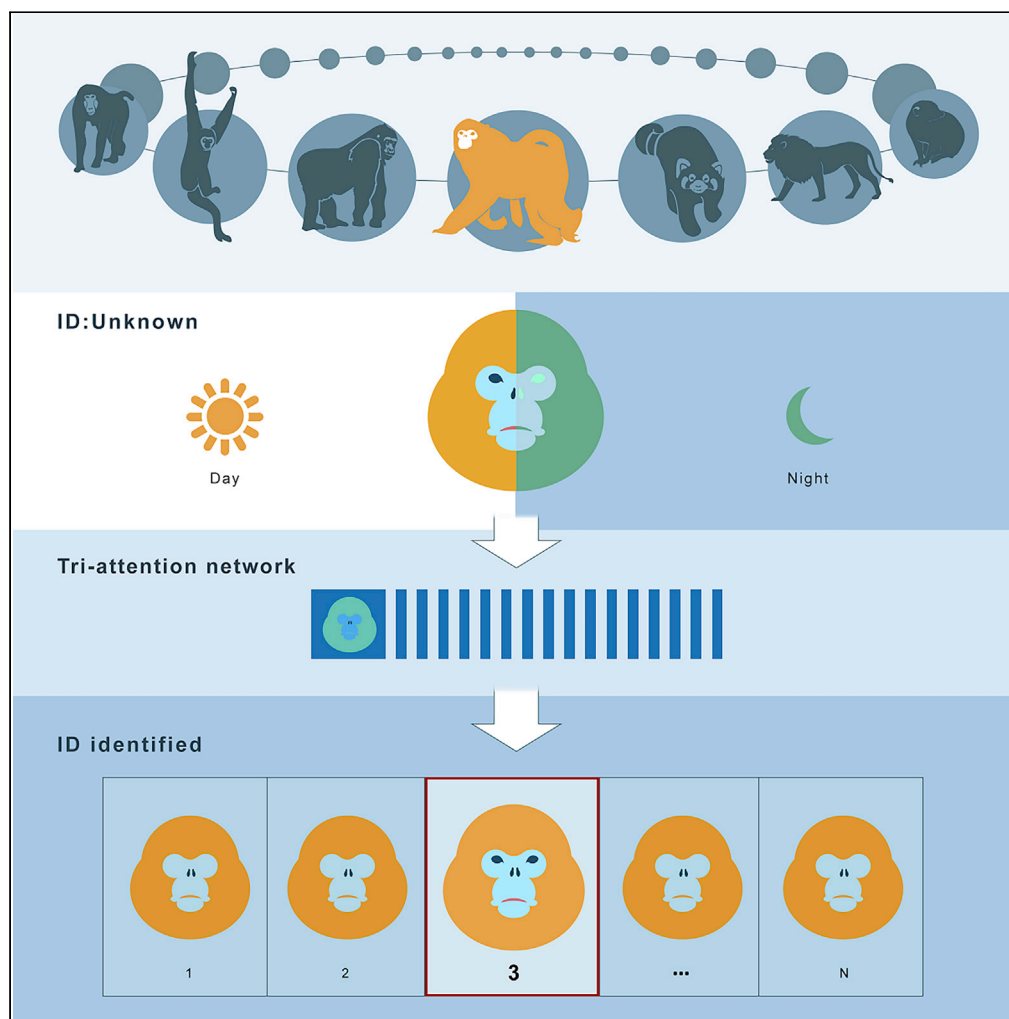


Article

Automatic Identification of Individual Primates with Deep Learning Techniques



Songtao Guo,
Pengfei Xu,
Qiguang Miao, ...,
Yewen Sun, Zhihui
Shi, Baoguo Li

songtaoguo@nwu.edu.cn

HIGHLIGHTS

The Tri-AI system can rapidly detect and identify individuals from videos and images

Tri-AI had an ID identification accuracy of 94% for 41 primates and 4 carnivores

The system could individually recognize 31 animals/s with images taken day or night

Systems like Tri-AI make around-the-clock monitoring and behavior analysis possible

Guo et al., iScience 23, 101412
August 21, 2020 © 2020 The
Authors.
[https://doi.org/10.1016/
j.isci.2020.101412](https://doi.org/10.1016/j.isci.2020.101412)

Article

Automatic Identification of Individual Primates with Deep Learning Techniques

Songtao Guo,^{1,11,12,*} Pengfei Xu,^{2,3,4,11} Qiguang Miao,^{5,6,11} Guofan Shao,⁷ Colin A. Chapman,^{1,8,9} Xiaojiang Chen,^{2,3,4} Gang He,¹ Dingyi Fang,^{2,3,4} He Zhang,¹ Yewen Sun,¹ Zhihui Shi,¹ and Baoguo Li^{1,10}

SUMMARY

The difficulty of obtaining reliable individual identification of animals has limited researcher's ability to obtain quantitative data to address important ecological, behavioral, and conservation questions. Traditional marking methods placed animals at undue risk. Machine learning approaches for identifying species through analysis of animal images has been proved to be successful. But for many questions, there needs a tool to identify not only species but also individuals. Here, we introduce a system developed specifically for automated face detection and individual identification with deep learning methods using both videos and still-framed images that can be reliably used for multiple species. The system was trained and tested with a dataset containing 102,399 images of 1,040 individuals across 41 primate species whose individual identity was known and 6,562 images of 91 individuals across four carnivore species. For primates, the system correctly identified individuals 94.1% of the time and could process 31 facial images per second.

INTRODUCTION

Answering many theoretical questions in ecology and conservation frequently requires the identification and monitoring of individual animals (Nathan, 2008). However, traditional marking methods are often costly and involve considerable risk to the animal, a risk that is typically unacceptable for endangered species (Fernandezduque et al., 2018). With the maturity of digital image acquisition and camera traps, it has become relatively easy to repeatedly capture images of animals; however, using these images to address many ecological questions requires accurate individual identification (Wang et al., 2013).

By employing image matching methods (Zeppezauer, 2013; Zhu et al., 2013; Chu and Liu, 2013; Finch and Murray, 2003) and machine learning (Loos and Ernst, 2013; Swanson et al., 2016; Nathan, 2008), researchers have accurately identified species from images using animal body surface characteristics, like colors (Zeppezauer, 2013; Zhu et al., 2013; Wichmann et al., 2010), shape (Chu and Liu, 2013; Tweed and Calway, 2002; Finch and Murray, 2003), and texture (Crouse et al., 2017). To identify individuals, however, the images of specific body parts are required (Norouzzadeh et al., 2018; Burghardt et al., 2004; Lahiri et al., 2011; Hiby et al., 2009; Karanth, 1995), and this has been done with penguin's abdomens (Burghardt et al., 2004), stripes of zebra (Lahiri et al., 2011) and tigers (Xu and Qi, 2008), and the unique spot and scar features on the backs of killer whales (Arzoumanian et al., 2005). Although helpful for specific studies, these methods are using species-specific traits and thus they cannot be used across species.

The objective of our study was to determine if animal facial images could be used as a universal part for individual detection and identification. Machine learning for facial recognition has been developed for animals. For example, Burghardt and Calic (2006) presented a method based on Haar-like features and AdaBoost algorithm to detect the lion's face and Ernst and Küblbeck (2011) extracted the features of the key facial points of chimpanzees for individual identification. Recently, Hou et al. (2020) used VGGNet for face recognition on 65,000 face images of 25 pandas and obtained an individual identification accuracy of 95%. Schofield et al. (2019) presented a deep convolutional neural network (CNN) approach for face detection, tracking, and recognition of wild chimpanzees from long-term video records in a 14-year dataset yielding 10 million face images from 23 individuals, and they obtained an overall accuracy of 92.5% for

¹Shaanxi Key Laboratory for Animal Conservation, School of Life Sciences, Northwest University, Xi'an 710069, China

²School of Information Sciences and Technology, Northwest University, Xi'an 710127, China

³Shaanxi International Joint Research Centre for the Battery-free Internet of Things, Xi'an, China

⁴Institute of Internet of Things, Northwest University, Xi'an, China

⁵School of Computer Science and Technology, Xidian University, Xi'an 710071, China

⁶Xi'an Key Laboratory of Big Data and Intelligent Vision, Xi'an 710071, China

⁷Department of Forestry and Natural Resources, Purdue University, West Lafayette, IN 47907, USA

⁸Department of Anthropology, Center for the Advanced Study of Human Paleobiology, George Washington University, Washington, DC 20037, USA

⁹School of Life Sciences, University of KwaZulu-Natal, Scottsville, Pietermaritzburg 3209, South Africa

¹⁰Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China

¹¹These authors contributed equally

¹²Lead Contact

*Correspondence.:
songtaoguo@nwu.edu.cn
<https://doi.org/10.1016/j.isci.2020.101412>



identity recognition and 96.2% for sex recognition. Such studies suggest that facial images could be a universal feature for individual detection and identification for mammals (Kumar et al., 2017).

Past studies have applied machine learning to recognize species; however, these methods use classical image segmentation and artificial features to perform coarse-grained species or individual identification and require high image quality, which results in that the recognition performances of these methods become poor if the image is occluded or the animal changes posture. As a result, this method is not practical for uncontrolled field conditions.

CNNs deal well with nonlinear problems due to their powerful feature extraction capabilities that provide new solutions for occlusion and posture changes of animals faces (Krizhevsky et al., 2012). CNNs have become an exciting and active research topic in the field of computer vision, and various derivative network models have been proposed, such as ZFNet (Zeiler and Fergus, 2014), VGG (Simonyan and Zisserman, 2015), GoogleNet (Szegedy et al., 2015), ResNet (He et al., 2016), and DenseNet (Huang et al., 2017). In the detection and identification of wildlife, the Log-Euclidean framework was used by Freytag et al. (2016) to improve the identification ability of CNN-based methods to predict the identity, age, and gender of chimpanzees. Norouzzadeh et al. (2018) used a variety of CNN frameworks and demonstrated that deep learning (DL) can automate species identification for 99.3% of 3.2 million wildlife images in the Serengeti and performed slightly better (96.6% accuracy) than teams of volunteers. Species or individual identification based on machine learning has achieved reasonably good results, especially when many images are available.

To accurately identify individuals using the same methods across multi-species, we designed a deep neural network with attention mechanism for individual identification of primates and other species. Subsequently, we developed and tested our system for automatic individual identification that uses DL techniques to identify body parts. We believe that this represents a breakthrough in software engineering to identify individual animals that will be very useful for biological research. We call our system Tri-AI, which is used for automated face detecting, identifying, and tracking. Here we report on our efforts to design the system Tri-AI for individual identification of 41 primate and 4 carnivore species (Figure 1).

RESULTS

Overview of the Experimental Results

We designed Tri-AI (related to Supplemental Information) to quickly detect, identify, and track individuals from videos or still-framed images of multiple species (related to Transparent Methods). To accomplish these tasks, the system Tri-AI needed to be fed in a great deal of data. Therefore, we established an Image Acquisition Standard for animals (related to Transparent Methods), and then obtained facial images of known individuals from images or videos to build an image set (related to Transparent Methods). This image set was used for face detection and individual identification separately. We trained and tested Tri-AI with an initial database that contained 102,399 images of 1,040 individuals whose identity was known, across 41 primate species and 6,562 images of 91 individuals across four carnivore species (related to Figure 1 and Table S1). Tri-AI achieved individual identification accuracy of 94.11% (mean \pm SE = 94.11% \pm 0.0024; Figure 1, Table S1) with an average detection and recognition time of 31 images per second. The identification accuracy for different species by the proposed Tri-attention network (related to Transparent Methods) in the system Tri-AI was obtained on the test image dataset (Table S1), and the accuracy of individual identification for golden snub-nosed monkey test images was 94.12% (related to Table S1). These images were obtained by mobile phones or single-lens reflex (SLR) cameras and have relatively high resolutions. The facial images are clear with little shielding. The images of most golden snub-nosed monkeys were captured on different days, whereas the images of other species were obtained in one day in zoos. Across the species examined, Tri-AI achieved correct detection accuracy (i.e., detection of the animal from the environment) of 91.1%–97.7%. Thus, our work suggests that the animals' faces can be used to identify individuals. The dataset and the related testing models have been released publicly at the database AFD (Animal Face Database): <http://dx.doi.org/10.17632/z3x59pv4bz.2>.

The three still-framed images in Figure 2A depict Tri-AI detection of the animal's facial area and subsequently identify the individual. If the individual is included in the existing database, Tri-AI identifies them and provides their identification number or name (Figures 2A, 2B1, and 2C1). If the individual is

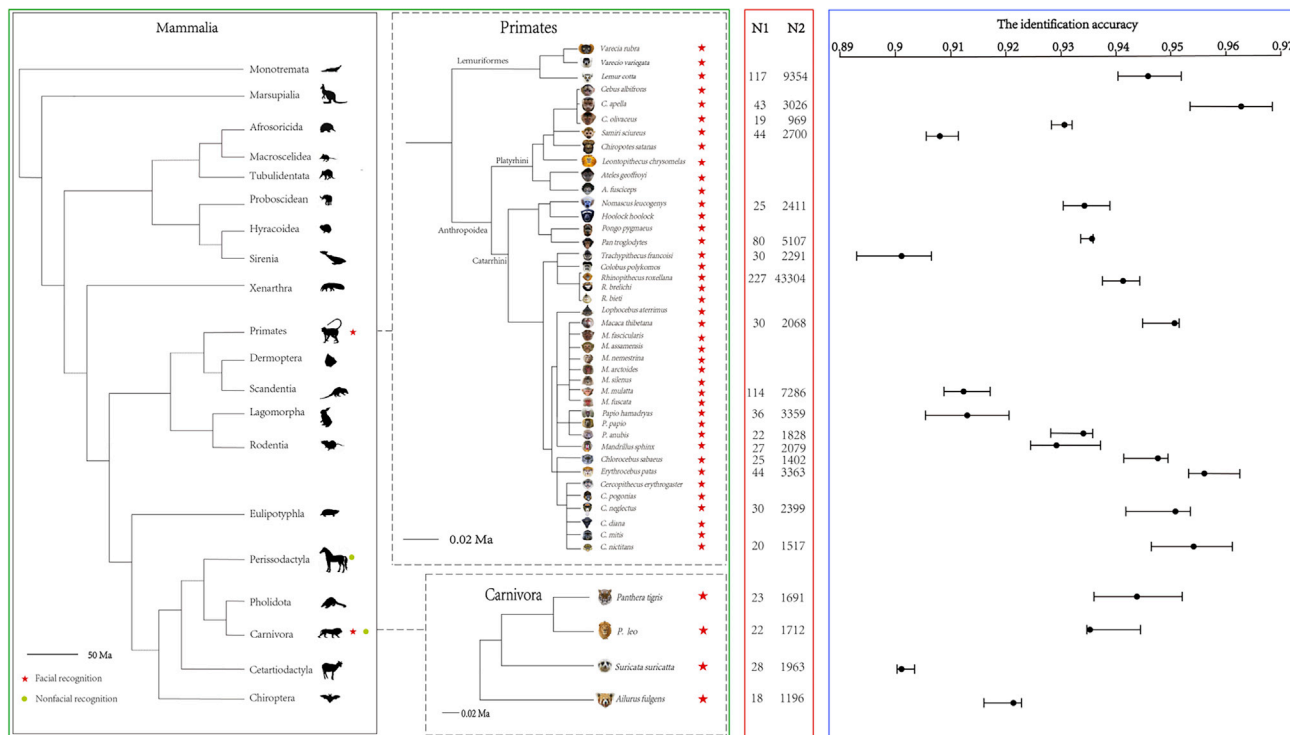


Figure 1. Application of Deep Learning in Animals and the Individual Identification Accuracy of 21 Species

Green dot, Non-facial (body) biometric character recognition; red stars, facial biometric recognition. Break line box: Tri-AI has successfully performed individual identification in the species in this study (in left green box). N1 is the number of the individuals, and N2 is the number of facial images for the corresponding species in our image dataset (in middle red box). For each species, we give the average, maximum, and minimum values of the accuracy from multiple tests (in right blue box).

not in the database, Tri-AI identifies the animal as a new individual and gives them a new number or name (Figures 2A, 2B2, and 2C2).

For videos, Tri-AI detects animal faces and identifies their identity frame by frame (Figures 2A–2C), and again if a new individual appeared in the videos, Tri-AI gives the animal a new ID or name (Figures 2A, 2B2, and 2C2). Meanwhile, all the facial images of the new individual can be collected to update the original image dataset. With video, more facial images of the new individual are collected. When the new individual has more than 60 facial images, the updated dataset is used to retrain Tri-attention network offline. If images of these individuals are collected in the future, the updated Tri-attention network will identify the individuals. Most importantly, this update will occur even if the observer has not been able to distinguish features to know the animal (Figure 2) and thus will help the researcher to learn new individuals. Our model identifies new individuals and collects the face images of the new individual according to the locations of the detected faces. Then the model is retrained offline on the update dataset with new individuals.

We obtained 98.70% accuracy in facial detection and 92.01% accuracy for individual identification with test videos of golden snub-nosed monkeys. Tri-AI needs to be trained with 60 or more images per individual to obtain stable identification accuracy (Figure 3). These results demonstrate that this AI software can reliably identify individuals across species using facial images (related to Tables S1 and S2). The promising performance of Tri-AI inspires that this approach can be used to automatically identify and track individuals of many wildlife species. Tri-AI can do real-time monitoring in the day or at night.

Details of Experimental Results

Face Detection for Images

Tri-AI can detect faces with a high level of accuracy. We randomly tested Tri-AI with 500 gray scale images of golden snub-nosed monkey and found that it successfully found 632 faces, but failed to detect a face in

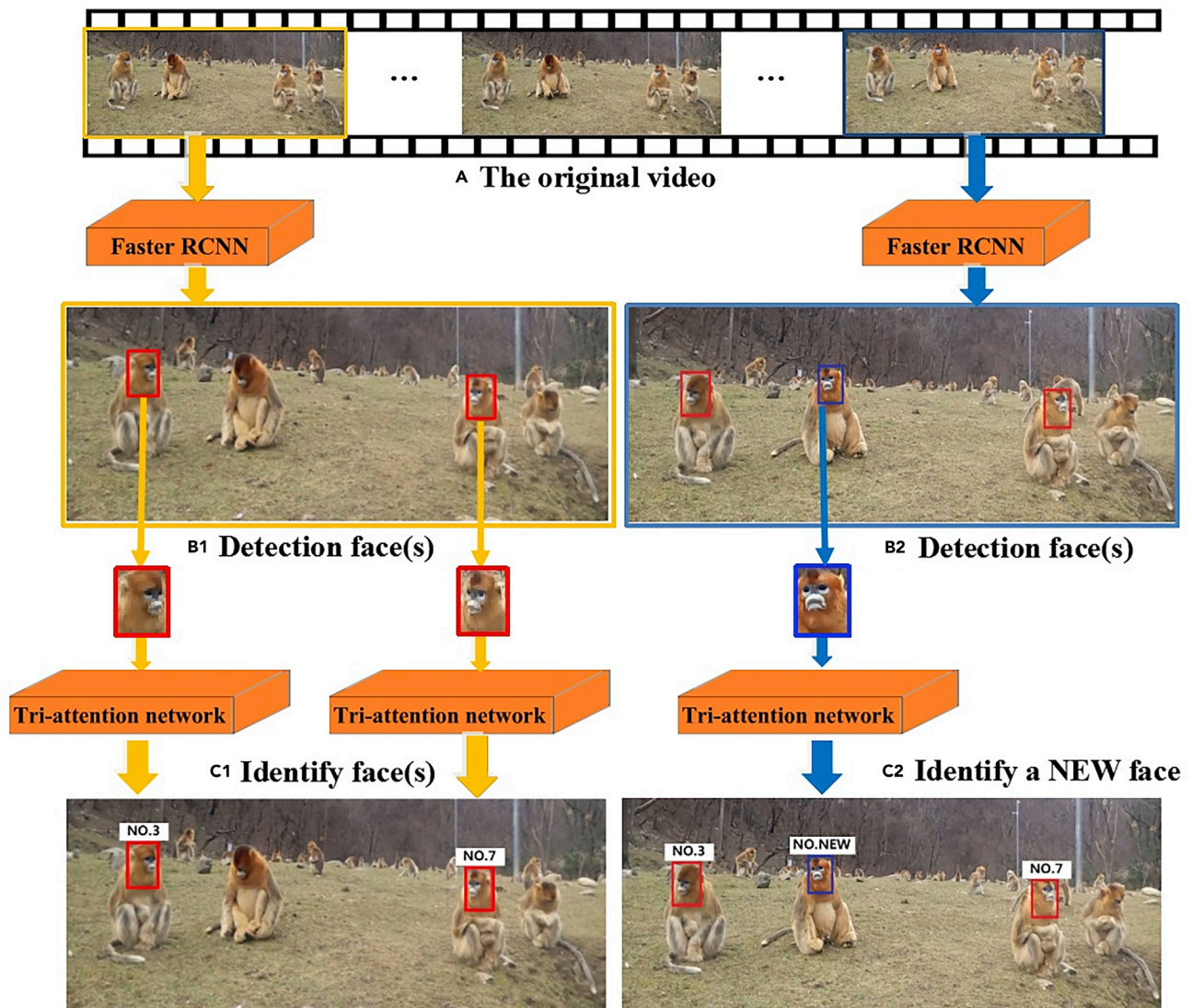


Figure 2. Face Detection and Identification of the Golden Snub-Nosed Monkeys by Using Deep Learning Methods in Tri-AI

(A) An original video has many frames, and the face areas of the monkeys must first be detected using Faster RCNN (Ren et al., 2015) from each frame.

(B1 and B2) The detected monkeys' faces are all marked in each frame by Faster RCNN and then are input to Tri-attention network for individual identification.

(C1) Tri-AI identifies and names all monkeys known in the database.

(C2) If Tri-AI finds a new monkey face, a new name is then given and automatically added to the database.

56 incidences (48 faces cannot be detected because they were shaded and all the faces that were falsely detected had only a partial face as judged by S.G.). We demonstrated that Faster RCNN (Figure S1) can be used to detect other species as long as it was provided with a sufficiently robust labeling file. We tested 1,500 images of three species: 500 images of golden snub-nosed monkeys, 500 images of Tibetan macaques (*Macaca thibetana*), and 500 images of individually known tigers (*Panthera tigris*). We checked the tested images one by one and obtained a high face detection level for each species: 91.13% for golden snub-nosed monkeys, 97.71% for Tibetan macaques, and 97.70% for tigers. The full precision-recall curve of Faster RCNN for face detection of golden snub-nosed monkeys is shown in Figure 4. It can be seen that the face detection performance for golden snub-nosed monkeys by Faster RCNN was high due to the upper right convex effect. We also provide samples of false face detection cases (Figure 5).

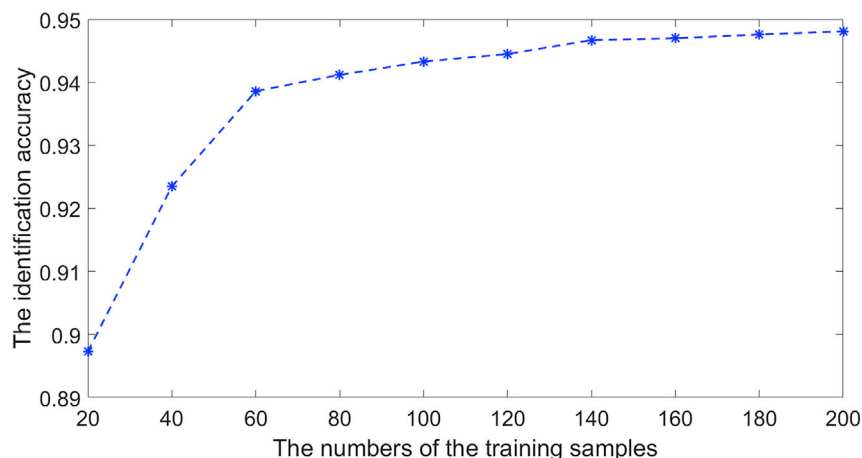


Figure 3. The Relationship between the Identification Accuracy and the Training Sample Numbers

The accuracy is improving with increasing in training samples for each individual, and the inflection point number of training images is around 60 for each individual of golden snub-nosed monkey.

In addition, Faster RCNN can be pre-trained by the labeled images of golden snub-nosed monkey, and then fine-tuned by small-scale labeled images of other species. For example, we pre-trained Faster RCNN with 1,200 labeled images of golden snub-nosed monkey, then fine-tuned the pre-trained model by 500 labeled images of Tibetan macaque, and we can get an average face detection accuracy of 96.00%, compared with 97.7% obtained by Faster RCNN trained on 1,200 labeled images of Tibetan macaque without pre-training.

Face Recognition for Images

We determined the identity of detected faces using Tri-AI and tested its ability to correctly assign individual identities. To do this, we used our image dataset of 41 primate species from wild and captive animals in China (102,399 facial images of 1,040 individuals) (related to [Transparent Methods](#)). We selected 17 primate species from the dataset that had more than 19 identified individuals and tested Tri-AI's accuracy using non-repetitive random sampling for five times. Tri-AI obtained an average identification accuracy of 93.58% (90.14%–96.27%; [Figure 1](#) and [Table S1](#)). This accuracy can be achieved rapidly as Tri-attention network ([Figures S2–S6](#)) accomplished individual identification of 1,000 facial images in 20 s. There were 24 species (No. 18–41, related to [Table S1](#)) in the dataset with fewer than 12 individuals. Despite having few individuals Tri-AI ([Figures S7–S16](#)) could identify species with an accuracy of 97.58% (a total of 7,883 facial images).

Training and Testing Samples

We randomly selected 60% of the facial images of all species as training samples, 10% of the images as validation samples, and the remaining 30% of the images as the tested samples. We detected the monkeys' facial images from all the captured images, and removed facial images with high similarities. To determine the samples size needed for training the Tri-AI system, we made a correlation analysis between the number of training samples and the identification accuracy. The identification accuracy was improving with increasing training sample numbers for each individual, whereas the inflection point number of training image was approximately 60 for reliably identifying an individual from its face images (see in [Figure 3](#)).

Face Detection and Individual Identification in Videos

We used 10 videos of golden snub-nosed monkeys as the test data with 22 individuals and evaluated the success rates of detection and identification frame by frame. The Tri-AI system obtained 98.70% face detection accuracy, 92.01% individual identification accuracy, and 87.03% identification accuracy for new individuals. For the individual identification accuracy, we used 864 frontal facial images in 10 videos; we found that 795 facial images could be correctly identified, but the remaining 69 facial images were incorrectly identified as wrong individuals.

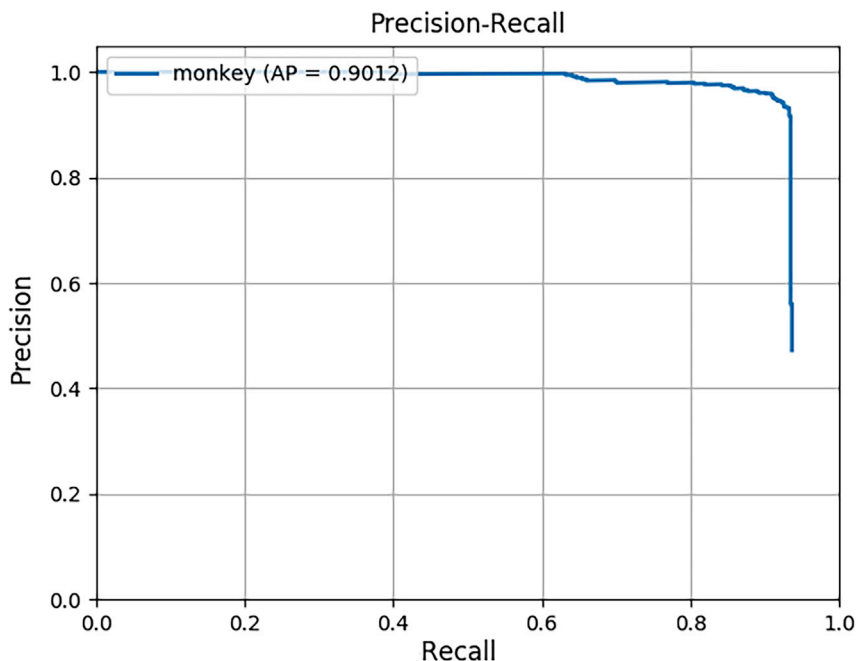


Figure 4. Precision-Recall Curves of Face Detection by Faster RCNN on Golden Snub-Nosed Monkeys

DISCUSSION

Our objective was to develop a system for identifications of individuals with deep network models (Tri-AI), which could quickly detect, identify, and track individuals from video or still-framed images from multiple species. To accomplish this, our system needed to be fed and learn tremendous prior knowledge. We built a dataset that contained 102,399 images of 1,040 individuals, whose identity was known, across all subfamilies of non-human primates and four carnivore species. Primates were chosen as the main testing species because advanced face identification techniques had been developed in human. Four carnivore species were chosen, as they do not have typical human-like faces.

Tri-AI can deal with color images and gray images (including night vision images [Figure S17]). As a result, Tri-AI can do real-time monitoring in the day or at night, meeting research needs of many research programs. The night vision images of golden snub-nosed monkeys and the individual recognition model have been also made publicly available in the database AFD (Animal Face Database): <http://dx.doi.org/10.17632/z3x59pv4bz.2>.

To extend the applications of Tri-AI, we examined Tri-AI on meerkats (*Suricata suricatta*), lions (*Panthera leo*), red panda (*Ailurus fulgens*), and tigers (*P. tigris*) (see in Figure 1), and Tri-AI had identification accuracy of 90.13%, 93.55%, 92.16%, and 94.38%, respectively. For the night vision images of golden snub-nosed monkeys, Tri-AI achieved an identification accuracy of 92.03% (related to Transparent Methods). This suggests that Tri-AI is generally effective for individual identification of a broad range of mammals.

The key to individual identification is to extract the effective features from the relatively fixed parts of a species, and the face is an effective and discriminating feature. Tri-AI has good performances in terms of face detection and identification for primates with their human-like faces and even for carnivores with their more divergent facial features. However, we have found that Tri-AI cannot obtain high identification accuracy for the animals with non-typical or odd face structure. For instance, the face of elephants with their long trunk is rather different and facial images of each individual have large variations due to their unstable and moving trunks, which leads to lower accuracy in elephant face recognition (Tri-attention network had an accuracy of 0.54 for 7 elephants with 459 images).

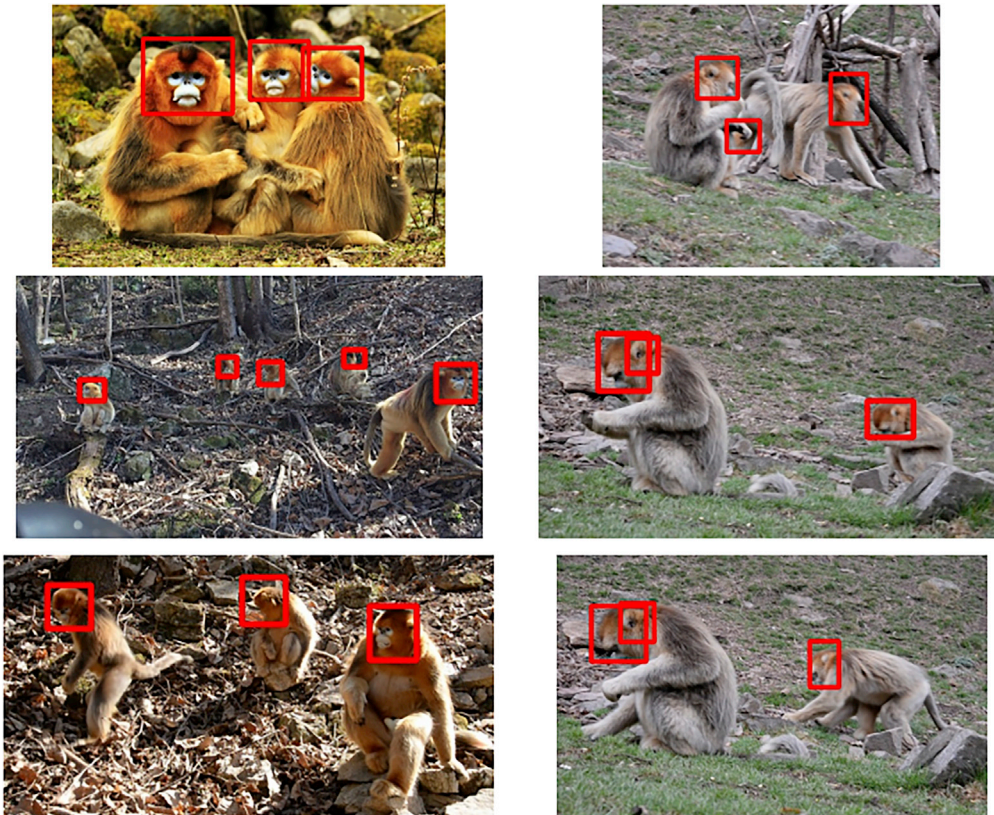


Figure 5. The Detected Results of Golden Snub-Nosed Monkeys by Faster RCNN

It is typically hard to capture the facial images of animals in the wild as they intend to avoid humans and their living environment often makes it very difficult to clearly capture their faces (related to [Transparent Methods](#)). Therefore, based on the existing commonly used image acquisition equipment, including cell phones, digital SLR cameras, and camera traps, we have developed basic strategies, parameters, and standards for capturing animal facial images needed to build the facial image dataset that can be widely applicable among species. Interspecific variations in the accuracy and efficiency of Tri-AI are caused by many factors, including the structures of network models, the number and the diversity of samples, and the quality of images.

Research into individual recognition has mainly focused on how to design a robust network model for complex data. The traditional machine learning methods use the artificial features and classifiers to identify individuals, and these methods can deal with the task of individual identification on relatively small datasets. Their ability to improve their performance is limited when the number of images becomes large. In contrast, DL methods can automatically extract the effective features from images and have superior abilities with large image datasets. These methods require large computational resources, and training DL methods is time-consuming. For individual identification, the CNN models are commonly used to extract the features from the whole images, but ignore some distinguishing parts, which may have significant distinguishing features. Therefore, we designed a new DL model with three channels of attention mechanism to improve the performance of individual recognition. Our Tri-attention network in the system Tri-AI achieved the individual identification accuracy of 94.11% at a speed of 50 images per second.

Nevertheless, when the image samples of each individual increase beyond approximately 60 images, the improvements in performance increase only slightly. Therefore, we attempted to develop individual recognition methods with higher recognition performance when the number of individuals is large and keeps increasing.

In general, the promising performance of Tri-AI provides us confidence that this approach can be used to automatically identify and geo-track individuals of many wildlife species. The operational use of this technology will improve wildlife research, conservation, and management.

Limitations of the Study

Although our database contains 102,399 images across 41 primate species and 6,562 images across 4 carnivore species, there were 24 primate species in the dataset with fewer than 12 individuals. Therefore, we need further to capture the images of more individuals of those species with fewer individuals in our dataset.

Resource Availability

Lead Contact

Songtao Guo (Email: songtaoguo@nwu.edu.cn) takes responsibility for the Lead Contact role.

Materials Availability

The codes of Tri-AI were written in C++ and Python and compiled using g++ on the operation system Ubuntu 14.04. Tri-AI can run on the workstation with Intel Xeon(R) CUP E5-2650 V4, Graphics: GeForce GTX 1080 8G, RAM: 64GB, and Storage: 1TB; Tri-AI also can run on the computer servers with higher hardware configuration.

Data and Code Availability

All the original animal facial images, the night vision images, the test videos, and the related testing models have been released publicly at the database AFD (Animal Face Database): <http://dx.doi.org/10.17632/z3x59pv4bz.2>.

SUPPLEMENTAL INFORMATION

Supplemental Information can be found online at <https://doi.org/10.1016/j.isci.2020.101412>.

METHODS

All methods can be found in the accompanying [Transparent Methods supplemental file](#).

ACKNOWLEDGMENTS

This work was supported by National Natural Science Foundation of China, 31872247, 31672301, 31270441, 31730104, 61973250; Strategic Priority Research Program of the Chinese Academy of Sciences, XDB31000000; Natural Science foundation of Shaanxi Province in China, 2018JC-022, 2016JZ009; National Key Programme of Research and Development, Ministry of Science and Technology, 2016YFC0503200; and the Shaanxi Science and Technology Innovation Team Support Project, 2018TD-026. We thank Prof. Patrick A. Zollner and Prof. Ruliang Pan for their comments and suggestions on the initial manuscript; we thank Shenzhen Safari Park and Shaanxi Academy of Forestry for their support; we thank Dr. Ya Wen for helping make the graphic abstract.

AUTHOR CONTRIBUTIONS

S.G. and P.X. designed the research and wrote the manuscript. G.S., C.A.C., B.L., Q.M., D.F., and X.C. contributed to the improvement of our ideas and to the revision of the manuscript. S.G., P.X., G.H., H.Z., Y.S., and Z.S. carried out the data collection and research experiments.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: October 14, 2019

Revised: June 19, 2020

Accepted: July 22, 2020

Published: August 21, 2020

REFERENCES

- Arzoumanian, Z., Holmberg, J., and Norman, B. (2005). An astronomical pattern-matching algorithm for computer-aided identification of whale sharks *Rhincodon typus*. *J. Appl. Ecol.* 42, 999–1011.
- Burghardt, T., Thomas, B., Barham, P.J., and Calic, J. (2004). Automated Visual Recognition of Individual African Penguins. Fifth International Penguin Conference, Ushuaia, Tierra del Fuego, Argentina, 1–10.
- Burghardt, T., and Calic, J. (2006). Analyzing animal behaviour in wildlife videos using face detection and tracking. *IEEE Proceedings-vision, Image Signal Process. (P-vis IMAGE Sign)* 153, 305–312.
- Chu, W., and Liu, F. (2013). An Approach of Animal Detection Based on Generalized Though Transform (ICCNCE (Atlantis Press)), pp. 117–120.
- Crouse, D., Jacobs, R.L., Richardson, Z., Klum, S., Jain, A.K., Baden, A.L., and Tecot, S.R. (2017). LemurFacID: a face recognition system to facilitate individual identification of lemurs. *BMC Zoolog.* 2, 1–14.
- Ernst, A., and Küblbeck, C. (2011). Fast Face Detection and Species Classification of African Great Apes. In 8th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS) (IEEE), pp. 279–284.
- Fernandezduque, M., Chapman, C.A., Glander, K.E., and Fernandezduque, E. (2018). Darting Primates: steps toward procedural and reporting standards. *Int. J. Primatol.* 39, 1009–1016.
- Finch, N., and Murray, P. (2003). Machine Vision Classification of Animals. 10th Annual Conference on Mechatronics and Machine Vision in Practice (MMVIP) (Perth, Australia), 9–11.
- Freytag, A., Rodner, E., Simon, M., Loos, A., Kuhl, H.S., and Denzler, J. (2016). Chimpanzee Faces in the Wild: Log-Euclidean Cnns for Predicting Identities and Attributes of Primates. In Proceedings of German Conference on Pattern Recognition (GCPR) (Springer), pp. 51–63.
- Hiby, L., Lovell, P., Patil, N., Kumar, N.S., Gopalaswamy, A.M., and Karanth, K.U. (2009). A tiger cannot change its stripes: using a three-dimensional model to match images of living tigers and tiger skins. *Biol. Lett.* 5, 383–386.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep Residual Learning for Image Recognition. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (IEEE), pp. 770–778.
- Huang, G., Liu, Z., Der Maaten, L.V., and Weinberger, K.Q. (2017). Densely Connected Convolutional Networks. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (IEEE), pp. 2261–2269.
- Hou, J., He, Y., Yang, H.B., Connor, T., Gao, J., Wang, Y.Y., Zeng, Y.C.H., Zhang, J.D., Huang, J.Y., Zheng, B.C.H., and Zhou, S. (2020). Identification of animal individuals using deep learning: a case study of giant panda. *Biol. Conserv.* 242, 108414.
- Karanth, K.U. (1995). Estimating tiger *Panthera tigris* populations from camera-trap data using capture-recapture models. *Biol. Conserv.* 71, 333–338.
- Kumar, S., Singh, S.K., Singh, R., and Singh, A. (2017). Animal Biometrics: Concepts and Recent Application. *Animal Biometrics* (Springer), pp. 1–20.
- Krizhevsky, A., Sutskever, I., and Hinton, G.E. (2012). Imagenet Classification with Deep Convolutional Neural Networks. In Neural Information Processing Systems (NIPS) (Curran Associates Inc), pp. 1097–1105.
- Lahiri, M., Tantipathananandh, C., Warungu, R., Rubenstein, D.I., and Berger-Wolf, T.Y. (2011). Biometric animal databases from field photographs: identification of individual zebra in the wild. In 1st ACM international conference on multimedia retrieval, 6 (ACM), pp. 1–8.
- Loos, A., and Ernst, A. (2013). An automated chimpanzee identification system using face detection and recognition. *EURASIP J. Image Vid.* 49, 1–17.
- Nathan, R. (2008). An emerging movement ecology paradigm. *Proc. Natl. Acad. Sci. U S A* 105, 19050–19051.
- Norouzzadeh, M.S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M.S., Packer, C., and Clune, J. (2018). Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proc. Natl. Acad. Sci. U S A* 115, E5716–E5725.
- Ren, S., He, K.M., Girshick, R., and Sun, J. (2015). Faster r-cnn: towards real-time object detection with region proposal networks. *Adv. Neural Inf. Process. Syst.* 91–99.
- Schofield, D., Nagrani, A., Zisserman, A., Hayashi, M., Matsuzawa, T., Biro, D., and Carvalho, S. (2019). Chimpanzee face recognition from videos in the wild using deep learning. *Sci. Adv.* 5, eaaw0736.
- Simonyan, K., and Zisserman, A. (2015). Very Deep Convolutional Networks for Large-Scale Image Recognition (ICLR), pp. 1–14.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2015). Going deeper with convolutions. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (IEEE), pp. 1–9.
- Swanson, A., Kosmala, M., Lintott, C., and Packer, C.A. (2016). Generalized approach for producing, quantifying, and validating citizen science data from wildlife images. *Conserv Biol.* 30, 520–531.
- Tweed, D., and Calway, A. (2002). Tracking Multiple Animals in Wildlife Footage. In IEEE International Conference on Pattern Recognition (ICPR) (IEEE), pp. 24–27.
- Wang, B.S., Wang, Z.L., and Lu, H. (2013). Facial similarity in Taihangshan macaques (*Macaca mulatta tcheliensis*) based on modular principal components analysis. *Acta Theriol. Sin.* 33, 232–237.
- Wichmann, F.A., Drewes, J., Rosas, P., and Gegenfurtner, K.R. (2010). Animal detection in natural scenes: critical features revisited. *J. Vis.* 10, 1–27.
- Xu, Q.J., and Qi, D.W. (2008). Parameters for texture feature of *Panthera tigris altaica* based on gray level co-occurrence matrix. *J. Northeast For. Univ.* 37, 125–128.
- Zeppelzauer, M. (2013). Automated detection of elephants in wildlife video. *EURASIP J. Image Vid.* 46, 1–23.
- Zhu, W., Drewes, J., and Gegenfurtner, K.R. (2013). Animal detection in natural images: effects of color and image database. *PLoS One* 8, e75816.
- Zeiler, M.D., and Fergus, R. (2014). Visualizing and Understanding Convolutional Networks (European conference on computer vision (ECCV)), pp. 818–833.

iScience, Volume 23

Supplemental Information

Automatic Identification of Individual Primates with Deep Learning Techniques

Songtao Guo, Pengfei Xu, Qiguang Miao, Guofan Shao, Colin A. Chapman, Xiaojiang Chen, Gang He, Dingyi Fang, He Zhang, Yewen Sun, Zhihui Shi, and Baoguo Li

Supplementary Information - Trasparent Methods

Table S2. The ablation experiments for Tri-attention network, related to Figure 1 and 2. The individual identification accuracies obtained by Object level attention model and Partial level attention model and Tri-attention network for 17 primate species with more than 19 individuals and four non-primate species.

Animal species	The channel of Object level attention model	The channel of Partial level attention model	Tri-attention network
Weeper Capuchin (<i>Cebus olivaceus</i>)	0.9010	0.7802	0.9305
Black-capped Capuchin (<i>Cebus apella</i>)	0.8904	0.9045	0.9627
Rhesus Macaque (<i>Macaca mulatta</i>)	0.9101	0.9322	0.9126
Common Chimpanzee (<i>Pan troglodytes</i>)	0.8357	0.8231	0.9358
Common Squirrel Monkey (<i>Saimiri sciureus</i>)	0.8928	0.8968	0.9083
Green Monkey (<i>Chlorocebus sabaeus</i>)	0.8449	0.8992	0.9477
Mandrill (<i>Mandrillus sphinx</i>)	0.8000	0.9182	0.9294
Golden Snub-nosed Monkey (<i>Rhinopithecus roxellana</i>)	0.8845	0.9340	0.9412
François' langur (<i>Trachypithecus francoisi</i>)	0.8571	0.8571	0.9014
Olive Baboon (<i>Papio anubis</i>)	0.9264	0.7341	0.9341
Ring-tailed Lemur (<i>Lemur catta</i>)	0.9193	0.7897	0.9457
De Brazza's Monkey (<i>Cercopithecus neglectus</i>)	0.9326	0.9299	0.9508
Patas Monkey (<i>Erythrocebus patas</i>)	0.9255	0.9184	0.9562
Greater Spot-nosed Monkey (<i>Cercopithecus nictitans</i>)	0.9014	0.8309	0.9543
Northern White-cheeked Gibbon (<i>Nomascus leucogenys</i>)	0.9132	0.9077	0.9342
Tibetan Macaque (<i>Macaca thibetana</i>)	0.9430	0.8134	0.9504

Crab-eating Macaque (<i>Macaca fascicularis</i>)	0.8965	0.7471	0.9337
Meerkat (<i>Suricata suricatta</i>)	0.9010	0.7796	0.9013
Lion (<i>Panthera leo</i>)	0.9192	0.6563	0.9355
Red Panda (<i>Ailurus fulgens</i>)	0.9170	0.8264	0.9216
Tiger (<i>Panthera tigris</i>)	0.9256	0.8807	0.9438

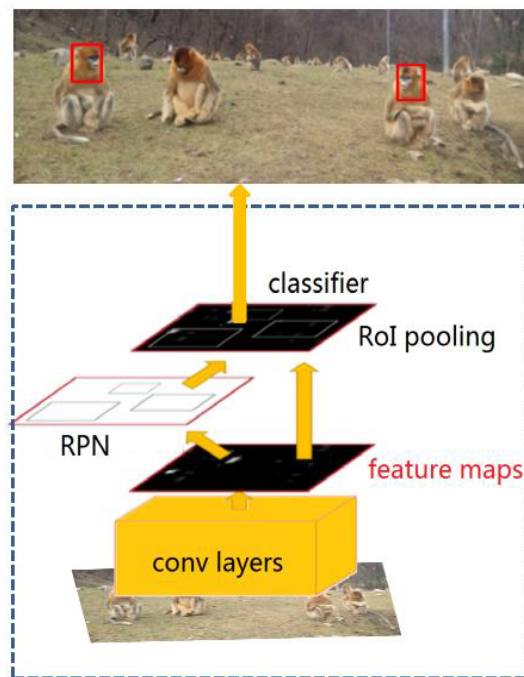


Figure S1. The face detection of golden snub-nosed monkeys by Faster RCNN, related to Figures 2, 4 and 5. We can detect the monkeys' faces when they show their faces. However, the other two individuals lower their heads in this moment, which results in Faster RCNN cannot find their faces.

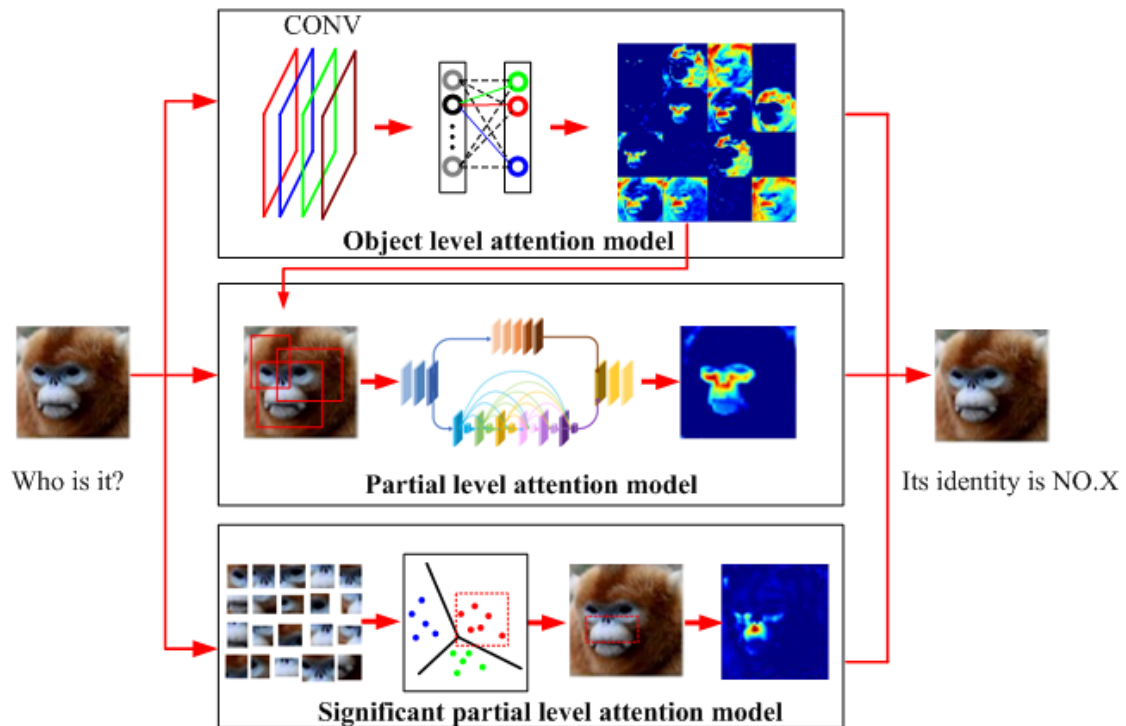


Figure S2. The framework of Tri-attention network, related to Figure 2. Object level attention model is used mainly to extract the features of the whole facial image. Significant partial level attention model pays attention to capture the local feature of a relatively fixed facial skin area. Partial level attention model focuses on the specific feature of a restricted smaller area, and different individuals would have their own specific facial areas.

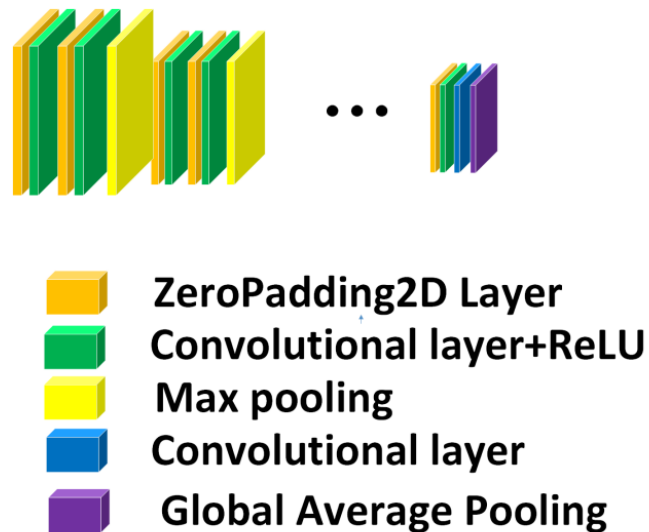


Figure S3. Object level attention model, related to Figure 2.

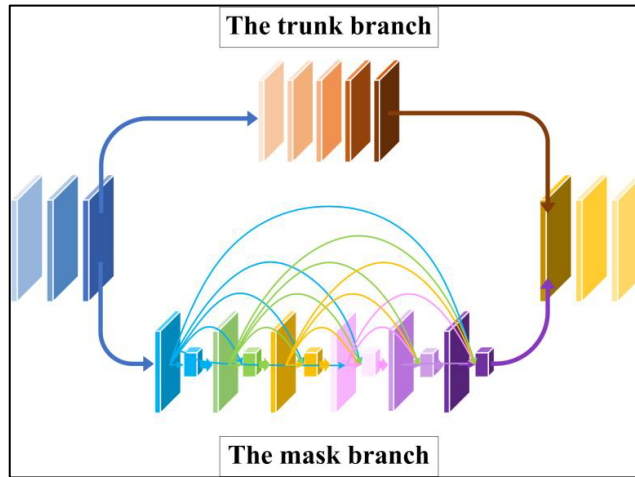


Figure S4. The structure of residual network with attention mechanism, related to Figure 2.

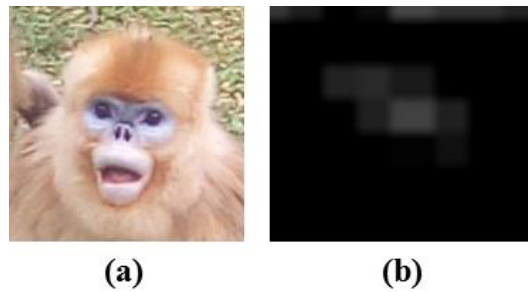


Figure S5. The feature map shows the model pays attention to feature extraction from the “skin area”, related to Figure 2. (a) The facial image of golden snub-nosed monkey. (b)The feature map for the attentional facial region of the “skin area”.

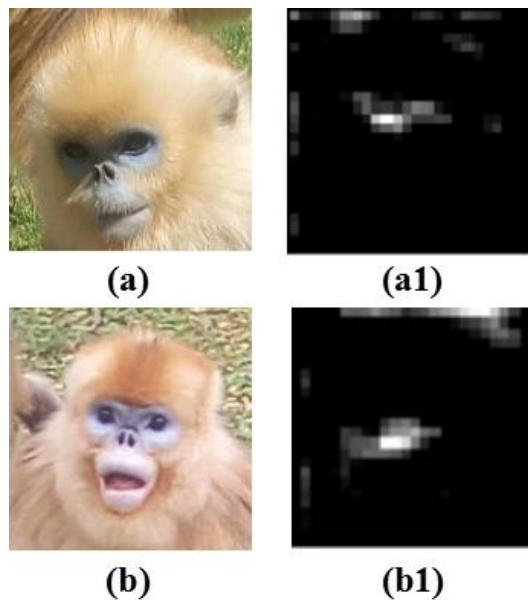


Figure S6. The original facial images of golden snub-nosed monkeys and the feature maps extracted by significant partial level attention model, related to Figure 2. (a) The facial image of golden snub-nosed monkey. (a1) The feature map for the attentional facial region. (b) The facial image of golden snub-nosed monkey. (b1) The feature map for the attentional facial region.

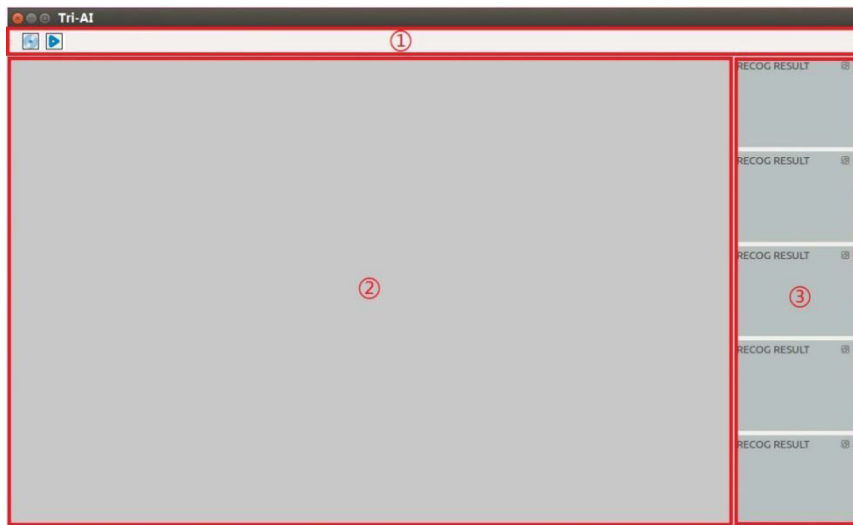


Figure S7. The main page of Tri-AI system, related to Figure 2.

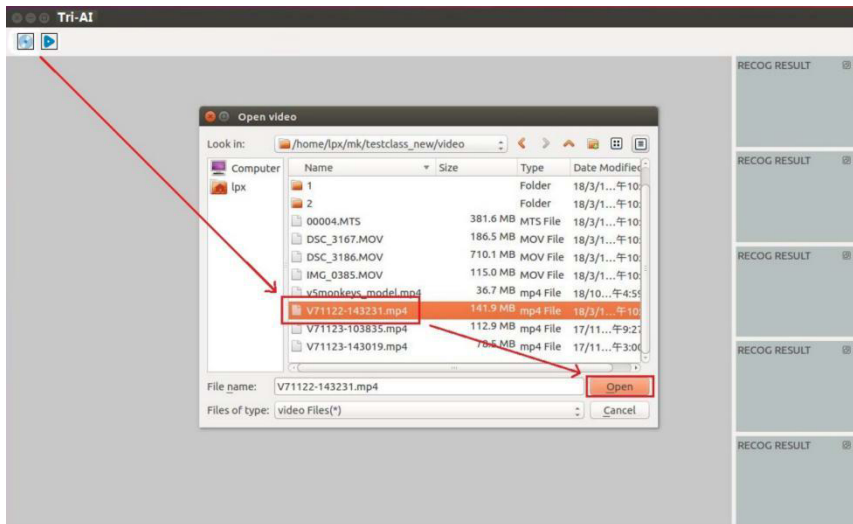


Figure S8. Select an image or a video, related to Figure 2.

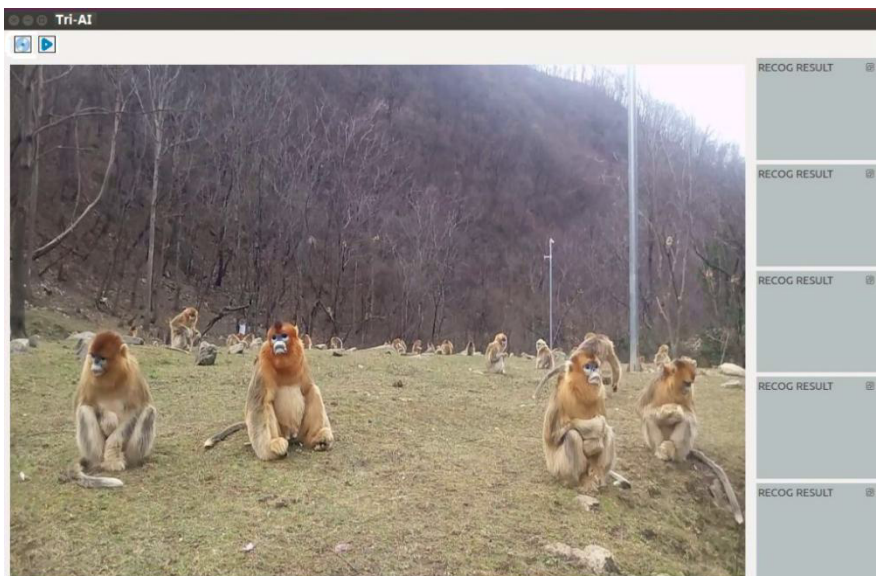


Figure S9. The first frame of a video shown in the area ②, related to Figure 2.

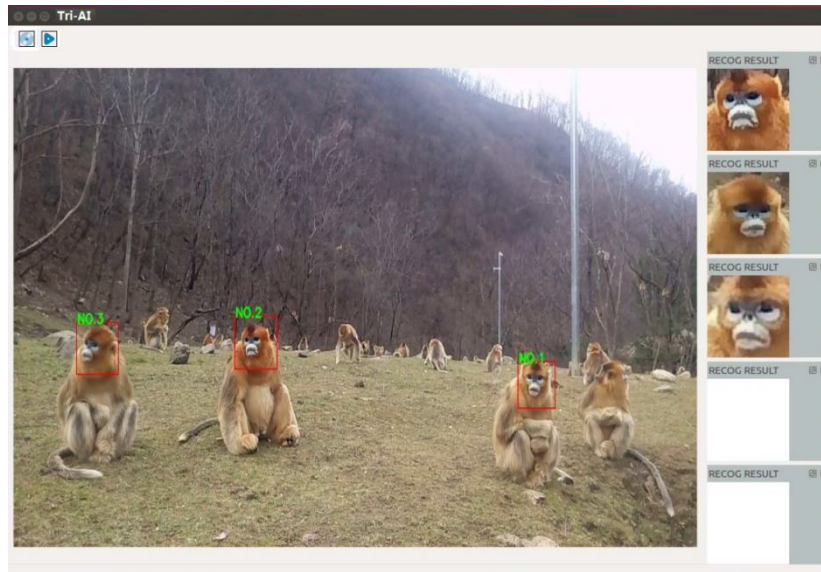


Figure S10. The detection and recognition results, related to Figure 2.

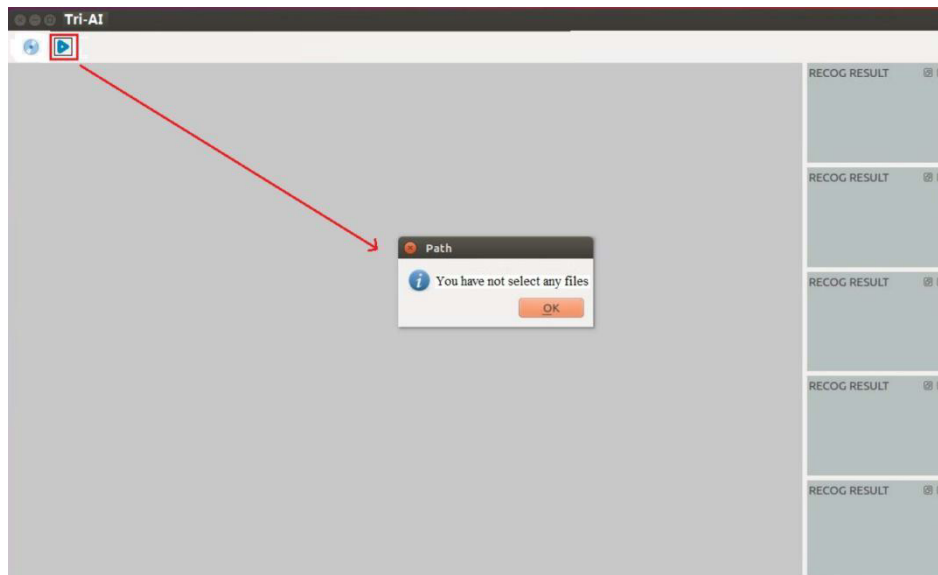


Figure S11. Exception handling, related to Figure 2.

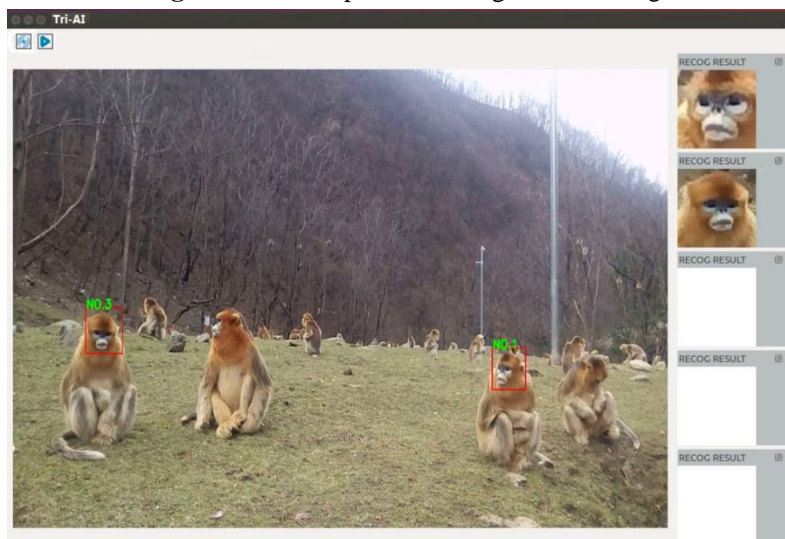


Figure S12. The results of system testing, related to Figure 2.

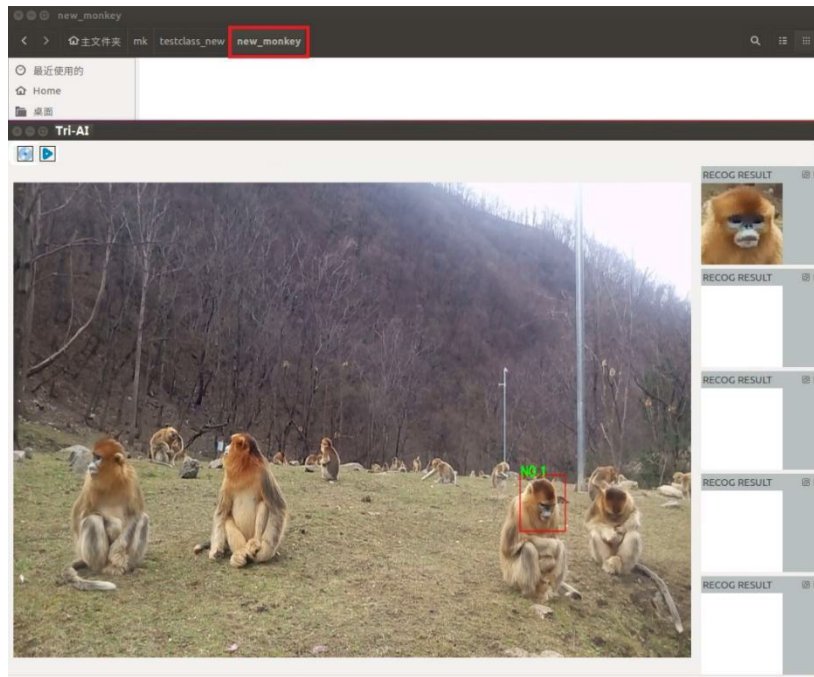


Figure S13. The recognition of a new individual. If the animals have their training facial images in the dataset, their identities can be recognized, related to Figure 2.

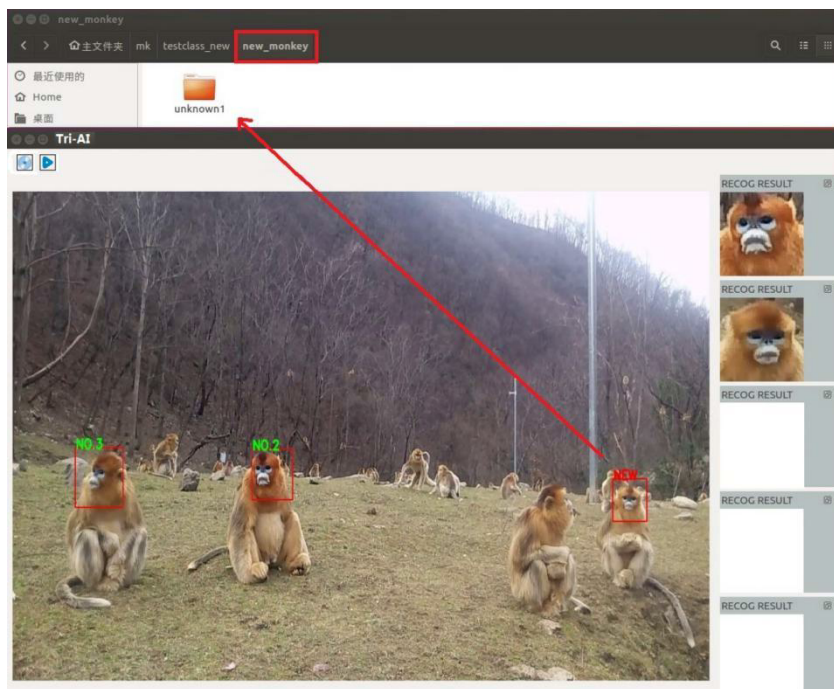


Figure S14. The recognition of a new individual. The new individual is identified, and its facial images are added to the dataset automatically, related to Figure 2.



Figure S15. Usable images of golden snub-nosed monkeys, related to Figure 1.

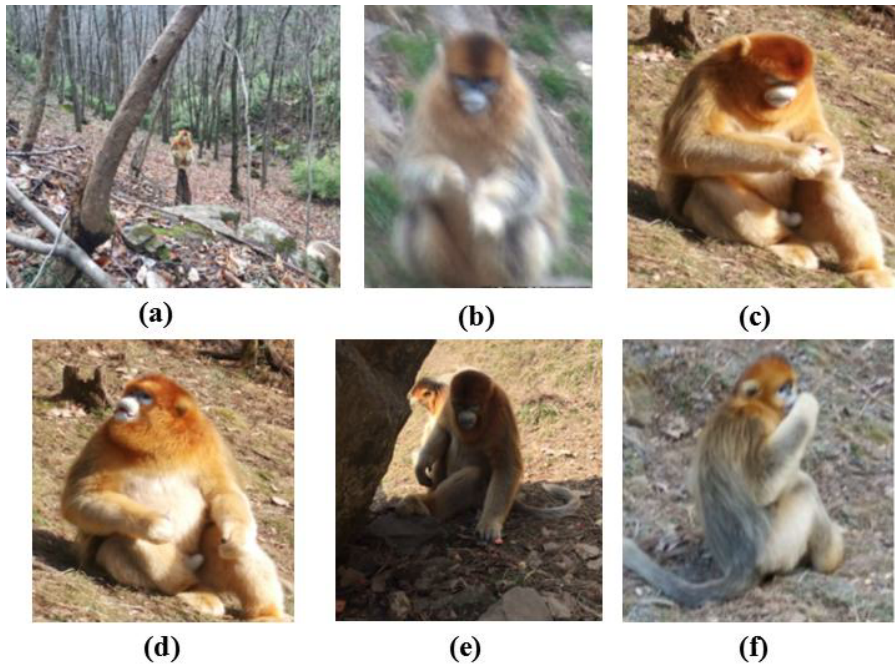


Figure S16. The unusable images under different conditions, related to Figure 1. (a) Too small face. (b) Motion blur. (c) Big angle of the face. (d) Big angle of the face. (e) Under shadow. (f) Covered by its hand.

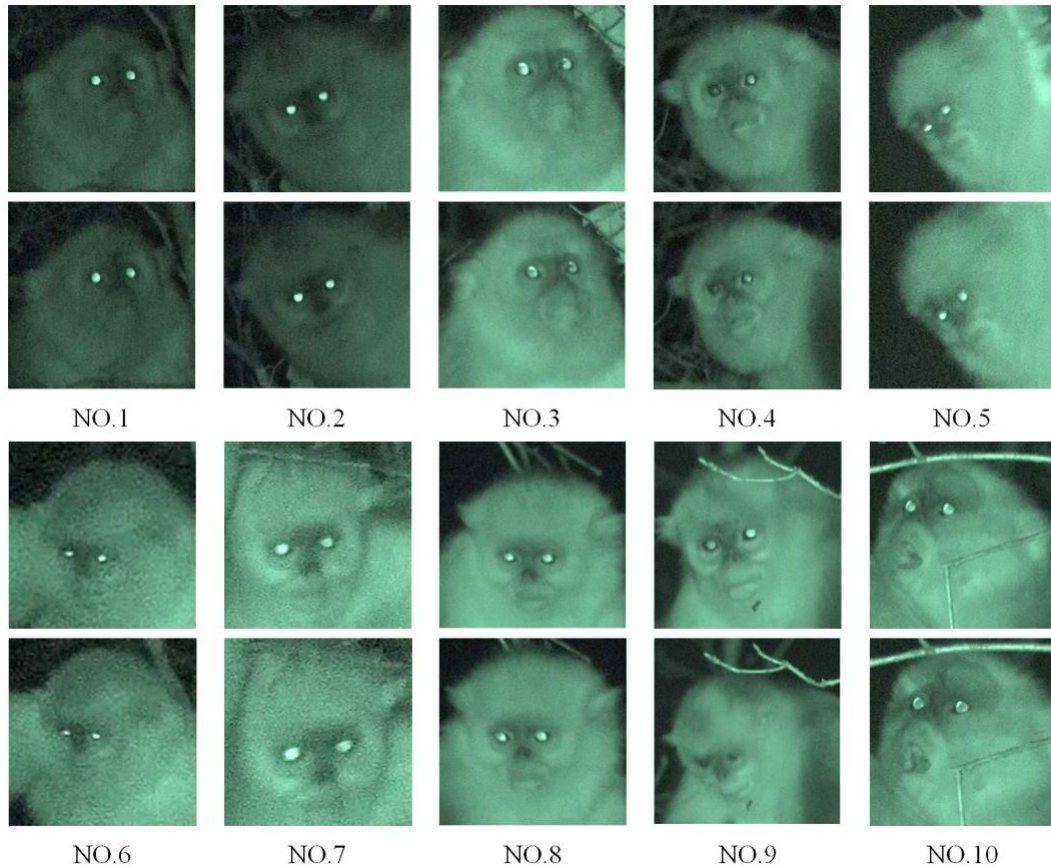


Figure S17. The night facial images of golden snub-nosed monkeys, related to Figures 2 and 3.

S1. Methods

We developed the system Tri-AI using deep learning techniques to automatically detect and identify individual animals by their faces captured from images or videos. Faster Region-based Convolutional Neural Network (RCNN) (Ren et al., 2015) was used to detect the animals' faces and these detected facial images were identified by Tri-attention network, which are described in the following two subsections.

Before Tri-AI could be used, we had to obtain training and testing images and conduct data preprocessing (Related to SI S4). The prepared dataset was used to train Faster RCNN and Tri-attention network to develop the models of face detection and identification. We tested the performances of Tri-AI with a large number of images (Table S1). We have explored the implications of animals' facial image recognition, and the face detection and identification of golden snub-nosed monkeys (*Rhinopithecus roxellana*) with Tri-AI was divided into two major steps.

(1) Face detection by Faster RCNN

To identify individual, we first detected the animals' faces from the images or videos with Faster

RCNN (Ren et al., 2015), which is an effective tool for object detection and has the advantages of simple labeling and high efficiency and accuracy. For object detection, Faster RCNN involves two modules. One module is a deep fully convolution network that provides the proposed regions of the faces, and the other module is a detector to identify the final face regions from the proposed regions. Faster RCNN is a single unified network that is particularly useful for detecting the faces of golden snub-nosed monkeys (Figure S1). After sending the images to a CNN model for feature extraction, the Region Proposal Network (RPN) generates the proposal regions based on the extracted features (Related to SI S6), and approximately 300 proposal regions were generated from each image. These regions were then mapped to the feature map of the last convolution layer of CNN. Finally, the Region of Interest (ROI) pooling layers enabled each ROI to generate a fixed-size feature map and the combination of classification and regression was performed using Softmax Loss and Smooth L1 Loss to locate the animals' faces.

Before detecting the animals' faces, the original parameters in Faster RCNN were trained with sufficient images of the animals. Our labeling task was accomplished by marking their faces on the images with the software of labeling. The trained Faster RCNN could be used to detect the animals' faces from images or videos. For example, we used 1,200 images of golden snub-nosed monkey for labeling to detect the monkeys' faces by Faster RCNN.

To verify the effectiveness of Faster RCNN, we made several test experiments. The trained Faster RCNN was used to detect the faces of golden snub-nosed monkeys, Tibetan macaques and tigers from their images, which were captured by cameras. For the test images, we randomly selected 500 images for each species, and counted the number of the detected faces from the detection results to obtain the detection accuracy of 91.1% for golden snub-nosed monkeys, 97.7% for Tibetan macaques, and 97.7% for tigers. Furthermore, ten videos of 22 golden snub-nosed monkeys were used as the test data to evaluate the success ratio of detection and identification. Then we checked the detection and recognition results frame by frame and obtain 98.70% of face detection accuracy. In addition, we used the full precision-recall curves to evaluate the Faster RCNN detector, and this resulted in some false detection cases.

We tracked animals' faces in the small areas around the locations of their detection in former frames. For a given video, Faster RCNN was used to detect the face regions on the first frame of the video, and then for the following frames. The face regions were detected by Faster RCNN only on

smaller image areas based on the former detected face locations.

(2) Individual identification of primates by Tri-attention network

For fine-grained recognition in the field of image recognition, the deep network models with attention mechanism would focus on a certain part of the image, and the extracted features may be the key features, which have the major differences from those of other classes. Thus, the performance of the classification algorithms can be improved. In this paper, a novel attention network model with three channels (Tri-attention network) was designed for the task of fine-grained individual recognition using primate facial images (Figure S2).

The Tri-attention network has three channels, and each channel has its own attention network models, including object level attention model to extract the global facial features, partial level attention model to focus on the facial region of interest, and significant partial level attention model to select a key smaller facial region for extracting the specific facial features. These three channels with different network structures extract three types of facial features, which are then combined for animal individual identification. Therefore, Tri-attention network focuses on the facial features in different levels simultaneously.

These three models constitute three channels of our Tri-attention network, and each channel pays attention to different areas of the primate facial images at different levels and extracts the corresponding features. The features extracted by each channel are fused by cascade operation and used for classification by softmax.

To prove the validity of the roles of each channel, we conducted ablation experiments on the Tri-attention network. While, the significant partial level attention model is used to extract the fine-grained features on only one small facial region, and these features need to be combined with the global facial feature for individual identification. Therefore, this single stream network is suitable for animal face recognition. We made individual identification experiments using object level attention model and partial level attention model and Tri-attention network. The results show different roles of different channels in Tri-attention network (Table S2). The Tri-attention network has improved performance improvements compared to each single channel.

(3) The Detailed Information of Tri-attention Network

In the Tri-attention network model, the object level attention model mainly deals with global features of the facial image. The partial level attention model is concerned with a relatively fixed local

area in the facial image, and the significant partial level attention model focuses on a specific smaller area selected from the facial images of each individual, with the locations of the areas being different for different individuals.

(A) Object level attention model

The features extracted by convolutional neural network (CNN) contain rich spatial information after multiple convolution and pooling layers. These feature maps need to be transformed into feature vectors by full connection layer. However, some effective information may be lost in this process, and the generalization ability of the model would be reduced. To solve this problem, the class activation mapping (CAM) strategy was used in object level attention model to reduce the impact of the fully connected layer on the loss of feature information. The global average pooling (GAP) layer was used in CAM instead of the fully connected layer in CNN, and the GAP layer calculated the average values of all the pixels in each feature map and converted all these average values into feature vectors for classification. The advantage of global average pooling layer is that there are no parameter settings, which means that the effective feature information can be preserved better, and the risk of overfitting be reduced.

Specifically, CAM was used to replace the fully connected layer with GAP. After GAP, the average value of each feature map in the last convolution layer was obtained. The final features were obtained by a weighted sum. The final extracted features were used as input in softmax for classification. We classified the numbers of regions which had an important impact for individual identification based on CAM. Due to the complex structure of Grad-CAM, we used the basic CAM here by considering the efficiency and complexity of the model.

The network structure of object level attention model has 12 convolution layers, 4 maximum pooling layers, and 1 GAP layer (Figure S3). The combination of the ZeroPadding2D layers and the convolutional layers can keep the sizes of the feature maps unchanged after convolution operations, while the sizes of the feature maps become a half of the original ones only after the maximum pooling operations. Correspondingly, the number of filters becomes twice as the former number after each pooling operation. Since the dimension of the output feature vector is related to the number of feature maps extracted by the previous convolution layer, we have added a convolution layer in front of the global average pooling layer, which ensures that the dimension of the resulting feature vector is consistent with the number of categories.

(B) Partial level attention model

Unlike human facial images, the primates general have more hair on their faces, and the shape and texture information of their hair are easily affected by many factors. However, the facial skin areas are relative invariant with respect to morphology and texture, which makes these skin areas more distinguishable. Therefore, partial level attention model mainly focuses on the facial skin area and minimizes the effects of hair and background on the recognition results. To achieve this, a residual network based on the attention mechanism was used in partial level attention model.

Different from object level attention model, the residual network with attention mechanism is a hierarchical structure. GAP is used to pay attention to extracting the feature maps of the entire facial images in object level attention model, while the residual network with attention mechanism used in partial level attention model is a stackable network structure, which can hierarchically pay attention to specific areas of the facial images. Our designed network model focuses on the skin area in primate facial images.

The structure of residual network with attention mechanism (ResNet-AM) is shown in Figure S4, and it mainly contains two branches, including the trunk branch and the mask branch. The trunk branch is a convolutional neural network consisting of three convolution layers and three ZeroPadding2D layers, and the output feature maps are $T_i(x)$. The mask branch processes the input feature maps to obtain the attention feature maps $M_i(x)$ with the same dimension as $T_i(x)$, and the weights of $M_i(x)$ are normalized. The final characteristics of local areas are expressed as:

$$H_i(x) = T_i(x) \times M_i(x)$$

The partial level attention model is used to focus on the feature extraction of the skin areas. Figure S5 shows the feature map in this model, and one can see this model mainly attends to "skin areas" for feature extraction. This partial level attention model is achieved by maximum pooling layer, which has two functions of reducing the dimensions of the extracted features and removing the redundant features of the images. After maximum pooling layer for down sampling and other layers, the obtained feature maps contain the important features from the original images for individual identification.

(C). Significant partial level attention model

Significant partial level attention model mainly focuses on the most significant area (the key facial

area) of the image, which consists of two steps: area selection and feature extraction. To generate the candidate areas for the facial images, Graph-Based segmentation (GBS) algorithm (Ohayon et al, 2013) are used to divide a facial image into a number of small candidate areas, and the color similarity S between the i^{th} and j^{th} small areas can be calculated with the following equation:

$$S = \sqrt{(r_i - r_j)^2 + (g_i - g_j)^2 + (b_i - b_j)^2}$$

The two areas with the most color similarity were then merged, and the two original sub-areas were removed. This process was performed until there were no original sub-areas left to be merged, so that we could obtain a set of candidate facial areas. For the set of candidate facial areas, we used a trained classification model to choose the key area for each individual, and we chose the merged facial areas with highest classification accuracy as the final key facial areas for the corresponding individual.

The significant partial level attention model can extract the highly separable features of key facial regions, which is reflected in the operation of key facial region selection from a large number of separated facial regions through a trained model (AlexNet) (Krizhevsky et al, 2012). Finally, the selected key facial region for each monkey's facial image was used for feature extraction. After the operation of convolution pooling, only the golden snub-nosed monkey 's eyes and nose were retained (Figure S6).

(D) Feature extraction from the key facial areas.

For the selected key facial areas, VGG16 network is used for feature extraction. VGG16 has 16 convolutional layers, the size of each convolution kernels is 3×3 , the convolution step is 1, and the number of convolution kernels increases from the initial 64 to 128, 256, and 512. The size of the convolution kernels in pooling layer is 2×2 . VGG16 has three fully connected layers. Since our Tri-attention network consists of three channels, the parameters of the full connection layer of VGG16 would inevitably affect the overall performance of the whole network. Therefore, the parameters of the three fully connection layers are set as 1024, 512, and 128. We used VGG16 rather than residual blocks in the significant partial level attention model to extract the key facial regions by considering the model efficiency.

S2. The operation instruction of Tri-AI

(1) Computer Requirements for Tri-AI

Tri-AI runs on a small workstation with Intel Xeon(R) CUP E5-2650 V4 (2.20GHZ \times 24), Graphics: GeForce GTX 1080 8g, RAM: 64GB, and Storage: 1TB. The codes of Tri-AI were written in C++ and

Python, and were compiled using g++. Its graphical interface was designed based on Qt, which was developed by the Qt Company in 1991 as an application development framework of cross-platform C++ graphical user interface. The supportive operation system was Ubuntu 14.04 (or an improved version of Ubuntu).

(2) Specific Operations of Tri-AI to Detect and Recognize Faces

Step 1: Start the application. Input the filename of our software in the command line, and press the enter button to start the software. This step intends to accomplish the tasks of animal face detection, identification, and tracking.

Step 2: The main interface will appear when our software finishes loading. The main interface has four areas (Figure S7). The area ① is functional area. Here, one can select the test images or videos and simply click on the “start” button to detect the monkeys’ faces from your selected images or videos. The area ② shows the selected images or videos and the corresponding detection and recognition results of monkeys’ faces. The area ③ presents the image of the face of the animal in the database that is most similar to the image that the software is being asked to identify.

Step 3: In the area ①, one can click on the first button to open an image file or a video file (Figure S8). When one clicks on “open” button, the selected image or video will appear in the area ② (Figure S9).

Step 4: In the area ①, one can click on “Start” button (the second button), and the system will start to detect the monkeys’ faces from the images or the videos and try to recognize who they are. If these monkeys show their frontal faces or profile faces not too far from directly ahead, the individuals will be identified and their corresponding information will be provided (e.g., IDs or names) and an image of each face will be shown in the area ③. These face images are selected from the database, and are the most similar to the monkeys’ faces the software is being asked to recognize (Figure S10).

(3) Exception handling

If one did not selected the test images or videos and simply click “Start”, a dialog box “You have not select any files” will pop up. In this case, click “OK” (Figure S11), and reselect a test image or video. One then can continue to make the following operations.

(4) System testing

Using golden-snub nosed monkeys as an example, when a monkey’s face is directed forward (frontal) or if the animal provides a profile face that is not too far off facing forward, Tri-AI will be

able to recognize this individual (Figure S12). Contrasting the same monkey presenting profile faces with different angles, the system can still get accurate recognition (Figure S10). If animals have their training facial images in the dataset, their identities can be recognized in the images or videos. Otherwise, they are identified as new individuals, and their facial images are added to the dataset automatically (Figures. S13 and S14).

S3. Basic Strategies, Parameters, and Standards for Capturing Animal Facial Images

Compared with acquiring human face images, it is much harder to capture animal facial images because both their living environment and behaviors are much uncontrollable. Therefore, we need to develop specific strategies to acquire different animals' facial images according to their habits and environment. In the process of capturing animals' facial images, it is desirable to use uniform norms and strategies to ensure the animals' facial images obtained within and between studies have high accuracy, regularity, and consistency. Based on our experience, we present guidelines for capturing facial images that are applicable to the theoretical and applied questions researchers are asking and the quality of images that can captured under often difficult field conditions.

(1) The basic strategies for capturing facial images:

(A) If there are few individuals in the population/group being sampled, potentially an observer can see all individuals at once and get independent facial images one by one. This is the ideal situation that can often be done in zoo settings or with very small groups in the wild. However, in most cases, there are many individuals in one area and they cannot be captured in a single image or a set of temporally close images. Therefore, it is necessary for at least two or more people to cooperate, while one taking pictures; others can track the individual's movement to ensure the same individual is not repeatedly photographed.

(B) If there are hundreds of individuals in the population/group or if it is difficult to see all individuals at once as they are somewhat cryptic, it is hard or impossible to separate them artificially, and difficult to distinguish them from memory without potentially months of field work. But if one can temporarily mark individuals, with a harmless color ink for example, then it is possible to mark animals after its photos have been taken. However, all images must be obtained within a short time frame, so that the mark does not fade.

(C) For most species, it will not be possible to mark the animals as they are difficult to see and avoid humans. However, for some group living animals that live in stable small groups, like ring-tailed

lemurs (*Lemur catta*), a number of observers can each select a single or a few animals or a particular area occupied by the group and can quickly take each individual's facial image at the same time. Once Tri-AI is used to assign individuals, it may be possible to study the images to learn the unique individual features.

(D) For species that are cryptic (e.g., many social carnivores) or avoid human observers (e.g., forest elephants (*Loxodonta cyclotis*), lowland gorillas (*Gorilla gorilla*), or giant forest hogs (*Hylochoerus meinertzhageni*)), but repeatedly come to specific locations, such as salt licks or waterholes, it may be possible to set multiple camera traps to simultaneously capture images of all or most of the group.

Above all, it is necessary to ensure that the repeated images, which are taken for the same individual, are unique individuals.

(2) Technical requirements for capturing primates' facial images

It is necessary to develop the basic parameters and standards to ensure that the collected images can meet the necessary requirements and we recommend the following specifications of the equipment, image quality, and quantity.

(A) Equipment: For most mammals, mobile phones take images of sufficient quality when the animals are less than 10 meters from the observer. However, if the animals are farther than 10 meters, high-quality single lens reflex (SLR) cameras which can take more than 5 images/sec images with over 20 million pixels are recommended.

(B) Image resolution: The primate faces in the images should have 50×50 or more pixels. This is not a difficult standard to meet with most cameras if the animal is relatively close.

(C) Image clarity: The primate faces should be clear, in focus, and not blurred caused by the animals' movement or camera person's hand tremor.

(D) The angles of the face: Images of the face should be taken within the face angles of 30° (or between -15° and +15°) and it is desirable to increase the diversity of the samples from different face angles.

(E) Light: We should avoid taking photos when these animals are in the strong light or under the shadow.

(F) Cover: It is not desirable to have large areas of the animals' faces covered, but if some parts of the face are obscure the images may still be useful.

(G) The number of the images: It is best to have at least one hundred facial images per individual.

(H) Mark the face image data: After the needed number of suitable images are taken, each image must be accurately labeled manually and the mark information should include species, individual identity, age, sex, etc..

(3) Examples of the captured images:

We illustrate the basic parameters, strategies, and standards for capturing facial images using golden snub-nosed monkey.

(A) Usable images: The facial images should be clear, have more than 50×50 pixels, have the face primarily facing forward (i.e., have small side angles), and have appropriate light exposure. The images shown in Figure S15 are suitable for use in our Tri-AI system.

(B) Unusable images. The following images cannot be used by the Tri-Ai System. The monkey's face in Figure S16 (a) is too small, (b) has motion blur, (c) the angles of the monkeys' face is inappropriate as the animal is looking down, (d) the angles of the monkeys' face is inappropriate as the animal is looking to the side, (e) the monkey's face is under shadow, and (f), the monkey's face is covered by its hand.

Using these basic parameters, strategies, parameters, and standards, we were able to capture the facial images of more than forty species of primates. For some of these species, we obtained more photographed more than 1,040 individuals (e.g., golden snub-nosed monkeys in Qinling Nature Reserve and *Macaca thibetanas* in Tangjiahe Nature Reserve). Also, we captured the images of many species of primates from 18 zoos, including those in Beijing, Shanghai, Dalian, Weihai, Qindao, Ningbo, and Shijiazhuang. The images of most golden snub-nosed monkey individuals were captured in a number of different days, while the images of other animals were obtained in single days in the zoos. We got 102,399 facial images of primates.

S4. Primate Facial Image Dataset

Establishing facial image dataset provides an important foundation for the current and future scientific research. To collect sufficient facial images for multiple wildlife individuals, we traveled to 18 zoos in 16 cities and 6 nature reserves in China and took images and videos of 1,040 individuals of 41 primate species, and selected 102,399 facial images.

The collections of these images were fraught with challenges, such animals living in a complex arboreal environment, animals moving very fast, individuals engaging in a wide range of activities,

light conditions changed frequently, and shadow often fell across the animal's face. To acquire these images and videos, we explored many approaches for image acquisition, image selection, facial image extraction, and labeling the images. Based on this experience, we make the following recommendations for future users.

(1) Image acquisition

The acquisition of animal facial images is much more difficult than capturing human facial images, as neither the environment nor the animal's behavior can be controlled. The user community of such approaches needs to develop standardized methods. The basic specifications and related requirements for the acquisition of primate facial images are explained in S3 above. Our early image-capturing work was done with the golden snub-nosed monkey of the Qinling Mountain in China and we collected facial images of 196 individuals. We also travelled to the Tangjiahe Nature Reserve and collected more than 2,000 facial images of 30 wild Tibetan macaques. Further, we captured images from 18 zoos in 16 cities including Chengdu, Taiyuan, Shijiazhuang, Beijing, Qinhuangdao, Dalian, Weihai, Qingdao, Jinan, Shanghai, Ningbo, Hangzhou, Suzhou, Wuxi, Nanjing, and Xi'an. In total, we obtained 102,399 facial images of 1040 individuals, from 41 species.

(2) Image data processing

After obtaining the images, the facial areas of each primate individual must be extracted from the images using manually screenshot methods or face detection methods, such as Faster RCNN. Primates are gregarious animals, which typically result in individual images showing multiple individuals. In this case, the most challenging task for us is individual recognition as many animals look so similar. Therefore, we had to carefully identify which animal the facial image was from. We did this by checking the position of each in all the images or repeatedly comparing the differences between them through visual interpretation. Therefore, most of the facial images were identified manually using screenshots to avoid data confusion. This manual method of face detection can only be applied to those images which have one or very few individuals. We found it best to have all the facial images in our dataset to be square, contain almost all the facial information of the individuals, and have as little background as possible.

(3) Construction of Primate Facial Image Datasets

When image acquisition and processing are completed, researchers need to build a dataset of facial images. For our dataset, we classified all the facial images by species and the facial images of

each individual were assigned a unique label. In our data set, there were differences in the number of individuals among species (e.g., golden snub-nosed monkey had 43,304 facial images of 227 individuals and each individual had an average of 191 facial images, *Cebus apella* had 3,026 facial images of 43 individuals, and the average number of facial images for each *C. apella* individual was 70). The average number of facial images for each individual in our dataset was 99. Finally, our image dataset had a total of 1,040 individuals and 102,399 facial images. A detailed list of all primate species, the individual numbers in each species and the total number of facial images of each species in our dataset are given in Table S1.

Our images of golden snub-nosed monkeys were captured in the wild. Thus there was greater variance in the images, compared to the other species where images were obtained from zoos, typically on single days. The dataset has been made publicly available at the database:

(AFD(Animal Face Database): <http://dx.doi.org/10.17632/z3x59pv4bz.2>).

S5. Identification of golden snub-nosed monkeys using taken at night

In total, 581 facial images taken of 24 golden snub-nosed monkeys at night were used for face recognition. Here, 60% were used as training samples, 10% for validation, and 30% as test samples. The identification accuracy of these night images was 92.03% and the night vision images for 10 golden snub-nosed monkeys are shown in Figure S17.

S6. A list of technical terms

CNN: A convolutional neural network (CNN or ConvNet) is a class of deep, feed-forward artificial neural networks that have successfully been applied to analyzing visual imagery. An early development of CNN for facial recognition was developed by Lawrence and colleagues (Krizhevsky et al, 2012), but it has been subsequently improved (Lecun et al, 2015).

RCNN: Regions with CNN features. RCNN is an object detection model based on the CNN network. It uses the selective search method to get 2,000 candidate boxes with different sizes, and then CNN network is used to extract the regional features for the classification of the objects and background.

Fast RCNN is a fast object detection model based on multi-task learning. In the training phase, Fast RCNN is 9 times faster than RCNN. During the testing phase, Fast RCNN is 213 times faster than RCNN (Girshick, 2015).

RPN: Region Proposal Network. RPN can obtain a series of object proposals from arbitrary images. It provides the suspected areas for object detection models.

Faster RCNN: This is a newer version of an object detection model of fast RCNN. The network structure mainly includes RPN and fast RCNN. RPN is used to select the suspected areas where the objects may exist in the image. The fast RCNN is used to identify whether these suspected areas actually are the objects.

ResNet: Residual Network. The residual network uses the convolutional layers to perform residual learning to solve the problem of performance degradation when the network goes deeper. The existing ResNet models have ResNet-20, ResNet-34, ResNet-51, ResNet-101, ResNet-152, and other improved related network models.

Shallow ResNet: Shallow ResNet is an improved deep residual network proposed in this paper. Shallow ResNet simplifies the network structure based on the traditional ResNet. The convolution layers are increased to form new types of residual blocks, which can improve feature learning ability. Shallow ResNet is a good solution to the problem that the traditional network models with fewer layers are hard to extract the deep features of the animal facial images, but the features extracted by deep network models lose more information.

Supplemental References

Girshick, R. (2015). Fast r-cnn. IEEE international conference on computer vision (ICCV) (IEEE, New York). 1440-1448.

Ren, S., He, K. M., Girshick, R., and Sun, J. (2015). Faster r-cnn: Towards real-time object detection with region proposal networks. Advances in Neural Information Processing Systems (NIPS). 91-99.

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. Conference on Neural Information Processing Systems (NIPS) (Curran Associates Inc). 1097-1105.

Lecun, Y., Bengio, Y., and Hinton, G. E. (2015). Deep learning. Nature, 521, 436-444.

Ohayon, S., Avni, O., Taylor, A. L., Perona, P., and Egnor, S. R. (2013). Automated multi-day tracking of marked mice for the analysis of social behaviour. J Neurosci Methods, 219, 10-19.