Objects in Context (MS-COCO) dataset [54] include Precision (P), Recall (R), and Mean Average Precision (mAP).

Precision (P) measures the proportion of correctly identified objects among all objects detected by the model and is calculated as

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Recall (R), on the other hand, quantifies the ability of the model to detect all relevant objects in the dataset and is computed as

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Additionally, Mean Average Precision (mAP) evaluates the detection accuracy across all classes. It's computed based on P and R for each class and then averaged to give an overall score. The Intersection over Union (IoU) metric is used to gauge the accuracy of object localization. It calculates the overlap between the ground truth bounding boxes ($b_g$) and the model's predicted bounding boxes ($b_{pred}$) as follows:

$$IoU = \frac{\text{Area}(b_{pred} \cap b_g)}{\text{Area}(b_{pred} \cup b_g)}$$

where $b_g$ represents the ground truth bounding box, and $b_{pred}$ denotes the bounding box predicted by the OD model. The IoU threshold acts as a filter to remove FP) bounding boxes with an IoU score below a certain threshold. This threshold determines the required accuracy for an object to be classified as detected or missed (e.g., IoU $\geq$ threshold). Different models may use various threshold values, such as 0.25, 0.5, or 0.75, during evaluation.

## IV. EXPERIMENTATION AND DISCUSSION

We have trained multiple OD models with different size variants of YOLOv8, YOLOv9 and YOLOv10 for benchmarking purposes. We split the dataset into 80% for training and 20% for testing. We have used default hyper-parameters suggested by the respective authors in their open-source codebase[13]. We used transfer learning to adapt a pre-trained model initially trained on the MS-COCO dataset and trained it on the SH17 dataset for 200 epochs. We followed [44] for training hyperparameters. We trained all models on the fixed image size of 640 instead of the original high-quality images due to the memory limitation by 2x NVIDIA RTX GPUs with a batch size of 128 for the nano models to 32 for bigger scale models. We also used a mosaic of 4 images along with horizontal flipping for data augmentation and Non-Max-Suppression (NMS) for post-processing the outputs of YOLOv8 and YOLOv9 model variants. We have consolidated the results of the model training in Table III evaluated on the separate test set. As mentioned, YOLOv9-e performs best among all variants with 70.9% mAP@50 and 48.7% mAP50-95. It has 58.1M parameters as compared to the next best YOLOv8-x, which has 68.2M parameters, effectively 15% fewer parameters. Furthermore, YOLOv9-c has comparable performance to YOLOv8-l with 32% fewer trainable parameters. YOLOv10-x has comparable performance with YOLOv9-c with 15% less trainable parameters, which significantly reduces the training and inferences time.

We present the training metrics of YOLOv9-e model, our best-performing model on SH17, in Fig. 3. The training of the model plateaus at around 170 epochs as the mAP50 does not improve after that. Table IV lists the class-wise accuracy of YOLOv9-e model on the randomly sampled test set. Hands class, having the most samples, performs well with 89.8% mAP. Tools and Foot class samples vary in types, safety shoes, slippers, joggers, and runners, while tools can be anything a person is working with, resulting in these classes showing low mAP.

These benchmarks further underscore the results reported in the [26], [27], [44] where YOLOv9 models perform comparably with reduced parameters due to the usage of PGI and GELAN. The YOLOv10 struggles on the small objects and SH17 consists of 52% objects covering less than 1% of the area(as explained in Section III-D). Figure 5 presents some visual differences between the models; YOLOv8-n, YOLOv8-m, YOLOv8-x, YOLOv9-e. For sample (2) in Fig. 5, YOLOv8-m and YOLOv8-x both made a FP prediction of the tool on the ground. For sample (3), YOLOv8-n model FP of face class. In row 5, all models failed to detect the tool in the worker's hand. The original labeled samples are mentioned in the Fig. 1.

### A. Generalization capability

We have verified the model's efficacy with the SH17 dataset on the cross-domain datasets. We selected the Pictor-PPE dataset [4] to validate the YOLOv9-e model. The dataset has only three classes: worker,

---

[13] https://github.com/ultralytics/ultralytics

TABLE III: Comparison of the models trained on SH17. **Bold** represents the best (a larger value is better).

| Model | Params | Images | Instances | P | R | mAP 50 | 50-95 |
|---|---|---|---|---|---|---|---|
| Yolo-8-n | 3.2 | 1620 | 15358 | 67.5 | 53.6 | 58.0 | 36.6 |
| Yolo-8-s | 11.2 | 1620 | 15358 | 81.5 | 55.7 | 63.7 | 41.7 |
| Yolo-8-m | 25.9 | 1620 | 15358 | 77.1 | 60.5 | 66.6 | 45.7 |
| Yolo-8-l | 43.7 | 1620 | 15358 | 76.7 | 62.9 | 68.0 | 47.0 |
| Yolo-8-x | 68.2 | 1620 | 15358 | 77.1 | 63.1 | 69.3 | 47.2 |
| Yolo-9-t | 2.0 | 1620 | 15358 | 75.0 | 52.6 | 58.5 | 37.5 |
| Yolo-9-s | 7.2 | 1620 | 15358 | 73.6 | 60.2 | 65.3 | 42.9 |
| Yolo-9-m | 20.1 | 1620 | 15358 | 77.4 | 62.0 | 68.6 | 46.5 |
| Yolo-9-c | 25.5 | 1620 | 15358 | 79.6 | 60.8 | 67.7 | 46.5 |
| Yolo-9-e | 58.1 | 1620 | 15358 | **81.0** | **65.0** | **70.9** | **48.7** |
| Yolo-10-n | 2.3 | 1620 | 15358 | 66.8 | 53.2 | 57.2 | 35.9 |
| Yolo-10-s | 7.2 | 1620 | 15358 | 75.8 | 57.0 | 62.7 | 40.9 |
| Yolo-10-m | 15.4 | 1620 | 15358 | 71.4 | 61.4 | 65.7 | 43.8 |
| Yolo-10-b | 19.1 | 1620 | 15358 | 77.7 | 59.1 | 65.8 | 45.1 |
| Yolo-10-l | 24.4 | 1620 | 15358 | 76.0 | 61.8 | 67.4 | 46.0 |
| Yolo-10-x | 29.5 | 1620 | 15358 | 76.8 | 62.8 | 67.8 | 46.7 |

TABLE IV: Class-wise accuracy of YOLOv9-e model.

| Class | Images | Instances | P | R | mAP 50 | 50-95 |
|---|---|---|---|---|---|---|
| all | 1620 | 15358 | 81.0 | 65.0 | 70.9 | 48.7 |
| hands | 1284 | 3212 | 91.4 | 83.9 | 89.8 | 64.8 |
| person | 1515 | 2734 | 90.9 | 89.2 | 92.1 | 77.9 |
| head | 1314 | 2427 | 94.8 | 89.1 | 93.5 | 74.3 |
| face | 1155 | 1855 | 96.0 | 88.1 | 93.8 | 73.8 |
| ear | 987 | 1612 | 91.2 | 75.5 | 84.3 | 55.0 |
| shoes | 320 | 956 | 79.2 | 62.8 | 70.8 | 43.2 |
| tool | 455 | 923 | 67.4 | 39.1 | 43.2 | 27.6 |
| gloves | 254 | 529 | 81.6 | 58.9 | 66.5 | 43.5 |
| glasses | 323 | 398 | 87.4 | 72.6 | 76.4 | 46.9 |
| helmet | 93 | 154 | 81.3 | 67.8 | 77.0 | 57.6 |
| face-mask | 75 | 151 | 88.8 | 73.2 | 75.5 | 49.2 |
| foot | 64 | 149 | 51.7 | 22.1 | 29.3 | 14.0 |
| safety-vest | 45 | 97 | 66.4 | 55.0 | 57.7 | 38.1 |
| ear-mufs | 38 | 49 | 79.6 | 46.9 | 57.1 | 40.5 |
| safety-suit | 28 | 45 | 65.7 | 53.3 | 58.5 | 38.3 |
| medical-suit | 30 | 43 | 86.0 | 65.1 | 68.5 | 40.6 |
| face-guard | 23 | 24 | 76.8 | 62.5 | 71.7 | 42.8 |

precision
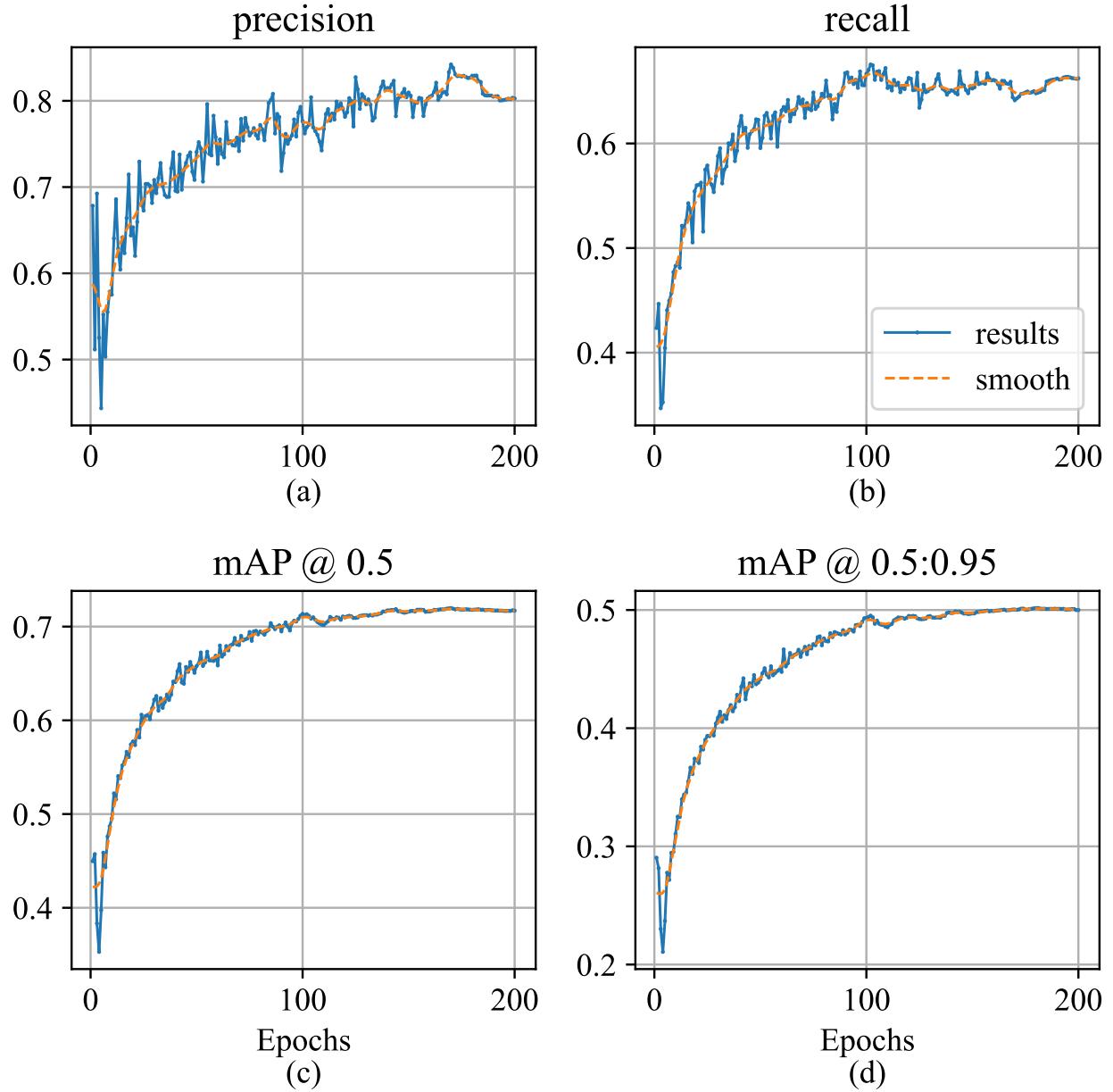
recall

mAP @ 0.5

mAP @ 0.5:0.95

Epochs
(c)

Epochs
(d)

Fig. 3: Training metrics of YOLOv9-e model that performs best among all.

hat, and vest, which we mapped to the Person, Helmet, and Safety-vest classes of SH17. We used all the publicly available images in the dataset for validation, and YOLOv9-e achieved 58.8% mAP, which shows that trained models can be used directly for PPE compliance in the manufacturing environment. SH17 has the least samples of safety-vest class; however, it still achieves 35.5% mAP50. The complete results are summarised in Table V. Figure 4 presents the confusion matrix of the complete Pictor-PPE as just the evaluation set to YOLOv9-e. It is evident that person class is only considered background in some instances, while the helmet is often mistaken with head class. We also show some visual results in the Fig. 6.
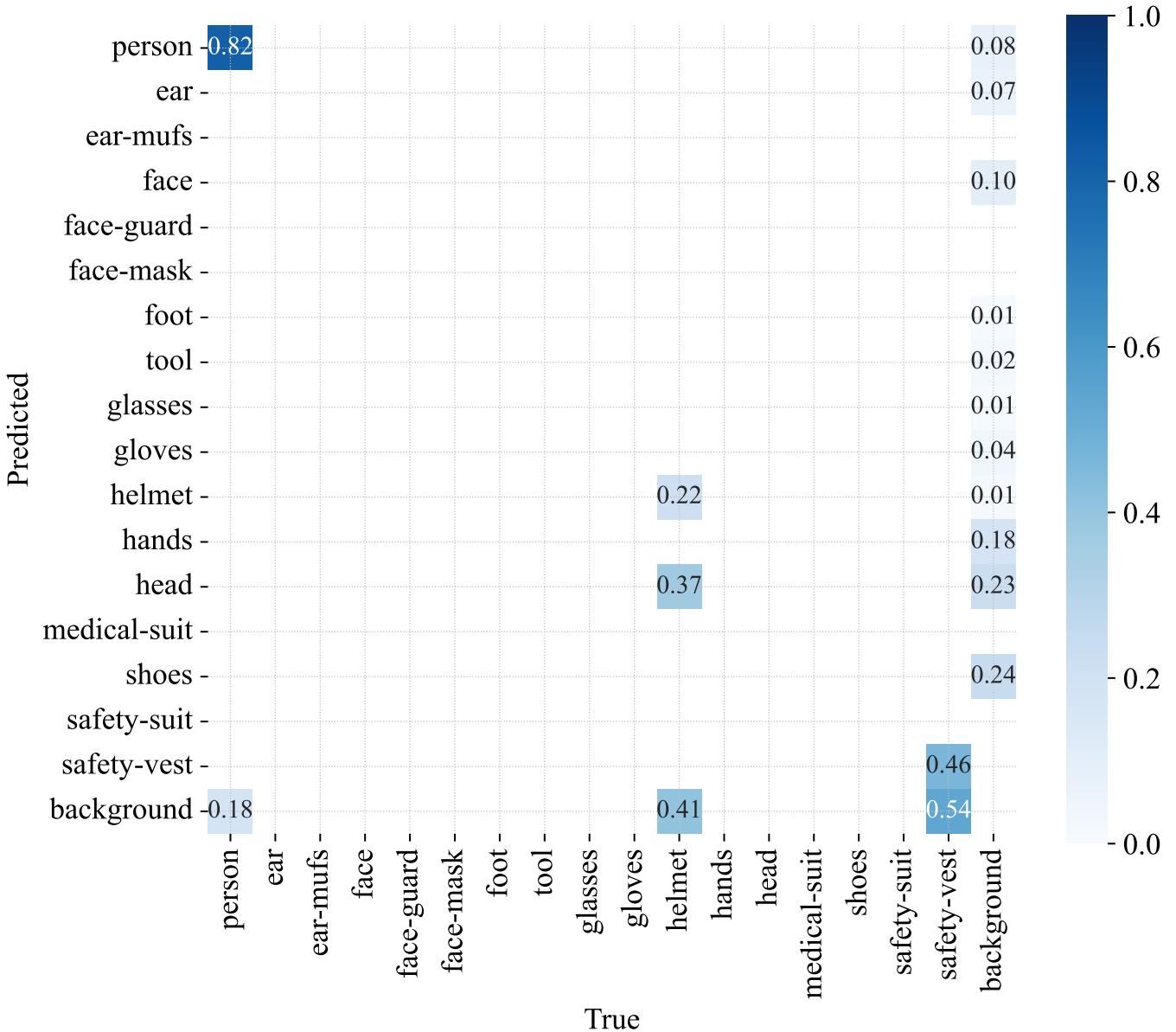
Fig. 4: Confusion Matrix of all the instances of Pictor-PEE dataset.

TABLE V: Class-wise accuracy of YOLOv9e model on the Pictor-PPE dataset [4].

| Class | Images | Instances | P | R | mAP | |
|---|---|---|---|---|---|---|
| | | | | | 50 | 50-95 |
| all | 654 | 3477 | 73.5 | 50.3 | 58.9 | 37.6 |
| person | 654 | 2080 | 83.6 | 80.7 | 85.5 | 61.6 |
| helmet | 451 | 1369 | 94.6 | 23.7 | 55.5 | 30.0 |
| safety-vest | 12 | 28 | 42.4 | 46.4 | 35.5 | 21.2 |

13
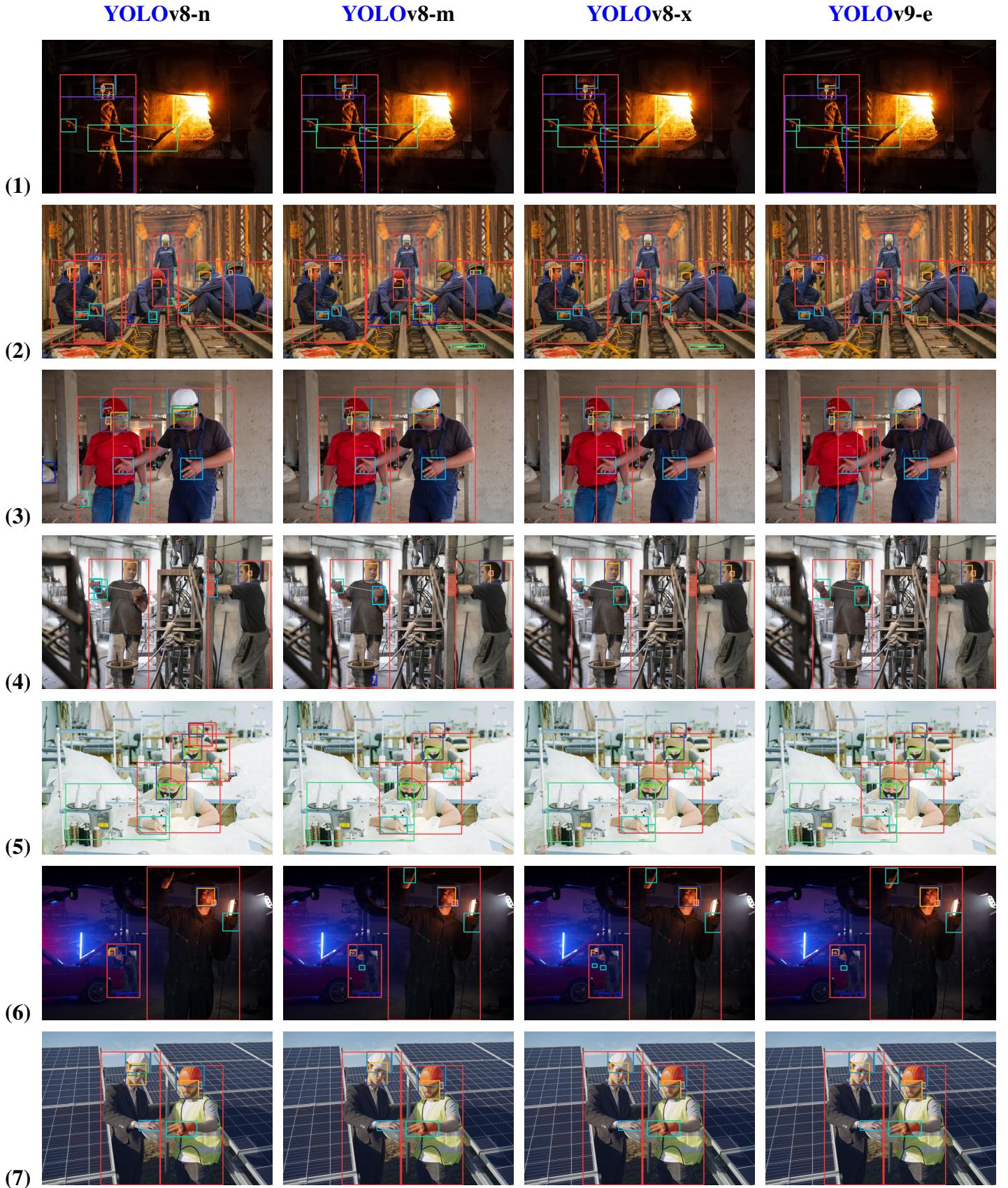


Fig. 5: Visual representation of the predictions made by different models; YOLOv8-n, YOLOv8-m, YOLOv8-x, YOLOv9-e. Best viewed online.

Fig. 6: Visual results of inference on Pictor-PPE dataset. Best viewed online.

## V. Conclusion

Human safety is a critical problem in the manufacturing environment, and in this study, we proposed a human safety dataset for ensuring the PPE compliance. We have open-sourced 8,099 images annotated in 17 different classes containing different PPE items and body parts. The dataset is publicly available for research and commercial purposes. To evaluate performance and provide some metrics on the dataset, we trained multiple models for benchmarking, among which YOLOv9-e performed best on the proposed dataset. We evaluated the performance of the trained model with the SH17 dataset on the Pictor-PPE dataset. We obtained satisfactory results that this model can be used in industrial environments to detect and ensure PPE compliance. For future work, further improvements in the performance of the minority classes with specialized models and custom training methodologies can be considered.

## Supplementary Materials

Supporting data is available at SH17 dataset GitHub repository.

## Authors Contributions

**Hafiz Mughees Ahmad:** Conceptualization, Methodology, Software, Validation, Data curation, Formal analysis, Visualization, Resources, Writing - Original draft preparation. **Afshin Rahimi:** Funding acquisition, Project administration, Software, Resources, Formal analysis, Supervision, Writing - Review & Editing.

## FUNDING

## INSTITUTIONAL REVIEW

Not applicable.

## DATA AVAILABILITY

The data presented in this study are available in SH17 dataset GitHub repository at https://github.com/ahmadmughees/sh17dataset.

## CONFLICT OF INTEREST

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## ACKNOWLEDGEMENT

## REFERENCES

[1] I. Sutton, "Chapter 6 - Personal protective equipment," in *Plant Design and Operations*, I. Sutton, Ed. Oxford: Gulf Publishing Company, Jan. 2015, pp. 127–137.

[2] Occupational Safety and Health Administration (OSHA), "Personal Protective Equipment," *U.S. Department of Labor*, 2004.

[3] V. Isailovic, A. Peulic, M. Djapan, M. Savkovic, and A. M. Vukicevic, "The compliance of head-mounted industrial PPE by using deep learning object detectors," *Scientific Reports*, vol. 12, no. 1, p. 16347, Sep. 2022.

[4] N. D. Nath, A. H. Behzadan, and S. G. Paal, "Deep learning for site safety: Real-time detection of personal protective equipment," *Automation in Construction*, vol. 112, p. 103085, Apr. 2020.

[5] F. Yu, X. Wang, J. Li, S. Wu, J. Zhang, and Z. Zeng, "Towards Complex Real-World Safety Factory Inspection: A High-Quality Dataset for Safety Clothing and Helmet Detection," Jun. 2023.

[6] A. M. Vukicevic, M. Djapan, V. Isailovic, D. Milasinovic, M. Savkovic, and P. Milosevic, "Generic compliance of industrial PPE by using deep learning techniques," *Safety Science*, vol. 148, p. 105646, Apr. 2022.

[7] B. Naticchia, M. Vaccarini, and A. Carbonari, "A monitoring system for real-time interference control on large construction sites," *Automation in Construction*, vol. 29, pp. 148–160, Jan. 2013.

[8] A. Kelm, L. Laußat, A. Meins-Becker, D. Platz, M. J. Khazaee, A. M. Costin, M. Helmus, and J. Teizer, "Mobile passive Radio Frequency Identification (RFID) portal for automated and rapid control of Personal Protective Equipment (PPE) on construction sites," *Automation in Construction*, vol. 36, pp. 38–52, Dec. 2013.

[9] C. Xiang, D. Yin, F. Song, Z. Yu, X. Jian, and H. Gong, "A Fast and Robust Safety Helmet Network Based on a Mutilscale Swin Transformer," *Buildings*, vol. 14, no. 3, p. 688, Mar. 2024.

[10] H. M. Ahmad, A. Rahimi, and K. Hayat, "Capacity Constraint Analysis Using Object Detection for Smart Manufacturing," Jan. 2024.

[11] J. Wu, N. Cai, W. Chen, H. Wang, and G. Wang, "Automatic detection of hardhats worn by construction personnel: A deep learning approach and benchmark dataset," *Automation in Construction*, vol. 106, p. 102894, Oct. 2019.

[12] njvisionpower, "Njvisionpower/Safety-Helmet-Wearing-Dataset," Apr. 2024.

[13] D. Peng, Z. Sun, Z. Chen, Z. Cai, L. Xie, and L. Jin, "Detecting Heads using Feature Refine Net and Cascaded Multi-scale Architecture," in *2018 24th International Conference on Pattern Recognition (ICPR)*, Aug. 2018, pp. 2528–2533.

[14] Z. Wang, Y. Wu, L. Yang, A. Thirunavukarasu, C. Evison, and Y. Zhao, "Fast Personal Protective Equipment Detection for Real Construction Sites Using Deep Learning Approaches," *Sensors*, vol. 21, no. 10, p. 3478, Jan. 2021.

[15] M.-E. Otgonbold, M. Gochoo, F. Alnajjar, L. Ali, T.-H. Tan, J.-W. Hsieh, and P.-Y. Chen, "SHEL5K: An Extended Dataset and Benchmarking for Safety Helmet Detection," *Sensors*, vol. 22, no. 6, p. 2315, Jan. 2022.

[16] "Safety Helmet Detection," https://www.kaggle.com/datasets/andrewmvd/hard-hat-detection.

[17] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jun. 2016, pp. 779–788.

[18] J. Redmon and A. Farhadi, "YOLO9000: Better, Faster, Stronger," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jul. 2017, pp. 6517–6525.

[19] ——, "Yolov3: An incremental improvement," in *arXiv Preprint arXiv:1804.02767*. ArXiv, 2018.

[20] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal Speed and Accuracy of Object Detection," Apr. 2020.

[21] G. Jocher, A. Stoken, J. Borovec, NanoCode012, ChristopherSTAN, L. Changyu, Laughing, tkianai, A. Hogan, lorenzomammana, yxNONG, AlexWang1900, L. Diaconu, Marc, wanghaoyang0106, ml5ah, Doug, F. Ingham, Frederik, Guilhen, Hatovix, J. Poznanski, J. Fang, L. Y. , changyu98, M. Wang, N. Gupta, O. Akhtar, PetrDvoracek, and P. Rai, "Ultralytics/yolov5: V3.1 - Bug Fixes and Performance Improvements," Zenodo, Oct. 2020.

[22] C. Li, L. Li, H. Jiang, K. Weng, Y. Geng, L. Li, Z. Ke, Q. Li, M. Cheng, W. Nie, Y. Li, B. Zhang, Y. Liang, L. Zhou, X. Xu, X. Chu, X. Wei, and X. Wei, "YOLOv6: A Single-Stage Object Detection Framework for Industrial Applications," Sep. 2022.

[23] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7464–7475.

[24] Z. Zheng, P. Wang, W. Liu, J. Li, R. Ye, and D. Ren, "Distance-IoU loss: Faster and better learning for bounding box regression," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 12 993–13 000.

[25] X. Li, W. Wang, L. Wu, S. Chen, X. Hu, J. Li, J. Tang, and J. Yang, "Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 002–21 012, 2020.

[26] C.-Y. Wang, I.-H. Yeh, and H.-Y. M. Liao, "YOLOv9: Learning What You Want to Learn Using Programmable Gradient Information," Feb. 2024.

[27] A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, and G. Ding, "YOLOv10: Real-Time End-to-End Object Detection," May 2024.

[28] Y. Chen, Q. Chen, Q. Hu, and J. Cheng, "DATE: Dual Assignment for End-to-End Fully Convolutional Object Detection," https://arxiv.org/abs/2211.13859v2, Nov. 2022.

[29] L. Xie, "Hardhat," Jan. 2019.

[30] J. Wu, N. Cai, W. Chen, H. Wang, and G. Wang, "Automatic detection of hardhats worn by construction personnel: A deep learning approach and benchmark dataset," *Automation in Construction*, vol. 106, p. 102894, Oct. 2019.

[31] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," in *ECCV*, 2016, vol. 9905, pp. 21–37.

[32] V. S. K. Delhi, R. Sankarlal, and A. Thomas, "Detection of Personal Protective Equipment (PPE) Compliance on Construction Site Using Computer Vision Based Deep Learning Techniques," *Frontiers in Built Environment*, vol. 6, 2020.

[33] Q.-H. Tran, T.-L. Le, and S.-H. Hoang, "A fully automated vision-based system for real-time personal protective detection and monitoring," *KICS Korea-Vietnam Int Jt Work Commun Inf Sci*, vol. 2019, no. 1, p. 6, 2019.

[34] S. Chen and K. Demachi, "A Vision-Based Approach for Ensuring Proper Use of Personal Protective Equipment (PPE) in Decommissioning of Fukushima Daiichi Nuclear Power Station," *Applied Sciences*, vol. 10, no. 15, p. 5129, Jan. 2020.

[35] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime Multi-person 2D Pose Estimation Using Part Affinity Fields," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Jul. 2017, pp. 1302–1310.

[36] F. Zhafran, E. S. Ningrum, M. N. Tamara, and E. Kusumawati, "Computer Vision System Based for Personal Protective Equipment Detection, by Using Convolutional Neural Network," in *2019 International Electronics Symposium (IES)*, Sep. 2019, pp. 516–521.

[37] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.

[38] M. Loey, G. Manogaran, M. H. N. Taha, and N. E. M. Khalifa, "Fighting against COVID-19: A novel deep learning model based on YOLO-v2 with ResNet-50 for medical face mask detection," *Sustainable Cities and Society*, vol. 65, p. 102600, Feb. 2021.

[39] P. Nagrath, R. Jain, A. Madan, R. Arora, P. Kataria, and J. Hemanth, "SSDMNV2: A real time DNN-based face mask detection system using single shot multibox detector and MobileNetV2," *Sustainable Cities and Society*, vol. 66, p. 102692, Mar. 2021.

[40] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted Residuals and Linear Bottlenecks," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, Jun. 2018, pp. 4510–4520.

[41] M. Ferdous and S. M. M. Ahsan, "PPE detector: A YOLO-based architecture to detect personal protective equipment (PPE) for construction sites," *PeerJ Computer Science*, vol. 8, p. e999, Jun. 2022.

[42] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "YOLOX: Exceeding YOLO Series in 2021," *arXiv:2107.08430 [cs]*, Aug. 2021.

[43] K. Kim, K. Kim, and S. Jeong, "Application of YOLO v5 and v8 for Recognition of Safety Risk Factors at Construction Sites," *Sustainability*, vol. 15, no. 20, p. 15179, Jan. 2023.

[44] G. Jocher, A. Chaurasia, and J. Qiu, "YOLO by Ultralytics," *URL: https://github. com/ultralytics/ultralytics*, 2023.

[45] L.-W. Lung and Y.-R. Wang, "Applying Deep Learning and Single Shot Detection in Construction Site Image Recognition," *Buildings*, vol. 13, no. 4, p. 1074, Apr. 2023.

[46] B. Di, L. Xiang, Y. Daoqing, and P. Kaimin, "MARA-YOLO: An Efficient Method for Multiclass Personal Protective Equipment Detection," *IEEE Access*, vol. 12, pp. 24 866–24 878, 2024.

[47] S. Han, W. Park, K. Jeong, T. Hong, and C. Koo, "Utilizing synthetic images to enhance the automated recognition of small-sized construction tools," *Automation in Construction*, vol. 163, p. 105415, Jul. 2024.

[48] R. Azizi, M. Koskinopoulou, and Y. Petillot, "Comparison of Machine Learning Approaches for Robust and Timely Detection of PPE in Construction Sites," *Robotics*, vol. 13, no. 2, p. 31, Feb. 2024.

[49] A. H. M. Rubaiyat, T. T. Toma, M. Kalantari-Khandani, S. A. Rahman, L. Chen, Y. Ye, and C. S. Pan, "Automatic Detection of Helmet Uses for Construction Safety," in *2016 IEEE/WIC/ACM International Conference on Web Intelligence Workshops (WIW)*, Oct. 2016, pp. 135–142.

[50] K. Li, X. Zhao, J. Bian, and M. Tan, "Automatic Safety Helmet Wearing Detection," in *2017 IEEE 7th Annual International Conference on CYBER Technology in Automation, Control, and Intelligent Systems (CYBER)*, Jul. 2017, pp. 617–622.

[51] B. E. Mneymneh, M. Abbas, and H. Khoury, "Vision-Based Framework for Intelligent Monitoring of Hardhat Wearing on Construction Sites," *Journal of Computing in Civil Engineering*, vol. 33, no. 2, p. 04018066, Mar. 2019.

[52] B. Balakreshnan, G. Richards, G. Nanda, H. Mao, R. Athinarayanan, and J. Zaccaria, "PPE Compliance Detection using Artificial Intelligence in Learning Factories," *Procedia Manufacturing*, vol. 45, pp. 277–282, Jan. 2020.

[53] Amazon Web Services, "Detecting personal protective equipment - Amazon Rekognition," https://docs.aws.amazon.com/rekognition/latest/dg/ppe-detection.html.

[54] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft COCO: Common Objects in Context," Feb. 2015.

APPENDIX A
SH17 DATASET

## A. Meta Data

We have also provided the collected meta-data of each image that can be used to build the dataset from the source or can be used as additional data during model training. The list of objects and their definitions in the meta data is mentioned in Table VI. We also provide scripts to build the data from the source in the SH17 dataset Github repository.

TABLE VI: Metadata Format for the Dataset

| Field | Description |
|---|---|
| Unique Identifier | A unique code assigned to each image for identification. |
| Width and Height | The dimensions of the image, specified in pixels. |
| URL | The web address where the image can be accessed. |
| Photographer Name | The name of the individual who captured the image. |
| Photographer URL | The web address leading to the photographer's portfolio or profile. |
| Photographer ID | A unique identifier for the photographer, used for reference or database purposes. |
| Average Color | The average color in the hexadecimal code. |
| Source | The origin or platform from which the image was obtained. |
| Liked | Indicates whether the image has been marked as liked or favored on the online platform, typically a boolean value. |
| Description | A summary or narrative about the image's content, context, or theme. |

## B. Person Demographics

In addition to the bounding boxes representing the Persons class, we have included additional tags based on visible features in the images. While these tags are primarily derived from observable characteristics, they may contain some inaccuracies, as visual features are not always the most reliable indicators. Nonetheless, we include these tags to identify and mitigate potential biases related to gender or ethnic representation. The tags encompass categories such as male, female, children, and various ethnic backgrounds, including Black, Brown, White, and Asian. Table VII provides the quantitative analysis of these tags.

TABLE VII: Quantitative analysis of the additional tags with Persons class.

| | Male | Female | Children | Total |
|---|---|---|---|---|
| White | 2432 | 2032 | 37 | 4501 |
| Black | 1098 | 776 | 7 | 1881 |
| Brown | 577 | 218 | 12 | 807 |
| Asian | 963 | 1272 | 52 | 2287 |
| Total | 5070 | 4298 | 108 | 9476 |

## C. Train Test split

We have used 80% in train and 20% in test data. We provide the separate list of training and test files in the SH17 dataset GitHub repository. Table VIII provides the complete list of instances in the training and test set.

TABLE VIII: Sorted list of instances in each class in Train and Test Set.

| Class | Train | Test |
|---|---|---|
| face-guard | 110 | 24 |
| medical-suit | 114 | 43 |
| safety-suit | 195 | 45 |
| ear-mufs | 269 | 49 |
| safety-vest | 433 | 97 |
| face-mask-medical | 519 | 151 |
| foot | 610 | 149 |
| helmet | 773 | 154 |
| glasses | 1547 | 398 |
| gloves | 2261 | 529 |
| shoes | 3604 | 956 |
| tools | 3724 | 923 |
| ear | 6118 | 1612 |
| face | 7095 | 1855 |
| head | 9558 | 2427 |
| person | 11068 | 2734 |
| hands | 12638 | 3212 |

**Hafiz Mughees Ahmad** completed his Bachelor's and Master's in Electrical Engineering from the Institute of Space Technology, Pakistan, in 2015 and 2018, respectively. He is currently pursuing a Ph.D. at the University of Windsor, Canada. Alongside his studies, he serves as a Deep Learning Engineer at IFIVEO CANADA INC. His previous roles include Research Associate at Istanbul Medipol University, Turkey, and Lecturer at the Institute of Space Technology, Pakistan. His research focuses on Computer Vision and Deep Learning, with applications in OD and real-time surveillance and monitoring in industrial manufacturing and production environments. He is a Graduate Student Member of IEEE.

**Afshin Rahimi** received his B.Sc. degree from the K. N. Toosi University of Technology, Tehran, Iran, in 2010, and the M.Sc. and Ph.D. degrees from Toronto Metropolitan University, Toronto, ON, Canada, in 2012, and 2017, respectively, in Aerospace Engineering. He was with Pratt & Whitney Canada from 2017 to 2018. He is currently an Associate Professor in the Department of Mechanical, Automotive, and Materials Engineering at the University of Windsor, Windsor, ON, Canada. Since 2010, he has been involved in various industrial research, technology development, and systems engineering projects/contracts related to the control and diagnostics of satellites, UAVs, and commercial aircraft subsystems. In recent years, he has also been involved with industrial automation and using technologies to boost manual labor work in industrial settings. He is a senior member of IEEE, a lifetime member of AIAA, and a technical member of the PHM Society.