

## Research Article

# DetReco: Object-Text Detection and Recognition Based on Deep Neural Network

Fan Zhang <sup>1,2</sup>, Jiaxing Luan <sup>1</sup>, Zhichao Xu,<sup>1</sup> and Wei Chen <sup>1</sup>

<sup>1</sup>School of Electrical and Information Engineering, China University of Mining and Technology (Beijing), Beijing Ding No. 11, Xueyuan Road, 100083 Beijing, China

<sup>2</sup>Institute of Intelligent Mining and Robotics, China University of Mining and Technology (Beijing), Beijing Ding No.11, Xueyuan Road, 100083 Beijing, China

Correspondence should be addressed to Fan Zhang; zf@cumtb.edu.cn

Received 13 April 2020; Accepted 12 June 2020; Published 14 July 2020

Guest Editor: Chi-Hua Chen

Copyright © 2020 Fan Zhang et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Deep learning-based object detection method has been applied in various fields, such as ITS (intelligent transportation systems) and ADS (autonomous driving systems). Meanwhile, text detection and recognition in different scenes have also attracted much attention and research effort. In this article, we propose a new object-text detection and recognition method termed “DetReco” to detect objects and texts and recognize the text contents. The proposed method is composed of object-text detection network and text recognition network. YOLOv3 is used as the algorithm for the object-text detection task and CRNN is employed to deal with the text recognition task. We combine the datasets of general objects and texts together to train the networks. At test time, the detection network detects various objects in an image. Then, the text images are passed to the text recognition network to derive the text contents. The experiments show that the proposed method achieves 78.3 mAP (mean Average Precision) for general objects and 72.8 AP (Average Precision) for texts in regard to detection performance. Furthermore, the proposed method is able to detect and recognize affine transformed or occluded texts with robustness. In addition, for the texts detected around general objects, the text contents can be used as the identifier to distinguish the object.

## 1. Introduction

Object detection [1, 2], as one of the most fundamental and challenging problems in computer vision, has received great attention in recent years. In the context of computer vision, object detection deals with the task of detecting instances of visual objects of specific classes such as humans, animals, and cars in digital images. It combines the cutting-edge technologies in many fields such as image processing, pattern recognition, automatic control, and artificial intelligence. Object detection is widely used in many fields including intelligent transportation systems [3, 4], advanced driver assistance systems (ADAS), and autonomous driving systems.

In intelligent traffic surveillance system [5], vehicle detection and recognition are a vital task. The automatic monitoring digital cameras take snapshots of passing

vehicles and other moving objects to provide valuable clues including license plate number, the vehicle type, and the driver’s facial image for authorities and other security departments. In recent years, autonomous cars and driverless vehicles have significantly changed the manner of transportation. Computer vision system is efficiently used in the development of ADAS. Sakhare et al. [6] have a detailed study of the vehicle detection in dynamic conditions. Yudin et al. [7] study vehicle detection in difficult areas with various architectures of deep neural networks [8].

In automated driving, detection and recognition of pedestrians, vehicles, traffic lights, and traffic signs [9] help avoid accidents and achieve safe driving. Collision avoidance systems are required for the driver to handle the emergence. Detecting pedestrians is essential for autonomous driving [10]. Zhang and Kim [11] propose a pedestrian detector which combines skip pooling from multiscale feature maps

and recurrent convolutional layers to detect pedestrians of different scales. Reliable traffic light detection and classification in urban environments are also crucial for automated driving [12, 13]. Kim et al. [12] develop a two-step method to detect traffic lights with SSD architecture. Lu et al. [14] utilize a visual attention model to detect traffic signals which is effective for the detection of small objects.

Object detection, which is the core of various intelligent transportation systems, has been a research hotspot in recent years. Meanwhile, the rapid development of deep learning has accelerated the development of object detection. Many deep learning based object detection techniques have led to giant breakthroughs and remarkable performance. Object detection can be divided into one-stage methods and two-stage methods. Object detection algorithm of two-stage methods usually involves two steps. Firstly, region proposals are obtained from the original image. Secondly, the classification and regression networks such as the R-CNN [15] (Regional Convolutional Neural Network) series are used to detect the region proposals. Object detection algorithm of one-stage method just needs one step. One-stage methods can accomplish the classification and bounding box regression tasks directly without finding the region proposals separately. Typical one-stage algorithms include SSD [16] (Single Shot Multibox Detector) and YOLO [17] (You Only Look Once).

R-CNN proposed by Ross B. Girshick uses selective search [18] method to perform ROI (Region of Interest) scaling and feature extraction on target images. Because R-CNN requires forward calculation for a large number of region candidates which may overlap each other, the speed of training and detection is very slow. Fast R-CNN [19] uses a feature extractor to extract the features of the entire image instead of extracting each image multiple times for each region proposal. Because Fast R-CNN does not extract features repeatedly, the processing time is significantly reduced. Faster R-CNN [20] uses a design similar to Fast R-CNN. Faster R-CNN replaces the selective search method with RPN (region proposal network), which solves the problem of excessive time overhead in generating ROI. The Faster R-CNN achieves high accuracy and detection speed to some extent, but it still cannot meet the real-time requirement.

Compared with Faster R-CNN, SSD has a significant advantage of detection speed. The network generates multiple feature maps at different scales. Then the classification and bounding box regression tasks are simultaneously done on multiscale feature maps. SSD is able to detect large objects effectively. YOLO is another one-stage method. It predicts bounding boxes and class probabilities of multiple objects simultaneously. However, different from the SSD algorithm, YOLO does not use multiscale feature maps for detection. Its generalization capability is poor for object with large scale variations compared with that of SSD. It leads to missed detection and low recognition accuracy. YOLOv2 [21] algorithm uses anchor mechanism which utilizes convolutional layers instead of fully connected layers as in YOLO to predict the bounding boxes. The disadvantage of using fully connected layer to predict bounding boxes is that the spatial

information of feature map is lost. However, the anchor mechanism directly predicts the bounding boxes on the feature map with convolutional layers. The spatial information of feature map is well preserved. Each feature point of the feature map corresponds to each grid of the original image. YOLOv2 improves the performance of the detection accuracy. YOLOv3 [22] algorithm adopts multiscale feature maps to predict bounding boxes. YOLOv3 uses FPN (Feature Pyramid Networks) concept which uses the output of the middle layers to merge with the output of the latter layers. The high-level features are passed to the low layers, so that small objects on low-level feature maps can be better detected. YOLOv3 has been greatly improved in regard to detection speed and accuracy.

The majority of the recent works related to deep neural networks has been devoted to detection or classification of object categories [23]. On top of that, another problem in computer vision that plays a vital role in intelligent transportation systems is the image-based text recognition. Text recognition aims to decode a sequence of labels from cropped text images.

The conventional methods recognize the text contents at character level. The characters of the text are segmented from the cropped text image. Then the segmented character regions are preprocessed and recognized. Different from the character-level recognition methods, recent text recognition methods do not require character segmentation in advance. One famous method is the multidigit number classification proposed by Goodfellow et al. [24], which is based on DCNN (deep convolutional neural network). The method requires selecting the maximum predictable sequence length in advance. This limits it to recognizing house number or license plate number whose length of texts is known beforehand. Another commonly used method is RNN (recurrent neural network) with CTC [25] (connectionist temporal classification). Shi et al. [26] and He et al. [27] propose RNN models to encode the features from the CNN and adopt CTC to decode the encoded sequence. The advantage of this method is that it can generate texts of any length. Furthermore, the nature of the Recurrent Neural Network determines that the model is able to learn the relationship between text and text temporal relations. Another type of method that does not require character segmentation of texts is attention mechanism. Lee and Osindero [28] use attention-based sequence-to-sequence structure to automatically focus on certain extracted CNN features and directly use text images to perform word string learning. This method implicitly learns character-level language models embodied in RNN. It is able to perform text recognition in unconstrained natural scenes.

Scene text recognition [29] in intelligent transportation systems has many applications, such as vehicle license plate recognition and road sign recognition. As an important part of intelligent transportation systems, vehicle license plate recognition is widely used in intelligent monitoring systems and parking systems. Automatic license plate recognition (ALPR) refers to the extraction of vehicle license plate information from an image or a sequence of images [30]. Chai

and Zuo [31] propose an automatic vehicle license plate recognition method which adopts edge detection algorithm in extraction and character segmentation and recognition. Chang et al. [32] use license plate recognition technology to track vehicle on the road in complex traffic conditions.

Object detection in applications refers to the detection under specific application scenarios, such as pedestrian detection, vehicle detection, and scene text detection. Text recognition in specific application scenarios can get more information from the objects on which the applications focus. In this paper, we propose a model which combines object-text detection and text recognition. The model is able to detect both texts and general objects simultaneously. The model combines object detection task and text detection task and recognizes the detected text contents. In addition, for the texts detected around general objects, the contents can be used as the identifier to distinguish the object. The method we propose can be applied to a wide range of applications in regard to intelligent transportation systems and has comprehensive capabilities of detection and recognition.

The contributions are summarized as follows:

- (1) We propose an object-text detection model for multiple objects which can simultaneously detect texts and general objects
- (2) We propose a text recognition framework that effectively combines text detection and recognition
- (3) The method we propose can detect multiple types of objects and instantiate the identities of the detected objects based on the identified text labels. The recognized text label is used as a valid identity of the object

## 2. Materials and Methods

The network structure in this paper consists of two parts: object-text detection network and the text recognition network. We use the YOLOv3 architecture which adopts a fully convolutional neural network [33] to detect objects and texts in real-scene images. The convolutional network is used to extract the features in multiple scales feature maps from the image. The classification and bounding box regression networks directly output the objectness score, the class of the object, and the coordinate offsets of the object at multiple feature maps. We use NMS [34] (nonmaximum suppression) to remove the redundant bounding boxes which have large overlap with the same object. We adopt a successful scene text recognition algorithm, CRNN [26] (Convolutional Recurrent Neural Network), in conjunction with object-text detection. According to the coordinates of the text type which are output from the object-text detection network, the text regions are cropped from the original image. A convolutional neural network is used to extract features from the text regions. The extracted feature maps need to be scaled to a uniform height with a fixed aspect ratio. We use the recurrent model to encode the feature sequences from the feature maps and CTC to decode the encoded sequence. The network structure we propose is shown in Figure 1.

*2.1. Architecture of the Object-Text Detection Network.* The backbone network adopts Darknet-53, which uses the former 52 layers without fully connected layer. The feature extraction network is a fully convolutional network. It is mainly composed of  $3 \times 3$  and  $1 \times 1$  convolution kernels and a large number of shortcut links with residual units [35–37]. The structure of the feature extraction network is shown in Figure 2. The network uses convolution kernel with stride instead of pooling layers to reduce the negative gradient effects brought by pooling. We also adopt a lot of data augmentation and batch normalization to avoid overfitting. In order to enhance the accuracy of the algorithm for small object detection, the network adopts upsampling and fusion methods which are similar to FPN [38] to implement the multiscale feature maps.

As shown in Figure 3, we assume the size of the input image to be  $416 \times 416$ . We extract three different scale feature maps from 26th, 43rd, and 52nd layers of the feature extraction networks in Figure 2. The scales of the extracted feature maps are  $13 \times 13$ ,  $26 \times 26$ , and  $52 \times 52$ . The feature fusion network outputs three different scale feature maps with upsampling and fusion. The top layer with a size of  $13 \times 13$  is concatenated with the  $26 \times 26$  feature map via one-time upsampling. Then it is concatenated with the  $52 \times 52$  feature map by upsampling twice. In this way, the high-level features from the top layer are passed to the low layers, which makes the model better at detecting small objects on low-level feature maps. Finally, the network generates three feature maps of different scales which are  $1/8$ ,  $1/16$ , and  $1/32$  of the original image.

The output layers in 3 different scales are also convolutional. In our experiments with our dataset which has twenty-one classes including twenty general categories and one text category, we predict 3 bounding boxes with different sizes at feature maps of each scale. The shape of the output tensor is  $N \times N[3 \times (4 + 1 + 21)]$ , where  $N$  is the scale of the feature map, 3 is the anchor boxes in features of different scales, 4 is the coordinate offsets of the bounding box, 1 is the objectness confidence prediction, and 21 is the object classes.

The network adopts the anchor-based mechanism. Each grid of the feature maps predicts 3 bounding boxes according to the anchor boxes of 3 different scales. There are in total 9 different scale anchor boxes which are generated from k-means clustering. The 9 clusters on the COCO dataset [39] are  $(10 \times 13)$ ,  $(16 \times 30)$ ,  $(33 \times 23)$ ,  $(30 \times 61)$ ,  $(62 \times 45)$ ,  $(59 \times 119)$ ,  $(116 \times 90)$ ,  $(156 \times 198)$ , and  $(373 \times 326)$ . The anchor boxes in different scale feature maps are shown in Figure 4.

The object-text detection network simultaneously predicts bounding boxes of texts and general objects conditioned on its input feature maps. At each grid of associated feature map, it outputs the objectness confidence, classification score, and coordinate offsets to its associated anchor boxes in a convolutional manner.

The object-text detection network adopts logistic regression to predict the bounding boxes and the objectness scoring on each anchor. Only the anchor with the highest objectness score is calculated. Each object can be detected by

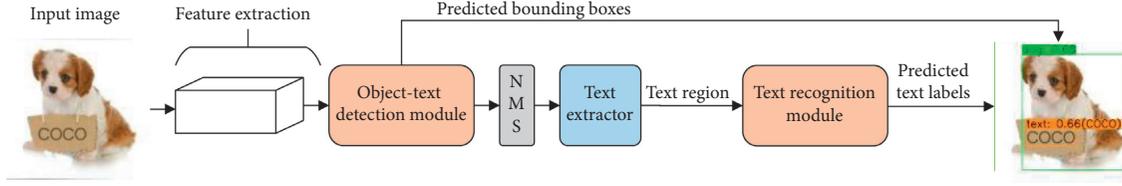


FIGURE 1: Overall architecture of the proposed network architecture. Feature maps are firstly extracted with convolutional layers. The object-text detection module is built on top of the feature maps to predict the bounding boxes of the texts and general objects. The NMS module is used to remove the redundant bounding boxes and retain the final positive bounding boxes. The text extractor extracts the text regions corresponding to the coordinates of the text bounding boxes which are the output of the object-text detection module. The text regions are then fed into the text recognition module.

only one anchor. This step is performed before prediction, which can remove unnecessary anchors and reduce the amount of calculation. In bounding box regression, the network outputs the coordinate offsets. The formula that converts offsets to bounding box coordinates is defined as

$$\begin{aligned} b_x &= \sigma(t_x) + c_x, \\ b_y &= \sigma(t_y) + c_y, \\ b_w &= p_w e^{t_w}, \\ b_h &= p_h e^{t_h}, \end{aligned} \quad (1)$$

where  $b_x, b_y, b_w,$  and  $b_h$  are the coordinates of the bounding box,  $c_x, c_y, p_w,$  and  $p_h$  are the coordinates of the anchor box, and  $\sigma(\cdot)$  represents the sigmoid function.

## 2.2. Loss Function of the Object-Text Detection Network.

Objectness confidence is the probability of predicting the existence of the object-text in anchor box. Objectness confidence loss adopts binary cross entropy. The objectness confidence loss function is defined as

$$\begin{aligned} L_{\text{conf}}(o, c) &= -\sum(o_i \ln(\hat{c}_i) + (1 - o_i) \ln(1 - \hat{c}_i)), \\ \hat{c}_i &= \text{Sigmoid}(c_i), \end{aligned} \quad (2)$$

where  $o_i \in \{0, 1\}$  represents the existence of the object-text in anchor box and  $\hat{c}_i$  represents the sigmoid probability of the existence of the bounding box.

Object-text classification score is the probability of the class which the object-text belongs to. The object-text class loss function is defined as

$$\begin{aligned} L_{\text{cls}}(O, C) &= -\sum_{i \in \text{Pos}} \sum_{j \in \text{cls}} (O_{ij} \ln(\hat{C}_{ij}) + (1 - O_{ij}) \ln(1 - \hat{C}_{ij})), \\ \hat{C}_{ij} &= \text{Sigmoid}(C_{ij}), \end{aligned} \quad (3)$$

where  $O_{ij} \in \{0, 1\}$  represents the existence of the object-texts' class  $j$  in anchor box  $i$  and  $\hat{C}_{ij}$  represents the sigmoid probability of the class  $j$  of the bounding box  $i$ .

Object-text detection model predicts the coordinate offsets between anchor boxes and bounding boxes. Equation (1) is used to convert the offsets to the coordinates of the bounding box. The object-text location loss adopts the GIoU [40] (Generalized Intersection over Union) method to

Layers	Filters	Feature maps	
Convolutional	32 3 * 3	416 * 416	
Convolutional	64 3 * 3 stride = 2	208 * 208	
Convolutional	32 1 * 1		×1
Convolutional	64 3 * 3		
Residual		208 * 208	
Convolutional	128 3 * 3 stride = 2	104 * 104	
Convolutional	64 1 * 1		×2
Convolutional	128 3 * 3		
Residual		104 * 104	
Convolutional	256 3 * 3 stride = 2	52 * 52	
Convolutional	128 1 * 1		×8
Convolutional	256 3 * 3		
Residual		52 * 52	
Convolutional	512 3 * 3 stride = 2	26 * 26	
Convolutional	256 1 * 1		×8
Convolutional	512 3 * 3		
Residual		26 * 26	
Convolutional	1024 3 * 3 stride = 2	13 * 13	
Convolutional	512 1 * 1		×4
Convolutional	1024 3 * 3		
Residual		13 * 13	

FIGURE 2: The feature extraction network. The backbone network of detection model uses the former 52 layers of the Darknet-53 without fully connected layer to extract features.

calculate the error between the bounding box and ground truth. The GIoU Loss Algorithm is defined as in Algorithm 1.

We use  $L_{\text{loc}} = \text{GIoU\_Loss}$  to form of the object-text location loss. The total loss function can be represented as

$$L_{\text{tol}} = \alpha L_{\text{conf}} + \beta L_{\text{cls}} + \gamma L_{\text{loc}}, \quad (4)$$

where  $\alpha, \beta,$  and  $\gamma$  are the weights of each loss. We empirically set  $\alpha = \beta = \gamma = 1$ .

**2.3. NMS Module.** The NMS module is applied to remove the redundant object-text bounding boxes detected from the same object. We adopt the NMS after the object-text detection on the object-text bounding boxes.

**2.4. Text Recognition.** After the object positions are detected from the object-text detection network, we pick out the text-type bounding boxes based on the text class. Firstly, the text

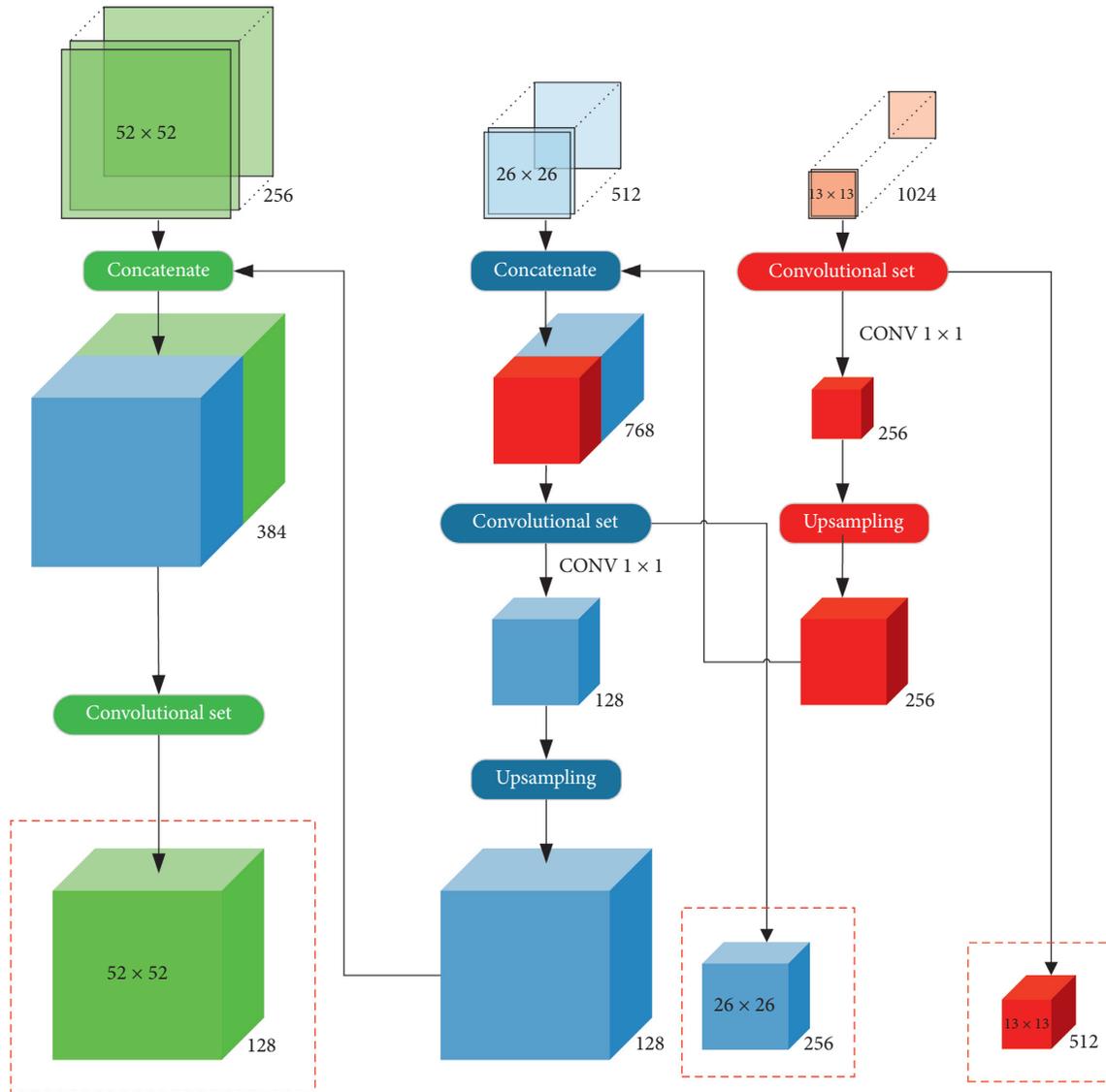


FIGURE 3: The feature fusion network. The generation of the 3 different scale feature maps uses upsampling and concatenation from convolutional layers. Firstly, the 3 different scale feature maps are extracted from the Darknet-53, respectively, in different convolutional layers. Secondly, the convolutional set which consists of a set of convolutional layers is used to reduce the channel of the  $13 \times 13$  feature map. Then the  $13 \times 13$  feature map is upsampled with stride 2 to concatenate with the  $26 \times 26$  feature map. The  $52 \times 52$  feature map concatenation process is similar to that of the  $13 \times 13$  feature map. Finally, the red dashed boxes are the extracted feature maps in 3 different scales.

extractor extracts the text regions corresponding to the coordinates of the text bounding boxes produced by the object-text detection module. Then the text recognition module preprocesses the extracted text regions by resizing them before they are fed into convolutional neural network. We scale the text regions to  $(32, 100, 3)$  with a fixed aspect ratio, where 32 is the fixed height, 100 is the maximum length, and 3 represents the number of the image channel. Finally, we use the scaled text region as the input of the convolutional layers.

We adopt the CRNN model as our text recognizer. Firstly, the convolutional layers extract the feature maps from the preprocessed text region. A sequence of feature vectors is extracted from left to right from the feature maps. Then each frame of the sequences which represents a vertical

region corresponding to the original text image becomes the input of the recurrent layers. The recurrent layers adopt the deep bidirectional LSTM [41] (long short-term memory) to encode the sequence of the feature vectors. Finally, we adopt CTC to predict the text label corresponding to the sequences from the recurrent layers.

### 3. Results and Discussion

*3.1. Experiment Setup.* The object-text detection network is trained with training images using Adam (adaptive moment estimation) [42]. We initialize the model with pre-trained weights on the COCO dataset. We divide the training process into two stages. In the first stage, we fix the backbone network and just train the classification and

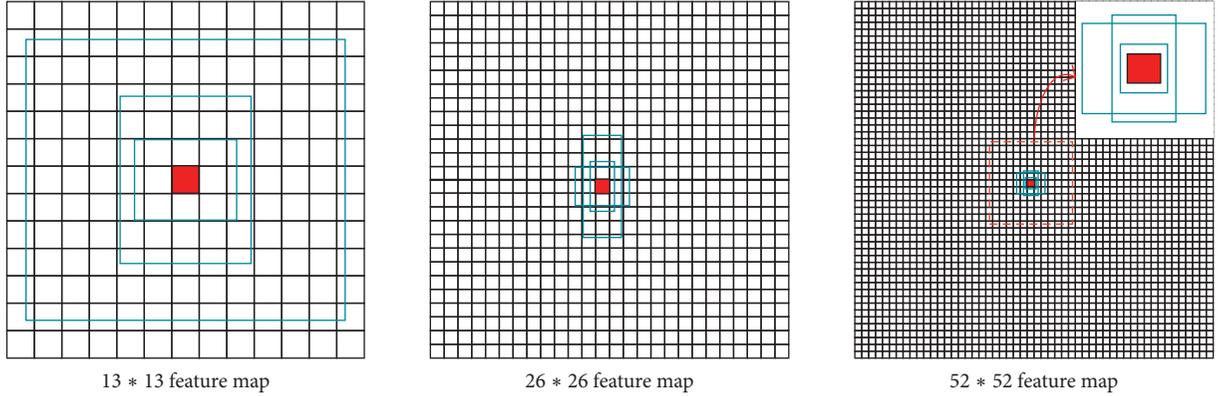


FIGURE 4: The anchor boxes in 3 scale feature maps.

regression network. In the second stage, we train the whole network.

When we train the object-text detection network, in the first two epochs in training, we adopt the method of gradually increasing the learning rate from low to high which is called “warmup stage” method. The network converges quickly with the large learning rate. And then, we need to stabilize the network with a low learning rate to avoid gradient oscillation. We adopt a cosine annealing strategy proposed by Loshchilov et al. [43]. At the  $i$ -th training step, the learning rate decays with a cosine annealing as follows:

$$\eta = \eta_{\text{end}} + 0.5(\eta_{\text{init}} - \eta_{\text{end}}) \left( 1 + \cos \left( \frac{T_{\text{cur}} - T_{\text{warm}}}{T_{\text{train}} - T_{\text{warm}}} \pi \right) \right), \quad (5)$$

where  $\eta_{\text{init}}$  is the initial value of the learning rate, which is set to  $10^{-4}$ ,  $\eta_{\text{end}}$  is the end value which we set to  $10^{-6}$ , and  $T_{\text{cur}}$  accounts for how many steps have been performed.  $T_{\text{train}}$  is the total steps during the training.  $T_{\text{warm}}$  represents the warmup steps in the first two epochs. The learning rate curve is shown in Figure 5.

The training algorithm of the object-text detection model is summarized as in Algorithm 2.

We use a CRNN model proposed by Shi et al. [26] as the text recognition network. The experiment uses a pretrained model trained on the synth90k dataset [44] to initialize the parameters of the text recognition model. We use NEOCR [45] dataset and SCUT FORU dataset to fine-tune the pretrained model. We set the training parameters as follows: The model training runs for 2000000 epochs. The batch size is 32. The initial learning rate is 0.01 with exponential decay of 0.1 every 500000 epochs. The experiment adopts gradient descent with momentum [46] to train the text recognition network. We set the parameter of momentum to 0.9.

The training algorithm of the text recognition model is summarized as in Algorithm 3.

**3.2. Dataset.** We evaluate the proposed method on four datasets: VOC 2007 [47], VOC2012 [48], ICDAR 2013 [49], and SCUT FORU DB. VOC2007 and VOC2012 are the datasets about object detection. ICDAR 2013 and SCUT

FORU DB are the datasets about text detection. We integrate them into a comprehensive dataset for detecting type-text object and general object simultaneously.

VOC2007 is the challenge to recognize objects from a number of visual object classes in realistic scenes. The database contains a total of 9963 annotated images. We use 5011 images as training set and 4952 images as testing set. There are twenty object classes in the dataset.

VOC2012 is the same challenge as VOC2007 which increases the size of the training set. There are 17125 training images in total. The testing set has not been released yet.

ICDAR 2013 is the Challenge 2 of ICDAR 2013 Robust Reading Competition, which contains horizontal texts. The dataset focuses on the reading of texts in real scenes. The images of the dataset refer to the text images focused around the text content of interest. The dataset consists of 229 training images and 233 testing images. Due to the fact that there are too few training images, we additionally use 1200 images from SCUT FORU training dataset.

SCUT FORU Database is released by the South China University of Technology. The dataset consists of Chinese2k and English2k. We only use the English2k dataset. The English2k dataset contains character annotations and word annotations. The characters of the dataset include 52 upper-lowercase letters and 10 Arabic numerals. The label format of the dataset is  $\{x, y, w, h, \text{label}\}$ .  $\{x, y\}$  are the top-left coordinates of the rectangular box.  $\{w, h\}$  are the width and height of the rectangular box.  $\{\text{label}\}$  is the word label of the text region. There are a total of 1715 images, of which 1200 are the training images and 515 are the testing images. The dataset has an average of 18.4 characters and 3.2 words per image.

COCO Dataset is a large-scale dataset for object detection, segmentation, and captioning. It contains more than 330K images and 200K labels. The COCO dataset has 80 object categories in total.

In the experiment, we integrate the datasets into a comprehensive dataset of 29265 images in total. There are 23565 training images and 5700 testing images. Since these datasets have different annotation formats, we need to convert them into a unified annotation format. The coordinates format of the annotation is defined as

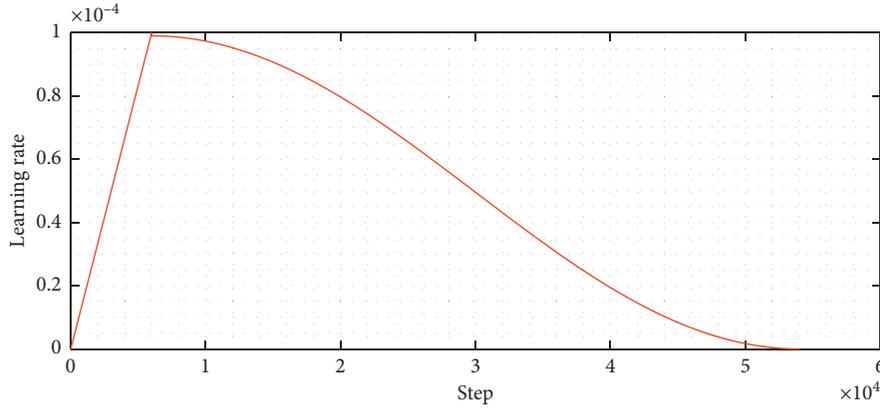


FIGURE 5: Learning rate schedule.

$\{x_{\min}, y_{\min}, x_{\max}, y_{\max}\}$ . We shuffle the combined dataset to feed into the model.

Text recognition network is performed with a CRNN model proposed by Shi et al. We use a pretrained model trained on the synth90k dataset and use NEOCR dataset and SCUT FORU dataset to fine-tune the pretrained model. The annotations in NEOCR dataset contain characters that are not in the English alphabet. We have modified the annotations by replacing the special characters to English letters that look similar. The text images in SCUT FORU dataset are cropped from the original images corresponding to the coordinates in annotations. The text images are resized to the size of  $(32 \times 100)$  before they are fed into the text recognition network.

**3.3. Evaluation Metrics.** We use mAP (mean Average Precision) as the measurement to evaluate the detection model performance. The mAP calculation is based on the following metrics [50]:

True Positives (TP): examples detected correctly with  $\text{IoU} \geq \text{threshold}$

False Positives (FP): negative examples detected incorrectly with  $\text{IoU} < \text{threshold}$

False Negatives (FN): the ground truths not detected

The threshold is usually set to 0.5, 0.75, or 0.95. In our evaluation, we set it to 0.5.

**Precision.** Precision is the percentage of correct positive predictions. The precision is defined as

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}. \quad (6)$$

**Recall.** Recall is the percentage of true positive detected among all relevant ground truths. The recall is defined as

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}. \quad (7)$$

**PR (Precision-Recall) Curve.** The PR curve is a good way to evaluate the performance of an object detector. The

precision and recall values of detected objects are plotted to get a PR curve. The area under the PR curve is called AP (Average Precision). The AP calculation is defined as

$$\text{AP} = \int P(R) dR, \quad (8)$$

where  $P(R)$  is the measured precision against recall. **mAP.** The mAP is the average of all categories of AP.

**3.4. Analysis of Experimental Results.** In order to verify the choice of YOLOv3 as the detection network in the proposed method, we compare the detection performance of different detection frameworks, namely, Fast R-CNN, Faster R-CNN, SSD, YOLO, YOLOv2, and YOLOv3. We as well compare different input size setups of YOLOv3. The results are shown in Table 1. All the detection frameworks being compared are trained on the VOC2007 and VOC2012 training datasets and the mAP is tested on the VOC2007 testing dataset. As can be seen in Table 1, the YOLOv3 framework with network input size of  $416 \times 416$  achieves the highest mAP among the frameworks being tested. Further, the YOLO series, be it YOLOv2 or YOLOv3, generally achieve higher mAP than other frameworks. It can therefore be concluded that the choice of the YOLOv3 framework in the proposed method is an optimized solution.

After we have confirmed the performance of the YOLOv3 in object detection, we further train it on the comprehensive dataset which is composed of the general object detection datasets of VOC2007 and VOC2012 and the text detection datasets of SCUT FORU and ICDAR2013. Then we test the performance of the frameworks on different testing datasets. The general objects detection testing dataset VOC2007 and the text detection datasets SCUT FORU and ICDAR2013 are used. We compare the performance of different detection frameworks with 3 categories out of the total 20 categories in the PASCAL VOC 2007 dataset. As shown in Table 2, the performance of YOLOv3 on the 3 categories is much better than other detection frameworks. We verify that the YOLOv3 has excellent performance on object detection. Our model achieves 70.0 mAP in the text

**Input:** The region of the  $GT$  and  $BB$  ( $GT, BB \subseteq S \in \mathbb{R}^n$ ), where  $S$  is the input image size.  
**Output:**  $GIoU$  Loss  
**Step 1.** Calculate the smallest enclosing region  $C, C \subseteq S \in \mathbb{R}^n$ ;  
**Step 2.**  $IoU = |GT \cap BB| / |GT \cup BB|$ ;  
**Step 3.**  $GIoU = IoU - (|C - (GT \cup BB)| / |C|)$ ;  
**Step 4.** Calculate the bounding box scale:  $gt\_scale = 2 - (GT/S)$ ;  
**Step 5.** Calculate the location loss function:  $GIoU\_Loss = O \cdot gt\_scale \cdot (1 - GIoU)$ , where  $O$  is the existence of the object-text in associated bounding box.

ALGORITHM 1: GIoU\_Loss algorithm.

**Input:** Parameter\_1: The training set  $X = \{x_1, x_2, x_3, \dots, x_i, \dots, x_N\}$ .  $N$  is the number of the batches.  $x_i$  is the  $i$ -th batch of the training set;  
Parameter\_2: The labels corresponding to the training images of which format is defined as:  $\{x_{\min}, y_{\min}, x_{\max}, y_{\max}, class\}$   
**Output:** Weights of the model  
**for**  $epoch$  **in**  $epochs$  **do**  
  **if**  $epoch \leq first\_stage\_epochs$  **then**  
     $initialize\_train\_weights \leftarrow classification\_regression\_weights$ ;  
  **else**  
     $initialize\_train\_weights \leftarrow backbone\_weights + classification\_regression\_weights$ ;  
  **for**  $batch$  **in**  $batches$  **do**  
    Predict the offsets, objectness and class:  
     $x_i, y_i \leftarrow preprocess(Input)$ ;  
     $\hat{y}_i \leftarrow forward(W, x_i)$ ;  
    Calculate the loss:  
     $l \leftarrow loss(\hat{y}_i, y_i)$ ;  
    Calculate the gradients:  
     $grad \leftarrow backward(l)$ ;  
    Update the model parameters:  
     $W \leftarrow update(grad)$ ;  
  **end for**  
**end for**

ALGORITHM 2: Detection network training algorithm.

detection task. We are not listing the text detection performance of the other methods because they do not feature text detection and recognition.

One may notice that the mAP of YOLOv3 in Table 2 is lower than that of Table 1. This is because we further train the YOLOv3 network on the text detection datasets. The detection of text objects reduces the mAP to a certain extent. In addition, the text object in datasets of VOC2007 and VOC2012 is not marked in the annotations. The detected texts in VOC2007 and VOC2012 will be seen as 'False Positives', thus the mAP would decrease.

Due to the imperfection of the comprehensive dataset which consists of general object datasets and text datasets, we improve the annotation information of the comprehensive dataset. We label the text objects in VOC2007 and VOC2012 and the general objects in ICDAR2013 and SCUT. This makes the comprehensive dataset of the object detection more accurate, reduces the false positive rate of the detection model in training and testing, and improves the detection accuracy as a whole. As can be seen from Table 2, the detection model used in the experiment has the highest detection accuracy on the YOLOv3 framework with the size

of  $544 \times 544$ . Table 3 compares the detection effect of the comprehensive dataset before and after the modification on the YOLOv3 framework with the size of  $544 \times 544$ . As shown in Table 3, the detection network on modified comprehensive dataset has higher accuracy on person and text objects than original dataset. The detection accuracy of the text object is significantly improved. The mAP on the modified comprehensive dataset has also improved.

### 3.5. Performance on Object-Text Detection and Recognition.

The model we propose performs two tasks: object-text detection and text recognition. The object-text detection network can detect general objects and text objects simultaneously. The text contents of the detected text regions from the detection network are recognized by the text recognition network. This section shows the detection and recognition results of test images in the experiment.

As shown in Figure 6, we mainly show some detection results of test images in transportation. The detection model can detect multiple objects in one image. It has good performance on both small objects and large objects. The text

```

Input: Parameter_1: The training set of text  $X = \{x_1, x_2, x_3, \dots, x_i, \dots, x_N\}$ .  $N$  is the number of the batches.  $x_i$  is the  $i$ -th batch of the training set;
Parameter_2: The text labels corresponding to the text training images:  $\{label\}$ 
Output: Weights of the model
for  $epoch$  in  $epochs$  do
   $initialize\_weights \leftarrow pretrained\_weights$ ;
  for  $batch$  in  $batches$  do
    Predict the recognition result:
     $x_i, y_i \leftarrow preprocess(Input)$ ;
     $\hat{y}_i \leftarrow forward(W, x_i)$ ;
    Calculate the CTC loss:
     $l \leftarrow ctc\_loss(\hat{y}_i, y_i)$ ;
    Calculate the gradients:
     $grad \leftarrow backward(l)$ ;
    Update the model parameters:
     $W \leftarrow update(grad)$ ;
  end for
end for

```

ALGORITHM 3: Recognition network training algorithm.



FIGURE 6: Examples of the object detection result on PACSAL VOC2007. We use rectangular boxes to represent the location of the detected objects. The different colours of the rectangular boxes represent the different categories.

detection dataset contains billboards, signboards, road sign, etc. Some texts exist in complex environments and they might be occluded. As shown in Figure 7, the detection model can detect the text in complex scenes. However, some text bounding boxes in images are not accurate enough, which may cause wrong recognition in texts. The object-text detection model we propose can simultaneously detect the text and general objects. Some detection examples are demonstrated in Figure 8.

The text recognition model can recognize the text contents of the text regions detected from the detection model. As shown in Figure 9, we demonstrate some examples of text recognition model on road sign. As shown in Figure 10, the text recognition model can recognize not

only the horizontal text, but also the affine distorted text. The affine distorted texts exist commonly due to the variations of the camera views. Yet the proposed model not only locates these texts, but also finds the contents of the texts.

Figure 11 gives a more application specific demonstration of the proposed object detection and text recognition model. In this scenario, information extracted by the text recognition module identifies the detected object. We use some cars images with plates as the proof of concept. The object-text detection model we have proposed can simultaneously detect the car and the plates on the car. Then the text recognition model recognizes the text contents on the plates.



FIGURE 7: Examples of the text detection result on SCUT FORU dataset. We use yellow solid boxes with overstriking to represent the location of the detected text.



FIGURE 8: Examples of object-text detection on the comprehensive dataset. The different colours of the rectangular boxes represent the different categories. The red dashed rectangular boxes represent the misdetections.



FIGURE 9: Examples of word spotting result. The yellow solid boxes with overstriking are the text locations. The words in parenthesis are the word recognition results. The red dashed boxes are the wrong word recognition.



FIGURE 10: Recognition examples of text with affine transforms. The yellow solid boxes with overstriking are the text locations. The words in yellow are the word recognition results. The red solid boxes are the misdetections.



FIGURE 11: Examples of object-text detection and recognition results. There are some images of car with plate. We use yellow box to represent the detection results of car. The yellow solid boxes with overstriking are the text locations. The words in yellow are the word recognition results.

TABLE 1: Detection performance on the PASCAL VOC 2007 testing dataset of different frameworks.

Detection frameworks	Trained on	mAP
Fast R-CNN [19]	VOC2007 + VOC2012	70.0
Faster R-CNN VGG-16 [20]	VOC2007 + VOC2012	73.2
Faster R-CNN ResNet [35]	VOC2007 + VOC2012	76.4
YOLO [17]	VOC2007 + VOC2012	63.4
SSD300 [16]	VOC2007 + VOC2012	74.3
SSD512 [16]	VOC2007 + VOC2012	76.8
YOLOv2 544 × 544 [21]	VOC2007 + VOC2012	78.6
YOLOv3 416 × 416	VOC2007 + VOC2012	<b>87.4</b>
YOLOv3 544 × 544	VOC2007 + VOC2012	86.8
YOLOv3 608 × 608	VOC2007 + VOC2012	86.1

TABLE 2: Performance of different detection frameworks on different testing datasets.

Detection frameworks	VOC2007				SCUT_FORU	ICDAR2013
	mAP	Car	Bus	Person	Text	Text
Fast R-CNN [19]	70.0	78.6	81.6	69.9	—	—
Faster R-CNN VGG-16 [20]	73.2	84.7	83.1	76.7	—	—
SSD300 [16]	74.3	84.2	83.0	76.2	—	—
SSD512 [16]	76.8	87.5	86.2	79.7	—	—
YOLOv3 416 × 416	77.9	90.4	85.4	85.7	70.0	64.2
YOLOv3 544 × 544	<b>78.3</b>	<b>91.5</b>	<b>89.4</b>	<b>86.8</b>	<b>70.0</b>	<b>70.0</b>
YOLOv3 608 × 608	77.9	90.4	85.4	85.7	70.0	64.2

TABLE 3: Detection performance between the comprehensive dataset and modified comprehensive dataset.

Detection frameworks	Dataset	mAP	Car	Bus	Person	Text
YOLOv3 544 × 544	Comprehensive dataset	76.9	<b>89.3</b>	<b>87.9</b>	83.5	66.0
	Comprehensive dataset (modified)	<b>77.2</b>	88.5	85.3	<b>83.9</b>	<b>72.8</b>

## 4. Conclusions

We present an object-text detection and recognition model in this article. The model not only detects the texts and general objects simultaneously but also recognizes the text contents inside the detected text bounding boxes. The method we have proposed combines both object detection and text recognition. In the applications of some scenarios, the recognized text contexts around the general objects are able to be used as the identifier to distinguish the object. The proposed method has potential in extensive applications, such as intelligent transportation systems and autonomous driving.

Possible directions for future research include the following:

- (1) Improving the dataset: this refers to adding more samples which contain both text and general object to train the network
- (2) Improving the detection network on the text detection: for example, the anchor box which is suitable for the text size can be used. We can use k-means to cluster on the dataset containing text objects to make the size of the generated anchor boxes more suitable for text
- (3) Optimizing the connection between the detection network and the recognition network: in our proposed model, the connection between detection and recognition network is the text region which is cropped from the original image corresponding to the coordinates of the detected text boxes. In order to optimize the connection, we can extract the feature map from the detection network as the input of the recognition network. The affine transformation is applied to the feature map extracted from detection network to fit the input size of recognition network. Thus, during backpropagation, the gradients can flow from the recognition network back to the detection network. The detection and recognition model can be regarded as an end-to-end model.

## Data Availability

The data used to support the findings of this study are available from the corresponding author upon request.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## Acknowledgments

This research was funded by the Foundation of the National Key Research and Development Program (grant number

2016YFC0801800), National Natural Science Foundation of China (grant number 51874300), National Natural Science Foundation of China and Shanxi Provincial People's Government Jointly Funded Project of China for Coal Base and Low Carbon (grant number U1510115), and the Open Research Fund of Key Laboratory of Wireless Sensor Network & Communication, Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences (grant numbers 20190902 and 20190913).

## References

- [1] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: a survey," pp. 1–40, 2019, <https://arxiv.org/pdf/1905.05055>.
- [2] Z.-Q. Zhao, P. Zheng, S.-T. Xu, and X. Wu, "Object detection with deep learning: a review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [3] C.-H. Chen, "A cell probe-based method for vehicle speed estimation," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. E103.A, no. 1, pp. 265–267, 2020.
- [4] C.-H. Chen, F.-J. Hwang, and H.-Y. Kung, "Travel time prediction system based on data clustering for waste collection vehicles," *IEICE Transactions on Information and Systems*, vol. E102.D, no. 7, pp. 1374–1383, 2019.
- [5] Y. Tang, C. Zhang, R. Gu, P. Li, and B. Yang, "Vehicle detection and recognition for intelligent traffic surveillance system," *Multimedia Tools and Applications*, vol. 76, no. 4, pp. 5817–5832, 2017.
- [6] K. V. Sakhare, T. Tewari, and V. Vyas, "Review of vehicle detection systems in advanced driver assistant systems," *Archives of Computational Methods in Engineering*, vol. 27, no. 2, pp. 591–610, 2020.
- [7] D. A. Yudin, A. Skrynnik, A. Krishtopik, I. Belkin, and A. I. Panov, "Object detection with deep neural networks for reinforcement learning in the task of autonomous vehicles path planning at the intersection," *Optical Memory and Neural Networks*, vol. 28, no. 4, pp. 283–295, 2019.
- [8] C.-H. Chen, F. Song, F.-J. Hwang, and L. Wu, "A probability density function generator based on neural networks," *Physica A: Statistical Mechanics and Its Applications*, vol. 541, p. 123344, 2020.
- [9] J. Zhang, M. Huang, X. Jin, and X. Li, "A real-time Chinese traffic sign detection algorithm based on modified YOLOv2," *Algorithms*, vol. 10, no. 4, p. 127, 2017.
- [10] H. C. Song, G. H. Lee, D.-S. Shim, and K. N. Choi, "Visual distinctiveness detection of pedestrian based on statistically weighting PLSA for intelligent systems," *International Journal of Control, Automation and Systems*, vol. 16, no. 2, pp. 815–822, 2018.
- [11] C. Zhang and J. Kim, "Multi-scale pedestrian detection using skip pooling and recurrent convolution," *Multimedia Tools and Applications*, vol. 78, no. 2, pp. 1719–1736, 2019.
- [12] J. Kim, H. Cho, M. Hwangbo, J. Choi, J. Canny, and Y. P. Kwon, "Deep traffic light detection for self-driving cars from a large-scale dataset," in *Proceedings of the 2018*

- 21st International Conference on Intelligent Transportation Systems (ITSC), pp. 280–285, Maui, HI, USA, November 2018.
- [13] K. Behrendt, L. Novak, and R. Botros, “A deep learning approach to traffic lights: detection, tracking, and classification,” in *Proceedings of the 2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1370–1377, Singapore, May–June 2017.
- [14] Y. Lu, J. Lu, S. Zhang, and P. Hall, “Traffic signal detection and classification in street views using an attention model,” *Computational Visual Media*, vol. 4, no. 3, pp. 253–266, 2018.
- [15] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 580–587, Columbus, OH, USA, June 2014.
- [16] W. Liu et al., “SSD: single shot MULTIBOX detector,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 9905, pp. 21–37, Springer, Cham, Switzerland, 2016.
- [17] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: unified, real-time object detection,” in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2016, pp. 779–788, Las Vegas, NV, USA, June–July 2016.
- [18] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders, “Selective search for object recognition,” *International Journal of Computer Vision*, vol. 104, no. 2, pp. 154–171, 2013.
- [19] R. Girshick, “Fast R-CNN,” in *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*, vol. 2015, pp. 1440–1448, Las Condes, Chile, December 2015.
- [20] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: towards real-time object detection with region proposal networks,” *Advances in Neural Information Processing Systems*, vol. 2015, pp. 91–99, 2015.
- [21] J. Redmon and A. Farhadi, “YOLO9000: better, faster, stronger,” in *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2017, pp. 6517–6525, Honolulu, HI, USA, July 2017.
- [22] J. Redmon and A. Farhadi, “YOLOv3: an incremental improvement,” 2018, <https://arxiv.org/abs/1804.02767>.
- [23] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [24] I. J. Goodfellow, Y. Bulatov, J. Ibarz, S. Arnaud, and V. Shet, “Multi-digit number recognition from street view imagery using deep convolutional neural networks,” in *Proceedings of the 2nd International Conference on Learning Representations*, Banff, Canada, April 2014.
- [25] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification,” in *Proceedings of the 23rd international conference on Machine Learning–ICML ’06*, pp. 369–376, Pittsburgh, PA, USA, June 2006.
- [26] B. Shi, X. Bai, and C. Yao, “An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition,” *IEEE Trans. Pattern Anal. Mach. Intell.* vol. 39, no. 11, pp. 2298–2304, 2017.
- [27] P. He, W. Huang, Y. Qiao, C. C. Loy, and X. Tang, “Reading scene text in deep convolutional sequences,” in *Proceedings of the 30th AAAI Conference on Artificial Intelligence. AAAI 2016*, pp. 3501–3508, Phoenix, Arizona, USA, February 2016.
- [28] C.-Y. Lee and S. Osindero, “Recursive recurrent nets with attention modeling for OCR in the wild,” in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2016, pp. 2231–2239, Las Vegas, NV, USA, June 2016.
- [29] H. Lin, P. Yang, and F. Zhang, “Review of scene text detection and recognition,” *Archives of Computational Methods in Engineering*, vol. 27, no. 2, pp. 433–454, 2020.
- [30] S. Du, M. Ibrahim, M. Shehata, and W. Badawy, “Automatic license plate recognition (ALPR): a state-of-the-art review,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 2, pp. 311–325, 2013.
- [31] D. Chai and Y. Zuo, “Extraction, segmentation and recognition of vehicle’s license plate numbers,” *Advances in Intelligent Systems and Computing*, vol. 887, pp. 724–732, 2019.
- [32] J.-K. Chang, S. Ryou, and H. Lim, “Real-time vehicle tracking mechanism with license plate recognition from road images,” *The Journal of Supercomputing*, vol. 65, no. 1, pp. 353–364, 2013.
- [33] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3431–3440, Boston, MA, USA, June 2015.
- [34] A. Neubeck and L. Van Gool, “Efficient non-maximum suppression,” in *Proceedings of the 18th International Conference on Pattern Recognition (ICPR’06)*, vol. 3, pp. 850–855, Hong Kong, China, August 2006.
- [35] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2016, pp. 770–778, Las Vegas, NV, USA, June 2016.
- [36] F. Zhang and Z. Xu, “A mine image reconstruction method based on residual neural network,” *Meitan Xuebao/Journal China Coal Soc.* vol. 44, no. 11, pp. 3614–3624, 2019.
- [37] F. Zhang, Z. Xu, W. Chen et al., “An image compression method for video surveillance system in underground mines based on residual networks and discrete wavelet transform,” *Electronics*, vol. 812 pages, 2019.
- [38] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature Pyramid networks for object detection,” in *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 2017, pp. 936–944, Honolulu, HI, USA, July 2017.
- [39] T. Y. Lin, M. Maire, S. Belongie et al., “Microsoft COCO: common objects in context,” in *Lecture Notes in Computer Science*, vol. 8693, pp. 740–755, Springer, Cham, Switzerland, 2014.
- [40] H. Rezaatofghi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, “Generalized intersection over union: a metric and a loss for bounding box regression,” 2019, <https://arxiv.org/abs/1902.09630>.
- [41] S. Hochreiter and J. Unger Schmidhuber, “Long shortterm memory,” *Neural Computation*, vol. 9, no. 8, Article ID 17351780, 1997.
- [42] D. P. Kingma, J. L. Ba, and “Adam,” “A method for stochastic optimization,” in *Proceedings of the 3rd International Conference on Learning Representations, ICLR 2015*, pp. 1–15, San Diego, CA, USA, May 2015.
- [43] I. Loshchilov and F. Hutter, “SGDR: stochastic gradient descent with warm restarts,” in *Proceedings of the 5th International Conference on Learning Representations ICLR 2017*, pp. 1–16, Toulon, France, April 2019.

- [44] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman, "Synthetic data and artificial neural networks for natural scene text recognition," pp. 1–10, 2014.
- [45] R. Nagy, A. Dicker, and K. Meyer-Wegener, "NEOCR: a configurable dataset for natural image text recognition," *Camera-Based Document Analysis and Recognition*, vol. 7139, pp. 150–163, 2012.
- [46] S. Ruder, "An overview of gradient descent optimization," pp. 1–14, 2016.
- [47] M. Everingham, L. van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (VOC) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [48] M. Everingham, S. M. Ali Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL visual object classes challenge: a retrospective," *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.
- [49] D. Karatzas, F. Shafait, S. Uchida et al., "ICDAR 2013 robust reading competition," in *Proceedings of the 12th International Conference on Document Analysis and Recognition ICDAR*, pp. 1484–1493, Barcelona, Spain, January 2013.
- [50] J. Davis and M. Goadrich, "The relationship between Precision-Recall and ROC curves," in *Proceedings of the 23rd International Conference on Machine Learning ICML 2006*, vol. 2006, pp. 233–240, Pittsburgh, PA, USA, June 2006.