

# A review of animal recognition in urban environments

**ChuBO He**

Central South university, Shaoshan South Road, Tianxin District, Changsha City,  
Hunan Province, Wenyuan Street

1605122983@qq.com

**Abstract.** This paper describes the importance of animal recognition in the context of increasing animal participation in daily life. The progress in this field is reviewed and the basic methods used in animal recognition are introduced. This paper collects some free and open source data sets for deep learning and introduces statistical indicators for beginners to train and evaluate their own models. This paper collected a statistical chart of the number of papers in the field over the past 60 years, and analyzed the changes and preferences in the field of animal recognition. Finally, the problems faced by the field of animal recognition and the possible future development direction and advantages are summarized. Aiming at the animal recognition in the urban environment, a multi-task network method is recommended to solve the problem of lack of animal recognition data and the resource occupation of neural network in the urban environment.

**Keywords:** animal recognition, deep learning, urban environment, multi-task network.

## 1. Introduction

In urban environments, pet identification plays an increasingly important role. With more and more pets living with people in cities, pet loss is becoming more frequent. In 2018, already 73.55 million Chinese living in cities had pets, among which 56.48 million people owned dogs or cats [1]. So a visual recognition system that helps people find and identify their pets in the city is crucial. Although there have been a lot of relevant studies, animal identification and detection is still not an easy thing, and there is still no single system that can effectively solve this problem. Traditional pet animal identification methods, such as permanent (Tattoo, Microchip), semi-permanent and temporary (RFID) methods, the security they provide for pets is not enough. Therefore, it is necessary to develop a robust biometric recognition system to identify individual pets. Animal visual recognition system can be a system based on pattern recognition. It extracts important feature sets from the biometric data of individuals input into the system, compares the differences between them and the preset feature sets, and makes predictions based on the comparison results. Animal visual recognition systems can take advantage of the differences and uniqueness of vocalizations, body dynamics and fur patterns to improve recognition efficiency.

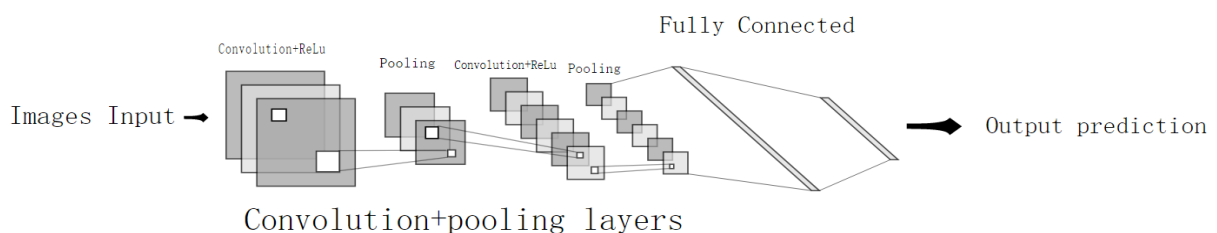
In 2013, Town, Marshall and Sethasien compared the individual similarities between species by enhancing the local contrast of images from the ventral side of Manta rays. The feature (i.e. spots) extraction algorithm they use is SIFT algorithm. This model was tested with only 720 images, although it took a lot of effort to get those 720 images. Town et al. achieved 51% accuracy and reported the difficulty of filming underwater [2]. At the same year, Loos and Ernst attempted to identify chimpanzee (Pan SPP.) faces, taking into account gradient changes in pixel values and features of

pixel grouping, to train a SVM classifier, which goal is make the linear decision boundary between categories reach the peak. Their test set was selected randomly by 20% from the data set. They accomplished an exactness of 84.0% and in the C-TAI information set, they accomplished 68.8% precision. Loos and Ernst look at this as a promising methodology, however it is restricted by the capacities of SVM classifiers [3].

In 2017, Hughes and Burghardt defined a locally naive Bayesian nearest neighbor nonlinear model for fin profile feature recognition. The model training set was a FinsScholl-2456 dataset (2456 images were included for 85 individuals), which was able to obtain correctly classified shark individuals with an accuracy of 82%. It was the best dorsal fin identification model at the time [4].

The algorithms and models mentioned above have the potential for improvement in performance, design approach, ease of implementation, and overall usability. The algorithm of Deep learning can solve these problems well.

One branch of deep learning research is to improve the performance of classification between unusual images, called fine-grained visual recognition. To identification the species of animal have become an active research area for fine-grained visual recognition, such as the classification of 675 similar moth species [5]. There are two mainstream approaches to improving model performance. The first is data enhancement, which is widely used to train image process networks, mainly by randomly flipping, cutting, rotating, blurring, moving and changing the color or light of images during each training step [6]. The second method is called "object localization", which classifies and predicts each unique part of the animal's body (i.e. head, back, ears, wings, etc.) separately, and finally concentrates on all combinations for final classification and prediction [7].



**Figure 1.** Example of neural network structure.

The deep learning methods described above all belong to the prediction results of only one animal class for each image. However, this does not apply to the recognition of camera trap images. To recognize different items (i.e., creatures), the specialists prepared object identifiers to disengage the pictures into regions going through CNN. Recently, three methods for object detection have become increasingly popular. The first is the Faster Region-CNN [8]. The second is YOLO, what separates the picture into lattices and predicts the picture through every framework cell in the organization utilizing a progression of predefined "anchors" connected with the order of anticipated shape and size [9]. Finally, there is a single multi-box detector, which is defined as a default set of boxes and adjusts these boxes to align objects of interest during training [10]. The target detection method is supplemented with a new evaluation metric called "joint crossing" (IOU), which is defined as the overlapping region of the actual region and the predicted region divided by the whole region containing the real and predicted regions [11].

In 2016, Freytag et al. used C-Zoo and C-TAI data sets to segment chimpanzee faces and train CNN architecture AlexNet. They reported a 92.0% and 75.7% improvement in accuracy compared to 84.0% and 68.8% of the original feature extraction method [12]. In 2017, Brust et al. trained

YOLO, after feature extraction, Brust et al. used the same program as Freytag et al. to train the same neural network, and finally tested it on a test set containing 500 images, achieving an accuracy rate of 90.8% [13].

This paper analyzes the hot spots and research directions in the field of animal recognition, aiming to inform readers of the existing problems and common methods in the field of animal recognition, and introduce readers to some common data sets and resources for future research. The following section will introduce common methods in the field of animal identification, and section 3 will introduce some available data set resources. The next part is the analysis of the existing research status. This paper makes statistics and analysis of the changes in the number of existing papers in the field of animal recognition, summarizes the focus and difficulties of animal recognition in the urban environment, and finally uses a very practical neural network framework in the urban environment.

## 2. Method

It can be seen in the literature that different groups of animals are studied by using various methods to recognize or classify the animal. However, there is not much significant research in this area. Studies have included the classification of cats and dogs, and surveillance of the wild lions face recognition can be used as an example, and for Marine fish, animal husbandry, as well as the identification of the rare animals, the number of related research is on the rise in recent years, but in general, the number is not much. In many of these studies, the methods used are varied and vary widely. Technologies include but are not limited to CNN, PCA, LDA, etc

### 2.1. Convolutional Neural Network, CNN

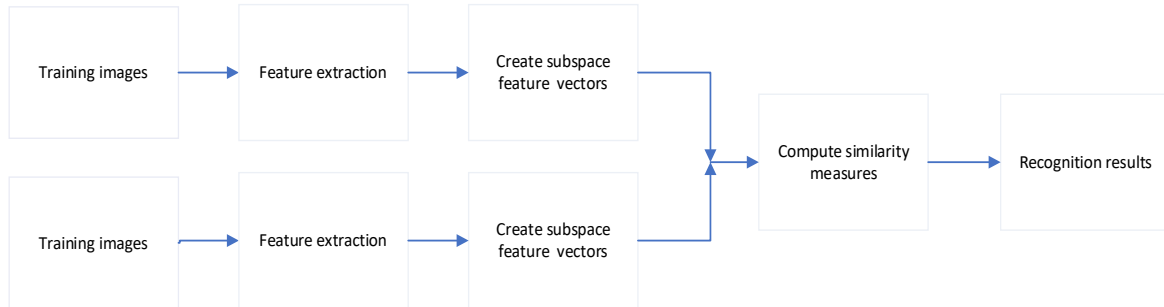
CNN is a layered perceptron suitable for deep learning algorithms to identify planar images. Each kernel in the different layers learns autonomously, without the need for predefined features. Compared with traditional techniques, it has the advantages of two-dimensional image recognition, local field perception, feature extraction, parameter sharing, hierarchical processing, image translation and rotation invariance recognition. CNN is a multi-stage pipeline network, which inputs images from one end and outputs predicted recognition probability from the other end. As displayed in the figure followed, CNN is basically a layer succession that can be partitioned into a few gatherings. Each layer comprises of a convolution layer in addition to a nonlinear activation function, typically a rectifying linear unit (ReLU), and a pooling layer, Mainly the maximum pooling layer. It closes with a few completely associated layers, the remainder of which is a result layer with expectations. Figure. 1 shows the classical structure of convolutional neural network.

What convolution does is extract features from inputs. In practice, CNN is learning the value of the convolution kernel in the training process. The more convolution kernels, the more features can be extracted. The extracted feature map size is controlled by three parameters: depth, which is the number of types of convolution kernels; Step length; Filled in. After each convolution operation, there will be a nonlinear operation. The reason for introducing nonlinear operations into convolutional neural networks is that most real-world problems that need to be learned are nonlinear. The purpose of pooling is to reduce the dimension of each feature graph as much as possible while preserving the most important information. The function of pooling can be summarized as follows: it makes the input dimension smaller and easier to operate; Reduce the parameters and computation in the network, so as to inhibit the overfitting. Enhance the network's robustness to small deformation, distortion and translation of the input image. The definition of Fully connected is that every neuron on the upper and lower levels connects. In addition to classification, adding full connection layer is also an effective method to learn nonlinear combination between features. Here it can be understood as feature classification extracted by pooling of paired convolution kernels.

### 2.2. Principal Component Analysis, PCA

Principal Component Analysis (PCA) is a common data Analysis method, which is often used for dimensionality reduction of high dimensional data and can be used to extract the main feature components of data. PCA not only transforms vector-based big data into unrelated vectors, but also transforms related vectors into small data and unrelated vectors, and transforms original data

accordingly. After that, the raw data is represented by the main components in the different dimensions. Characteristic names obtained by representing different sizes. Among them, the variance of the first principal component is the largest, and the variance of the other basic components will decrease successively.



**Figure 2.** The structure of LDA.

PCA can be summarized as a linear projection in the direction of the least reconstruction error (maximum variance). Finding the strongest pattern is the most important goal in PCA.

### 2.3.LDA

Linear Discriminant Analysis (LDA) is also called Fisher Linear Discriminant (FLD). Unlike PCA, LDA is supervised learning. The basic idea of LDA is to project high-dimensional pattern samples into the optimal discriminant vector space, so as to extract classification information and compress the dimension of feature space. After the projection, the maximum inter-class distance and the minimum intra-class distance of pattern samples are guaranteed in the new subspace, that is, the pattern has the best separability in this space. Therefore, it is an effective feature extraction method. Using this method, the interclass dispersion matrix of the projected pattern sample can be maximized and the intra-class dispersion matrix can be minimized. The structure of LDA is shown in Figure 2.

## 3. Dataset

It has been proved that when the training set and test set come from two different data sets, the performance of the algorithm will be significantly worse than when they come from different partitions of the same data set. Therefore, the current research in the field of image recognition largely needs to test the performance of new algorithms by using benchmark data sets. This fact runs counter to most machine learning systems, in which the more training samples, the better the classification results should be. The same is true in animal recognition.

Although there are many image data sets in the field of animal recognition, there is no standard publicly available object recognition benchmark. Therefore, many researchers choose to create a data set by themselves, possibly because of the particularity of research methods, which require images to participate in training to have the same size and pixel. Next, a few data sets are useful.

### 3.1.Snapshot Serengeti

This dataset is based on the deployment of 225 Camera traps across 1 [13], 125 square kilometres of Serengeti National Park in Tanzania to assess the spatio-temporal dynamics between species. The camera has been operating continuously since 2010 and has taken a total of 99,241 days and 1.2 million photos as of 2013. As users view images, they record species, numbers of individuals, related behaviors and the presence of young birds. This provides a unified classification and raw image resource for researchers studying the dynamics of multiple species in complete ecosystems, as well as a great boon for researchers in deep learning and image recognition. More than 28,000 registered users contributed 10.8 million categories, saving image recognition workers a lot of time.



**Figure 3.** Example of the image.

### 3.2. *Dogs vs Cats*

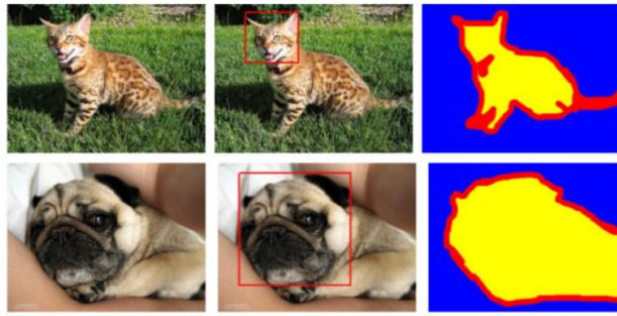
The dog vs cat dataset comes from a 2013 contest on Kaggle [14]. The dataset was reconstructed from a larger dataset of 3 million hand-annotated photos of cats and dogs, which is developed by Petfinder.com in partnership with Microsoft. The training set consists of 12,500 pictures of cats and dogs labeled as cat and dog, with a total of 25,000 pictures. The pictures are in 24-bit JPG format, namely RGB three-channel images, with different sizes. The test set consists of 12,500 cat or dog images, unlabeled, also in 24-bit JPG format, RGB three-channel images of different sizes. Pierre Sermanet won the competition, achieving about 98.914% classification accuracy on 70% of the subsamples in the test dataset. His method was used in a later paper [15].



**Figure 4.** Give a typical example.

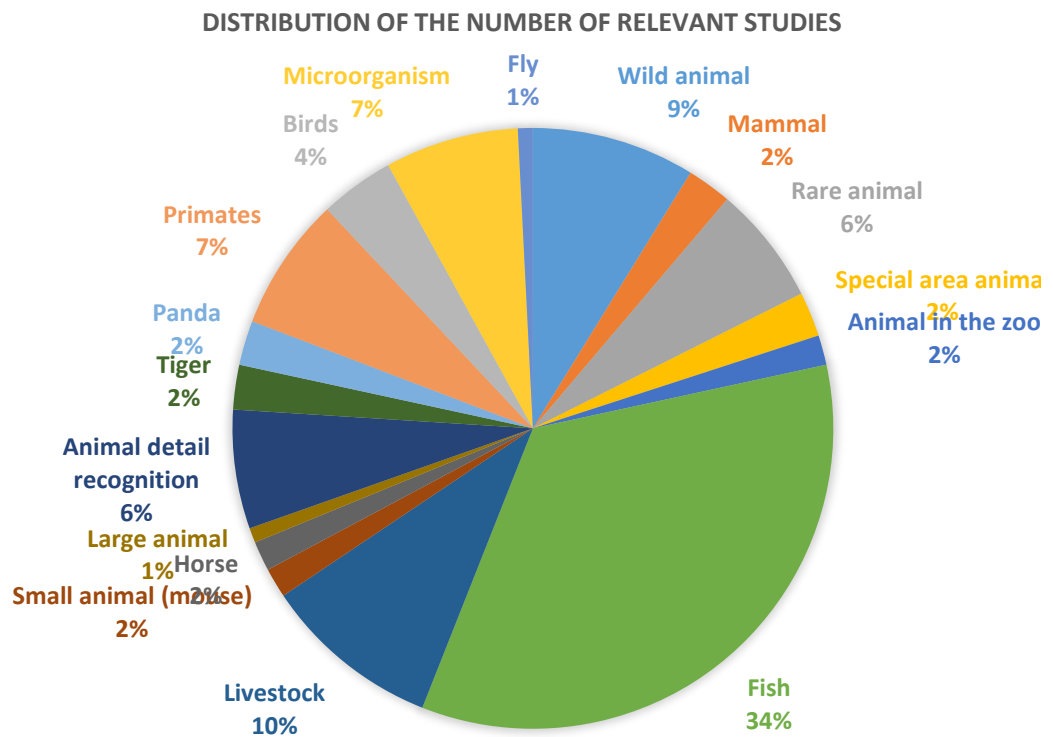
### 3.3. *The Oxford-IIIT Pet Dataset*

This is a publicly available pet dataset provided by the University of Oxford [16], which containing 37 categories, each with approximately 200 images. The images vary greatly in proportion, posture and lighting. All images have category, header ROI, and pixel-level correlation labels that can be used for three-image segmentation. Each image in the dataset is annotated with : (a) specific species and specific breed names; (b) Bounding box around animal heads (ROI); And (c) pixel-level foreground background segmentation (Trimap).



**Figure 5.** An example of annotations.

### 3.4. Animal-10



**Figure 6.** Distribution of the number of relevant studies.

Animal-10 is a data set for 10 animal categories: dog [17], cat, horse, chicken, Spider, cow, butterfly, sheep, elephant, and squirrel. The kit contains 28,000 medium quality animal images, ideal for testing animal recognition or classification tasks. All pictures were gathered from Google Images and have been checked by human. At the same time, there is incorrect data to simulate real situations (for example, images taken by application users). This dataset is provided by individual users and made public on Kaggle.

### 3.5. Stanford Dogs Dataset

The Stanford Dog Dataset contains 120 dog images from around the world, images and annotations from ImageNet, recommended for fine-grained image classification. The dataset includes 120 categories of dogs with 20,580 images and Boundingboxes and corresponding breed annotations.

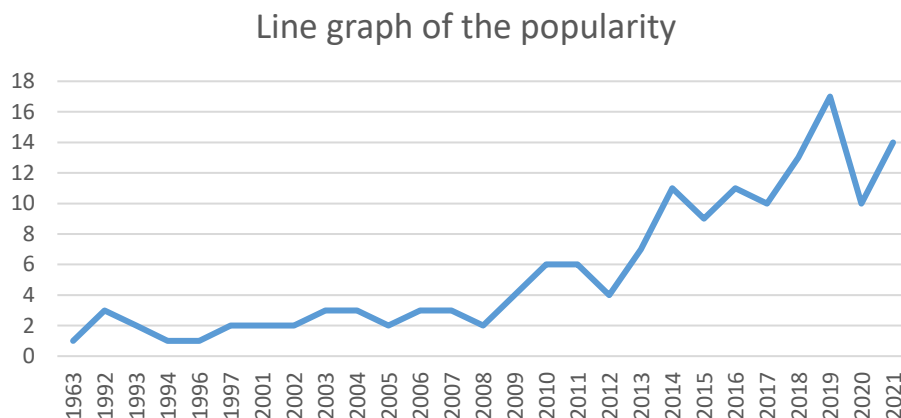
By the way, If we want to evaluate how good the trained model is, the evaluation indexes of animal classification include Accuracy, Precision, Recall, F1-score, Specificity and Sensitivity, confusion

matrix, ROC curve and AUC , The official DOCUMENT of SKLEARN Model 《Evaluation: Quantifying the Quality of Predictionsprovides 》 made a very detailed descriptions and implementations of the predictions for binary classification, multi-classification and multi-label problems.

## 4. Review

### 4.1.Statistic

Lately, the field of animal recognition and classification has become increasingly popular, and there are more and more related studies. The following table shows the heat analysis in this field (data from IEEE) (Fig.6 and Fig.7) .



**Figure 7.** Line chart of number of papers over time.

As you can see, the number of papers is increasing year by year. It is worth mentioning that most researchers do not only focus on single animals, and only a few taxonomic studies on single species have been conducted

### 4.2.Challenge

Although authentication and recognition usually use the same classification algorithm, they are two different application scenarios. To better understand the task and difficulty of animal detection and classification, the following factors must be considered as they may significantly affect the performance of animal detection and classification models.

**Lighting and other image acquisition conditions:** The quality of the images used may be affected by lighting variations, light distribution, intensity, or camera characteristics such as sensor response and lens.

**Illusion:** Due to the specific nature of the urban environment, parts of the animal image may be obscured by other objects to confuse the prediction model.

In fact, animal recognition has its own particularity, unlike most object recognition. Pictures of animals are often similar. Research into other object categories has focused on how to enable computers to distinguish objects that are easy to distinguish with the naked eye, such as airplanes and cats. The most well known global benchmarks (e.g., Caltech-101 [18], Caltech-256 [19], PASCAL VOC [20]) include many item classes, the majority of which are outwardly unique.Indeed, even in the bigger data set (ImageNet) , classes are characterized in light of significant level ontologies, so visual similarities between them are rare, if any, by accident. But if you have to recognize or classify the creature, it seems inevitable.Although animal recognition also requires the problem of distinguishing between dogs and cats of different breeds, as well as cats and cats of the same breed, this is an example of a challenging fine-grained object classification. The difficulty is that these may only be distinguished by inconspicuous phenotypic subtleties, which, due to the highly morphed



nature of these animals' bodies and certain environments (such as cities), it is hard to detect and record by surveillance systems, let alone identify.

At present, fine-grained image classification basically adopts deep learning method and achieves good results. In particular, it very well may be generally separated into the accompanying classifications:

While other methods are difficult to obtain differentiated local details for fine-grained classification, Deep Convolutional Neural Network (DCNN) can solve this problem well

- Location-based recognition method: firstly, the discriminating parts are found, and then feature extraction and classification are carried out. This method can be divided into strong supervision and weak supervision.

- The method based on network integration is as follows: multiple DCNNs are used to discriminate similar features in fine-grained recognition.

- High-order coding method of convolution features: CNN features are transformed into higher-order features and then classified, mainly including Fisher Vector, bilinear model and kernel fusion.

#### *4.3.Application direction*

Animal biometrics are very useful in field research. The system uses remote audio and video recordings to quickly locate species, and significantly increases the number of predetermined identification objects the system can identify, rather than just locating them in time and space

Christian Kublbeck and Andreas Ernst used face detection algorithms to detect the facial features of African apes, chimpanzees and gorillas. Classification rates varied with the quality and visibility of orangutans' facial images, achieving a classification rate of 89-97% under optimal conditions. The system, which can be applied to real-time data from remote cameras, with limited computing resources, the time needed to detect species from the input videos can be reduced significantly, from hundreds of hours of thousands of video clips to hours or even minutes. The system will output a list of details such as the species name, date, time, and location. As a result, a lot of video film can be handled routinely for research purposes (For example, detecting the presence of great apes, estimating visit rates, monitoring whether populations have migratory trends or calculating how much habitat overlap chimpanzees and gorillas have.

The limitation of manual data classification is solved. It can be of great help to ecological researchers because of the unique output of animal biometrics and its extraction of some key features. The definition of animal appearance attributes helps to more realistically discover, distinguish and identify species, individuals and their behavior and form. Unlike human observation, this work avoids variation and error due to human observer subjectivity, skill, or stubbornness to past experience.

Use animal biometrics to analyze behavior. Zheyuan Wang et al. used support vector machine (SVM) to divide the depth image into five behaviors of stationary [21], walking, grooming, feeding and rotation through the recognition of rodent behaviors. But automated animal biometrics applications explicitly extract behavior that is intrinsically linked to a controlled environment. Focus on less complex tasks, such as swarm or swarm movement of bees, fish, or birds. Unique motor activities, such as the gait of quadrupeds [21], generate enough motor features to determine whether an animal exists from the video and classify behaviors such as walking, trotting or tracking [21]. There are many good algorithms in this field, but there is still some way to go before a system can be assembled to detect complex audio-visual movements and behaviors

It's a great way to relieve the stress of having a large pet population in the city. Pradeepa Jeyaraj mentioned in the article that animal biometrics can well meet the demand for fine-grained animal classification in the urban environment [22], which is different from the fuzzy classification provided by intelligence in the past, and can well solve the serious threat to human society posed by the excessive number of stray dogs and cats.

#### *4.4.Solution: new practical framework of multi-task neural network*

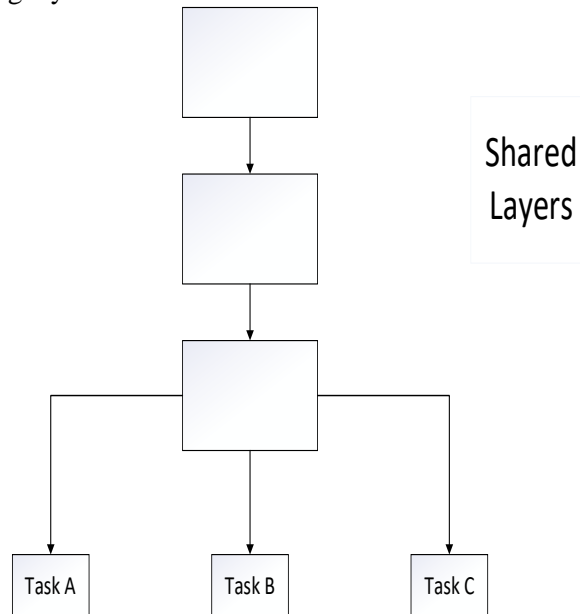
In practice, the actual situation of change always ask us to customize the neural network to meet the need, but a neural network can meet most of the actual scene takes up a lot of computing resources, and requires a lot of data sets to train, it is lack of neural network training, and in be used actually



solve single task in the scene, There's a lot of redundancy in calling the whole neural network. In addition, in the real world, the tasks we need to face are often interrelated, and the common idea is to train a network model that can have different branches solve different tasks on the basis of sharing some characteristics and parameters. Multi-task network aims to improve generalization ability and processing efficiency through joint learning between related tasks. So we recommend multi-head neural network to solve this kind of problem.

Under the logic of multi-task, different branch layers are usually used to solve a specific task. The common multi-task network structure is shown in the figure below, which first passes through some Sharing layers and then divides into specific layers to solve specific tasks. Hard-parameter Sharing is the mainstream of multi-task model.

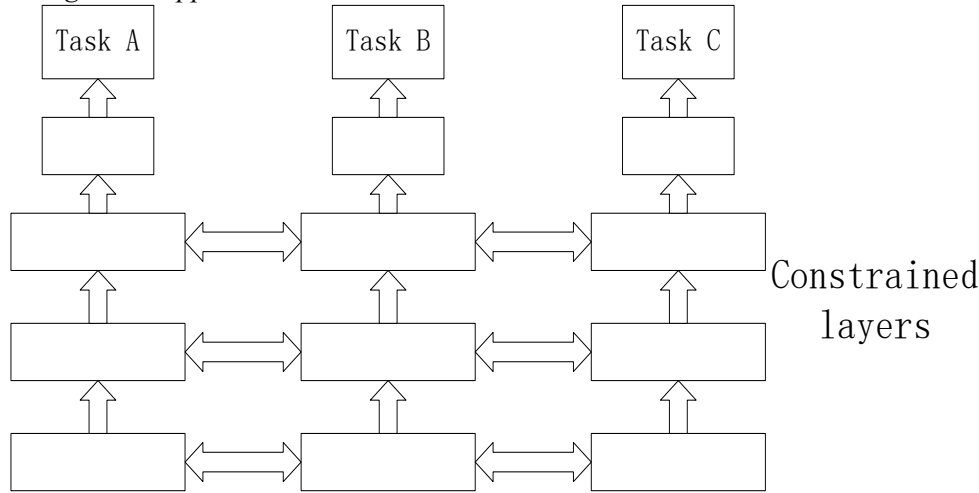
However, the structure of such parameter sharing is not fixed and varies from task to task, and the optimal number of sharing layers is different for different tasks. This conclusion is based on the article cross-Stitch Networks for Multi-Task Learning by CVPR2016. The author once conducted such an experiment, listing all possible sharing layers for different tasks and comparing the performance differences between sharing layers.



**Figure 8.** The structure of Hard-parameter Sharing.

The following figure shows the shared soft parameters. The parameters of each task are different, and in quite a few scenarios, we cannot directly use the exact same parameters of other tasks, only some parameters similar to this task. Soft parameter sharing is usually given in a regularized form, You can use L2, trace norm, etc.

#### 4.5. Advantages and application scenarios



**Figure 9.** The structure of Soft-parameter Sharing

Each task has more or less sample noise, which may differ from task to task, and eventually multiple task learning can cancel out some of the noise (similar to bagging's idea that noise from different tasks exists in all directions and eventually averages to zero).

For some noisy tasks or insufficient training samples with high dimensions, the model may not be able to learn relevant features.

Some characteristics may be difficult to learn in the primary task (such as only having a high order of correlation, or being suppressed by other factors), but they are good to learn in the auxiliary task. These features can be learned through ancillary tasks such as Hints (Predicting important Features).

By learning a large enough hypothesis space, you can perform better on certain new tasks in the future (solve cold priming), provided that these tasks are homogeneous.

Constrain the model as a form of regularity. The so-called Revolution Bias. Alleviates the complexity of overfitting and reduces the Rademacher complexity of the model.

## 5. Conclusion

This paper mainly provides beginners in the field of animal recognition with the practical background of the subject and some cutting-edge research, and also provides some basic methods of image recognition processing. This article has collected some open source and free data set resources to avoid spending too much time building data sets. At the same time, some statistical indicators are given for beginners to evaluate their own models. Finally, this paper proposes a solution of multi-task network to solve the problems of insufficient and chaotic data and redundant neural network resources in the urban background.

## References

- [1] 2019-2025 "China Pet Industry Market Analysis Research and Development Trend Research Report" released by Zhiyan Consulting
- [2] Town C, Marshall A, Sethasathien N. Manta M atcher: automated photographic identification of manta rays using keypoint features[J]. Ecology and evolution, 2013, 3(7): 1902-1914.
- [3] Loos A, Ernst A. An automated chimpanzee identification system using face detection and recognition[J]. EURASIP Journal on Image and Video Processing, 2013, 2013(1): 1-17.
- [4] Hughes B, Burghardt T. Automated visual fin identification of individual great white sharks[J]. International Journal of Computer Vision, 2017, 122(3): 542-557.
- [5] Rodner E, Simon M, Fisher R B, et al. Fine-grained recognition in the noisy wild: Sensitivity analysis of convolutional neural networks approaches[J]. arXiv preprint arXiv:1610.06756, 2016.
- [6] Souri Y, Kasaei S. Fast bird part localization for fine-grained categorization[C]//Proc. 3rd

- Workshop Fine-Grained Vis. Categorization (FGVC3) CVPR. 2015.
- [7] S. Ren, K. He, R. Girshick and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137-1149, 1 June 2017, doi: 10.1109/TPAMI.2016.2577031.
  - [8] Redmon J, Divvala S, Girshick R, et al. You only look once: Unified, real-time object detection[C]//*Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016: 779-788.
  - [9] Liu W, Anguelov D, Erhan D, et al. Ssd: Single shot multibox detector[C]//*European conference on computer vision*. Springer, Cham, 2016: 21-37.
  - [10] Nowozin S. Optimal decisions from probabilistic models: the intersection-over-union case[C]//*Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014: 548-555.
  - [11] Taigman Y, Yang M, Ranzato M A, et al. Deepface: Closing the gap to human-level performance in face verification[C]//*Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014: 1701-1708.
  - [12] Freytag A, Rodner E, Simon M, et al. Chimpanzee faces in the wild: Log-euclidean CNNs for predicting identities and attributes of primates[C]//*German Conference on Pattern Recognition*. Springer, Cham, 2016: 51-63.
  - [13] Swanson A, Kosmala M, Lintott C, et al. Snapshot Serengeti, high-frequency annotated camera trap images of 40 mammalian species in an African savanna[J]. *Scientific data*, 2015, 2(1): 1-14.
  - [14] Parkhi O M, Vedaldi A, Zisserman A, et al. Cats and dogs[C]//*2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012: 3498-3505.
  - [15] Sermanet P, Eigen D, Zhang X, et al. Overfeat: Integrated recognition, localization and detection using convolutional networks[J]. *arXiv preprint arXiv:1312.6229*, 2013.
  - [16] Renwick R L. Smart Cat-Door System (SCDS)[J]. 2020.
  - [17] Gupta S N, Brown N B. Adjusting for Bias with Procedural Data[J]. *arXiv preprint arXiv:2204.01108*, 2022.
  - [18] Fei-Fei, Li, Rob Fergus, and Pietro Perona. "Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories." 2004 conference on computer vision and pattern recognition workshop. IEEE, 2004.
  - [19] Griffin, Gregory, Alex Holub, and Pietro Perona. "Caltech-256 object category dataset." (2007).
  - [20] Vicente, Sara, et al. "Reconstructing pascal voc." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014.
  - [21] Wang Z, Mirbozorgi S A, Ghovanloo M. Towards a kinect-based behavior recognition and analysis system for small animals[C]//*2015 IEEE Biomedical Circuits and Systems Conference (BioCAS)*. IEEE, 2015: 1-4.
  - [22] Jeyaraj P, Aponso A. A Review of Techniques for Image Classification to Enhance Online Animal Adoption Speed[C]//*Proceedings of the 2020 12th International Conference on Computer and Automation Engineering*. 2020: 114-118.