

GB-YOLO: a lightweight model for mask detection on IoT system

Cheng Huang, Yishen Liu, Zihan Xia, Siyang Jiang, Hao Tian, and Yongbin Yu

Abstract—The combination of the Internet of Things and deep learning technology is usually accompanied by many problems, such as limited bandwidth and computing resources. IoT combined with deep learning often causes system freezes or delays due to limited computing resources. Upgrading the hardware equipment of the IoT system requires a large economic cost, but using a lightweight deep learning model can reduce the consumption of hardware resources to adapt to the actual scene. In this paper, we combine IoT technology and propose a lightweight deep learning model, GB-YOLO, to assist people in mask detection, vehicle counting, and target tracking, which does not take up too many computing resources. We deployed GB-YOLO on the server side, and completed the training in the container. The weight file after training was deployed in Docker, and then combined with Kubernetes to get the final experimental results. The resulting graph can be displayed by opening a browser at the edge node and entering the relevant IP address. Users can also perform certain operations on the results in the front end of the browser, such as drawing a horizontal line in the road to complete the local vehicle count. These operations are also fed back to the server for interaction with developers. For GB-YOLO, the recognition speed and accuracy are faster than before. At the same time, compared with the previous version, the model itself requires less storage space and is easier to deploy, making the model easier to implement in the operation of edge nodes. Theoretical analysis and experimental results verify the feasibility and superiority of the proposed method.

Index Terms—internet of things, artificial intelligence, deep learning, YOLOv5, object detection, object tracking

I. INTRODUCTION

WITH increase in speed and bandwidth of internet, Internet of Things (IoT) is pushing the market and the life of people to a new level and knocking on the door with new opportunities for interactions with other fields, such as 5G [1], artificial intelligence [2] and so on. For IoT itself, it is an extended network based on the Internet. It combines various information sensing devices with the network to form a huge network, so that people can operate and manage

Cheng Huang and Yishen Liu are the joint first authors. Both of them are from the Chinese University of Hong Kong, Department of Information Engineering and the Jockey Club School of Public Health and Primary Care, and are currently studying for a master's degree.

Zihan Xia is currently pursuing a master's degree at the University of Electronic Science and Technology of China, and he is about to go to UC San Diego for a Ph.D. Siyang Jiang is currently pursuing the M.S. degree at National Taiwan University, Taiwan, China. He is going to go to The Chinese University of Hong Kong in September 2022 to pursue a Ph.D. Hao Tian is from the Chinese University of Hong Kong, Department of Information Engineering, and are currently studying for a master's degree. Hao Tian is about to go to The Hong Kong Polytechnic University for a Ph.D.

Yongbin Yu is currently an Associate Professor with the School of Information and Software Engineering, UESTC and he is corresponding Author of this paper.

by issuing instructions to the devices of this network. For example, our home appliances can be connected to a server through IoT technology, and then users can use their mobile phones to command the server to control these furniture [3]. But with more and more devices adding IoT systems, IoT's central processing platform will face a large amount of data processing. How to deal with these huge amounts of data, or how to achieve efficient scheduling in limited computing resources to prevent congestion, has become a hot research direction. This means that running such a large amount of data requires powerful hardware. At the same time, better strategies must be developed to prevent conflict. But in many real cases, we often get limited resources, such as no powerful GPU and CPU, insufficient human resources, and imperfect allocation strategies, so the existing IoT systems will perform unsatisfactorily. For example, when researchers use cameras to count vehicles, a single node can run and compute without overloading the core server. However, if multiple edge nodes start running, transmitting video, etc., the server may be overloaded. At the same time, human resources will also increase to ensure and maintain the normal operation of the system. Based on such situations, developing a lightweight and easy-to-deploy 'tool' to assist people in serving IoT systems can not only reduce human resource consumption, but also help the system process large amounts of data more efficiently.

To make IoT application be used in real-world scenarios, many researchers combine artificial intelligence technology with IoT [4] [5] [6]. Artificial intelligence can operate and maintain the IoT system like a human. For example, it can help the communication of the IoT system to allocate bandwidth reasonably, and it can also complete some monitoring on the system: vehicle counting, mask detection, etc. Different tasks often require different artificial intelligence models, which also involves detailed branches of AI. For artificial intelligence itself, there are many branches, such as machine learning, deep learning, reinforcement learning, etc. Many of these branches can also intersect with each other, such as deep reinforcement learning. Deep learning is to learn the inherent laws and representation levels of sample data, and the information obtained during the learning process is of great help to the interpretation of the data. Its ultimate goal is to enable machines to have the ability to analyze and learn like humans, and to recognize data such as words, images, and sounds, like image recognition [7][8] and segmentation [9][10], speech recognition [11][12], target tracking [13][14], etc.. This technologies can be used in some special application scenarios, such as during the new crown epidemic, people need to wear masks. Identifying whether people are wearing masks can help the government

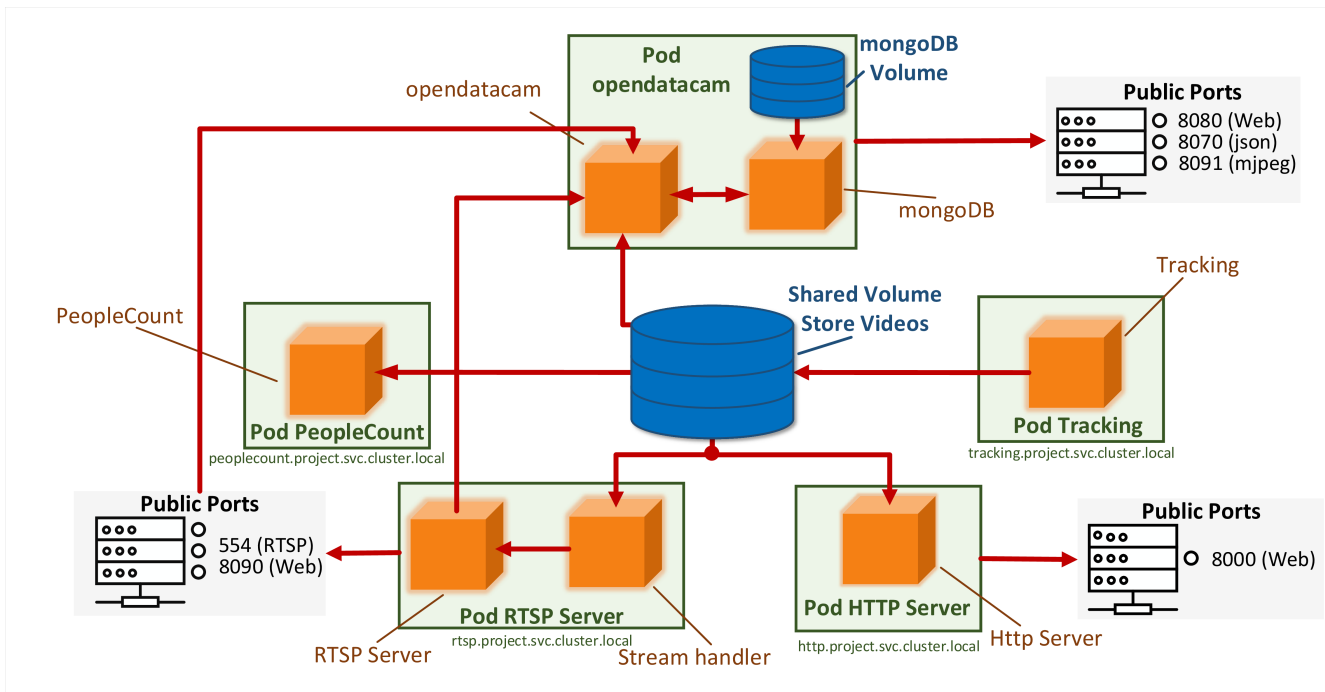


Fig. 1. Overview diagram of our approach.

implement effective epidemic prevention measures. Epidemic prevention personnel can capture pictures of people by calling the cameras in the IoT system, and then use AI for mask identification.

In this paper, we combine the YOLOv5 model [15] with BiFPN [16] and Ghost-Net [17] to propose a lightweight model, GB-YOLO, that can be used for dense object detection. This model is lightweight, which means that it does not require as much computational resources as larger models. Such a model does not preempt too many IoT system resources to keep itself running. This ensures that the system is not easily overloaded if multiple edge nodes start computing at the same time. At the same time, we combined Ghost-Net to further reduce the redundant feature calculation in the image, which can be faster in speed. This can enhance the real-time performance of the system, reduce delays, and collect the latest mask wearing data at the first time. Finally, we introduce BiFPN to replace the PANs network, in order to extract features bidirectionally. Doing so allows our model to have better detection performance than the original model when dealing with dense targets and small targets. On the basis again, we deploy this model on the IoT platform system, as shown in Fig. 1, which can effectively assist people in operations such as vehicle counting and mask detection. We collect videos from *Pod Tracking* and save it in *Share Volume Store Videos*. GB-YOLO model is stored in *Pod opendatacam*. *mongoDB* is a database based on distributed file storage. The *Pod RTSP Server* is a real-time streaming protocol-based server, and its purpose is to transmit the processed data to the client. *Pod HTTP Server* is bound to an IP address and port number and listens for incoming TCP connections from clients on this address. Finally, you can open the web according to the corresponding port number to view the results.

The main contributions of this article are summarized as follows.

- First, we use the Ghost convolution to replace the convolution in YOLOv5 to make the model become a lightweight one, which can reduce redundant features during model processing and speed up the model.
- Second, we modified the PANet layer in YOLOv5 to EfficientDet-BiFPN to realize the bidirectional fusion of top-down and bottom-up deep and shallow features, and enhance the transfer of feature information between different network layers.
- Third, we combine the improved YOLOv5 model with IoT technology to make the system intelligent, which can assist in people counting and data collection.

II. RELATED WORK

The combination of IoT and AI manifests in many fields. Meng Chen and his team proposed an animal motion tracking system based on wireless AI-powered IoT sensors (AIIS). Such a tracking system can monitor and track small animal motions by analyzing the collected accelerations and angular velocities. They used the IoT system to collect information from mice by deploying microchips on them. Then, a machine learning algorithm was developed to process the behavioral information of the mice [18]. Jiwei Zhang and his team propose DRLTrack, a framework for target tracking with a collaborative deep reinforcement learning (C-DRL) in Edge-IoT with the aim to obtain two major objectives: the QoT and resource-efficient network performance, which defines the Edge-AI-aided collaborative tracking through the DQN [19]. Anirudh S. Chakravarthy and his team propose a novel and robust segmentation framework (DroneSegNet) for UAV

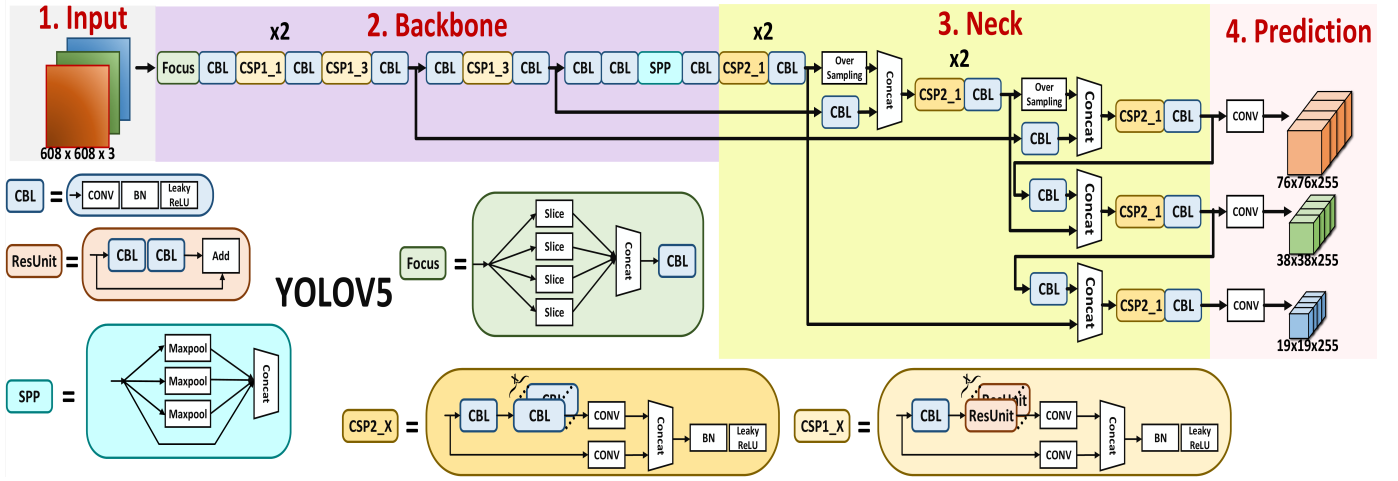


Fig. 2. The structure of YOLOv5.

captured images. The UAVs integrated within IoT framework are capable of delivering IoT services, which can advance the capability for unmanned operations [20][21]. Tao Bai and his team propose a novel framework, called AI-GAN, where a generator, a discriminator, and an attacker are trained jointly, which can be used for testing and evaluating the robustness of IoT systems [22]. Muhammed Golec and his team have proposed a security- and privacy-based lightweight framework called iFaaSBus, which uses the concept of IoT, machine learning (ML), and function as a service (FaaS) or serverless computing to diagnose the COVID-19 disease and manages resources automatically to enable dynamic scalability [23]. In these studies on the combination of AI and IoT, AI is based on IoT. Through the IoT system, data can be continuously supplemented, and training can be completed on this basis. At the same time, the final model after AI training is also used to assist the IoT system. The movement of AI from the cloud to the edge solves the bandwidth and security issues that hinder IoT development. The massive data brought by IoT feeds back AI and complements each other.

Because of this, we can find the most suitable AI tools under the IoT system in a special environment according to the principle of special needs of special scenarios, so as to maximize the effectiveness of both parties. For vehicle counting and mask detection, models in the field of computer vision can assist people well. Many object detection models, such as R-CNN [24] and SSD [25], have shown good results in real-life scenarios, such as face recognition, trajectory tracking, etc. These models perform very well with sufficient training iterations and datasets. Among them, R-CNN also derives more advanced models such as Fast R-CNN [26] and Faster R-CNN [27]. However, these models have some problems: the model file after training becomes larger; the model itself has a complex structure; the model is not easy to be deployed, etc. In the case of limited computing power and computing resources, it will take a lot of time to complete the training, and due to the limited computing power, delays and offsets may occur when performing some real-time tasks. Therefore, a model that is lightweight and easy to deploy

without sacrificing accuracy and speed can easily function in this environment, which are YOLOs (You Only Look Once and its plus versions) [15][28][29][30][31]. The new version of YOLO has made some improvements on the previous version. The v5 version only adds some new improvement ideas on the basis of the v4 version, so that its speed and accuracy have been greatly improved, and the size of the model has been greatly improved and also, much smaller than before. The features of the v5 version of the model just apply to our application in the limited computing situation. Not only that, in order to maximize its effectiveness under limited computing power, we optimized the convolutional neural layer in the v5 version to combine it with Ghost-Net [17]. This can prevent the v5 model from using the redundant features in the feature extraction layer to complete the training, which will cause the extension of training time and fitting problems. At the same time, we modified the PANet layer in YOLOv5 to BiFPN [16] to realize the two-way fusion of top-down and bottom-up deep and shallow features, enhance the transfer of feature information between different network layers, and significantly improve the detection accuracy of the YOLOv5 algorithm. Better detection performance.

III. METHODOLOGY

In this section, we introduce YOLOv5, Ghost-Net and BiFPN, IoT systems and our method, GB-YOLO we propose which consists of YOLOv5, Ghost-Net and BiFPN.

A. YOLOv5

As shown in Fig. 2, the network structure of YOLOv5 is divided into four parts: Input, Backbone, Neck, and Head.

The input end mainly includes three parts: Mosaic data enhancement, image size processing and adaptive anchor box calculation. Mosaic data enhancement combines four images to achieve the effect of enriching the background of images [15][32]; the sizes of images processing adaptively adds the least black borders to the original images of different lengths and widths, and uniformly scales them to the standard size; the

adaptive anchor frame is calculated at the initial anchor. On the basis of the frame, the output prediction frame is compared with the real frame, the gap is calculated and then updated in the reverse direction, and the parameters are continuously iterated to obtain the most suitable anchor frame value.

Backbone mainly includes BottleneckCSP [33] and Focus modules [15][34]. The BottleneckCSP module greatly reduces the amount of computation while enhancing the learning performance of the entire convolutional neural network; the Focus module slices the image, expands the input channel by 4 times, and obtains the down-sampling feature map through a convolution. Downsampling reduces computation and increases speed.

Neck adopts the structure of combining FPN [35] and PANs [36], which combines the conventional FPN layer with the bottom-up feature pyramid, and fuses the extracted semantic features and position features. The model obtains richer feature information.

Head outputs a vector with the class probability of the target object, the object score, and the location of the bounding box for that object. The detection network consists of three detection layers, and feature maps of different sizes are used to detect target objects of different sizes. Each detection layer outputs the corresponding vector, and finally generates the predicted bounding box and category of the object in the original image and labels it.

The previous v4 version, the structure is roughly the same as v5, but not as good as v5 in some details.

- Input: in the model training stage, some improvement ideas are proposed, including Mosaic data enhancement [32], adaptive anchor box calculation, and adaptive image scaling.
- Backbone: v5 integrates some new ideas in other detection algorithms, like Focus structure [34] and CSP structure [33].
- Neck network: the target detection network often inserts some layers between the BackBone and the last Head output layer, and the FPN+PAN [35][36] structure is added to Yolov5.
- Head output layer: the anchor box mechanism of the output layer is the same as that of YOLOv4. The main improvements are the loss function GIOU_Loss [37] during training, and the DIOU_nms [38] for prediction box screening.

B. Ghost-Net

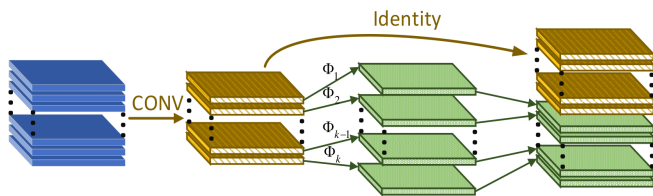


Fig. 3. The Ghost module.

Sufficient or redundant information in feature layers can always ensure a comprehensive understanding of the input data, and there are many similarities between feature layers, which are like ghosts of each other. This is also the origin of the Ghost network name. The core of it is that there is no

need to generate these redundant feature maps with a large number of FLOPs and parameters. It is because the feature maps are partially similar, so there is no need to deconvolute them all, and directly use part of them as repeated redundant information to process [17].

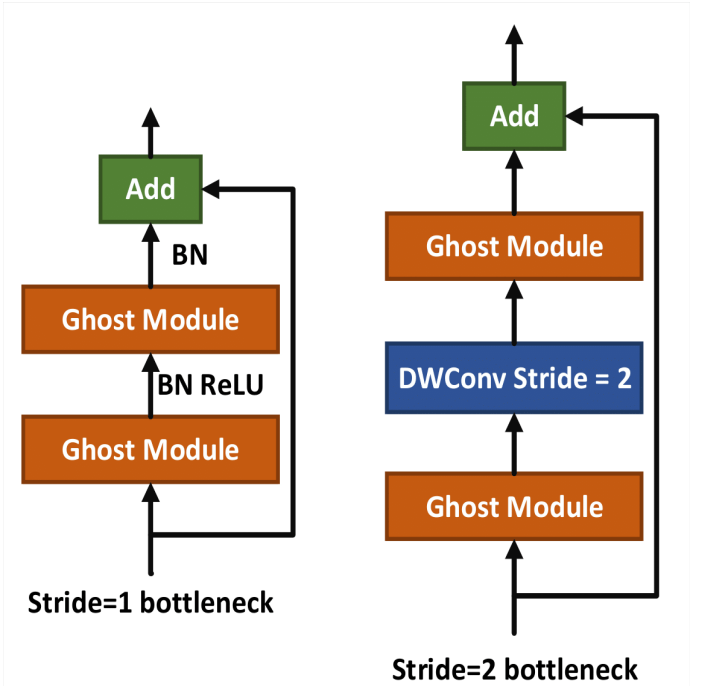


Fig. 4. Ghost bottleneck with different parameters.

As shown in Fig. 3, the Module in GhostNet is divided into two steps to obtain the same number of feature maps as ordinary convolutions: i. smaller amount of convolution. ii. cheap operations. In this way, the feature maps of the upper and lower parts are phantoms of each other. The first step is to apply a traditional convolution operation to the input tensor to obtain a tensor of yellow feature maps from light to dark in the Fig. 3, also known as intrinsic features. Then, the linear operation represented by the ϕ function (Depthwise Separable Convolution) [39][40] is applied to the feature map to generate a red feature map tensor from light to dark in the figure, and finally stacked together as the output of the Ghost Module, which is used to generate a corresponding Ghost feature map for a series of single-channel feature maps.

Set the width of the input tensor to be w , the height to be h , the number of channels to be c , the number of channels of the output tensor to be n , and the size of the standard convolution kernel to be $k*k$. The *conv* here should not contain nonlinear activation functions. The output of feature maps in the Ghost Module can be divided into n/s groups, each of which is s feature maps. The s feature maps contain the calculation result of a standard convolution kernel and the result obtained by $s-1$ linear transformations. The operation of adding offsets is ignored, and the approximation task $k*k$ is equal to $d*d$. Then the ratio of computational complexity between ordinary convolution and Ghost-Net is as follows:

$$\begin{aligned}
 r_s &= \frac{n \cdot h' \cdot w' \cdot c \cdot k \cdot k}{\frac{n}{s} \cdot h' \cdot w' \cdot c \cdot k \cdot k + (s-1) \cdot \frac{n}{s} \cdot h' \cdot w' \cdot d \cdot d} \\
 &= \frac{c \cdot k \cdot k}{\frac{1}{s} \cdot c \cdot k \cdot k + \frac{s-1}{s} \cdot d \cdot d} \approx \frac{s \cdot c}{s + c - 1} \approx s
 \end{aligned} \quad (1)$$

As shown in (1), using the Ghost Module to transform the standard convolution kernel can reduce the calculation amount and the number of parameters of the model by approximately s times.

Based on advantages of Ghost Modules, the author introduces the Ghost bottleneck (G-bneck) specially designed for small CNNs, as shown in Fig. 4. Borrowing from MobileNetV2, ReLU is not used after the second Ghost module, and other layers apply BN and ReLU after each layer. For efficiency reasons, the initial convolution in the Ghost module is a point convolution.

And this structure mainly adopts the common design pattern of channel enlargement and reduction. It first passes through a Ghost Module that enlarges the number of channels, and then connects to a Ghost module that reduces the number of channels, so that the number of channels before and after remains unchanged, so as to perform a channel-by-channel addition operation with the original tensor connected through the shortcut. There are two types of bottlenecks of this type. The structure on the right side in Fig. 4 is based on the structure on the left, and a Depthwise Conv with stride=2 is added. Correspondingly, in order to be able to add directly, the shortcut branch needs to perform a downsampling.

C. BiFPN

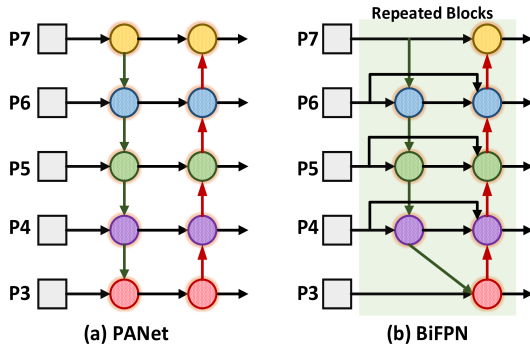


Fig. 5. Structural comparison of PANs and BiFPN.

PANs [36] is an improved FPN algorithm [35] which adds a bottom-up channel, but it brings a problem, that is, a large amount of computation, and BiFPN [16] improves it on the basis of FPN, adds edges with contextual information to the original FPN module, and multiplies each edge with a corresponding weight. As shown in Fig. 5, compared with PANs, BiFPN adds residual links to enhance the representation ability of features; moreover, it removes nodes with a single input edge, because the nodes on the input edge do not perform feature fusion, so they have information comparison Less, there is no contribution to the final fusion, on the contrary, removal can also reduce the amount of calculation. At the same

time, it also performs weight fusion, that is, adding a weight to each scale feature of the fusion to adjust the contribution of each scale. Among them, Fast-softmax is proposed to improve the detection speed.

Since different input features have different resolutions, their contributions to output features are usually unequal. So let each input add extra weight and let the network understand the importance of each input feature. M. Tan and his team have considered three ways to solve the problem: i. Unbounded fusion; ii. Fusion based on softmax; iii. Fast normalized fusion.

As shown in (2), w_i is a learnable weight that can represent vectors, scalars and multidimensional tensors. However, since the scalar weights are unbounded, the training will be unstable. Based on this, weight normalization is used to limit the value range.

$$O = \sum_i w_i \cdot I_i \quad (2)$$

Also, as shown in (3), (2) is combined with softmax yields (3).

$$O = \sum_i \frac{e^{w_i}}{\sum_j e^{w_j}} \cdot I_i \quad (3)$$

(3) Apply softmax to each weight such that all weights are normalized to the 0 to 1 range to represent their input importance. And then, (3) adopts a fast fusion method, which is (4), and its value is passed through the Relu function to ensure numerical stability, and the value of its normalized weight is also between 0 and 1.

$$O = \sum_i \frac{w_j}{\epsilon + \sum_j w_j} \cdot I_i \quad (4)$$

In general, Bi-FPN is equivalent to giving different weights to each layer for fusion, making the network pay more attention to important layers, and reducing the node connections of some unnecessary layers.

D. IoT Platform

The establishment of the entire IoT system and how to deploy and operate on the edge node and Kubemetes platform are shown in Fig. 6. We combine the Dockerfile with the dependency file to create a docker file, and then upload the file to the docker hub. Furthermore, we store the YOLOv5 configuration file and the improved configuration file in config.json, and participate in the development of Kubemetes [41] together with docker hub through configmap. Then run it in opendatacampod, and apply for video data from mongoDB pod through port 27017. Finally, we can open a browser through port: 8070, 8080, or 8090 to view the output results.

E. GB-YOLO

As shown in Fig. 7, The general structure of the improved YOLOv5 algorithm has not changed. We replaced the original CSP structure with Ghost Bottleneck, which can greatly optimize the parameter scale and computing resource consumption of the network without affecting the detection accuracy of network pre-training. There are many other ways to reduce

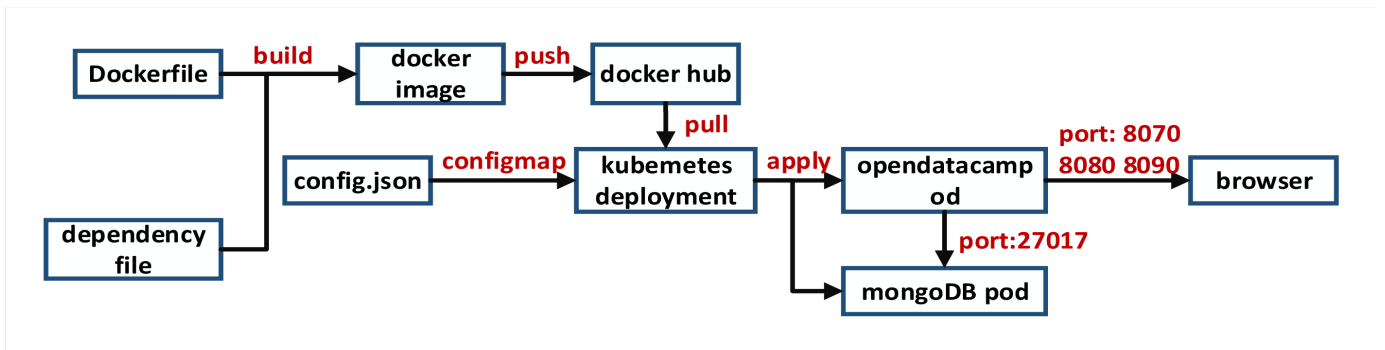


Fig. 6. Flowchart of the IoT system we use.

redundant parameter computation and redundant features, such as atrous convolution [42]. The receptive field of atrous convolution is very large. When the amount of parameters is certain, ordinary convolution can only extract small blocks of features, while atrous convolution can increase the hole rate to make more overlapping sampling areas on the input feature map for each sampling, so as to obtain denser characteristic response. Atrous convolution can be used when the network layer needs a large receptive field, but the computing resources are limited and the number or size of the convolution kernel cannot be increased. But this also brings two problems: not all pixels participate in the calculation, so the obtained features are discontinuous, which is defective in pixel-level detection; the information using a large dilation rate is only effective for some large object segmentation, while for small objects it does not help. When our model detects masks, some of the masks have a very small proportion of the screen, or the video is blurry and the resolution is not high. These situations will lead to difficulties in sampling and thus cannot be accurately identified. Therefore, atrous convolution is not suitable for this scenario.

And also, we replaced the Concat unit with BiFPN, which can realize the two-way fusion of top-down and bottom-up deep and shallow features, and enhance the transfer of feature information between different network layers. Because the resolutions of the feature maps of different layers are different, their contributions to the fused features are different, so weights are introduced in the feature fusion stage. Unlike PANs, PANs adds a bottom-up channel based on FPN, which

makes PANs have better accuracy than FPN, but requires more parameters and computation. If the node with only one input edge and output edge in the PANs is removed, and if the input and output nodes are at the same level, an extra edge is added to fuse more features without increasing the cost. As shown in (2), there is no limit on weight, which may cause unstable training. In (2), the calculation of softmax is slow, which is especially prominent in the front end. In order to ensure that the weight is greater than 0, the relu function is used before the weight. In order to ensure that the weight is greater than 0, the relu function is used before the weight. (3) Similar accuracy to (2) but can be 30% faster. Based on it, BiFPN achieves similar accuracy to repeated FPN+PANet, but uses far fewer parameters and FLOPs. With additional weighted feature fusion, our BiFPN further achieves the best accuracy with fewer parameters and FLOPs.

And then, as shown in Fig. 1, we collect videos from *Pod Tracking* and save it in *Share Volume Store Videos*. GB-YOLO model is stored in *Pod opendatacam*. *mongoDB* is a database based on distributed file storage. The *Pod RTSP Server* is a real-time streaming protocol-based server, and its purpose is to transmit the processed data to the client. *Pod HTTP Server* is bound to an IP address and port number and listens for incoming TCP connections from clients on this address. Finally, you can open the web according to the corresponding port number to view the results.

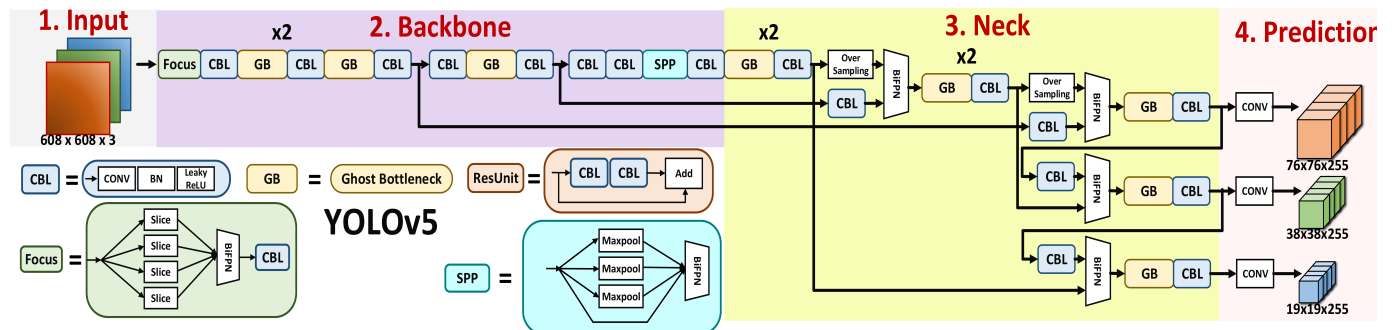


Fig. 7. The structure of our proposed method.

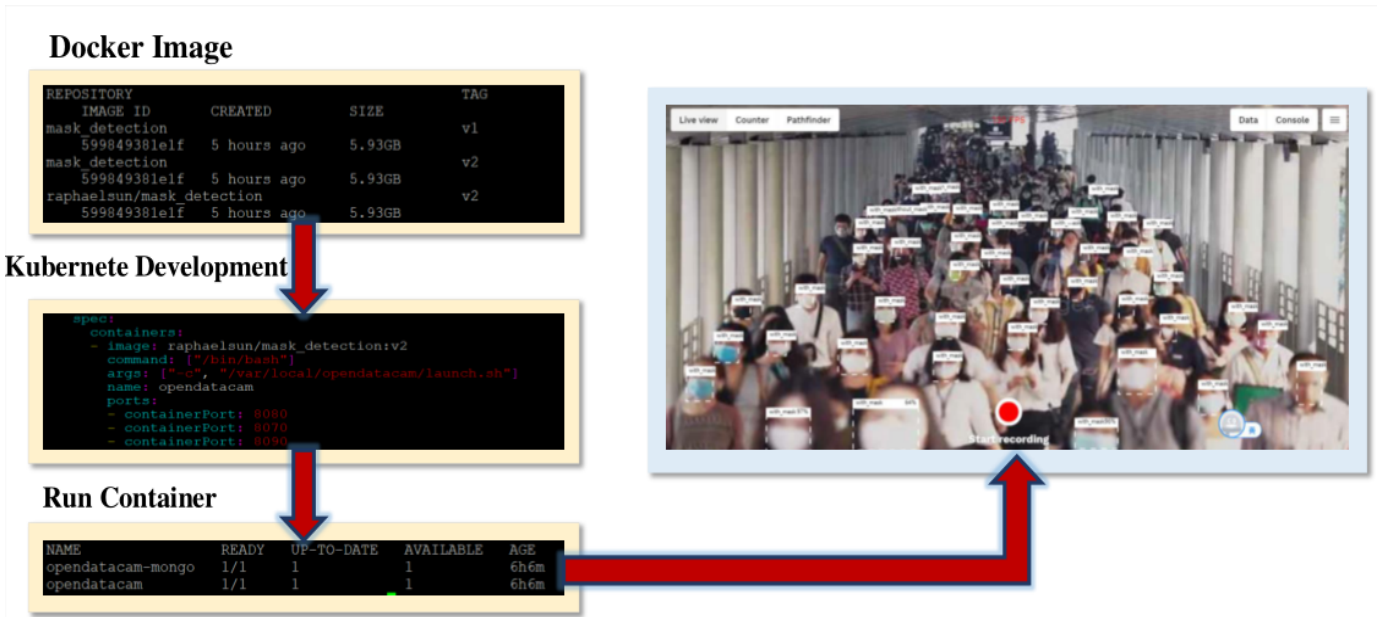


Fig. 8. The experimental result of our method.

IV. HARDWARE EQUIPMENT AND DATA PREPARATION

A. Hardware Equipment

Our GB-YOLO model and deployment of the iot system use the Ubuntu18.04 operating system, 16GB of memory, and a weak GPU acceleration. Several related servers are on this system.

The Docker weight file and the Kubernetes deployment are both on the personal computer. The final viewing of the results via the web was also performed on a personal computer.

B. Dataset

Our dataset comes from RMFD(Real-World Masked Face Dataset) [43]. The dataset is mainly divided into two parts: real

face mask data and simulated mask face data. Among them, the real mask face recognition dataset contains 5,000 mask faces and 90,000 normal faces of 525 people. The simulated mask face recognition dataset includes a simulated mask face dataset of 10,000 people and 500,000 faces. These images can be annotated with Labelme software to generate labels, and then divide the dataset. We use a ratio of 8:2 to divide the training and testing datasets in this article.

C. Performance Measurements

For the measurement of detection or recognition performance of machine learning model results, confusion metric is widely used, it is a performance-based metric, and the widely used metrics for model evaluation are discussed as follows:



Fig. 9. Vehicle tracking and its count.

1. *True Positive (TP)*: In attack detection, the *TP* indicates that Class A is correctly identified as belonging to Class A.

2. *True Negative (TN)*: This matrix indicates that Class A is correctly identified as not belonging to Class A.

3. *False Positive (FP)*: It indicates that Class A is not correctly identified as belonging to Class A.

4. *False Negative (FN)*: It indicates that Class A is not correctly identified as not belonging to Class A.

However, using the metrics above, different measures can be made to better evaluate the model. For accurate detection, the classifier minimizes the values of the FP and FN metrics. However, the selected metrics used in this article are detailed below.

Accuracy: In mask detection, it can be described as Whether or not wearing a mask is a false detection that does not match the real situation. However, using performance measurement metrics, the accuracy can be defined mathematically as

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

And we only used the above given equation for the proposed technique performance evaluation

V. EXPERIMENTAL RESULTS AND ANALYSIS

A. Experimental Results

As shown in Fig. 8, this is the final result of our experiment after combining the IoT system with the improved YOLOv5 model. We opened the IP address and its port output by the server on the personal computer, and the results showed the real-time detection of masks. The experimental results will not be displayed at the beginning, because the CPU is used for computing, and it will take a while to appear.

TABLE I
COMPARISON OF MASK DETECTION EFFECTS OF DIFFERENT MODELS

Model	mask	un_mask	FPS
YOLOv4	95.12%	96.03%	74
YOLOv5	96.62%	97.87%	107
Faster R-CNN	95.78%	95.31%	71
SSD	95.98%	96.63%	96
GB-YOLO	97.08%	97.85%	128

We first deployed the relevant Docker Image on the server, then combined it with Kubernetes, then we ran the container through Kubernetes, and finally opened the corresponding port through the web to see the result.

Not only that, we also extend our method to other domains such as living objects and their statistics, vehicle counting, etc, which is shown in Fig. 9. In this web server, we have also developed many other functions, such as recording, scribing, etc. which can assist people to count data and measure traffic. At the same time, the data will also be stored in the background, which is convenient for secondary reference.

At the same time, we also tested GB-YOLO separately, as shown in Fig. 10 and compared it with other models, as shown in Table. I. We compare v4, v5, SSD, Faster R-CNN and our model. These models are all run on the original operating system, and the same dataset is used for experiments. Among



Fig. 10. The experimental result of our improved model.

them, the test of the number of frames is mainly based on the video of mask detection, and the rest of the video such as vehicle counting is not included in it.

B. System Performance

We recorded the operation of each part of the system when the IoT system was running. When the system starts to run the IoT system and activates a node, the GPU resources consumed are not large, and as the mask detection begins, the GPU and CPU speed up. Also always filled. However, compared with the previous v5 model, the detection effect can make better use of gpu resources. In the same case of full load, our method is faster and has a higher frame rate. Under the same limited computing power, our improved model can support detection of more objects.

And we need less computing power for the same detection effect, which allows our model to independently detect other

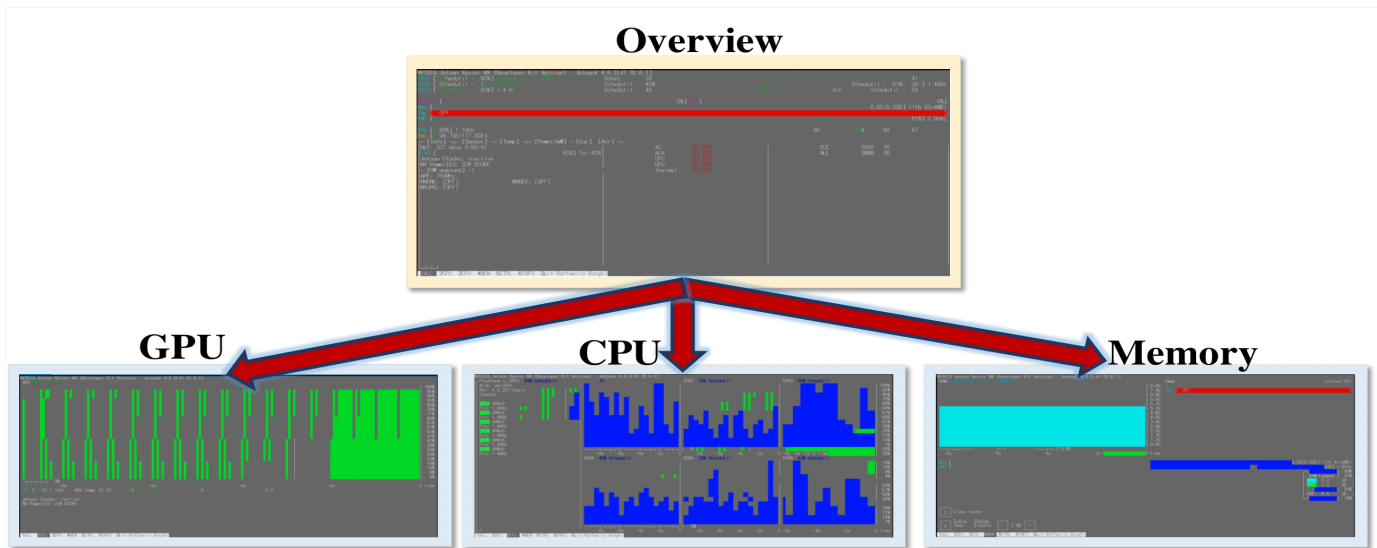


Fig. 11. The system performance of our method.

objects concurrently or run two improved models at the same time.

C. Analysis of Results

Compared with the previous target detection algorithm, the YOLOv5 model itself is lightweight, and the speed is much faster than the previous model, and the accuracy is not inferior. On this basis, we combine it with Ghost-Net and BiFPN, reduce its convolutional layer parameters and redundant features and increase bidirectional adoption, so that the model can spend less computing power and obtain more effective features. Correspondingly, its accuracy has been greatly improved compared to the previous model, and the number of detected frames has also increased. And our IoT system can support the operation of such a lightweight model without being overloaded when less computing power is allocated. On this basis, the IoT system can also run other projects at the same time. The system itself plus a lightweight model, easy to deploy and quick to clone. If a large model such as Faster R-CNN is deployed, the system can run, but the number of edge nodes in parallel and the computing power allocated will be reduced. At the same time, the load on the system itself will also increase, making it unable to play its full effect.

VI. CONCLUSION

In this paper, we successfully combine an IoT system and an improved YOLOv5 model for mask detection in daily life. Our method can well assist government personnel and anti-epidemic workers to count the wearing of masks, so as to formulate effective anti-epidemic policies. We use the lightweight model YOLOv5 and combine Ghost-Net and BiFPN on it, which enables it to exert more computing power with less computing resources. The improved YOLOv5, GB-YOLO, combined with our deployed IoT system can better perform concurrent operations, and different edge nodes can

run different projects at the same time, such as vehicle counting, person tracking, mask detection, etc. Furthermore, this IoT system combined with YOLOv5 is easy to expand and deploy, and it is very easy to clone on the server. Its lightweight setting can help people maximize efficiency under the condition of limited computing power.

VII. FUTURE WORK

During the experiment, we discovered several interesting phenomena, such as Fig. 12 and Fig. 13.

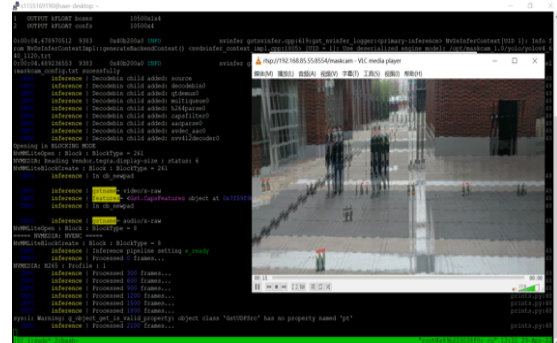


Fig. 12. The system performance of our method.

We found that if the environmental parameters required by deep learning do not meet the requirements of the platform, it will be incompatible with imagination. This can lead to some unintended consequences, such as bokeh in the image. In this case, although the wearing of the mask can be detected on the server side, when we observe the edge nodes, we find that the pixels are extremely blurred, as shown in Fig. 12. Although there are a few frames where you can see the mask being detected on the web, it will not be very helpful for human inspection. If you look up the mask location parameters of the server, it is difficult for non-professionals to understand the meaning, which will make it difficult for people to directly understand.



Fig. 13. The system performance of our method.

For improving the model, we also found that if a person is wearing a mask and is not wearing it, but there are other distractors, such as a hand blocking it, there will be a superimposed state of wearing and not wearing a mask, as shown in Fig. 13. In this case, we will analyze from the inside of the neural network. And on this basis, a model for mask detection is being specially developed for the future.

REFERENCES

[1] X. Liu and X. Zhang, "Rate and Energy Efficiency Improvements for 5G-Based IoT With Simultaneous Transfer," in *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 5971-5980, Aug. 2019.

[2] G. Xu et al., "TT-SVD: An Efficient Sparse Decision-Making Model With Two-Way Trust Recommendation in the AI-Enabled IoT Systems," in *IEEE Internet of Things Journal*, vol. 8, no. 12, pp. 9559-9567, 15 June 2021.

[3] J. Hwang, L. Nkenyereye, N. Sung, J. Kim and J. Song, "IoT Service Slicing and Task Offloading for Edge Computing," in *IEEE Internet of Things Journal*, vol. 8, no. 14, pp. 11526-11547, 15 July 2021.

[4] F. Hussain, R. Hussain, S. A. Hassan and E. Hossain, "Machine Learning in IoT Security: Current Solutions and Future Challenges," in *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 1686-1721, thirdquarter 2020.

[5] B. Yuan, J. Wang, P. Wu and X. Qing, "IoT Malware Classification Based on Lightweight Convolutional Neural Networks," in *IEEE Internet of Things Journal*, vol. 9, no. 5, pp. 3770-3783, 1 March 2022.

[6] B. Mao, Y. Kawamoto and N. Kato, "AI-Based Joint Optimization of QoS and Security for 6G Energy Harvesting Internet of Things," in *IEEE Internet of Things Journal*, vol. 7, no. 8, pp. 7032-7042, Aug. 2020.

[7] R. He, J. Cao, L. Song, Z. Sun and T. Tan, "Adversarial Cross-Spectral Face Completion for NIR-VIS Face Recognition," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 5, pp. 1025-1037, 1 May 2020.

[8] C. Fu, X. Wu, Y. Hu, H. Huang and R. He, "DVG-Face: Dual Variational Generation for Heterogeneous Face Recognition," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 6, pp. 2938-2952, 1 June 2022.

[9] K. He, G. Gkioxari, P. Dollár and R. Girshick, "Mask R-CNN," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 386-397, 1 Feb. 2020.

[10] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh and J. Liang, "UNet++: Redesigning Skip Connections to Exploit Multiscale Features in Image Segmentation," in *IEEE Transactions on Medical Imaging*, vol. 39, no. 6, pp. 1856-1867, June 2020.

[11] C. Fan, J. Yi, J. Tao, Z. Tian, B. Liu and Z. Wen, "Gated Recurrent Fusion With Joint Training Framework for Robust End-to-End Speech Recognition," in *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 198-209, 2021.

[12] L. Chai, J. Du, Q. -F. Liu and C. -H. Lee, "A Cross-Entropy-Guided Measure (CEGM) for Assessing Speech Recognition Performance and Optimizing DNN-Based Speech Enhancement," in *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 106-117, 2021.

[13] Z. Zhang, B. Zhong, S. Zhang, Z. Tang, X. Liu and Z. Zhang, "Distractor-Aware Fast Tracking via Dynamic Convolutions and MOT Philosophy," 2021 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 1024-1033.

[14] N. Wang, W. Zhou, J. Wang and H. Li, "Transformer Meets Tracker: Exploiting Temporal Context for Robust Visual Tracking," 2021 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 1571-1580.

[15] YOLOv5, <https://github.com/ultralytics/yolov5>.

[16] M. Tan, R. Pang and Q. V. Le, "EfficientDet: Scalable and Efficient Object Detection," 2020 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 10778-10787.

[17] K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu and C. Xu, "GhostNet: More Features From Cheap Operations," 2020 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1577-1586.

[18] M. Chen et al., "Wireless AI-Powered IoT Sensors for Laboratory Mice Behavior Recognition," in *IEEE Internet of Things Journal*, vol. 9, no. 3, pp. 1899-1912, 1 Feb. 2022.

[19] J. Zhang, M. Z. A. Bhuiyan, X. Yang, A. K. Singh, D. F. Hsu and E. Luo, "Trustworthy Target Tracking With Collaborative Deep Reinforcement Learning in EdgeAI-Aided IoT," in *IEEE Transactions on Industrial Informatics*, vol. 18, no. 2, pp. 1301-1309, Feb. 2022.

[20] A. S. Chakravarthy, S. Sinha, P. Narang, M. Mandal, V. Chamola and F. R. Yu, "DroneSegNet: Robust Aerial Semantic Segmentation for UAV-Based IoT Applications," in *IEEE Transactions on Vehicular Technology*, vol. 71, no. 4, pp. 4277-4286, April 2022.

[21] N. Hossein Motlagh, T. Taleb and O. Arouk, "Low-Altitude Unmanned Aerial Vehicles-Based Internet of Things Services: Comprehensive Survey and Future Perspectives," in *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 899-922, Dec. 2016.

[22] T. Bai et al., "Toward Efficiently Evaluating the Robustness of Deep Neural Networks in IoT Systems: A GAN-Based Method," in *IEEE Internet of Things Journal*, vol. 9, no. 3, pp. 1875-1884, 1 Feb. 2022.

[23] M. Golec, R. Ozturac, Z. Pooranian, S. S. Gill and R. Buyya, "iFaaSBus: A Security- and Privacy-Based Lightweight Framework for Serverless Computing Using IoT and Machine Learning," in *IEEE Transactions on Industrial Informatics*, vol. 18, no. 5, pp. 3522-3529, May 2022.

[24] R. Girshick, J. Donahue, T. Darrell and J. Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation," 2014 *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580-587.

[25] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C. Y. Fu, et al., "SSD: Single shot multibox detector," *European Conference on Computer Vision*, pp. 21-37, 2016.

[26] R. Girshick, "Fast R-CNN," 2015 *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1440-1448.

[27] S. Ren, K. He, R. Girshick and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137-1149, 1 June 2017.

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

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

[30] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement" in *arXiv:1804.02767*, 2018.

[31] Alexey Bochkovskiy, Chien-Yao Wang, Hong-Yuan Mark Liao, "YOLOv4: Optimal speed and accuracy of object detection," 2020.

[32] I. Lai and W. Tsai, "Secret-Fragment-Visible Mosaic Image-A New Computer Art and Its Application to Information Hiding," in *IEEE Transactions on Information Forensics and Security*, vol. 6, no. 3, pp. 936-945, Sept. 2011.

[33] C. -Y. Wang, H. -Y. Mark Liao, Y. -H. Wu, P. -Y. Chen, J. -W. Hsieh and I. -H. Yeh, "CSPNet: A New Backbone that can Enhance Learning Capability of CNN," 2020 *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020, pp. 1571-1580.

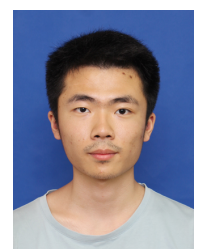
- [34] Z. Li et al., "Toward Efficient Safety Helmet Detection Based on YoloV5 With Hierarchical Positive Sample Selection and Box Density Filtering," in *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1-14, 2022, Art no. 2508314.
- [35] T. -Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan and S. Belongie, "Feature Pyramid Networks for Object Detection," 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 936-944, doi: 10.1109/CVPR.2017.106.
- [36] S. Liu, L. Qi, H. Qin, J. Shi and J. Jia, "Path Aggregation Network for Instance Segmentation," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 8759-8768.
- [37] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, Silvio Savarese, "Generalized Intersection over Union: A Metric and A Loss for Bounding Box Regression," arXiv: 1902.09630, 2019.
- [38] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, Dongwei Ren, "Distance-IoU Loss: Faster and Better Learning for Bounding Box Regression," in arXiv: 1911.08287, 2019.
- [39] L. Sifre and S. Mallat, "Rigid-motion scattering for texture classification", arXiv:1403.1687 [cs], Mar. 2014.
- [40] A. Howard et al., "Searching for MobileNetV3," 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 1314-1324.
- [41] Kubernetes, <https://kubernetes.io/>.
- [42] L. -C. Chen, G. Papandreou, I. Kokkinos, K. Murphy and A. L. Yuille, "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834-848, 1 April 2018.
- [43] Real-World Masked Face Dataset, <https://github.com/X-zhangyang/Real-World-Masked-Face-Dataset>.



Cheng Huang Cheng Huang is currently pursuing the M.S. degree at The Chinese University of Hong Kong, Hong Kong, China. He received his B.Eng degree in Software Engineering from University of Electronic Science and Technology of China, Chengdu, China in 2020. He was a research assistant of Department of Electronic Engineering in Tsinghua University, Beijing, China, in 2021. His research interests include artificial intelligence and deep learning. He will pursue a Ph.D. in the future.



Yishen Liu Yishen Liu is currently pursuing the M.S. degree in Public Health at the Chinese University of Hong Kong. In 2021, she received her B.S. degree in Chinese Medicine in Guangzhou University of Chinese Medicine, China. Her current research interests include in health systems strengthening, disease control and chinese medicine development. She is working as a research assistant in Perking University, ShenZhen Graduate School, School of Chemical Biology and Biotechnology now.



Zihan Xia Zihan Xia is currently pursuing the M.S. degree at the University of Electronic Science and Technology of China, Chengdu, China. He received his B.Eng degree in Communication Engineering from University of Electronic Science and Technology of China in 2021. He was selected as a member of the Yingcai Honors College in 2018. His research interests include deep learning and energy-efficient circuits. He is going to go to The University California San Diego in September 2022 to pursue a Ph.D.



Siyang Jiang Siyang Jiang received his M.S. degree in Electrical Engineering from National Taiwan University in 2021. He received his B.S. degree in Information & Computation Science from Shaanxi Normal University in 2019. His research interests include few-shot learning, high-performance computing and embedding system. He is going to The Chinese University of Hong Kong for pursuing a Ph. D. degree.



Hao Tian Tian Hao is currently pursuing the M.S. degree in Information Engineering at the Chinese University of Hong Kong. In 2021, he received his B.S. degree in computer science and technology and B.Eng degree in Apparel Design and Engineering from Taiyuan University of Technology, China. His current research interests include smart wearable devices, AIoT, and computer vision. He is going to go to Hong Kong Polytechnic University in September 2022 to pursue a Ph.D.



Yongbin Yu Yongbin Yu was born in Sichuan, China, in 1975. He received the Ph.D. degree from the University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2008. He visited the University of Michigan at Ann Arbor, Ann Arbor, MI, USA, in 2013, and the University of California at Santa Barbara, Santa Barbara, CA, USA, in 2016. He worked as the Guest Deputy Director with the Department of Big Data Industry, Sichuan Provincial Economic and Information Commission, in 2018. He is currently an Associate

Professor with the School of Information and Software Engineering, UESTC. His research interests include memristor-based neural network, swarm intelligence, natural language processing, and big data. Dr. Yu won the First Prize of Science and Technology Award of Tibet Autonomous Region in 2018.