

A multitask model for realtime fish detection and segmentation based on YOLOv5

QinLi Liu ^{Equal first author, 1}, **Xinyao Gong** ^{Equal first author, 1}, **Jiao Li** ¹, **Hongjie Wang** ¹, **Ran Liu** ¹, **Dan Liu** ¹, **Ruoran Zhou** ¹, **Tianyu Xie** ¹, **Ruijie Fu** ¹, **Xuliang Duan** ^{Corresp. 1}

¹ College of Information Engineering, Sichuan Agricultural University, Ya'an, Sichuan, China

Corresponding Author: Xuliang Duan
Email address: duanxuliang@sicau.edu.cn

The accuracy of fish farming and real-time monitoring are essential to the development of “intelligent” fish farming. Although the existing instance segmentation networks (such as Maskrcnn) can detect and segment the fish, most of them are not effective in real-time monitoring. In order to improve the accuracy of fish image segmentation and promote the accurate and intelligent development of fish farming industry, this paper uses YOLOv5 as the backbone network and object detection branch, combined with semantic segmentation head for real-time fish detection and segmentation. The experiments show that the object detection precision can reach 95.4% and the semantic segmentation accuracy can reach 98.5% with the algorithm structure proposed in this paper, based on the golden crucian carp dataset, and 116.6 FPS can be achieved on RTX3060. On the publicly available dataset PASCAL VOC 2007, the object detection precision is 73.8%, the semantic segmentation accuracy is 84.3%, and the speed is up to 120 FPS on RTX3060.

A multitask model for realtime fish detection and segmentation based on YOLOv5

Qinli Liu^{1,*}, Xinyao Gong^{1,*}, Jiao Li¹, Hongjie Wang¹, Ran Liu¹, Dan Liu¹, Ruoran Zhou¹, Tianyu Xie¹, Ruijie Fu¹, Xuliang Duan¹

¹ College of Information Engineering, Sichuan Agricultural University, Ya'an, Sichuan, China;

* These authors contributed equally to this work.

Corresponding Author:

Xuliang Duan¹

Xinkang Road, Ya'an, Sichuan, 625014, China

Email address: duanxuliang@sicau.edu.cn

Abstract

The accuracy of fish farming and real-time monitoring are essential to the development of “intelligent” fish farming. Although the existing instance segmentation networks (such as Maskrcnn) can detect and segment the fish, most of them are not effective in real-time monitoring. In order to improve the accuracy of fish image segmentation and promote the accurate and intelligent development of fish farming industry, this paper uses YOLOv5 as the backbone network and object detection branch, combined with semantic segmentation head for real-time fish detection and segmentation. The experiments show that the object detection precision can reach 95.4% and the semantic segmentation accuracy can reach 98.5% with the algorithm structure proposed in this paper, based on the golden crucian carp dataset, and 116.6 FPS can be achieved on RTX3060. On the publicly available dataset PASCAL VOC 2007, the object detection precision is 73.8%, the semantic segmentation accuracy is 84.3%, and the speed is up to 120 FPS on RTX3060.

Keywords: realtime monitoring, golden crucian carp dataset, multi-task, object detection, semantic segmentation, YOLOv5

Introduction

Humans have known how to catch fish since ancient times, and now the demand for fish is increasing year by year. The nutritional value of fish is extremely high and fish is easy to absorb, which is popular with people. In the process of the aquaculture, managers need to gain information about the lives of fish, shrimp, shellfish and other aquatic organisms, specifically, species, behavior identification and biomass estimation (total weight of fish and shrimp in specific waters). Among them, the biomass estimation is the total weight of fish and shrimps in a specific water. The information on the biomass of fish, shrimp and other aquatic organisms at

various growing stages is crucial, because managers need to optimize feeding requirements and make effective decisions based on this information.

The conventional estimation of fish biomass mainly adopts manual fishing and weighing, which not only consumes human resources and decreases efficiency, but also may have adverse effects on the growth of fish. Since there is a relationship between the weight of organisms, and their body length and image area, the weight can be estimated indirectly through a deep learning-based image segmentation method to predict daily feed intake of aquatic organisms. It can also monitor the growth rate of aquatic organisms, control breeding density, determine optimal harvesting time, and ensure optimal utilization of facility investment. Due to the reliance on visual pattern matching in recognition process, it is convenient and least invasive compared to traditional manual monitoring. Non-invasive techniques offer significant advantages in terms of cost, safety and convenience. Therefore, using computer vision methods for intelligent detection of fisheries is an inevitable trend in the development of fish farming industry in modern society [1].

There are mainly four tasks in computer vision regarding image recognition: Classification, Location, Detection, and Segmentation. Classification focuses on “What broad category of object is in this photograph” [2], and is concerned with the overall content of the image. It is mainly divided into two types of tasks: dichotomous and multiclassification tasks, and is widely used in agriculture, medicine, soil, etc. [3~5]. The main task of localization is to find where the object is, generally in the form of a bounding box to circle the location of the object [6]. Unlike classification, detection is to find out which objects are in the photograph and obtain their category information and location information, which is a combination of classification and location [2]. Segmentation mainly focuses on pixels, solving the problem of “What pixels belong to the object in the image” [2]. Segmentation includes semantic segmentation and instance segmentation. The semantic segmentation is to determine whether pixels go well with the object without distinguishing different instances of the same category, while the instance segmentation needs to determine different instances of the same category on the basis of semantic segmentation, which is a combination of object detection and semantic segmentation.

Object detection is a collection of classification and regression, mainly aiming to find out all the targets of interest in an image and determine their categories and locations. Early object detection algorithms were based on handcrafted features. The representative ones are Viola Jones Detectors [7], HOG Detectors [8] and Deformable Part based Model (DPM) [9]. These algorithms have laid an important foundation for later object detection algorithms. The development of object detection algorithms was limited due to the saturation of handcrafted features. However, with the advent of Convolutional Neural Networks, object detection begins developing unprecedentedly, and object detection algorithms based on deep learning has become mainstream. Deep learning-based object detection algorithms are mainly divided into two categories, One-Stage and Two-Stage. Common algorithms for One-Stage include OverFeat [6], YOLOv2 [10], YOLOv3 [11], YOLOv4 [12], SSD [13], and RetinaNet [14], while for Two-Stage there are R-CNN [15], SPP-Net [16], Fast R-CNN [17], and Faster R-CNN [18], etc.

Currently, these algorithms are widely used in many fields such as security, military, transportation, and daily life, etc. [19~22]. Among the algorithms for object detection, One-Stage's object detection algorithm enables a completely single training of shared features and rapid acceleration, while guaranteeing accuracy. For example, YOLOv1 enables real-time object detection with the basic network running at 45 frames per second and processes streaming videos in real time with a latency of less than 25 ms [23].

As a quintessential computer vision problem, semantic segmentation involves taking some raw data (e.g., planar images) as input and converting them into masks with highlighted regions of interest. Semantic segmentation is mainly divided into standard semantic segmentation and instance-aware semantic segmentation, commonly using FCN [24], OCRNet [25], Deeplabv3+ [26], and UPerNet [27]. These algorithms are mainly used in self-driving cars, geological detection, facial segmentation, precision agriculture, clothing classification, and other fields [28–32].

Instance segmentation is currently a comparatively mature technique in the field of real-time scene understanding and image information processing [33]. It is the localization of instances in an image using a object detection algorithm, and then on the basis of semantic segmentation, the target objects in different localization frames are further segmented into specific objects of classified categories. In contrast to the bounding box of object detection, instance segmentation can be accurate to the edges of objects, while compared with semantic segmentation, it needs to annotate different individuals of the same object on the graph. The existing instance segmentation networks, such as Mask R-CNN [34] and Faster R-CNN [18], have the high rate of accuracy, however, they fail to process details such as segmented edges meticulously and have a low speed.

All stages of fish growth and development are susceptible to various external factors, possibly resulting in poor growth or even death of fish, which brings serious losses to fish farmers. The focus of current intelligent fishery research is the way to analyze and understand the state of fish growth process in a timely and effective manner. Therefore, our research will focus on the real-time nature of the model, so that we can later deploy it on cell phones or monitoring devices to be applied to fish farms for real-time monitoring of fish growth process. In recent years, researches on fish image detection segmentation based on deep learning have attracted great attention. Chen G et al [35] proposed an automatic fish classification system based on deep learning, Akgul T et al [36] proposed fish detection in turbid underwater, Xiu L et al [37] proposed a deep and lightweight network for detecting fish, Wang J H et al [38] proposed a detection of abnormal behaviors of underwater fish using artificial intelligence techniques, Funkuralshdaifat N F et al [39] proposed a deep learning framework for segmentation of fish with underwater videos, KM Knausgrd et al [40] proposed a method for detecting and classifying temperate fish, Piyadarsan Parida and Nilamani Bhoi (2018) [41] proposed a hybrid transition region color image segmentation method based on dual transition region extraction for fish image segmentation applications. Xiangxiao Lei (2018) et al [42] proposed an image segmentation method based on equivalent 3D entropy and artificial fish population

optimization algorithm, which is more efficient than the traditional 3D entropy method and the equivalent 3D entropy method. Kazim Raza and Song Hong (2020) [43] proposed a quick and accurate method for fish detection based on the improved YOLO-v3 model and migration learning. Its average mean Accuracy Precision (mAP) of 87.56%, which increased from 87.17% to 91.30%, compared with the experimental analysis of the YOLOv3 model and the improved model. Kewei Cai et al [44] (2020) proposed an improved YOLOv3 fish detection model based on MobileNetv1 as the backbone, which combined YOLOv3 with MobileNetv1 for real detection of farmed fish. C. S. Arvind et al [45] used Mask R-CNN model along with GOTURN tracking algorithm for fish detection and tracking, showing that image multi-region parallel processing and tracking are accurate with F1 score of 0.91 at 16 frames per second inland. Nevertheless, most methods mentioned above suffer from slow segmentation and low speed, and also cannot be better monitored in real time. Thus, we use a multi-task network YOLOv5 as the backbone network and object detection branch, combined with semantic segmentation head optimized by GhostC3 module, improving the speed and accuracy of the model in general.

In short, our main contributions are:

1. Construct a golden crucian carp segmentation dataset: Because of little public dataset of fish segmentation, we have built a new golden crucian carp dataset, containing 640 semantic segmentation images of 10 golden crucian carps.
2. Give comparison to the existing mainstream instance segmentation networks, object detection networks and semantic segmentation networks.
3. Detect and segment golden crucian carp based on YOLOv5. We use a multitasking model to solve the problem of monitoring fish and acquiring fish segmentation images in real time, which perform real-time detection and segmentation while ensuring high accuracy. At the same time, we use Ghost module to further optimize the model and reduce the number of parameters of the model. Compared with most of the networks framed by instance partitioning, the FPS of our model is higher and can reach the standard of real-time monitoring.

Materials & Methods

Multi-task learning model is learning multiple tasks together, and the efficiency and quality of studying each task can be improved by learning the connections and differences between different tasks. Different from traditional single-task learning model, which seeks to use a specific model [46] to accomplish the task, multi-task learning model is more effective to improve the generalization of models. In addition, computational repetition can be minimized, inference speed can be increased, and memory utilization can be reduced by sharing layers across multiple tasks. [47]

In this paper, our dataset is divided into two types of annotations, and our model also includes 2 tasks: object detection and semantic segmentation of golden crucian carp. These two share the backbone based on convolutional neural network, and the shared layer improves the feature extraction ability of golden crucian carp and obtains better performance and

generalization. The architecture of our model is shown in Figure 1, and the specific experimental procedure and methodology will be pointed out below.

2.1. Acquisition of Materials

Golden crucian carp is a kind of fish with strong physique, strong resistance and wide eating habits. It is easy to raise and doesn't need good care. At the same time, the shape of golden crucian carp is similar to crucian carp and grass carp, so the method we proposed can also be well promoted to the cultivation of crucian carp and grass carp in the future. To sum up, we took the golden crucian carp as the experimental object. During the experiment, we collected pictures of fish in a non-intrusive way, without any impact on the fish itself.

We purchased 10 easily distinguishable golden crucian carps in the ornamental fish market and used a transparent tank, 80 cm long, 30 cm wide and 60 cm high, as the breeding environment. The size of the purchased tank and the number of fish is reasonable, so as to ensure that the data collected will not be too concentrated and easy to distinguish. The 10 golden crucian carps were raised in a transparent square tank, provided with sufficient oxygen, and fed with common ornamental fish feed. At the same time, in order to simulate a real environment in the water, we added different doses of water conditioner DEBAO and nitrifying bacteria YUECAI to the initial water environment. Finally, it formed a stable and balanced water environment and achieved the optimal effect.

We used the DJI pocket2 camera to capture golden crucian carps from different angles and distances. In order to ensure the reliability of the experimental data, the images were taken randomly, taking into account the effects of sunlight as well as indoor lighting and other factors. When collecting the dataset, we placed the camera about 15-20 cm away from the front of the fish tank and adjusted the camera angle to ensure that we could get a complete figure of the inside of the fish tank without too much background outside the fish tank, which not only ensure the diversity of the collected images, but also enhance the adaptability of subsequent models to various environments [1]. In the deep learning model, the quality of the image has a crucial impact on the training results of the model. The higher the quality of the image annotation is, the better the performance of the model may be. Images without golden crucian carp are useless images for data labeling work, and images with golden crucian carp heavily obscured by aquatic plants and images with turbid water will make labeling work difficult. Therefore, for the captured images, we manually screened and removed some images of poor quality (such as images of golden crucian carp not captured, images of golden crucian carp heavily occluded by aquatic plants and images with turbid water).

Finally, we obtained a dataset consisting of 1858 images from 10 golden crucian carps. There are 2 categories of datasets, which contain 1858 general frame object detection datasets and 640 segmentation datasets.

We raised golden crucian carps and photographed them, and our experiments met the ethical requirements for animal welfare, and we obtained the affidavit of approval of animal ethical and welfare approved by the Sichuan Agricultural University IACUC, numbered

20200054. After the data collection was completed, we released the ten golden crucian carps and placed them in an ornamental fish pond on campus, where they could be fed and cared for by relevant professionals.

The above dataset was annotated by 10 annotators under the guidance of professionals. The annotation work is divided into two main parts. In the first part, the object detection dataset was produced. Labellmg is a graphical image annotation tool written in Python language, often used for the production of object detection datasets. The annotations are stored as xml files in PASCAL VOC format, and YOLO format is also supported.. Therefore, we use labellmg to frame the target, as shown in Figure 2, and the annotation is done by framing golden crucian carps' bodies with a horizontal box to ensure that each fish is completely framed into the corresponding box.

For the second part, the segmented dataset is created. Labelme, an image annotation tool developed by the Computer Science and Artificial Intelligence Laboratory (CSAIL) of MIT, can be used to annotate images in polygons, polylines, points and other forms, which is a common tool for image annotation in segmentation tasks. Therefore, we use labelme to label golden crucian carp and generate json format for the corresponding image. As shown in Figure 3, the body outline of each golden crucian carp in the image is framed with dense points, and the number of golden crucian carp in each image varies. Finally we obtained the annotation results of about 6000 golden crucian carps, and these datasets were used for the subsequent experiments of semantic segmentation and instance segmentation.

In the subsequent experiments, we divided both the object detection and segmentation datasets into a training dataset, a validation dataset, and a test dataset, all in the ratio of 8:1:1.

2.2. Data enhancement for object detection

YOLOV5 passes each batch of training data through the data loader and enhances the training data at the same time. The data loader performs three types of data enhancement: scaling, color space adjustment and mosaic enhancement. For the training of the object detection YOLOv5 model, we used various tricks such as Fliplrud, MixUp, Mosaic, and HSV_Aug for data enhancement on the golden crucian carp dataset, analyzing and comparing the effect of different combinations of tricks on the golden crucian carp dataset, and selecting the most suitable tricks.

2.2.1. Fliplrud

Fish are moving freely in the water and their positions and poses are diverse. The dataset should be contained more information of fish to improve the recognition ability of the model. Therefore, we set up two methods respectively: random horizontal flip and random vertical flip, and the execution probability is 0.5. As shown in Figure 4, fliplr is flipping the image 180 degrees from left to right or right to left, and flipud is flipping 180 degrees from top down or bottom up.

2.2.2. MixUp

MixUp is an algorithm for mixed-class enhancement of images used in computer vision. It can mix images between different classes, specifically by fusing the features and labels of two samples to expand the training dataset.

We took the fusion ratio λ to obey the beta (α, β) distribution (where the α value range is generally $[0.1, 0.4]$, and the λ value range is between $[0, 1]$). As shown in Equation 1, we mix two random samples, and for each input batch image, together with the randomly selected image, it is fused with the randomly selected image, separately for the image itself and the corresponding label according to the fusion ratio λ , to obtain the mixed tensor Mixed_Batch. Among them, the principle of MixUp actually uses linear interpolation.

$$\text{Mixed_batch} = \lambda * \text{batch}_1 + (1 - \lambda) * \text{batch}_2 \quad (1)$$

2.2.3. Mosaic

Mosaic randomly selects four images from the dataset at a time. After randomly cropping, randomly flipping, and color gamut changes for each image, the four images are placed in the order of top left, top right, bottom left, and bottom right and superimposed to form a new image. As shown in Figure 5, this is the image obtained after Mosaic operation during the training process. Mosaic not only enriches the background of the dataset images, but also expands the dataset and improves the robustness of the model.

2.2.4. HSV_Aug

RGB is the most common color space that we often come into contact with. The golden crucian carp dataset image we've collected is an RGB image. Images acquired in natural environments are often disturbed by factors such as lighting, occlusion, shadows, etc. For such changes in brightness, the three components in the RGB color space are highly correlated with it. In the HSV color space, the change of light can be directly expressed through the change of brightness, which is closer to people's perception of color and more suitable for image processing. Therefore, we converted the values of the R, G, and B components of the golden crucian carp data into HSV components according to the following Equations (2–4), thereby obtaining an HSV image. In this way, the characteristics of the image can be expressed more intuitively, and the effect can be significantly enhanced.

$$V = \max(R, G, B) / 255 \quad (2)$$

$$S = (\max(R, G, B) - \min(R, G, B)) / (\text{float})\max(R, G, B) / 255 \quad (3)$$

$$H = \begin{cases} 0^\circ, & \text{if } \max(R, G, B) == \min(R, G, B); \\ 60 * \frac{G - B}{\max(R, G, B) - \min(R, G, B)} + 0, & \text{if } \max(R, G, B) == R \text{ and } G \geq B; \\ 60 * \frac{B - R}{\max(R, G, B) - \min(R, G, B)} + 120, & \text{if } \max(R, G, B) == G; \\ 60 * \frac{R - G}{\max(R, G, B) - \min(R, G, B)} + 240, & \text{if } \max(R, G, B) == B; \\ 60 * \frac{G - B}{\max(R, G, B) - \min(R, G, B)} + 360, & \text{if } \max(R, G, B) == R \text{ and } G < B \end{cases} \quad (4)$$

2.2.5. FocalLoss

FocalLoss is a kind of dynamic scaling cross-entropy loss based on dichotomous cross-entropy, which focuses on hard-to-distinguish samples by a dynamic scaling factor to reduce the weight of easy-to-distinguish samples in the training process. FocalLoss mainly solves the problem of one-stage object detection cases with extreme imbalance between foreground and background classes during training. The formula is as follows:

$$FL(p_t) = -(1 - p_t)^\gamma \log(p_t), \quad (5)$$

where γ is a constant in the range of values $[0, 5]$, and the formula is as follows:

$$P_t = \begin{cases} p & \text{if } y = 1 \\ 1-p & \text{otherwise,} \end{cases} \quad (6)$$

where the values of y are 1 and -1, representing the foreground and background, respectively, and the values of p are 0~1, representing the probability that the model prediction belongs to the foreground.

By using the above data enhancement techniques, we are able to diversify our dataset, reduce overfitting of the data, and effectively improve the generalization ability of the model.

2.3. Methods of detection and estimation

We used the golden crucian carp dataset to test some of the current mainstream instance segmentation networks, and the experimental results are shown in Table 1.

mAP (mean Average Precision) is an important measure of the object detection algorithm, which is calculated by the comprehensive weighted average of the average accuracy rate (AP) of all class detection. We calculate the mAP for IoU at 0.5, and the process requires a definite integral to the Precision-Recall curve to find the area. The [mAP@0.5](#) can be calculated as

$$mAP@0.5 = \frac{1}{n} \sum_{i=0}^{n-1} \int_0^1 p(r_i) dr, \quad (7)$$

where n is the number of categories and $p(r_i)$ is a Precision-Recall curve function formed by using the model's Recall as the abscissa and Precision as the ordinate.

FPS refers to the number of images that can be processed by the network model in 1s.

As is seen from Table 1, although these instance segmentation networks have high precision, FPS is low and real-time monitoring is poor. According to the principle of instance segmentation, instance segmentation can be regarded as a combination of object detection and semantic segmentation. Moreover, the current object detection networks (such as Yolov4s, CenterNet, etc.) are fast and have a good real-time detection effect. Therefore, we try to use the combination of object detection and semantic segmentation instead of instance segmentation to achieve real-time detection and segmentation.

In a word, our main methods in this paper are as follows: using the Yolov5 network as the backbone network and object detection head, adding semantic segmentation head in Yolov5, using Ghost module to further optimize the semantic segmentation head, replacing the activation function.

2.3.1. Object detection

In this paper, we primarily use the YOLOv5 network, version 5.0, as the backbone network and object detection branch, and the full name of YOLO is you only look once, which means that you only need to browse once to identify the category and location of objects in the figure.

The YOLOv5 model was publicly released by Ultralytics on June 9, 2020, and their code and experimental results are open-sourced at <https://github.com/ultralytics/yolov5/tree/v5.0>. The YOLOv5 model is based on the improved YOLOv3 model, with YOLOv5s, YOLOv5m, YOLOv5l, and YOLOv5x models. In terms of model volume, the number of parameters for YOLOv5s, YOLOv5m, YOLOv5l, and YOLOv5x are 7.3M, 21.4M, 47.0M, and 87.7M, respectively. The detection speed of YOLOv5s is faster than all the other three models, especially 2 times faster than YOLOv5x. In general, the smaller the model is, the faster the detection speed is, and it's also easier to deploy on embedded devices, so we choose YOLOv5s for our subsequent experiments.

YOLOv5 is a One-Stage network, which has improved accuracy and speed compared with other networks before YOLOv5 in the YOLO series, and its network flow is shown in Figure 6.

The network structure of YOLOv5 mainly is composed of Backbone, Neck, and Head, where YOLOv5 uses CSPDarknet as Backbone to extract abundant information features from the input image. Neck utilizes PANet structure, and Neck is mainly used to generate feature pyramids. YOLOv5 uses the same Head as YOLOv3, which is a 1*1 convolutional structure with three sets of outputs.

Either Adam or SGD optimizer can be chosen in YOLOv5, and we use SGD optimizer.

For the CSP structure in YOLOv5, YOLOv5 uses two kinds of CSP structures, the first one mainly uses CSP1_X in Backbone, where the Bottleneck is using the Resunit structure, and the second one is using CSP2_X in Neck, where the Bottleneck does not use the Resunit structure.

YOLOv5 has significant advantages shown as follows:

- Compared with the Darknet framework used in YOLOv4, using the Pytorch framework is very user-friendly and easy to train your own dataset, and implement into production
- It's easy to read code, integration of a large number of computer vision techniques, very conducive to learning and learning
- Not only is it easy to configure the environment, but also the model training is rapid, and batch inference can produce real-time results
- It's able to effectively inference directly on individual images, batch images, video and even webcam port inputs
- The Pytorch weights can be easily converted to ONNX format for Android, which can then be converted to OPENCV format, or to IOS format via CoreML for direct deployment to mobile applications.
- The speed of YOLOv5s object recognition is impressively up to 140FPS, and the using experience is very excellent

2.3.2. Semantic segmentation

In this paper, we add a segmentation head to YOLOv5 network for the segmentation task. The segmentation head takes the output of the 16th layer of YOLOv5 as input, and the channel

configuration is 512. It is mainly composed of C3, SPP and Up Sampling modules. Our segmentation head network refers to the design idea of Deeplabv3+ [26], PSPNet [48] and other networks. The specific network structure is shown in Figure 7.

SPP [16] is the abbreviation of spatial pyramid pooling, which can effectively avoid image distortion caused by image region clipping and scaling operations.

The Up Sampling here uses bilinear upsampling, which is one of the interpolation algorithms and is an extension of linear interpolation. It uses four real existing pixel values around the target point in the original image to jointly determine a pixel value in the target map, and the core idea is to perform linear interpolation once in each of the two directions.

As shown in Figure 8, we need to ask for the pixel values of point P. We know the coordinates of Q11, Q21, Q12, Q22, and P. The pixel values of Q11, Q21, Q12, and Q22 are also known. So first use the single linear interpolation about X to calculate the pixel values of R1 and R2 respectively, as shown in Equations 8 and 9.

$$f(x,y_1) \approx \frac{x_2 - x}{x_2 - x_1} f(Q_{11}) + \frac{x - x_1}{x_2 - x_1} f(Q_{21}) \quad (8)$$

$$f(x,y_2) \approx \frac{x_2 - x}{x_2 - x_1} f(Q_{12}) + \frac{x - x_1}{x_2 - x_1} f(Q_{22}) \quad (9)$$

The letters $f(Q_{11})$, $f(Q_{12})$, $f(Q_{21})$, $f(Q_{22})$, x_1 , x_2 , and x in the equation on the right are already known, and the derived $f(x,y_1)$ and $f(x,y_2)$ are the pixel values of R1 and R2.

Then the pixel values of point P are calculated with single linear interpolation about the y direction to obtain Equation 10.

$$f(x,y) \approx \frac{y_2 - y}{y_2 - y_1} f(x,y_1) + \frac{y - y_1}{y_2 - y_1} f(x,y_2) \quad (10)$$

The letters y_1 , y_2 , and y are known in the equation on the right. $f(x,y_1)$ and $f(x,y_2)$ are the values of R1 and R2 pixels found in the previous equation.

The C3 module in the segmentation head is a feature extraction module designed by YOLOv5 based on CSPNet (Cross Stage Partial Network) (Wang et al., 2020), in which the stacking of multiple Bottleneck layers causes significant overhead.

In addition, affected by the GhostNet [49] function of the lightweight network architecture, we replaced the C3 module with the GhostC3 module. The GhostC3 module structure is shown in Figure 9. The main difference between GhostC3 and C3 is that it replaces the Bottleneck layer in C3 with the GhostBottleneck layer, which makes GhostC3 use fewer parameters, but has higher accuracy than the original. In our test, we found that after replacing C3 module with GhostC3 module, the overall parameters used by the model changed from 7731361B to 7577985B, a decrease of 2%.

GhostC3 mainly consists of Conv, concat, and GhostBottleneck modules, and the main components of GhostBottleneck module are GhostConv. Among them, the structure of GhostConv is shown in Figure 10, C1 and C2 are the number of input and output channels, half of the output feature map comes from one regular convolution, and the other half is generated by

5×5 Depthwise convolution on the result of the first one. Compared with the original convolutional layer, GhostConv is able to achieve the same or even more efficient feature extraction with less complexity. In CNN networks, it is common to have a lot of convolutional computations with redundant intermediate feature maps, while GhostConvolution forces the network to learn useful features from half of the convolutional kernels.

The structure of the GhostBottleneck layer is similar to the residual part of the Resnet. Figure 11 shows the structure for a step size of 1, which consists of two Ghost convolutional stacks with 1×1 kernels and residual connections. The GhostBottleneck layer with a step size of 1 can be used to increase the number of channels of the input feature map and to expand the processing for later operations.

Figure 12 shows the structure of GhostBottleneck with a step size of 2. Depthwise convolution is applied to downsample the feature map between GhostConv. To ensure that the features in the residual branch have the same dimension as the original features, the Depthwise convolution is used for downsampling and the original 1×1 convolution is used to change the number of channels. The GhostBottleneck layer with a step size of 2 enables feature extraction by stacked GhostConv when downsampling the feature map and lifting the channel. At the same time, it mitigates the gradient disappearance problem by residual concatenation.

The GhostBottleneck structure approach can not only effectively reduce the model parameters and computation, but also improves the detection efficiency of the model by optimizing the feature map with the Ghost module. Therefore, we replaced the original C3 module with the GhostC3 module.

2.3.3. Activation function improvement

The YOLOv5 model provides a variety of activation functions, such as ReLU, Swish, etc. The activation functions can increase the nonlinear factors. The linear model is not sufficiently expressive, and a nonlinear function is introduced as the excitation function so that the deep neural network is more expressive.

Therefore, to further improve the accuracy of our model, for the activation function, we replaced the Swish activation function used Conv of the YOLOv5 model with the HardSwish activation function.

The Swish activation function (as shown in Figure 13) replaces the ReLU activation function, which can significantly improve the accuracy of the neural network, and the specific definition of the Swish activation function is shown in Equation 11.

$$f(x) = x * \text{sigmoid}(x) \quad (11)$$

Although this non-linearity improves the precision, the sigmoid function is composed of exponentials, which are much more computationally expensive on mobile devices. The sigmoid activation function can be fitted with the segmented linear function HardSigmoid, as shown in Equation 12.

$$\text{Hardsigmoid}(x) = \begin{cases} 0, & x \leq -3 \\ 1, & x \geq 3 \\ \frac{x}{6} + \frac{1}{2}, & \text{otherwise} \end{cases} \quad (12)$$

As a consequence, replacing sigmoid with Hardsigmoid can greatly reduce the cost of operations, which led to the birth of HardSwish (As shown in Figure 14), as specified in the Equation 13.

HardSwish(x)

$$= x * \text{Hardsigmoid}(x) = x * \frac{\text{ReLU6}(x + 3)}{6} = x * \begin{cases} 1, & x \geq 3 \\ \frac{x}{6} + \frac{1}{2}, & -3 < x < 3 \\ 0, & x \leq -3 \end{cases} \quad (13)$$

The derivative of this function with respect to x is:

$$\text{HardSwish}'(x) = \begin{cases} 1, & x \geq 3 \\ \frac{x}{3} + \frac{1}{2}, & -3 < x < 3 \\ 0, & x \leq -3 \end{cases} \quad (14)$$

Results

Our experiments are divided into three main steps.

First, we validate multiple object detection models, and then select the network with better performance as the baseline to be used as the model in the first step.

Second, we add a segmentation head to the baseline model and perform ablation experiment on it to select the most suitable tricks.

Third, we further optimize the proposed model and compare the results with those produced by other models to demonstrate the competitive advantage of our approach.

Through stepwise experiments, we can ensure that each step of our identification of golden crucian carp is a local optimal solution, thus an effective model is built for the detection and segmentation of golden crucian carp.

3.1. Model selection

3.1.1. Object detection model selection

Currently, the mainstream target detection models are CenterNet, YOLOv4, YOLOv5, EfficientDet, and RatinaNet, and we trained these five target detection models on the existing golden crucian carp dataset, and the results of these model tests are listed in Table 2. To ensure the accuracy of comparison between the results of each model above, we uniformly conducted the experiments on RTX 5000 with uniform settings of parameters: image size of 640x640, optimizer of Adamw, initial learning rate set to 0.001, and epoch of 300.

In the above experiments, we choose Precision, Recall, F1-score, mAP@0.5, mAP@0.5:0.95 and Inference@batch_size 1 as the performance evaluation indexes for model effectiveness, and we consider these indexes together to determine the selection of our final

model. Through Table 2, we find that CenterNet has the highest precision of 95.21% on the golden crucian carp dataset, but its Recall value is lower compared to other networks, only reaching 92.48%, while YOLOv5s Recall reaches 95.38%. Therefore, we reconcile Precision and Recall for the mean. And we summed Precision and Recall to obtain the F1-score, which is exactly the same. But in terms of mAP, YOLOv5s is better. At the same time, from the overall point of view, YOLOv5s reaches the top of each index compared with the rest of the network models, so we finally choose YOLOv5s.

In the experimental prediction process, the following four cases occur for the golden crucian carp dataset: (a) correctly identify the golden crucian carp, indicated by TP (b) incorrectly identified other things as golden crucian carp, denoted by FP (c) did not identify golden crucian carp, denoted by FN (d) non- golden crucian carp things were not identified as golden crucian carp, denoted by TN. To obtain the specific values of these three cases, the calculation of IoU is required, which is the degree of overlap between the predicted boxes and the boxes marked in the original figure, as shown in Figure 15, which is calculated as in Equation 15 (where G is the real box and D is the predicted box), and in general, we write down the number of detected boxes with IoU greater than 0.5 as TP, while the data of detected boxes whose values are less than or equal to 0.5 are written down as FP.

$$IoU = \frac{G \cap D}{G \cup D} \quad (15)$$

Using TP,FP and FN, we can calculate the above four metrics in the following way.

As shown in Equation 16, Precision is the ratio of the number of correct predictions in the forecast results.

$$Precision = \frac{TP}{TP + FP} \quad (16)$$

Recall denotes the probability of being detected in all positive samples during the prediction of experimental results, which in this experiment is the probability of being detected in all pictures containing golden crucian carp, and is calculated as in Equation 17.

$$Recall = \frac{TP}{TP + FN} \quad (17)$$

The F1-score can well distinguish the strengths and weaknesses of the algorithm and is the summed average of Precision and Recall, as shown in Equation 18.

$$F1 - score = 2 \times \frac{Recall \times Precision}{Recall + Precision} \quad (18)$$

mAP@0.5:0.95 represents the average mAP at different IoU thresholds (from 0.5 to 0.95 in steps 0.05).

Inference@batch_size 1 indicates the inference time required for a picture.

The experimental results show that YOLOv5s and CenterNet show relatively better performance, with F1-scores reaching 0.94 and mAP@0.5 around 0.95. However, the Inference@batch_size 1 of YOLOv5 is smaller than CenterNet ,and takes less inference time to

process the same photo, so we chose YOLOv5s as the backbone network for subsequent experiments.

3.1.2. Selection of data enhancement tricks

We added a semantic segmentation head in YOLOv5 for subsequent experiments. To further improve the effectiveness of our model and enhance its generalization ability, we used different tricks to process the model and further ablation experiments were conducted to extract the best model training configuration. Table 3 and Table 4 show the effects of using tricks such as HSV_Aug, FocalLoss, Mosaic, MixUp and the evaluated metrics afterwards. This symbol, “√” indicates that the technique was used and “×” indicates that it was not used.

Considering the connection and impact in different data enhancement strategies (e.g., MixUp [50] is implemented on top of Mosaic, thus MixUp is only used as the Mosaic strategy is enabled), we finally selected only 12 combinations as shown in Table 3 for our experiments.

The results of the experiments are shown in Table 4, and the results show that some tricks decrease the accuracy of our network, while some tricks increase the accuracy of the network.

For the semantic segmentation model, mIoU is its standard metric. The formula of mIoU is shown in Equation 19, which is labeled as Seg mIoU in Table 4. In our experiment, the golden crucian carp dataset is mainly divided into two categories, fish and background, so k is 1.

$$mIoU = \frac{1}{k+1} \sum_{i=0}^k \frac{TP}{FN + FP + TP} \quad (19)$$

seg ACC refers to the pixel accuracy, and its corresponding formula is shown in Equation 20.

$$seg\ ACC = (TP + TN) / (TP + TN + FP + FN) \quad (20)$$

FocalLoss is often used to suppress background classes in unbalanced data and target detection[51]. As can be seen in the results of Table 4 for both group 2 and 3, and group 7 and 10, the accuracy of the model after FocalLoss processing decreases with the golden crucian carp dataset and the corresponding settings in this paper. We speculate that there is no significant imbalance in the space allocation between background and object classes in the data, causing FocalLoss not to work.

All of our detectors were experimented using pre-trained weights. The experimental results of groups 3 and 4 can clearly see that Mosaic has a certain effect on improving the precision of the model, and the experimental results of groups 6 and 7 can also show that Mosaic can improve the performance of the model, so we decided to use Mosaic. Looking at the experiments in groups 8 and 9, we can see that MixUp has little effect on model precision, recall and other parameters. Group 5 added MixUp to group 4, and the precision of the model decreased, and group 11 added MixUp to group 10, and the performance of the model also decreased, so we decided not to use MixUp. Similarly, the experiments in groups 7 and 8 led us to reject FlipLrud. Therefore, we chose the tricks of group 8, that is, we only used HSV_Aug and Mosaic for subsequent experiments.

3.2. Parameter setting

In the Yolov5s model with segmentation head, we finally used the 8th set of tricks as in Table 4, and we set the hyperparameters of HSV color model as H (hue): 0.015, S (saturation): 0.7, V (brightness): 0.4, and mosaic as 1. We set the image input size as 640*640, epochs We set the input size to 640*640, epochs to 300, SGD optimizer and initial learning rate to 0.0015, and put the custom golden crucian carp dataset with batch size of 2 into the network for training, and our experimental environment for training is ubuntu 20.04 with RTX3060 graphics card.

3.3. Experimental results

Our final model uses YOLOv5s version 5.0 as the backbone network, and adds a segmentation head with Ghost module optimization, and replaces the original Swish activation function with the HardSwish activation function, which is relatively less computationally expensive. The final object detection accuracy is 0.954, the recall is 0.930, and the speed is 116.6 FPS on RTX3060.

The experimental results after improving the semantic segmentation head and replacing the activation function are shown in Table 5.

According to Equation 16, we can see that precision refers to the probability of True positive in all samples that are predicted to be correct, mainly focusing on "finding precision"; According to Equation 17, we can see that recall refers to the probability that True positive accounts for the sample that is actually True, and focuses mainly on "finding the full". If you just look at recall or precision, you may go to extremes without knowing it. We expect the results to be both accurate and complete, so we choose F1-score, which combines precision and recall, as the evaluation metric. We tend to choose the model with a higher F1-score.

After replacing the C3 module in segmentation head with the GhostC3 module, the F1-score and FPS values of the model were improved compared with the results of the initial version of the segmentation head model. After replacing the Swish activation function with the HardSwish activation function on the segmentation head model improved with the Ghost module, the F1-score of the model increase to more than 0.94 and the FPS increase by 0.6, which is relatively obvious. Therefore, we will use the model of the Ghost module that has optimized the segmentation head and the activation function is HardSwish as our final model.

As shown in Table 6, our multi-task model shows advantages over other models in both detection and segmentation tasks. In particular, in terms of speed, the FPS of our model is much higher than the other models. To further demonstrate the effectiveness of our model, we conduct experiments on the publicly available dataset PASCAL VOC 2007. We divide the dataset into training set, validation set and test set in the ratio of 8:1:1, and the size of the training images is set to 512*512. Finally, our object detection accuracy is 0.738, semantic segmentation accuracy is 0.843, and FPS is 120. We use our final model to make predictions for the golden crucian carp videos and take four images from it, as shown in Figure 16.

From this prediction result, we can see the video prediction effect of our model: high accuracy of object detection and high coincidence of semantic segmentation with the object to be detected.

Discussion

The current development trend shows that instance segmentation technology has a large potential for application in smart fisheries. However, the problems of instance segmentation, such as low speed and inferior real-time detection, have contributed to the poor results of its application. In this paper, we study a multi-task model based on YOLOv5, add a segmentation head for semantic segmentation based on the original YOLOv5 framework, improve the segmentation head by Ghost module, and replace the activation function of the model to further improve the overall accuracy of the model. Based on this, our discussions are as follows:

4.1. Contribution to smart fisheries

Fishery is an indispensable part of China's national economy. With the continuous improvement of people's living standards, the process of fish farming has received wide social attention. At all the stages of fish growth and development are susceptible to a variety of external factors, possibly resulting in poor fish growth or even death, thus causing serious losses to farmers. Therefore, the major focus of current research on smart fisheries is how to analyze and understand the state of fish growth process in a timely and effective manner. To address this problem, researchers have proposed many effective methods. Among them, real-time monitoring of fish using deep learning techniques is by far the most effective method. In this paper, we add the optimized segmentation head of Ghost module to the YOLOv5s object detection model to achieve the function of real-time detection and segmentation, and experimentally prove that our method greatly improves the detection speed in real-time monitoring, and the object detection accuracy and semantic segmentation accuracy are above 0.95. Therefore, our method can be applied to the smart fishery aspect, which is important for the monitoring activities of fish growth and development process.

4.2. Limitations and future work

At present, our model mainly utilizes a object detection algorithm with a semantic segmentation algorithm, and the segmentation head technique is used to effectively combine the two algorithms to form a new network model that guarantees the detection of the object while being able to segment it. Our experiments test on the existing golden crucian carp dataset, and the final results have achieved an accuracy of more than 95% and an FPS of 116.6, showing excellent results. However, our experiments still have some limitations. Our dataset is only for the golden crucian carp, and we have not yet experimented the effect of the model on the dataset of common fish such as grass carp and crucian carp. In the future, we will work on the expansion of the original golden crucian carp dataset to diversify the fish species to ensure that the dataset can contain more fish with generalizability, and we will also keep on expanding the quantities of each type of fish to form a large-scale fish dataset. On this foundation, we will conduct further experiments and optimization on the basis of this model to ensure the universality of the model. Our ultimately purpose is to implement the model in practice, and we intend to develop it into a

system that can estimate the biomass of fish and continuously optimize and expand its functions to effectively promote the intelligent development of fishery farming.

Conclusions

In this paper, we first detected and segmented golden crucian carps by the general instance segmentation algorithm. Although the experimental results were comparatively good in terms of accuracy, each FPS was low and it was difficult to reach the standard of real-time monitoring. Subsequently, we added the Ghost module optimized semantic segmentation head to the YOLOv5 model and replaced the activation function, and finally propose our model. It is experimentally confirmed that our model can effectively identify grass carp and perform real-time segmentation, with object detection accuracy up to 95.4% , semantic segmentation accuracy up to 98.5%, and a speed of 116.6 FPS on RTX3060 , which can effectively meet the requirements related to real-time detection.

In the future, our research will be further developed in the following aspects.
1. In this paper, the survival environment of our fish is a fish tank, but in the actual farm, it may be more complex. For this reason, we will consider further optimization of the experimental environment to ensure a closer fit to the farm's culture environment and increase the robustness of the model.

2. Grass carp is a unique fish species and an important freshwater economic fish in China. Considering its mature farming technology, low-cost input, little difficulty in management , high survival rate and high economic value, we choose grass carp as the research target. However, the characteristics of grass carp are not obvious enough compared with golden crucian carp. For this reason, we intend to add some hardware such as one-way pipes to assist us in collecting more detailed and comprehensive data sampling of grass carp to avoid the problem of poor model training results caused by insufficient features.

3. We will further extend our features to perform fish pose estimation, behavior recognition, target detection of fish diseases and semantic segmentation of fish diseases based on semantic segmentation, in order to form a completely intelligent fish farming system and improve the foundation for future application in actual farms.

Acknowledgements

We are grateful to Prof.Duan and other teachers for their support and assistance. It's their help that we can carry out research so smoothly during our college life.

Data Availability

The following information was supplied regarding data availability :

The model is available at <https://doi.org/10.5281/zenodo.7413306>

The golden crucian carp detection dataset is available at <https://doi.org/10.5281/zenodo.7413354>

The golden crucian carp segmentation dataset is available at

<https://doi.org/10.5281/zenodo.7410176>

References

1. Lin, Bin, Kailin Jiang, Zhiqi Xu, Feiyi Li, Jiao Li, Chaoli Mou, Xinyao Gong, and Xuliang Duan. "Feasibility Research on Fish Pose Estimation Based on Rotating Box Object Detection." *Fishes* 6, no. 4 (2021): 65.
2. Brownlee J. *Deep learning for computer vision: image classification, object detection, and face recognition in python*[M]. Machine Learning Mastery, 2019.
3. Ashraf T, Khan Y N. Weed density classification in rice crop using computer vision[J]. *Computers and Electronics in Agriculture*, 2020, 175: 105590.
4. Liaqat, Amna, Muhammad A. Khan, Muhammad Sharif, Mamta Mittal, Tanzila Saba, K. Suresh Manic, and Feras NH Al Attar. "Gastric tract infections detection and classification from wireless capsule endoscopy using computer vision techniques: A review." *Current Medical Imaging* 16, no. 10 (2020): 1229-1242.
5. Srivastava P, Shukla A, Bansal A. A comprehensive review on soil classification using deep learning and computer vision techniques[J]. *Multimedia Tools and Applications*, 2021, 80(10): 14887-14914.
6. Sermanet, Pierre, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. "Overfeat: Integrated recognition, localization and detection using convolutional networks." *arXiv preprint arXiv:1312.6229* (2013).
7. Viola P, Jones M. Rapid object detection using a boosted cascade of simple features[C]//*Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001. Ieee*, 2001, 1: I-I.
8. Dalal N, Triggs B. Histograms of oriented gradients for human detection[C]//*2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05). Ieee*, 2005, 1: 886-893.
9. Felzenszwalb P, McAllester D, Ramanan D. A discriminatively trained, multiscale, deformable part model[C]//*2008 IEEE conference on computer vision and pattern recognition. Ieee*, 2008: 1-8.
10. Redmon J, Farhadi A. YOLO9000: better, faster, stronger[C]//*Proceedings of the IEEE conference on computer vision and pattern recognition. 2017: 7263-7271*.
11. Redmon J, Farhadi A. Yolov3: An incremental improvement[J]. *arXiv preprint arXiv:1804.02767*, 2018.
12. Bochkovskiy A, Wang C Y, Liao H Y M. Yolov4: Optimal speed and accuracy of object detection[J]. *arXiv preprint arXiv:2004.10934*, 2020.
13. Liu, Wei, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. "Ssd: Single shot multibox detector." In *European conference on computer vision*, pp. 21-37. Springer, Cham, 2016.

14. Lin, T. Y., P. Goyal, R. Girshick, K. He, and P. Dollár. "Focal loss for dense object detection. Proceedings of the IEEE International Conference on Computer Vision." (2017): 2980-2988.
15. Girshick, Ross, Jeff Donahue, Trevor Darrell, and Jitendra Malik. "Rich feature hierarchies for accurate object detection and semantic segmentation." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 580-587. 2014.
16. He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Spatial pyramid pooling in deep convolutional networks for visual recognition." IEEE transactions on pattern analysis and machine intelligence 37, no. 9 (2015): 1904-1916.
17. Girshick R. Fast r-cnn[C]//Proceedings of the IEEE international conference on computer vision. 2015: 1440-1448.
18. Ren, Shaoqing, Kaiming He, Ross Girshick, and Jian Sun. "Faster r-cnn: Towards real-time object detection with region proposal networks." Advances in neural information processing systems 28 (2015).
19. Saikia, Surajit, Eduardo Fidalgo, Enrique Alegre, and Laura Fernández-Robles. "Object detection for crime scene evidence analysis using deep learning." In International Conference on Image Analysis and Processing, pp. 14-24. Springer, Cham, 2017.
20. Janakiramaiah, B., G. Kalyani, A. Karuna, L. V. Prasad, and M. Krishna. "Military object detection in defense using multi-level capsule networks." Soft Computing (2021): 1-15.
21. Li, Yanfen, Hanxiang Wang, L. Minh Dang, Tan N. Nguyen, Dongil Han, Ahyun Lee, Insung Jang, and Hyeonjoon Moon. "A deep learning-based hybrid framework for object detection and recognition in autonomous driving." IEEE Access 8 (2020): 194228-194239.
22. Jiang, Qingsheng, Dapeng Tan, Yanbiao Li, Shiming Ji, Chaopeng Cai, and Qiming Zheng. "Object detection and classification of metal polishing shaft surface defects based on convolutional neural network deep learning." Applied Sciences 10, no. 1 (2019): 87.
23. Redmon, Joseph, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 779-788. 2016.
24. Long J, Shelhamer E, Darrell T. Fully convolutional networks for semantic segmentation[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2015: 3431-3440.
25. Yuan, Yuhui, Xiaokang Chen, Xilin Chen, and Jingdong Wang. "Segmentation transformer: Object-contextual representations for semantic segmentation." arXiv preprint arXiv:1909.11065 (2019).
26. Chen, Liang-Chieh, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. "Encoder-decoder with atrous separable convolution for semantic image segmentation." In Proceedings of the European conference on computer vision (ECCV), pp. 801-818. 2018.
27. Xiao, Tete, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. "Unified perceptual parsing for scene understanding." In Proceedings of the European conference on computer vision (ECCV), pp. 418-434. 2018.

28. Masood, Sharjeel, Fawad Ahmed, Suliman A. Alsuhibany, Yazeed Yasin Ghadi, M. Y. Siyal, Harish Kumar, Khyber Khan, and Jawad Ahmad. "A deep learning-based semantic segmentation architecture for autonomous driving applications." *Wireless Communications and Mobile Computing* 2022 (2022).
29. Hu, Guang, Zhengwang Hu, Jiangping Liu, Fei Cheng, and Daicheng Peng. "Seismic Fault Interpretation Using Deep Learning-Based Semantic Segmentation Method." *IEEE Geoscience and Remote Sensing Letters* (2020).
30. Khan K, Mauro M, Leonardi R. Multi-class semantic segmentation of faces[C]//2015 IEEE International Conference on Image Processing (ICIP). IEEE, 2015: 827-831.
31. Du, Zhenrong, Jianyu Yang, Cong Ou, and Tingting Zhang. "Smallholder crop area mapped with a semantic segmentation deep learning method." *Remote Sensing* 11, no. 7 (2019): 888.
32. de Souza Inácio A, Brilhador A, Lopes H S. Semantic segmentation of clothes in the context of soft biometrics using deep learning methods[C]//Anais do 14o Congresso Brasileiro de Inteligência Computacional. Curitiba, PR: ABRICOM. 2019: 1-7.
33. LIANG, Xin-yu, Xi-kun LIN, Ji-chuan QUAN, and Kai-hong XIAO. "Research on the Progress of Image Instance Segmentation Based on Deep Learning." *ACTA ELECTRONICA SINICA* 48, no. 12 (2020): 2476.
34. He, Kaiming, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask r-cnn." In *Proceedings of the IEEE international conference on computer vision*, pp. 2961-2969. 2017.
35. Chen G, Sun P, Shang Y. Automatic fish classification system using deep learning[C]//2017 IEEE 29th International Conference on Tools with Artificial Intelligence (ICTAI). IEEE, 2017: 24-29.
36. Akgül T, Çalik N, Töreyn B U. Deep Learning-Based Fish Detection in Turbid Underwater Images[C]//2020 28th Signal Processing and Communications Applications Conference (SIU). IEEE, 2020: 1-4.
37. Li X, Tang Y, Gao T. Deep but lightweight neural networks for fish detection[C]//OCEANS 2017-Aberdeen. IEEE, 2017: 1-5.
38. Wang, Jung-Hua, Shih-Kai Lee, Yi-Chung Lai, Cheng-Chun Lin, Ting-Yuan Wang, Ying-Ren Lin, Te-Hua Hsu, Chang-Wen Huang, and Chung-Ping Chiang. "Anomalous behaviors detection for underwater fish using ai techniques." *IEEE Access* 8 (2020): 224372-224382.
39. Alshdaifat N F F, Talib A Z, Osman M A. Improved deep learning framework for fish segmentation in underwater videos[J]. *Ecological Informatics*, 2020, 59: 101121.
40. Knausgård, Kristian Muri, Arne Wiklund, Tonje Knutsen Sjørdalen, Kim Tallaksen Halvorsen, Alf Ring Kleiven, Lei Jiao, and Morten Goodwin. "Temperate fish detection and classification: a deep learning based approach." *Applied Intelligence* 52, no. 6 (2022): 6988-7001.
41. Parida P, Bhoi N. Dual transition region extraction based colour image segmentation: Application to fish image segmentation[J]. *Global Journal of Computer Science and Technology*, 2018.

42. Lei X, Ouyang H, Xu L. Image segmentation based on equivalent three-dimensional entropy method and artificial fish swarm optimization algorithm[J]. Optical Engineering, 2018, 57(10): 103106.
43. Raza K, Song H. Fast and accurate fish detection design with improved YOLO-v3 model and transfer learning[J]. International Journal of Advanced Computer Science and Applications, 2020, 11(2).
44. Cai, Kewei, Xinying Miao, Wei Wang, Hongshuai Pang, Ying Liu, and Jinyan Song. "A modified YOLOv3 model for fish detection based on MobileNetv1 as backbone." Aquacultural Engineering 91 (2020): 102117.
45. Arvind, C. S., R. Prajwal, Prithvi Narayana Bhat, A. Sreedevi, and K. N. Prabhudeva. "Fish detection and tracking in pisciculture environment using deep instance segmentation." In TENCON 2019-2019 IEEE Region 10 Conference (TENCON), pp. 778-783. IEEE, 2019.
46. Ma, Jiaqi, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. "Modeling task relationships in multi-task learning with multi-gate mixture-of-experts." In Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining, pp. 1930-1939. 2018.
47. Lu, Kangjie, Jingwen Huang, Junhui Li, Jiashun Zhou, Xianliang Chen, and Yunfei Liu. "MTL-FFDET: A Multi-Task Learning-Based Model for Forest Fire Detection." Forests 13, no. 9 (2022): 1448.
48. Zhao, Hengshuang, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. "Pyramid scene parsing network." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2881-2890. 2017.
49. Han, Kai, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. "Ghostnet: More features from cheap operations." In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 1580-1589. 2020.
50. Zhang, Hongyi, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. "mixup: Beyond empirical risk minimization." arXiv preprint arXiv:1710.09412 (2017).
51. Li, Danyang, Houcheng Su, Kailin Jiang, Dan Liu, and Xuliang Duan. "Fish Face Identification Based on Rotated Object Detection: Dataset and Exploration." Fishes 7, no. 5 (2022): 219.

Table 1 (on next page)

Results of the instance segmentation networks.

1 **Table 1 :**
2 **Results of the instance segmentation networks.**

3

| Model | mAP@0.5 | FPS |
|--------------------|---------|------|
| Mask R-CNN | 0.887 | 10.6 |
| YOLACT | 0.887 | 23.0 |
| Cascade Mask R-CNN | 0.916 | 8.8 |

4

Table 2(on next page)

Comparison of Object Detection Models.

Table 2 :
Comparison of Object Detection Models.

| Model | Precision | Recall | F1-score | mAP@0.5 | mAP@0.5:0.95 | Inference @batch_size 1 (ms) |
|--------------|-----------|--------|----------|---------|--------------|------------------------------|
| CenterNet | 95.21% | 92.48% | 0.94 | 94.96% | 56.38% | 32 |
| Yolov4s | 84.24% | 94.42% | 0.89 | 95.28% | 52.75% | 10 |
| Yolov5s | 92.39% | 95.38% | 0.94 | 95.38% | 58.31% | 8 |
| EfficientDet | 88.14% | 91.91% | 0.90 | 95.19% | 53.43% | 128 |
| RatinaNet | 88.16% | 93.21% | 0.91 | 96.16% | 57.29% | 48 |

¹ Refers from [1]

Table 3(on next page)

YOLOv5s + segmentation head with different data enhancement tricks.

Table 3 :
YOLOv5s + segmentation head with different data enhancement tricks.

| HSV_Aug | FocalLoss | Mosaic | MixUp | Fliplrud | Group |
|---------|-----------|--------|-------|----------|-------|
| × | × | × | × | × | 1 |
| √ | × | × | × | × | 2 |
| √ | √ | × | × | × | 3 |
| √ | √ | √ | × | × | 4 |
| √ | √ | √ | √ | × | 5 |
| √ | √ | × | × | √ | 6 |
| √ | √ | √ | × | √ | 7 |
| √ | × | √ | × | × | 8 |
| √ | × | √ | √ | × | 9 |
| √ | × | √ | × | √ | 10 |
| √ | × | √ | √ | √ | 11 |
| √ | √ | √ | √ | √ | 12 |

Table 4(on next page)

The experiment results of different data enhancement tricks.

Table 4 :
The experiment results of different data enhancement tricks.

| Group | Precision | Recall | F1-score | mAP@0.5 | Seg Acc | Seg mIoU | FPS |
|-------|-----------|--------|----------|---------|---------|----------|-------|
| 1 | 0.904 | 0.931 | 0.917 | 0.955 | 0.981 | 0.978 | 117.1 |
| 2 | 0.905 | 0.935 | 0.920 | 0.963 | 0.980 | 0.977 | 115.0 |
| 3 | 0.888 | 0.924 | 0.906 | 0.947 | 0.975 | 0.972 | 114.4 |
| 4 | 0.922 | 0.939 | 0.930 | 0.966 | 0.985 | 0.983 | 114.8 |
| 5 | 0.904 | 0.924 | 0.914 | 0.957 | 0.984 | 0.982 | 113.8 |
| 6 | 0.897 | 0.918 | 0.907 | 0.942 | 0.985 | 0.983 | 118.7 |
| 7 | 0.899 | 0.937 | 0.918 | 0.955 | 0.982 | 0.980 | 113.4 |
| 8 | 0.954 | 0.917 | 0.935 | 0.976 | 0.985 | 0.984 | 115.1 |
| 9 | 0.952 | 0.924 | 0.938 | 0.974 | 0.986 | 0.984 | 116.2 |
| 10 | 0.935 | 0.924 | 0.929 | 0.964 | 0.984 | 0.982 | 115.3 |
| 11 | 0.893 | 0.910 | 0.901 | 0.943 | 0.986 | 0.984 | 114.3 |
| 12 | 0.900 | 0.936 | 0.918 | 0.954 | 0.986 | 0.984 | 113.3 |

Table 5 (on next page)

Experimental results after improving the semantic segmentation head and replacing the activation function.

YOLOv5s+segmentation head(C3) is Group 8 of Table 4, YOLOv5s+segmentation head(GhostC3) means replace C3 in the segmentation head with GhostC3, YOLOv5s +segmentation head(GhostC3)+HardSwish means replacing C3 in the segmentation head with GhostC3 and replacing the original Swish activation function with the HardSwish activation function.

Table 5:
Experimental results after improving the semantic segmentation head and replacing the
activation function. YOLOv5s+segmentation head(C3) is Group 8 of Table 4,
YOLOv5s+segmentation head(GhostC3) means replace C3 in the segmentation head with
GhostC3, YOLOv5s +segmentation head(GhostC3)+HardSwish means replacing C3 in the
segmentation head with GhostC3 and replacing the original Swish activation function with
the HardSwish activation function.

| Group | Precision | Recall | F1-score | mAP@0.5 | Seg Acc | Seg mIoU | FPS |
|--|-----------|--------|----------|---------|------------|-------------|-------|
| YOLOv5s+segmentation head(C3) | 0.954 | 0.917 | 0.935 | 0.976 | 0.985 | 0.984 | 115.1 |
| YOLOv5s+segmentation head(GhostC3) | 0.952 | 0.927 | 0.939 | 0.975 | 0.985 | 0.983 | 116.0 |
| YOLOv5s +segmentation head(GhostC3)+ HardSwish | 0.954 | 0.930 | 0.942 | 0.976 | 0.985 | 0.983 | 116.6 |

Table 6(on next page)

Comparison results of different models in segmentation task and detection task.

Table 6 :
Comparison results of different models in segmentation task and detection task.

| Model | mAP@0.5 | Seg Acc | Seg mIoU | FPS |
|--|---------|---------|----------|-------|
| FCN | \ | 0.982 | 0.965 | 8.6 |
| OCRNet | \ | 0.985 | 0.967 | 10.1 |
| Deeplabv3+ | \ | 0.981 | 0.966 | 8.0 |
| UPerNet | \ | 0.982 | 0.966 | 8.6 |
| ANN | \ | 0.982 | 0.965 | 7.0 |
| CcNet | \ | 0.982 | 0.965 | 6.9 |
| DANet | \ | 0.982 | 0.966 | 8.0 |
| PSPNet | \ | 0.981 | 0.965 | 9.1 |
| Mask R-CNN | 0.887 | \ | \ | 10.6 |
| YOLOACT | 0.887 | \ | \ | 23.0 |
| Cascade Mask R-CNN | 0.916 | \ | \ | 8.8 |
| YOLOv5s+segmentation head(GhostC3)+HardSwish (Our method) | 0.976 | 0.985 | 0.983 | 116.6 |

Figure 1

The network structure of our model. It contains two tasks, which are the detection task and the segmentation task.

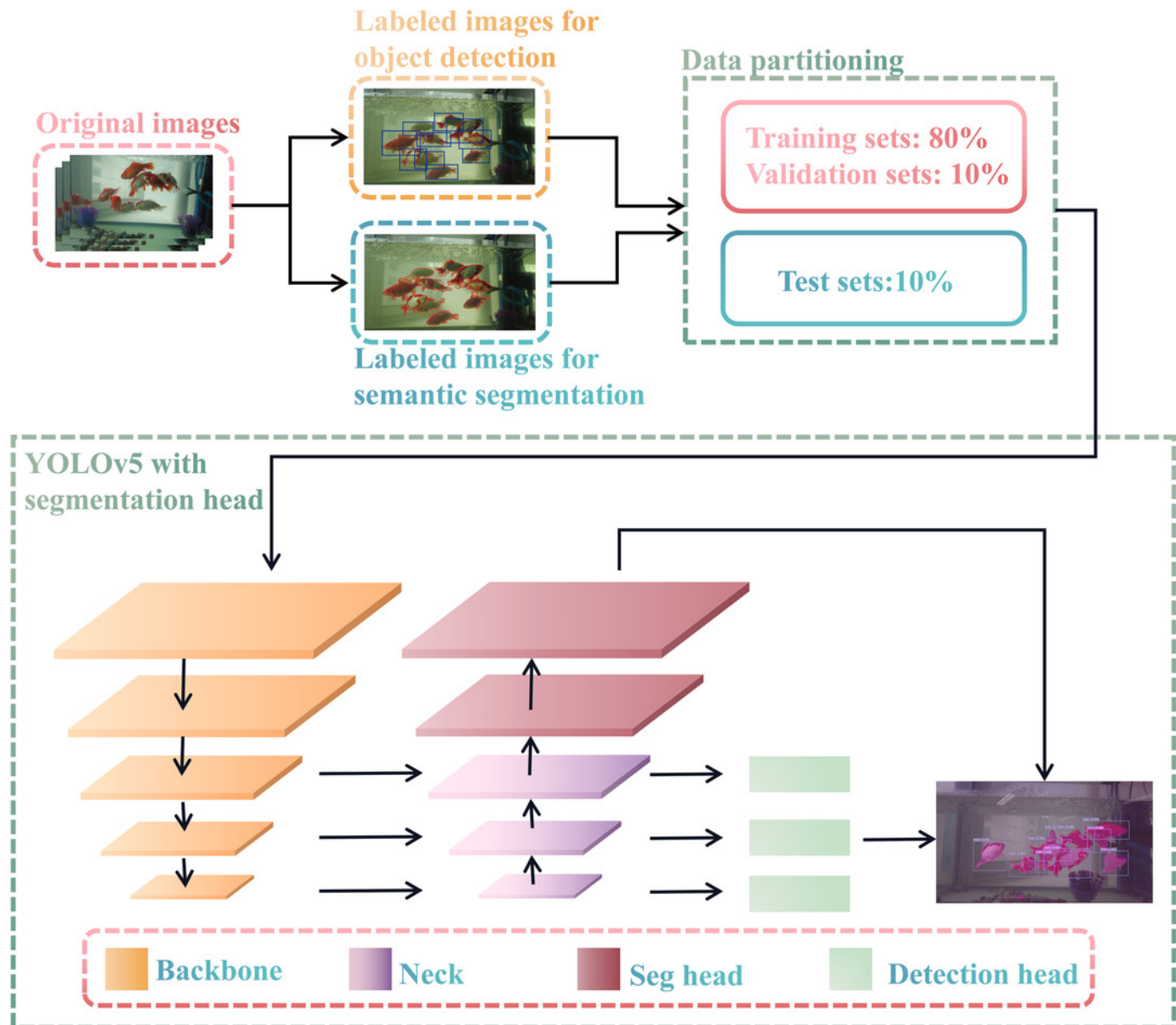


Figure 2

Annotation of object detection dataset.

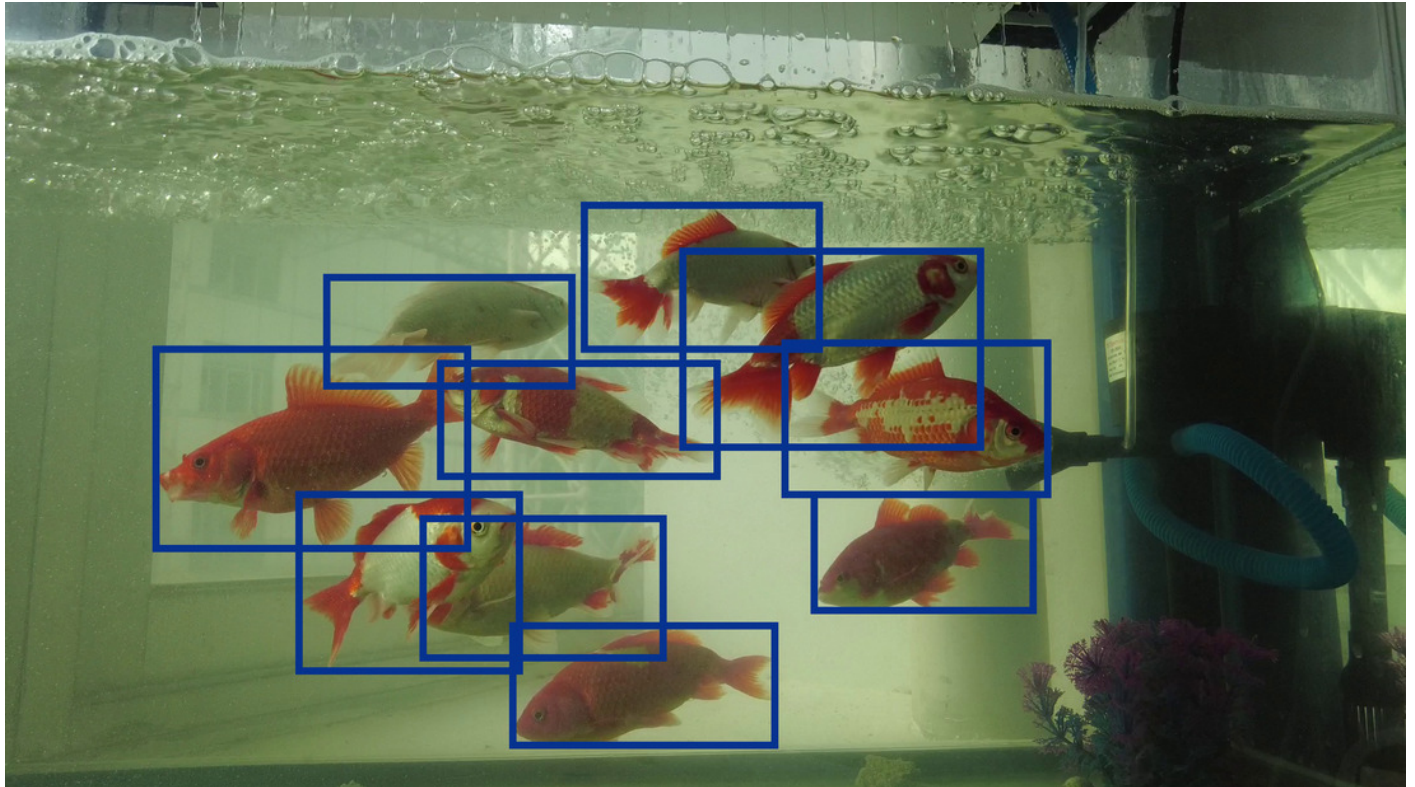


Figure 3

Annotation of segmentation dataset.



Figure 4

Figure (a) is a picture in our dataset, Figure (b) is a picture of Figure (a) flipped left and right, Figure (c) is a picture of Figure (a) flipped up and down, and Figure (d) is a picture of Figure (a) flipped up, down, left and right.



(a)



(b)



(c)



(d)

Figure 5

Training images after mosaic operation.

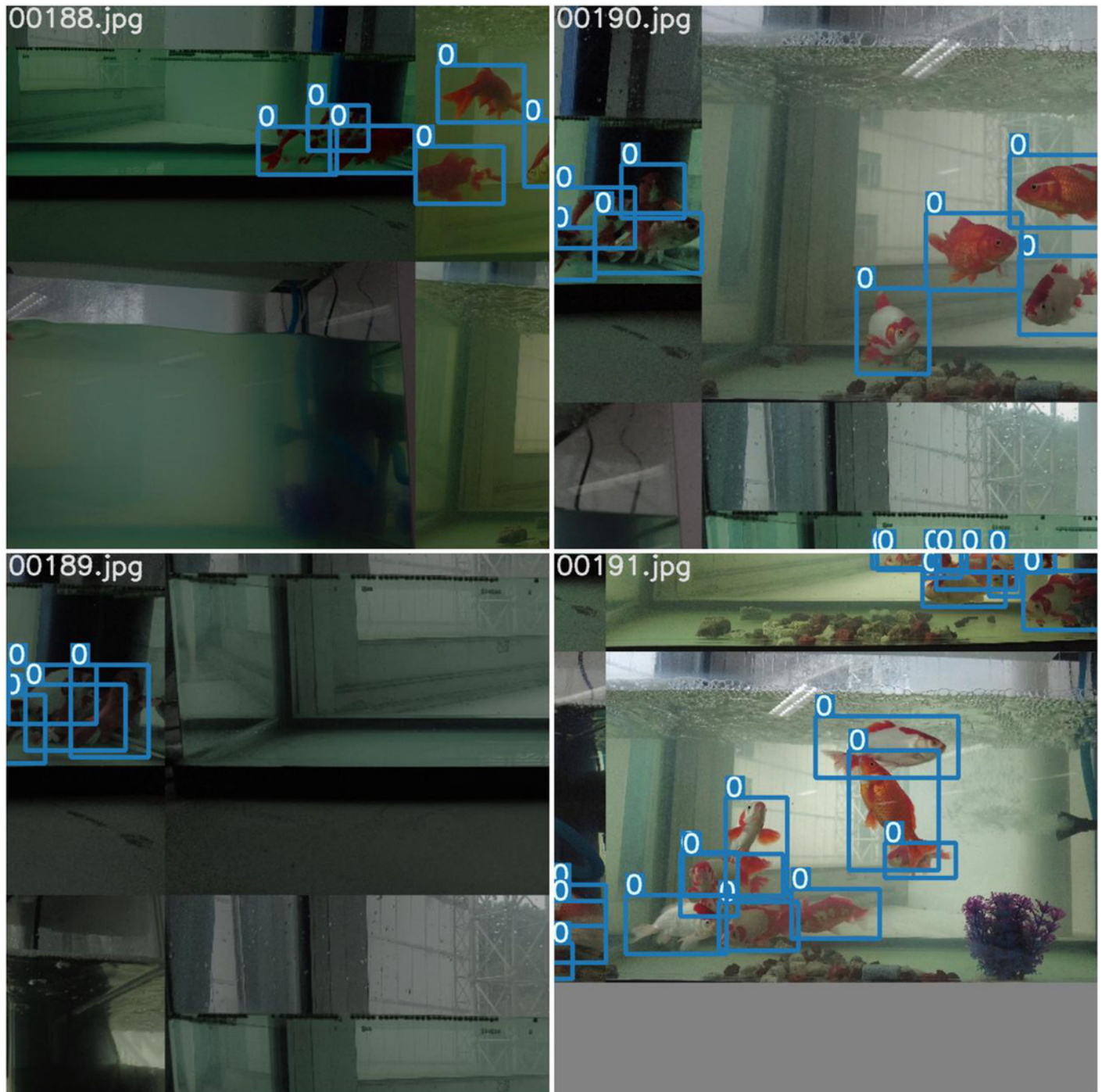
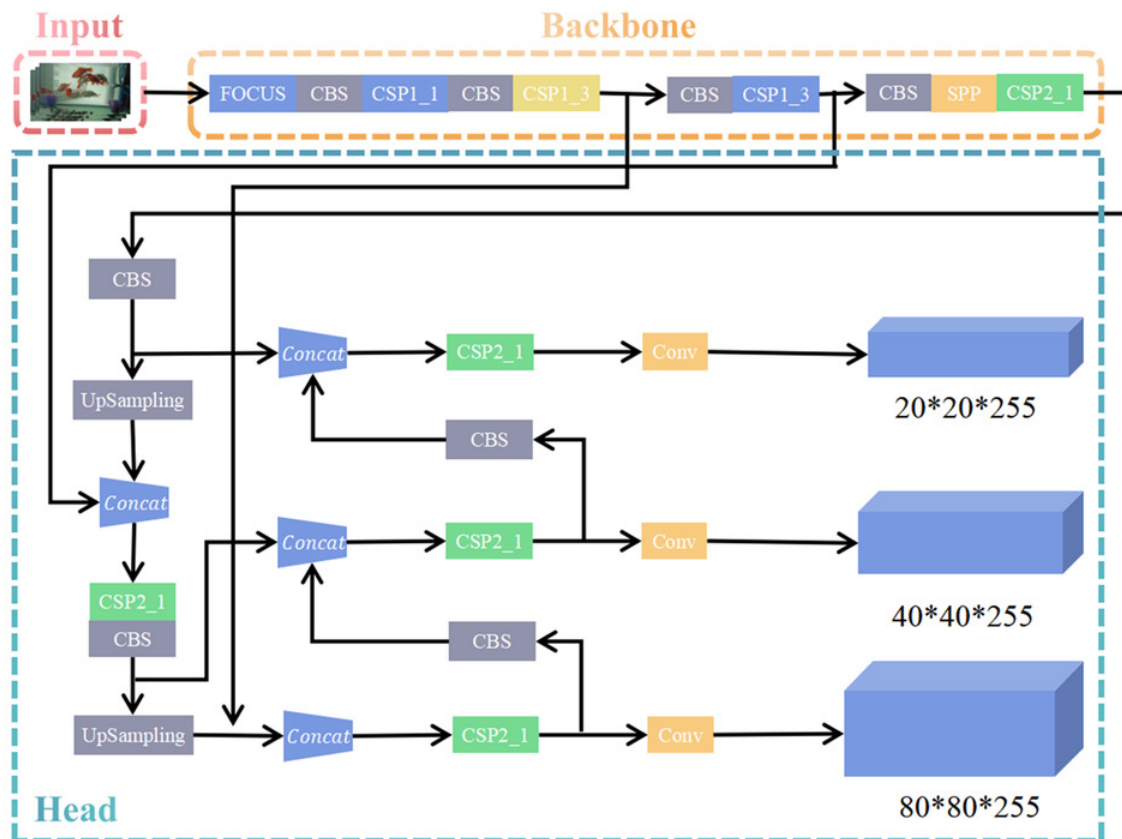


Figure 6

The network structure and application of YOLOv5s-5.0.

The CBS component in the figure is composed of the Convolutional layer + BatchNormalization + SiLU activation function. The Resunit component draws on the residual structure in the Resnet network and can play a role in building a deeper network. The CSP1_X and CSP2_X component draws on the CSPNet network structure. The FOCUS component is to slice the data, which can play a role in the down-sampling operation without information loss. The SPP component adopts the maximum pooling method of 1×1 , 5×5 , 9×9 , and 13×13 for multi-scale fusion.



Basic components

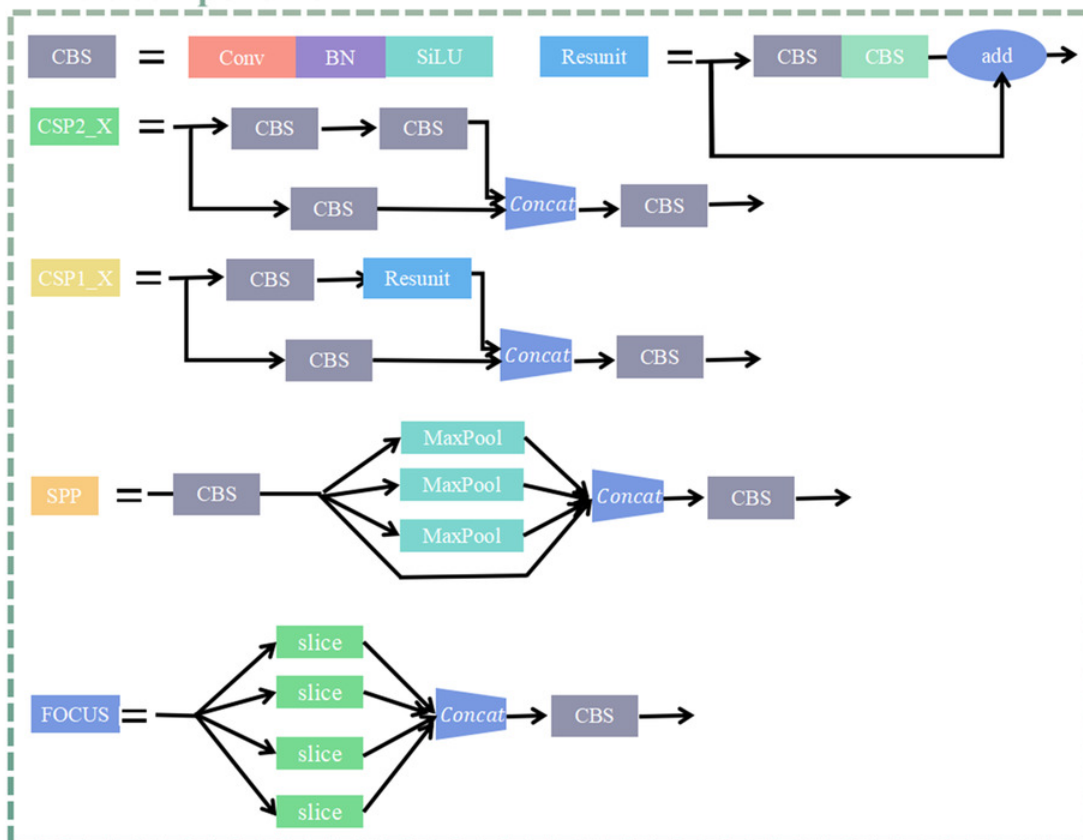


Figure 7

Structure of segmentation head.

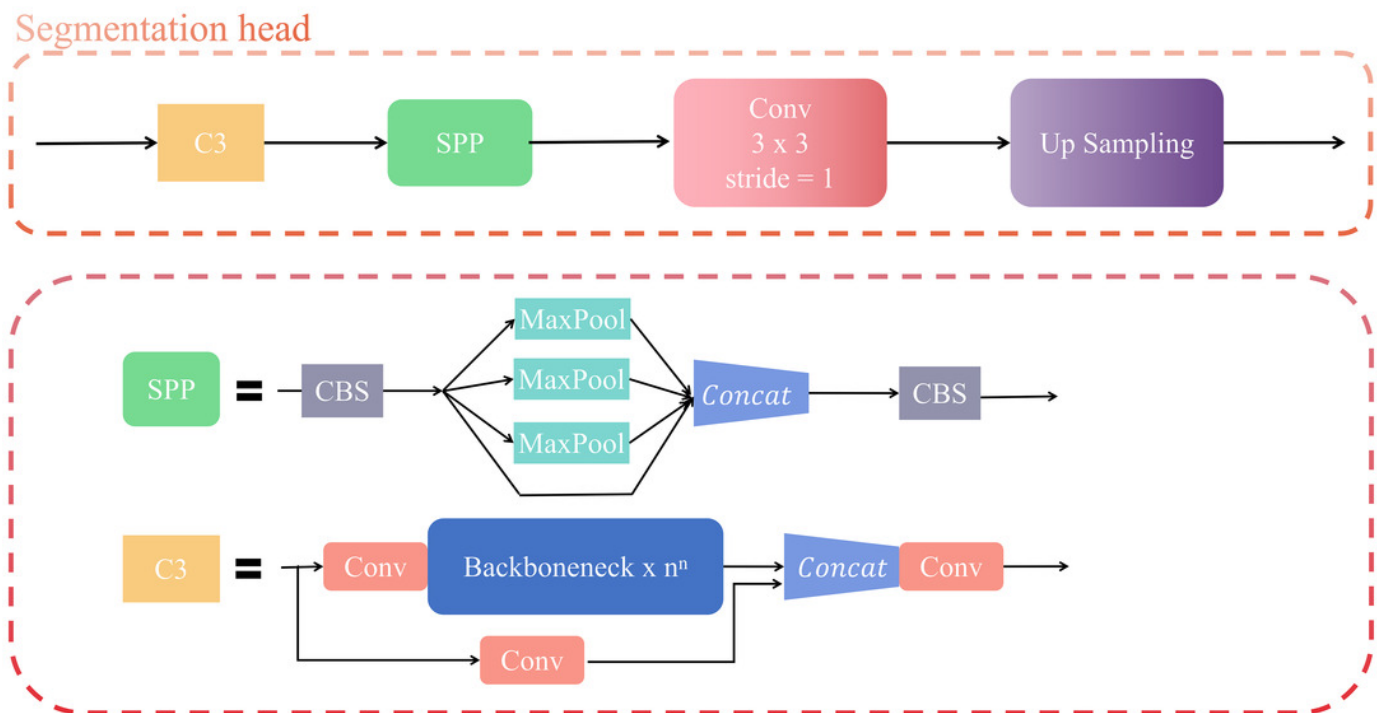


Figure 8

Example diagram of bilinear interpolation algorithm.

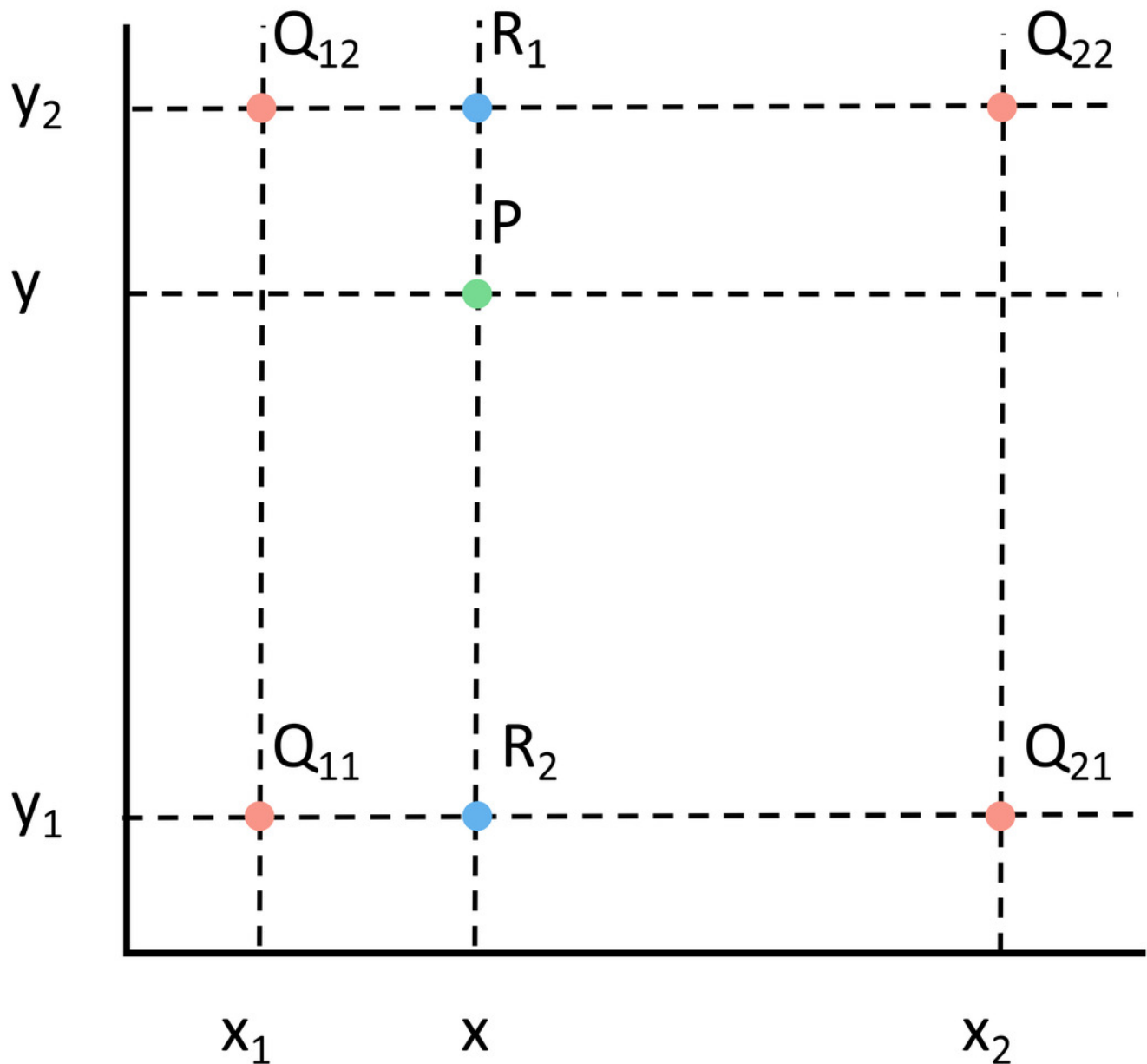


Figure 9

The structure of GhostC3 module.

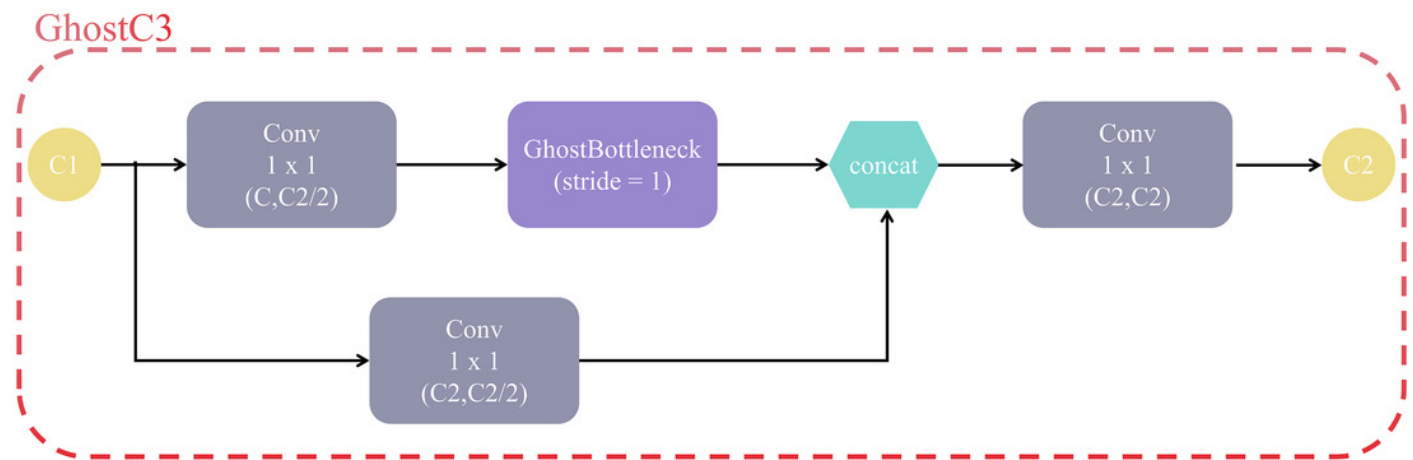


Figure 10

The structure of GhostConv module.

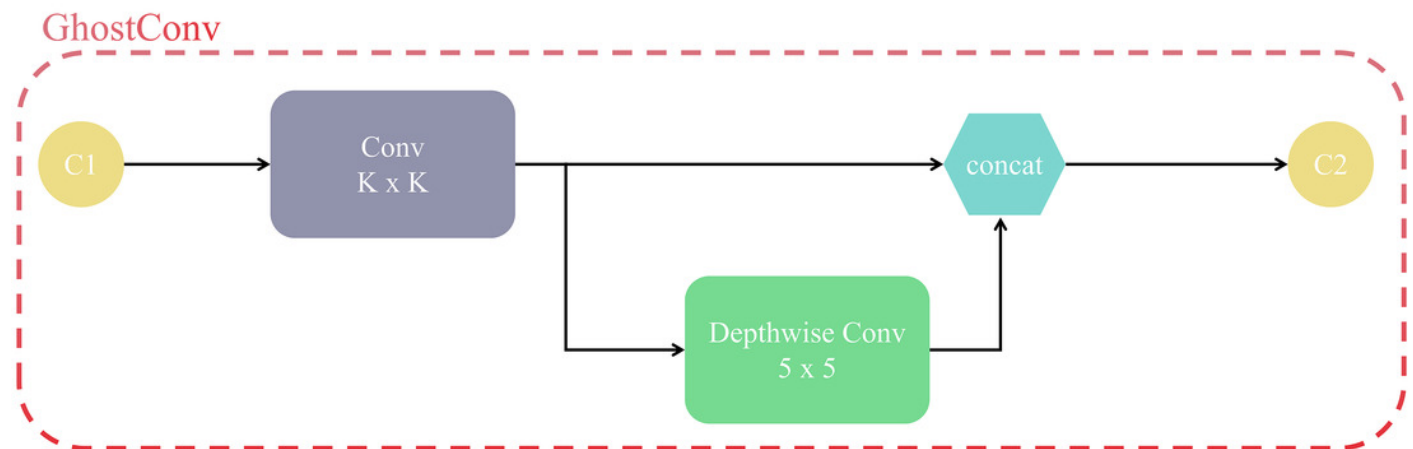


Figure 11

The structure of GhostBottleneck(stride=1).

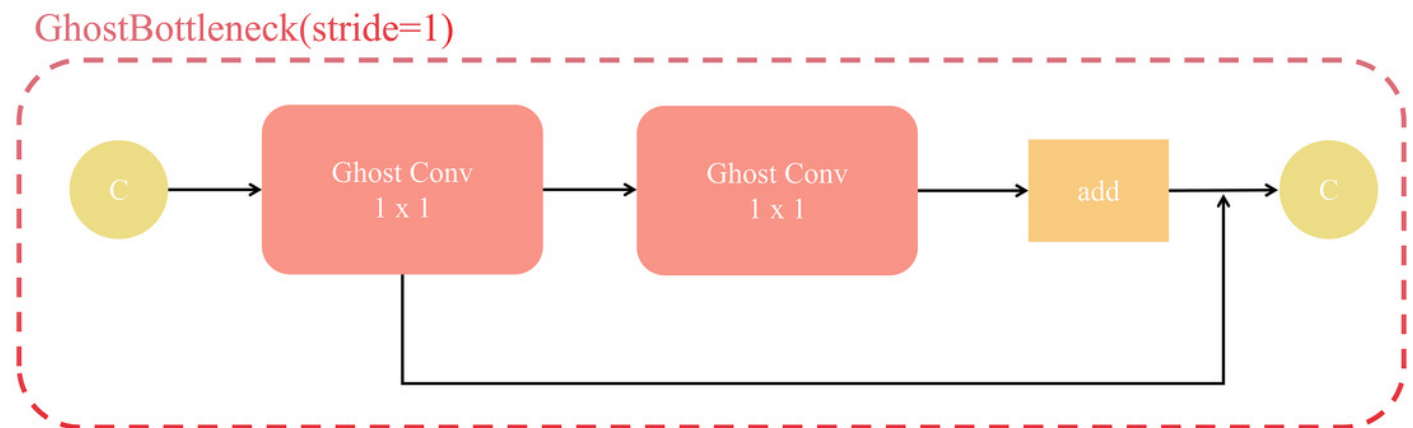


Figure 12

The structure of GhostBottleneck(stride=2).

GhostBottleneck(stride=2)

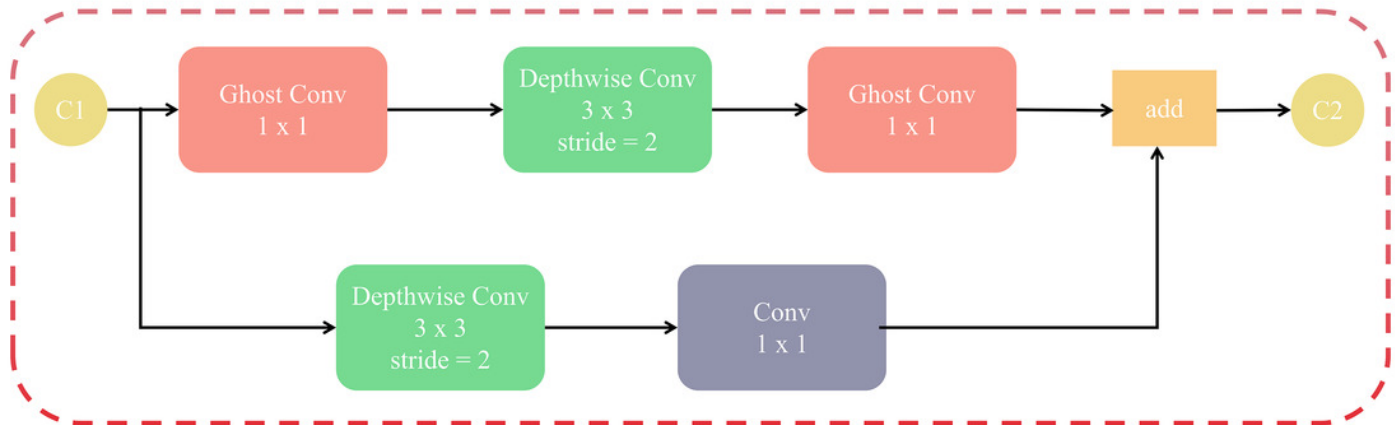


Figure 13

Swish function diagram.

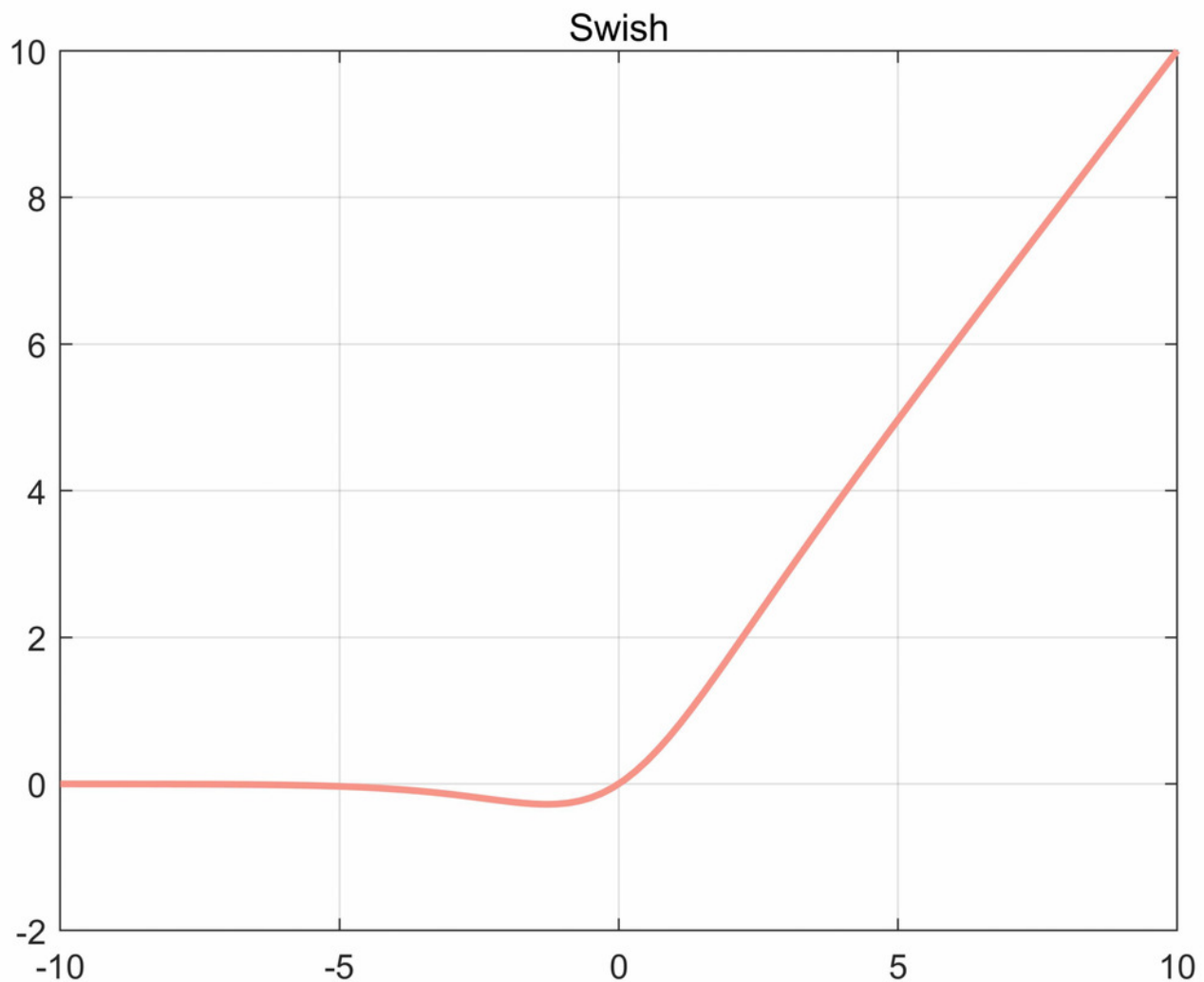


Figure 14

HardSwish function diagram.

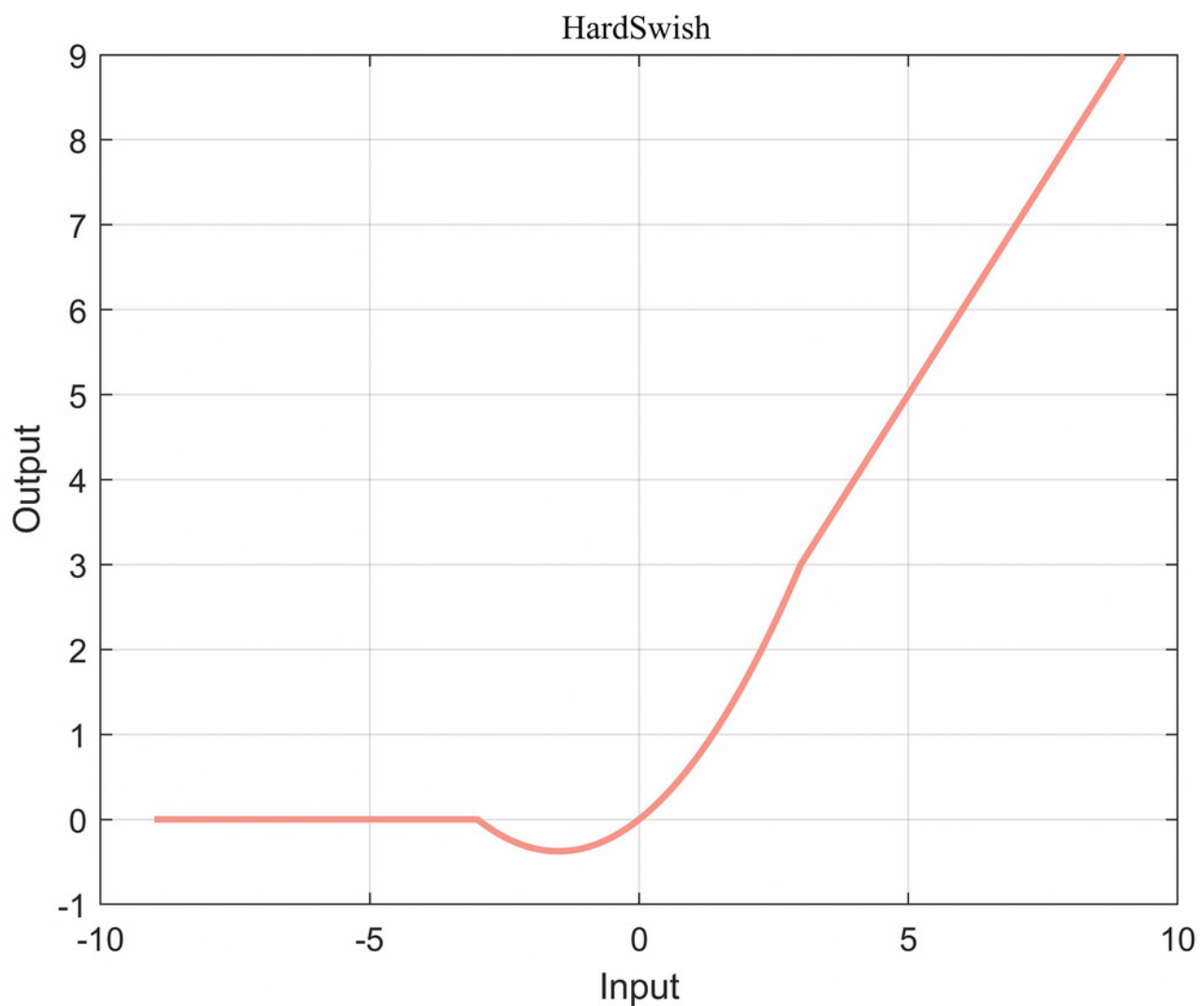


Figure 15

Visual representation of IoU.

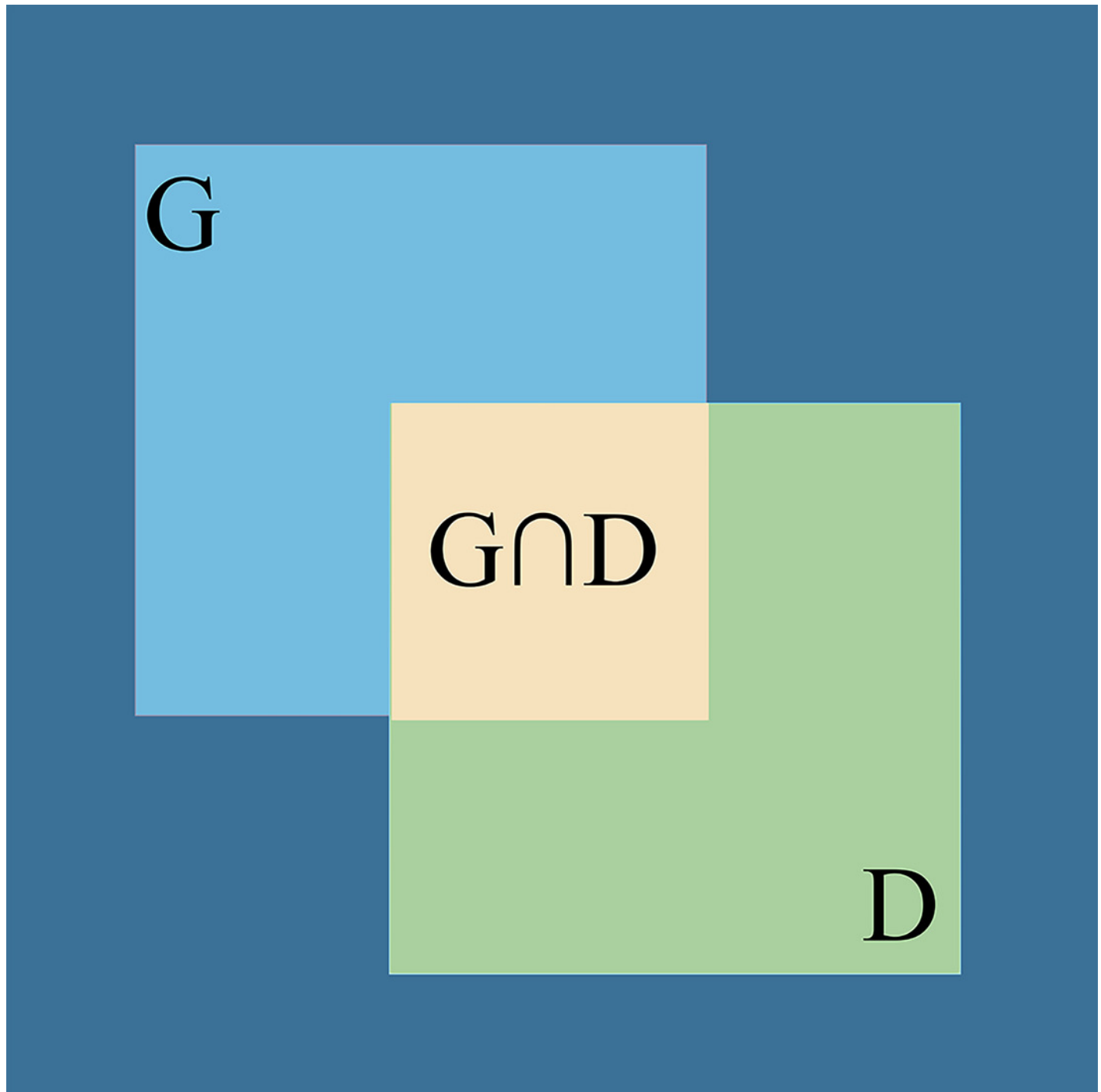


Figure 16

The results of our model for detecting and segmenting the golden crucian carp.

