# JAST

# A study of duck detection using deep neural network based on RetinaNet model in smart farming

Jeyoung Lee, Hochul Kang*

*Department of Digital Media, The Catholic University of Korea, Bucheon 14662, Korea*

## Abstract

In a duck cage, ducks are placed in various states. In particular, if a duck is overturned and falls or dies, it will adversely affect the growing environment. In order to prevent the foregoing, it was necessary to continuously manage the cage for duck growth. This study proposes a method using an object detection algorithm to improve the foregoing. Object detection refers to the work to perform classification and localization of all objects present in the image when an input image is given. To use an object detection algorithm in a duck cage, data to be used for learning should be made and the data should be augmented to secure enough data to learn from. In addition, the time required for object detection and the accuracy of object detection are important. The study collected, processed, and augmented image data for a total of two years in 2021 and 2022 from the duck cage. Based on the objects that must be detected, the data collected as such were divided at a ratio of 9 : 1, and learning and verification were performed. The final results were visually confirmed using images different from the images used for learning. The proposed method is expected to be used for minimizing human resources in the growing process in duck cages and making the duck cages into smart farms.

**Keywords:** Duck detection, Duck farming, Smart farming, Object detection, Deep neural network, Computer vision

**ORCID**
Jeyoung Lee
https://orcid.org/0000-0002-1464-7839
Hochul Kang
https://orcid.org/0000-0002-7733-2287

## INTRODUCTION

In a duck cage, ducks are placed in various states. In particular, if a duck is overturned and falls or a duck is dead during growth, a person must make the duck stand up or collect the duck. To that end, it was necessary for humans to continuously manage the cage during the growing process of ducks. In order to improve the foregoing, this study proposes a method to use an object detection algorithm to utilize a robot in a duck cage to observe ducks to check if any duck fell or died and make any duck fell stand up and collect any duck dead. According to Zaidi et al. [1], object detection means the work to classifying and localize all objects present in the image when an input image is given. Object detection algorithms can be largely divided into one-stage methods and two-stage methods, and each method has advantages and disadvantages. The one-stage method is faster but less accurate. Data are necessary to train artificial intelligence (AI) algorithms. In particular, a lot of processed data is required to use an object detection algorithm. However, there is no processed public data about the state of ducks in a duck cage environment. Therefore, in order to detect objects in the duck cage, it was necessary to

**Availability of data and material**
Upon reasonable request, the datasets of this study can be available from the corresponding author.

**Authors' contributions**
Conceptualization: Lee J, Kang H.
Data curation: Lee J, Kang H.
Formal analysis: Lee J, Kang H
Methodology: Lee J, Kang H.
Software: Lee J, Kang H.
Validation: Lee J, Kang H.
Investigation: Lee J, Kang H.
Writing - original draft: Lee J, Kang H.
Writing - review & editing: Lee J, Kang H.

**Ethics approval and consent to participate**
This article does not require IRB/IACUC approval because there are no human and animal participants.

firsthand collect, process, and augment data. This study collected, processed, and augmented image data from a duck cage for a total of two years of 2021 and 2022. The data collected as such will be discussed again in Materials and Methods. Finally, among the one-stage algorithms, RetinaNet [2] was used for learning and experiment. Unlike published data, data collected firsthand have many limitations. In particular, problems of the limited number of data and the imbalance of the correct answer to the data often occur. RetinaNet [2] is the most common algorithm that enables solving the imbalance problem of correct answers in collected data. By utilizing RetinaNet, it is possible to solve the bias of learning models created by the problems of imbalance of correct answers in data caused by relatively insufficient data collection.

This study is closely related to object detection in smart farms. Gikunda and Jouandeau [3] and Dhanya et al. [4] collected and investigated cases where artificial intelligence was used in relation to smart farms. Dhanya et al. [4] state that the agricultural industry is going through a process of rapid digital transformation and that technology is being made more powerful by state-of-the-art approaches such as artificial intelligence technology. Sa et al. [5] proposes a DeepFruits model that finds about five kinds of fruits, such as sweet pepper and rockmelon, in a greenhouse using Faster Regions with Convolutional Neural Networks (R-CNN) [6]. Bargoti and Underwood [7] propose a method for finding apples, mangos, and almonds in an orchard by applying the DeepFruits [5] network. Sørensen et al. [8] propose a method for finding thistles that cause loss in crop yield using DenseNet [9] based on aerial photographs of crops. Albuquerque et al. [10] studies a method for identifying water in a watering machine based on Mask R-CNN [11] in image frames captured by an unmanned aerial vehicle (UAV). Osorio et al. [12] compared and analyzed Mask R-CNN [11], SVMs [13], and YOLOv3 [14] for methods to detect weeds in lettuce crops. Riekert et al. [15] conducted a study on a method to find a pig's position using Faster R-CNN [6]. Tedesco-Oliveira et al. [16] applied Faster R-CNN [6] and SSD [17] to study the development of an automated system for predicting cotton yields from color images acquired with a simple mobile device. Zhou et al. [18] compared various back-bone networks of SSD [17] to conduct a study on a method to find kiwi fruit in real time. Tang et al. [19] propose a method of applying object detection to detect the distribution and precise shape of center pivot irrigation systems. Shojaeipour et al. [20] applied two-stage YOLOv3 [14]-ResNet50 [21] to study a method for detecting the mouth region of a cow from a cow face image dataset for livestock welfare and management. Syed-Ab-Rahman et al. [22] propose an end-to-end anchor-based model to detect and classify citrus disease states.

Based on this, our paper analyzes the method of directly collecting, processing, and augmenting data for object detection on the state of ducks in a duck cage, and the application and the results of application of object detection algorithms. In order to check whether learning is successfully carried out using the collected data, the data are divided at a ratio of 9:1 based on the objects that must be detected and are learned and verified. As for the evaluation, the average precision (AP) is measured using the separated data for evaluation, and the final result is visually checked using images different from the images used for learning. The proposed method is expected to be used for minimizing human resources in the growing process in duck cages and making smart farms.

## MATERIALS AND METHODS

### Data collection

Data collection and generation is one of the most important and time-consuming tasks in any field of artificial intelligence. In this study, the data necessary for object detection are largely the video data of ducks in the duck cage, the bounding boxes that specify the locations of ducks by image frame, and the state class labels. However, there are no studies similar to this or it is not a common

situation. That is, there is no public data. Therefore, this study proceeds from the data collection stage. When raising ducks in duck cages, ducks are not raised from eggs. Generally, baby ducks hatched from eggs are brought to a duck cage and raised, and all are delivered after a certain age. This is a characteristic of broiler ducks, and because of this characteristic, it is difficult to secure a large amount of data. However, deep learning requires a large amount of data in various types. To solve this problem, this study received image data directly from the duck cage over two years, 2021 and 2022, and uses techniques such as data augmentation. When receiving video data, the main point of view is whether the video has an appropriate height that can be used in real situations and whether the duck states are sufficiently diverse. An example of the video data provided is as shown in Figs. 1 and 2.
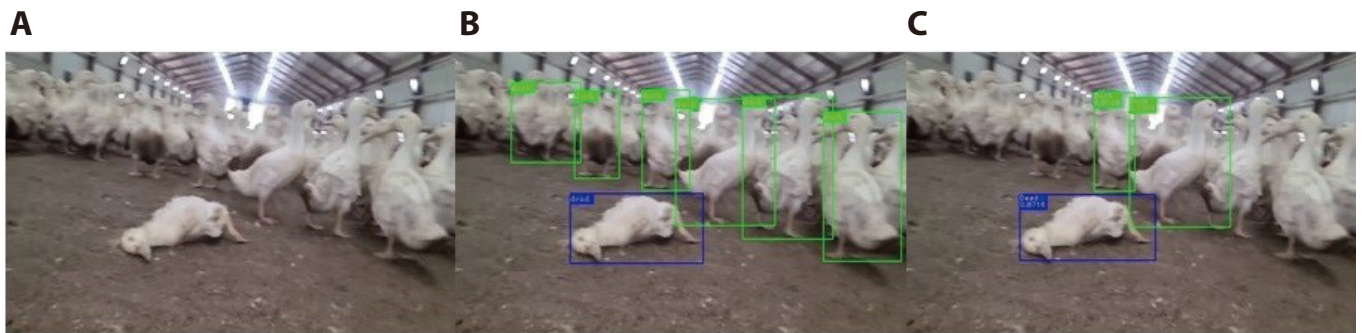


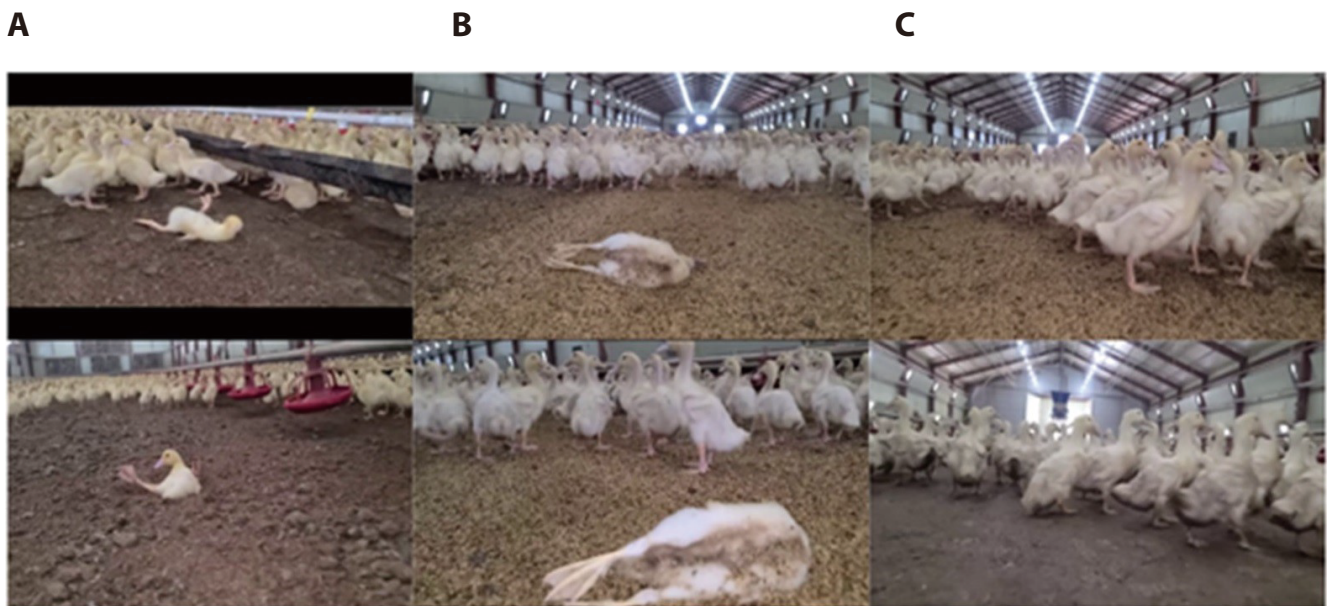**Fig. 1. An example of the detection result from the video.** (A) input Image, (B) ground truth, (C) our detection result.



**Fig. 2. An example of the state of the duck in images.** (A) slap, (B) dead, (C) normal.

## Data labeling

The training images are extracted from the video as frames at the duck farm in 2021, and the bounding boxing and class labeling are carried out directly by human hands. There are three states where ducks can exist in the image: normal, fallen, and dead. In this case, as the length of the video increases, the number of frames becomes too large. As a result, the differences between the images between the frames of the video are not large, and as the video moves, frames where it is difficult to recognize the shapes of the ducks occur. In addition, when humans firsthand create labels, as the number of images increases, the problem of taking longer time also occurs. That is, taking and using all image frames is not good for learning and only increases the data generation time. In order to solve this problem, this study selected only one image per 5 to 10 frames, and labeled the 1,285 images selected as such first. Duck cages raise large numbers of ducks. Therefore, when labeling an image for object detection, there is a problem that the number of ducks is excessively large, and ducks are dense. To solve this problem, it is necessary to clarify criteria when creating labels and to establish common rules. In this study, labels are created based on the duck in the frontmost of the image. In addition, only those ducks whose face, body, tail, and feet are clearly identified are identified in the normal state. The characteristics of the dataset created are examined with the labels and images created with the rule. Some problems were found due to the labeling results of the 2021 data. The ratios of dead ducks and fallen ducks in the data are overwhelmingly insufficient. This study solves this problem in three methods. First, we added more data which is provided in 2022 for improving the performance of the detection, and apply it to train. Second, we solved the problem by augmenting insufficient data using a data augmentation technique. Finally, the focal loss proposed in RetinaNet [2] is used. Focal loss was proposed to solve the class imbalance problem. The problem that humans firsthand carry out labeling one by one occurs. If labeling is carried out by humans, there is the problem that a long time is taken, and the stability of the label cannot be guaranteed. To solve the foregoing problems, the object detection model was first trained using the 2021 data. Thereafter, using the model, an automatic labeling program was created. An example of automatic labeling program is shown in Fig. 3. Based on the program, the 2022 duck cage image data provided later were extracted by image frame, and thereafter, labeling was carried out first using an automatic labeling program. Finally, the labeling was inspected and corrected by humans to save time and improve stability. As such, 2,852 images and labels were finally created. An example of a label created as such is shown in Fig. 4.

## Dataset

The number of data sets finally created is 2,852. The average size of the image is 1,748.30 and 999.94 for the width and height, respectively, and the total numbers of normal ducks, fallen ducks, and dead ducks in all images are 10,461, 1,208, and 381, respectively. The maximum number of normal ducks, fallen ducks, and dead ducks in one image is 24, 1, and 1, respectively. Ducks in all states may or may not exist. Also, ducks in various states may appear simultaneously. The ratios of one duck object to image are 0.056, 0.053, and 0.082, respectively. Ducks in most states appear evenly throughout the image, but dead ducks always appear below the halfway of the image (Table 1).

## RetinaNet training

The purpose of this study is to find duck objects in the duck cage in real time. There are many similar object detection algorithms. However, as a characteristic of the collected datasets, the ratio of fallen ducks and dead ducks is overwhelmingly lower than that of normal ducks. This problem is called the state imbalance problem. To solve this problem, this study uses RetinaNet [2]. RetinaNet [2] has the advantage that the backbone model and the region proposal network can be freely
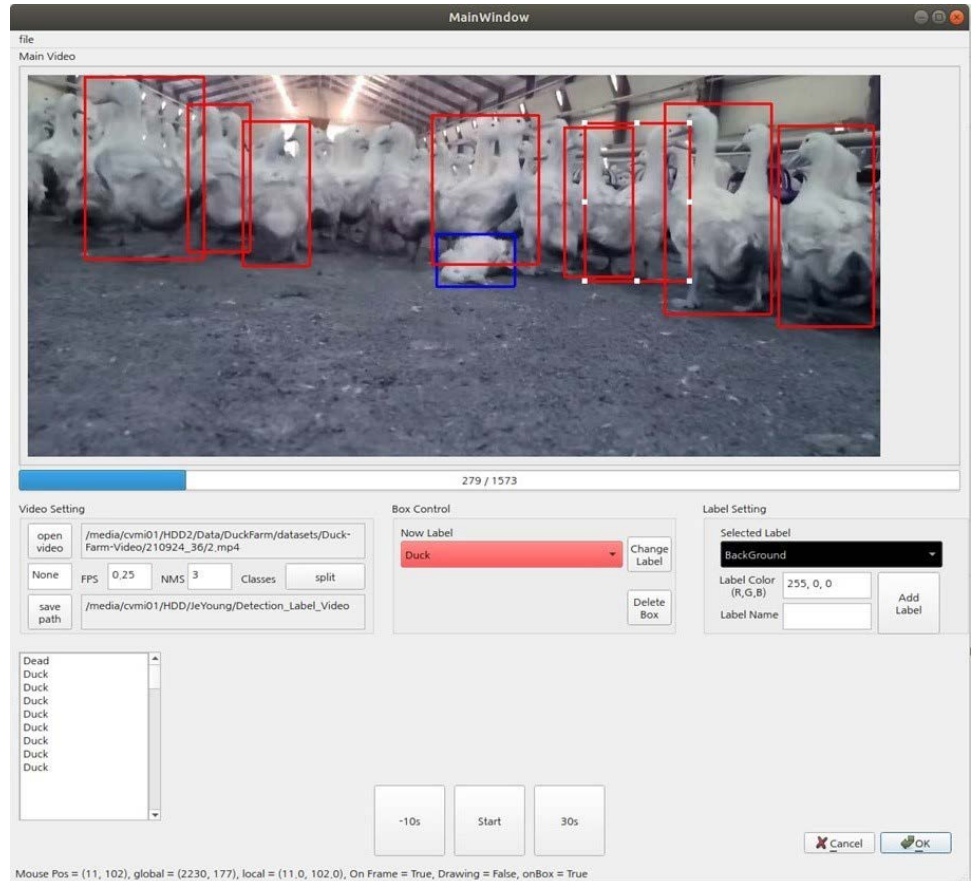
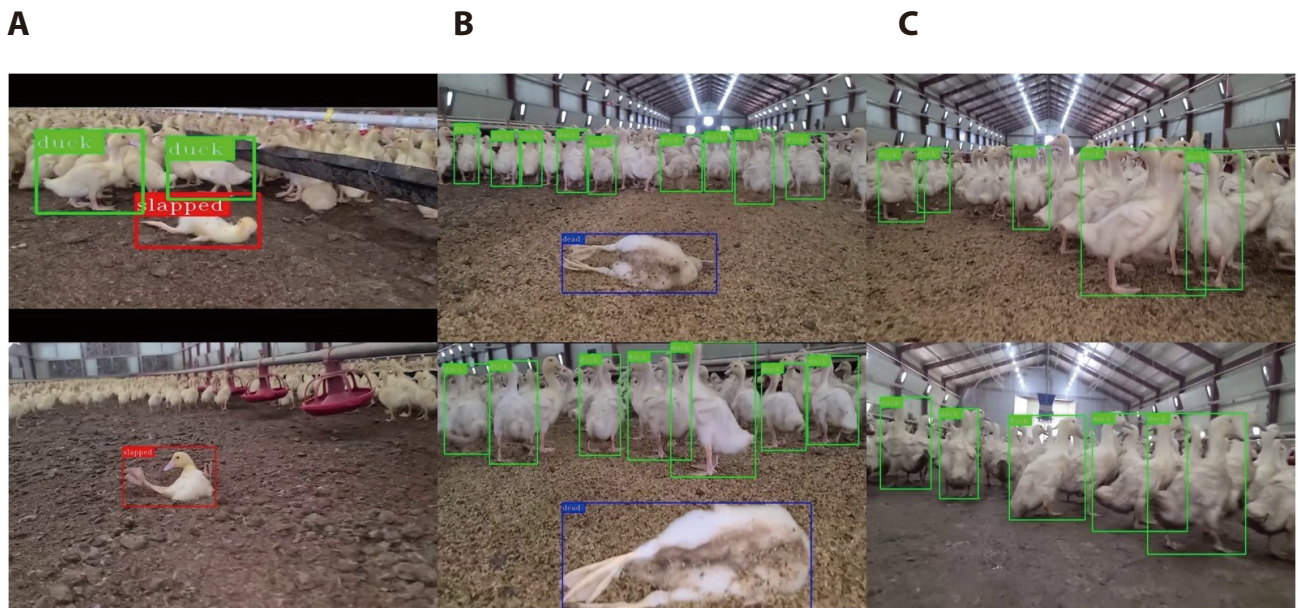**Fig. 3.** A screenshot of our auto labeling program.



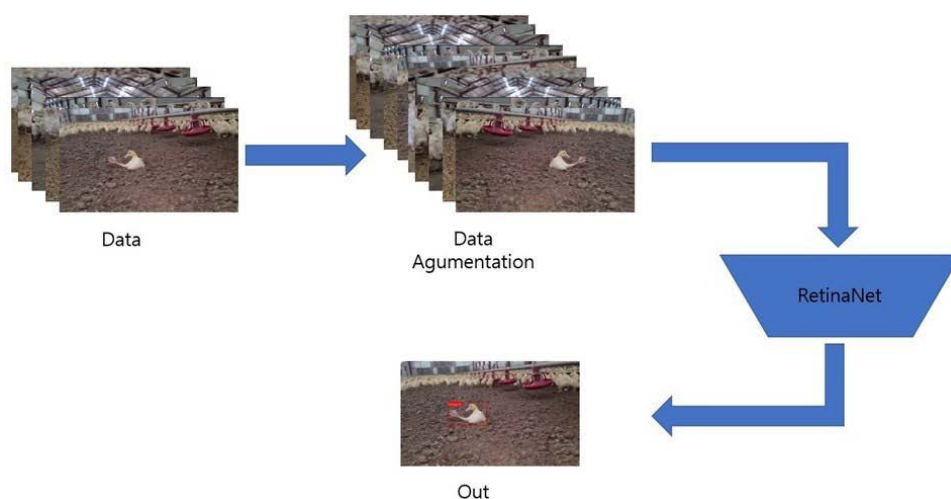**Fig. 4.** An example of labeled images. (A) slap, (B) dead, (C) normal.

**Table 1. Dataset Information**

| | Total number | Maximum number | Minimum number | Average region rate | Minimum top left x | Minimum top left y | Maximum top left x | Maximum top left y |
|---|---|---|---|---|---|---|---|---|
| Duck | 10,461 | 24 | 0 | 0.0563 | 0.00 | 0.00 | 1,818.65 | 896.53 |
| Slap | 1,208 | 1 | 0 | 0.0531 | 0.00 | 0.00 | 1,380.64 | 850.54 |
| Dead | 381 | 1 | 0 | 0.0825 | 0.00 | 97.48 | 1,611.73 | 832.36 |

changed. In addition, it is easy to apply new datasets because many studies have been conducted. Furthermore, the introduction of the focal loss solves the problem of state imbalance to some extent. The focal loss is an extended version of the cross entropy loss that reduces the weights of easy examples and focuses learning on difficult examples. Finally, real-time object detection is possible because it is a one-stage model. Therefore, RetinaNet [2] is used as the basic model of this study. A figure of the learning pipeline using RetinaNet [2] is as shown in Fig. 5.

### Data augmentation

The more the data used in deep learning, the better the deep learning. However, the total number of data used in this study is 2852. Many studies try to obtain more data for learning. However, when it is difficult to secure additional data, data are increased through data augmentation. This study augments data before using the data for learning. The techniques used in that case are brightness conversion, contrast conversion, saturation conversion, rotation, random resize, and flip. For brightness, contrast, and saturation conversions, values between 0.9 and 1.1 are randomly applied based on the image value. In the case of rotation, values between –20 degrees and 20 degrees are applied according to the characteristics of the image. Flip is applied left and right, and the application probability is 0.5. For random resize, a length of one of 640, 672, 704, 736, 768, and 800 is selected based on the length of the shortest side, and the length of the longest side is up to 1,333. Finally, each technique is applied independently of the other. That is, several techniques may be applied at the same time, or none may be applied. Fig. 6 is an example of an image to which augmentation was applied.



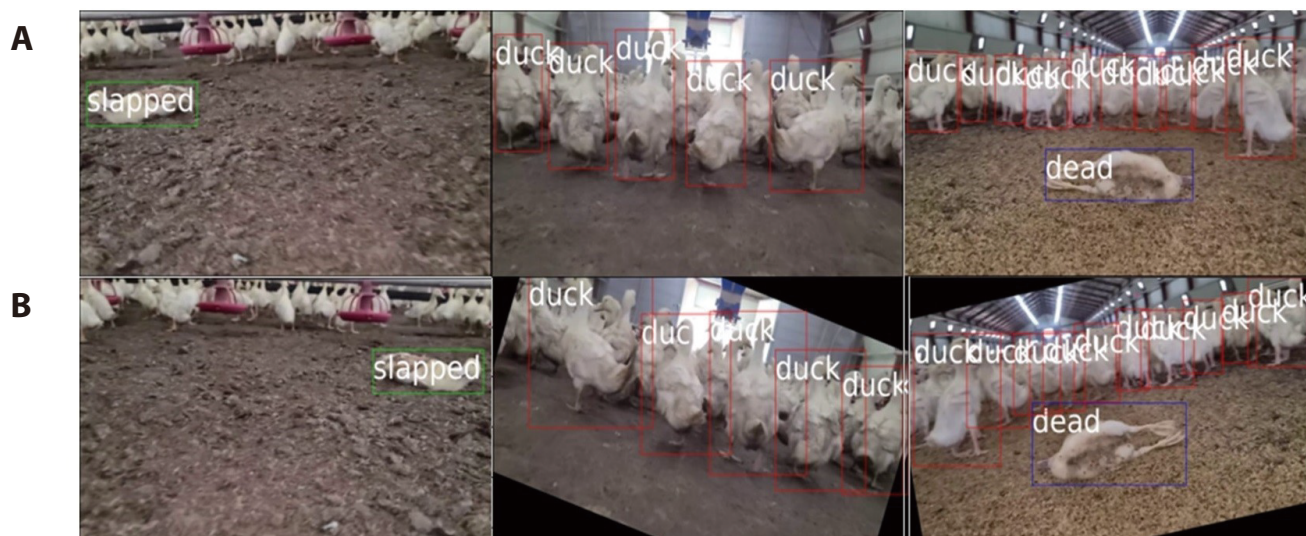**Fig. 5. The training pipeline of our duck detection algorithm.**

**Fig. 6. An example of data augmentation.** (A) original images, (B) augmented images.

## Fine tuning

Fine-tuning is a method used to train one's own model based on an existing model that has been trained. Many deep learning approaches use fine-tuning to achieve a task. In this study too, the RetinaNet [2] model pretrained using the Common Objects in Context (COCO) dataset is fined-tuned and learned. There are two models prepared for fine-tuning, 1x model and 3x model, which will be used depending on the schedule. Heet al. [23] questioned fine-tuning and studied a new way of learning. They introduce training scheduling techniques, batch normalization, and methods that do not use fine-tuning. According to them, a learning schedule to search the COCO Dataset once based on the COCO Dataset is defined as a 1x schedule. That is, the prepared 1x pretraining model means a model that searches the COCO dataset once, carries out 90,000 iterations, and has learning rates reduced to 1/10 at 60 k and 80 k. The 3x pretraining model is a model that searches the COCO dataset twice, caries out 270,000 iterations, and has the learning rate reduced to 1/10 every 210 k and 250 k. In this study, both models are used for learning and the results are compared thereafter.

## Train details

For learning and validation, the data are divided into train data and validation data at a ratio of 9:1. When dividing the data, the data are divided based on classes so that the data can be divided fairly by class. In addition, a total of three models are learned: a model to which data augmentation was not applied, a model to which data augmentation was partially applied, and a model to which data augmentation was fully applied. As for the model to which data augmentation was partially applied, it was found that the model to which only random resize and random flip were applied as elements found during learning performed better. Details can be found in Result Section. The basic RetinaNet [2] used in learning is a combination of ResNet50 [21] and Feature Pyramid Network (FPN) [24]. In addition, two models trained on the COCO dataset were prepared. We fine-tune from the two prepared models. In this case, focal loss is used as the loss and Stochastic Gradient Descent (SGD) is used as the optimizer. The basic learning rate is 1e-3, and the warm-up scheduler and the step scheduler are used as the learning schedulers. Therefore, the learning rate is

first warmed up to 1,000 iterations. The step scheduler reduces the basic learning rate by 1e-1 each at the last iterations, 5,000 and 6,000 iterations. The batch size is 16 and the iteration is 7,000. One RTX 3090 was used for learning, and the time taken for the learning was about 2 hours.

## RESULTS

The most commonly used value to measure performance in object detection is AP. In short, AP means the percentage of correct answers in the predicted boxes. AP is again divided into AP50, AP75, etc. according to the ratio of intersection over union (IoU) according to the degree of overlap between the predicted box and the correct answer box. AP means the average accuracy measurement method for all ratios of IoU, which increases by 0.05 from 0.5 to 0.95, AP50 means when IoU is greater than 0.5, and AP75 means when IoU is greater than 0.75. In this study, how accurate the combination of basic ResNet50 [21] and FPN [24] is checked for each AP according to the pretraining model and whether augmentation is carried out. Table 2 is a table of measurement of AP for 270 pieces of validation data. Table 3 is the result of measurement of AP by class for the same validation data.

According to Tables 2 and 3, it can be seen that the performance of the 3x model is basically higher than that of the 1x model. In addition, the performance of the model to which only random resize and flipping were applied is superior to that of the model to which full augmentation was applied for validation data. It can be seen that excessive augmentation does not help the validation performance because the number of validation data is small, and the images are mainly those images with angles and shapes similar to those of the learning images. However, this is far from generalization, which is the goal of learning. Therefore, the validation data are augmented through flipping and rotation to generate 2,770 validation data, and more general performance is measured thereafter. The results are in Tables 4 and 5 below.

**Table 2.** Duck detection RetinaNet result

| Backbone | Scheduler | Augmentation | AP | AP50 | AP75 |
|---|---|---|---|---|---|
| Resnet50-FPN | 1x | None | 73.969 | 97.035 | 87.633 |
| Resnet50-FPN | 3x | None | 74.630 | 97.046 | 88.686 |
| Resnet50-FPN | 1x | Part | 79.599 | 98.060 | 91.569 |
| Resnet50-FPN | 3x | Part | 79.797 | 98.023 | 91.569 |
| Resnet50-FPN | 1x | All | 66.286 | 97.788 | 81.559 |
| Resnet50-FPN | 3x | All | 67.101 | 97.711 | 84.954 |

AP, average precision; FPN, Feature Pyramid Network.

**Table 3.** Duck detection RetinaNet result by class

| Backbone | Scheduler | Augmentation | Duck | Slap | Dead |
|---|---|---|---|---|---|
| Resnet50-FPN | 1x | None | 62.291 | 76.794 | 82.821 |
| Resnet50-FPN | 3x | None | 61.985 | 79.549 | 82.357 |
| Resnet50-FPN | 1x | Part | 68.187 | 84.082 | 86.527 |
| Resnet50-FPN | 3x | Part | 68.467 | 85.362 | 85.563 |
| Resnet50-FPN | 1x | All | 58.852 | 72.208 | 67.797 |
| Resnet50-FPN | 3x | All | 59.518 | 72.910 | 68.876 |

AP, average precision; FPN, Feature Pyramid Network.

**Table 4.** Duck detection augmentation validation data RetinaNet result

| Backbone | Scheduler | Augmentation | AP | AP50 | AP75 |
|---|---|---|---|---|---|
| Resnet50-FPN | 1x | None | 34.413 | 86.847 | 16.997 |
| Resnet50-FPN | 3x | None | 33.917 | 86.609 | 16.562 |
| Resnet50-FPN | 1x | Part | 37.432 | 91.314 | 20.255 |
| Resnet50-FPN | 3x | Part | 37.340 | 90.682 | 19.787 |
| Resnet50-FPN | 1x | All | 70.984 | 97.182 | 88.584 |
| Resnet50-FPN | 3x | All | 70.784 | 97.361 | 89.745 |

AP, average precision; FPN, Feature Pyramid Network.

**Table 5.** Duck detection augmentation validation data RetinaNet result by class

| Backbone | Scheduler | Augmentation | Duck | Slap | Dead |
|---|---|---|---|---|---|
| Resnet50-FPN | 1x | None | 32.513 | 38.990 | 31.737 |
| Resnet50-FPN | 3x | None | 32.061 | 37.264 | 32.426 |
| Resnet50-FPN | 1x | Part | 37.408 | 41.936 | 32.953 |
| Resnet50-FPN | 3x | Part | 37.499 | 41.278 | 33.244 |
| Resnet50-FPN | 1x | All | 62.786 | 76.781 | 73.386 |
| Resnet50-FPN | 3x | All | 64.474 | 76.871 | 71.007 |

FPN, Feature Pyramid Network.

Through the results in Table 4 and Table 5, it can be seen that the generalization performance of the model to which full augmentation was applied is better. Therefore, in this study, the test is conducted using a model to which full augmentation is applied. In addition, between the 1x model and the 3x model, the 3x model generally has better performance. However, in the present evaluation, the average AP performance of the 1x model was shown to be better. Since the AP75 performance of the 3x model was better, the 3x model was used and applied to images different from the images used for learning and evaluation. Because the images to which the models were applied as such have no information of the actual objects, it was checked with eyes whether the images were searched well. The results checked with the eyes are as shown in Figs. 7, 8, and 9.

In addition, the average inference time per one image for all models is within 0.003 seconds. This shows that the inference time of this model is short and effective. Therefore, the model can be used for real-time detection.

## DISCUSSION

This study collected and defined anomalous object detection datasets for making a smart farm for anomalous duck detection in a duck cage environment. Thereafter, using the datasets, learning and evaluation were caried out utilizing RetinaNet, a one-stage network. Finally, for good results, image augmentation, warm-up scheduler, etc. were used for comparison to explore the best algorithm between basic ResNet50 and FPN models. The datasets defined through the foregoing were shown to be usable and basic model guidelines were established. However, there are some limitations. First, the backbone network was not changed. In the case of object detection, the performance varies greatly depending on the size of the backbone network and the method of the region of interest network. If the size of the backbone model is increased, the accuracy will increase. However, due to the definition of the problem that objects should be detected in real time, a search process to find a network of an appropriate size is necessary. Second, a method that uses an object detection

**Fig. 7. The result of the detection for normal ducks.** (A) input image, (B) ground truth, (C) our detection result.



**Fig. 8. The result of the detection for slapped ducks.** (A) input image, (B) ground truth, (C) our detection result.
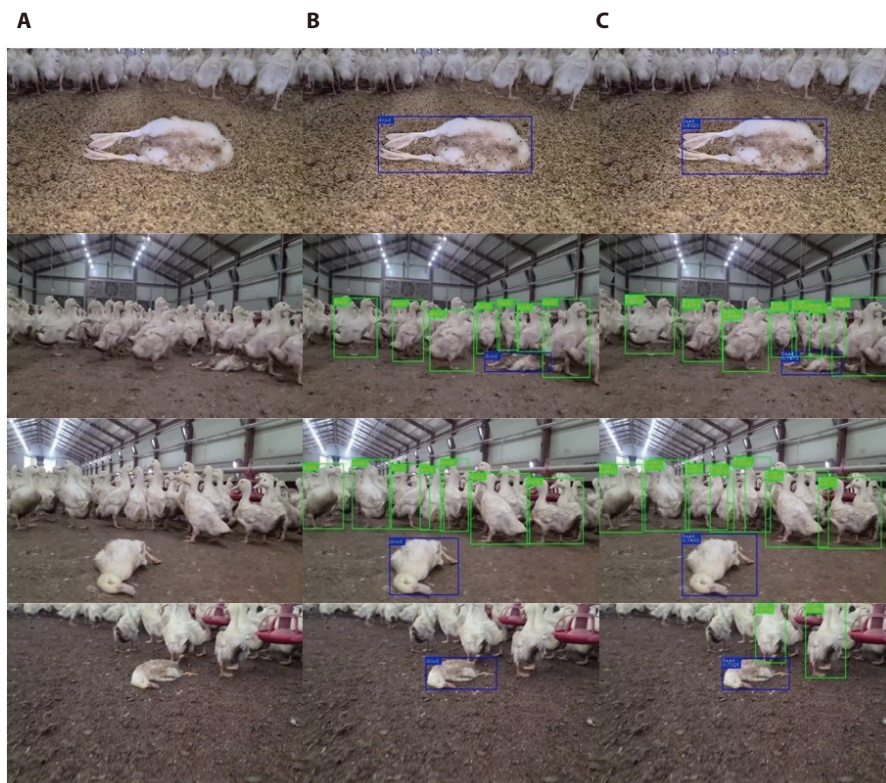
**Fig. 9. The result of the detection for dead ducks.** (A) input image, (B) ground truth, (C) our detection result.

model other than RetinaNet is necessary. RetinaNet is a network that has been studied a lot and has characteristics suitable for solving our problems, but it is also an old model. This means that experiments should be caried out on other models that advanced RetinaNet while retaining the features. Finally, research on the improvement of a new network tailored to the datasets is needed. Currently, we applied our datasets based on a famous model and focused on exploring how well it performs. A study like this is also a study, and through this, we showed that our problem definition is solvable and that our datasets can be used well in a general model. However, this does not mean that general models published well fit our datasets. Research on new models that fit the characteristics of our datasets is also needed. All of these limitations will be addressed in the future based on this study by utilizing and developing the insights found in this study.

## REFERENCES

1. Zaidi SSA, Ansari MS, Aslam A, Kanwal N, Asghar M, Lee B. A survey of modern deep learning based object detection models. Digit Signal Process. 2022;126:103514. https://doi.org/10.1016/j.dsp.2022.103514
2. Lin TY, Goyal P, Girshick R, He K, Dollár P. Focal loss for dense object detection. In: Proceedings of the IEEE International Conference on Computer Vision; 2017; Venice, Italy. p. 2980-8.
3. Gikunda PK, Jouandeau N. State-of-the-art convolutional neural networks for smart farms: a review. In: Arai K, Bhatia R, Kapoor S, editors. Intelligent Computing-Proceedings of the Computing Conference. Chanm: Springer; 2019 p. 763-75.

4. Dhanya VG, Subeesh A, Kushwaha N, Vishwakarma DK, Kumar TN, Ritika G, et al. Deep learning based computer vision approaches for smart agricultural applications. Artif Intell Agric. 2022;6:211-29. https://doi.org/10.1016/j.aiia.2022.09.007

5. Sa I, Ge Z, Dayoub F, Upcroft B, Perez T, McCool C. Deepfruits: a fruit detection system using deep neural networks. Sensors. 2016;16:1222. https://doi.org/10.3390/s16081222

6. Ren S, He K, Girshick R, Sun J. Faster R-CNN: towards real-time object detection with region proposal networks. In: Advances in neural information processing systems 28: 29th Annual conference on neural information processing systems 2015; 2015; Montreal. p. 91-9.

7. Bargoti S, Underwood J. Deep fruit detection in orchards. In: 2017 IEEE International Conference on Robotics and Automation (ICRA); 2017; Singapore. p. 3626-33.

8. Sørensen RA, Rasmussen J, Nielsen J, Jørgensen RN. Thistle detection using convolutional neural networks. In: EFITA WCCA 2017 Conference; 2017; Montpellier, France.

9. Zhu Y, Newsam S. DenseNet for dense flow. In: 2017 IEEE International Conference on Image Processing (ICIP); 2017; Beijing, China. p. 790-4.

10. Albuquerque CKG, Polimante S, Torre-Neto A, Prati RC. Water spray detection for smart irrigation systems with mask R-CNN and UAV footage. In: 2020 IEEE International Workshop on Metrology for Agriculture and Forestry (MetroAgriFor); 2020; Trento, Italy. p. 236-40.

11. He K, Gkioxari G, Dollar P, Girshick R. Mask R-CNN. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV); 2017; Venice, Italy. p. 2961-9.

12. Osorio K, Puerto A, Pedraza C, Jamaica D, Rodríguez L. A deep learning approach for weed detection in lettuce crops using multispectral images. AgriEngineering. 2020;2:471-88. https://doi.org/10.3390/agriengineering2030032

13. Noble WS. What is a support vector machine? Nat Biotechnol. 2006;24:1565-7. https://doi.org/10.1038/nbt1206-1565

14. Redmon J, Farhadi A. Yolov3: an incremental improvement. arXiv:180402767 [Preprint]. 2018 [cited 2023 May 17]. https://doi.org/10.48550/arXiv.1804.02767

15. Riekert M, Klein A, Adrion F, Hoffmann C, Gallmann E. Automatically detecting pig position and posture by 2D camera imaging and deep learning. Comput Electron Agric. 2020;174:105391. https://doi.org/10.1016/j.compag.2020.105391

16. Tedesco-Oliveira D, da Silva RP, Maldonado W Jr, Zerbato C. Convolutional neural networks in predicting cotton yield from images of commercial fields. Comput Electron Agric. 2020;171:105307. https://doi.org/10.1016/j.compag.2020.105307

17. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu CY, et al. SSD: single shot multibox detector. In: European Conference on Computer Vision. Cham: Springer; 2016. p. 21-37.

18. Zhou Z, Song Z, Fu L, Gao F, Li R, Cui Y. Real-time kiwifruit detection in orchard using deep learning on AndroidTM smartphones for yield estimation. Comput Electron Agric. 2020;179:105856. https://doi.org/10.1016/j.compag.2020.105856

19. Tang J, Arvor D, Corpetti T, Tang P. Mapping center pivot irrigation systems in the southern Amazon from Sentinel-2 images. Water. 2021;13:298. https://doi.org/10.3390/w13030298

20. Shojaeipour A, Falzon G, Kwan P, Hadavi N, Cowley FC, Paul D. Automated muzzle detection and biometric identification via few-shot deep transfer learning of mixed breed cattle. Agronomy. 2021;11:2365. https://doi.org/10.3390/agronomy11112365

21. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2016; Las Vegas, NV. p. 770-8.

22. Syed-Ab-Rahman SF, Hesamian MH, Prasad M. Citrus disease detection and classification

using end-to-end anchor-based deep learning model. Appl Intell. 2022;52:927-38. https://doi.org/10.1007/s10489-021-02452-w

23. He K, Girshick R, Dollar P. Rethinking ImageNet pre-training. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV); 2019; Seoul. p. 4918-27.

24. Lin TY, Dollár P, Girshick R, He K, Hariharan B, Belongie S. Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition; 2017; Honolulu, HI. p. 2117-25.