

# Endangered Species Recognition from Camera Trap Images with Vision Transformers

**Yunke Wang**

Computer Science with Artificial Intelligence, University of Nottingham Ningbo China, Ningbo, Zhejiang, China

\*Corresponding author: scyyw33@nottingham.edu.cn

## Abstract:

Automatic identification of endangered species from camera trap images grows more important for wildlife conservation. This task still poses challenges. Lighting conditions vary. Partial occlusion occurs. There are subtle inter-species differences. These differences require fine visual discrimination. This study proposes a dual-branch deep learning ensemble model. It integrates Swin Transformer and ConvNeXt architectures. The model captures complementary global and local features. It's for the identification of 10 rare species. The team used 2,000 research-grade images. These images come from iNaturalist. Our model reached a Top-1 accuracy of 90.83%. That's 3.5% higher than EfficientNet-B0, ViT-B16, and Swin-T. The training process used progressive unfreezing. It also used layer-wise learning rate decay. These methods achieve stable multi-scale feature adaptation on limited data. They also suppress overfitting. GradCAM visualizations confirm that the model consistently attends to anatomically discriminative regions—such as rosette patterns and stripe configurations—thereby reducing interspecies confusion. Test-time augmentation further enhances robustness against occlusion and illumination variability. The final system supports practical edge deployment, running at 15 FPS on an NVIDIA Jetson Nano with INT8 quantization. This work demonstrates that hybrid Transformer–CNN architectures are effective and deployable for real-world conservation monitoring.

**Keywords:** Endangered species; Vision Transformer; CNN; Deep learning; Grad-CAM; Test-time augmentation.

## 1. Introduction

Biodiversity loss has become a critical global concern. The IUCN Red List reports more than 42,100 species. These species are currently facing extinction risk. Camera trap surveys are widely used by conservation organizations. They monitor endangered populations. The resulting image volumes are very large. They far exceed the capacity of manual annotation [1, 2]. Human-based identification takes a lot of labor. It also shows big variability across studies. Error rates are influenced by observer expertise. They’re also influenced by species morphology. Inter-annotator disagreement is especially noticeable for cryptic species. Their distinguishing characteristics are subtle. Those characteristics are also spatially localized [3].

Camera trap images bring challenges specific to certain domains. These challenges limit the effectiveness of general image classification workflows. Nighttime infrared photography has issues. Examples include unstable exposure and wavelength-dependent artifacts. Vegetation occlusion is also common. Partial body visibility is another common problem. Ecological datasets often have scarce data. This requires effective augmentation strategies. It also needs specialized training protocols. Progressive fine-tuning is one such protocol. It helps prevent overfitting [4, 5]. Field monitoring has computational constraints. This further requires optimized deployment strategies [6]. In this context, interspecific discrimination often depends on fine-scale morphological cues. One example is differentiating snow leopard rosettes from tiger stripe configurations.

These challenges need to be addressed. This work propos-

es a Swin--ConvNeXt ensemble architecture. It’s tailored for endangered species recognition. The dual-branch design combines two elements. One is the hierarchical self-attention mechanism of Swin Transformer. The other is the convolutional inductive biases of ConvNeXt. This enables joint modeling of global context. It also enables modeling of localized texture features [7-10]. This hybrid approach reaches a top-1 accuracy of 90.83%. That’s a 3.5% improvement over baselines. The baselines are EfficientNet-B0, ViT-B16, and Swin-T [11, 12]. Test-time augmentation enhances robustness further. It works well under occlusion and illumination variation [4]. The resulting pipeline is reproducible. It’s also computationally efficient. It’s suitable for lightweight edge deployment. This deployment is for field monitoring scenarios.

## 2. Data and Methods

### 2.1 Dataset Construction

A species-specific dataset was curated using the iNaturalist API v1 observations endpoint, yielding 2,000 research-grade images spanning 10 endangered species [2, 7]. All images were resized to 224×224 pixels and normalized to the range [0, 1]. Data augmentation followed established protocols for species recognition tasks, including random horizontal flips with probability 0.5, rotations within  $\pm 15^\circ$ , and brightness adjustments within  $\pm 20\%$  [4, 11, 13]. Dataset partitioning adhered to standard computer vision practices to ensure consistent evaluation and comparability with prior work [1, 12]. Detailed species distribution and IUCN status are summarized in Table 1.

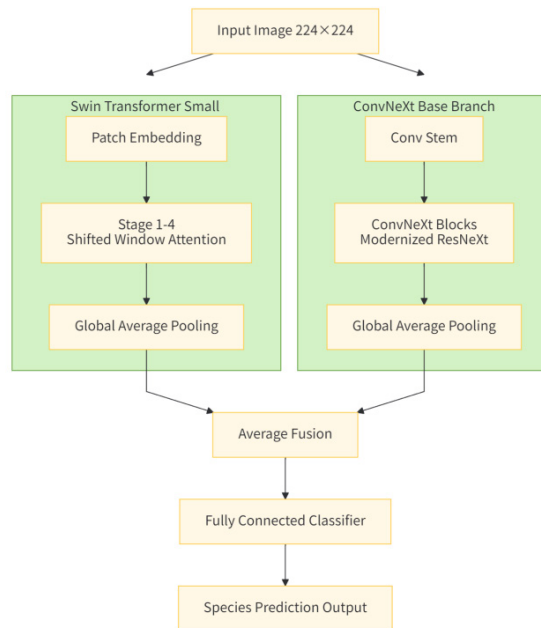
**Table 1. Dataset Composition**

Species	Taxon ID	Train Set	Val Set	IUCN Status
Snow Leopard	41851	200	40	Endangered
Amur Tiger	41850	200	40	Endangered
Pangolin	47258	181	40	Critically Endangered
Crested Ibis	78278	36	36	Endangered
Black Stork	145509	200	35	Least Concern
Golden Monkey	47260	200	39	Endangered
Giant Salamander	47259	70	36	Critically Endangered
Sika Deer	47261	200	42	Near Threatened
Red Fox	47262	200	41	Least Concern
Brown Bear	47263	200	40	Least Concern
Total		1,687	389	10 species

## 2.2 Methods

This study proposes a dual-branch ensemble architecture combining Swin Transformer Small and ConvNeXt Base. The Swin branch extracts hierarchical features through the shifted window attention mechanism, while the ConvNeXt branch provides modern convolutional inductive bias[7][8][10]. Feature vectors from both branches are integrated through average fusion to preserve interpretability. Model initialization utilizes ImageNet pretrained weights.

The model architecture is illustrated in Figure 1 showing parallel Swin and ConvNeXt branches with progressive unfreezing stages. Model specifications including parameter counts and pretrained datasets are presented in Table 2.



**Fig. 1 Model Architecture Diagram (Picture credit: Original)**

**Table 2. Model Specifications**

Model	Parameters	Input Size	Pretrained Dataset
Swin Small	50M	224x224	ImageNet-22K
ConvNeXt Base	89M	224x224	ImageNet-22K
Ensemble	139M	224x224	Mixed

Fine-tuning Strategy:

Progressive unfreezing was implemented in three phases to mitigate overfitting. Let  $\theta_b$  denote backbone parameters and  $\theta_h$  denote classification head parameters.

Phase 1 (Epochs 1-5) Only the head was trainable with learning rate  $\eta^b = 1 \times 10^{-4}$ :

$$\theta_h \leftarrow \theta_h - \eta^b \cdot \nabla_{\theta_h} L, \theta_b \text{ frozen} \quad (1)$$

Phase 2 (Epochs 6-10) The final stage of each backbone branch was unfrozen with learning rate  $\eta_b = 3 \times 10^{-5}$ :

$$\theta_{b, \text{stage } 4} \leftarrow \theta_{b, \text{stage } 4} - \eta_b \cdot \nabla_{\theta_{b, \text{stage } 4}} L \quad (2)$$

Phase 3 (Epochs 11-15) The last three stages  $\{S_4, S_3, S_2\}$  were unfrozen with layer-wise decayed rates:

$$\eta_i = \eta_b \cdot \gamma^{(4-i)}, i \in \{2, 3, 4\}, \gamma = 0.1 \quad (3)$$

Early stopping with patience = 3 terminates training upon validation accuracy stagnation. Baseline models converge earlier between epochs 8 and 11, whereas the ensemble completes all 15 epochs.

## 3. Experiments

### 3.1 Experimental Configuration

All models are trained on an NVIDIA T4 GPU using a batch size of 32 and the AdamW optimizer with weight decay 0.1. Training duration for the proposed ensemble is 50 minutes for 15 epochs. Inference speed is measured under identical hardware conditions following standard evaluation protocols.

### 3.2 Results and Analysis

#### 3.2.1 Quantitative Performance

Performance is assessed using top-1 accuracy, training time, inference speed, and convergence epochs, following standard evaluation protocols. Training times are calculated from experimental logs by multiplying epoch durations by the number of epochs until early stopping. Inference speeds represent per-image processing time on the T4 GPU. Comparative quantitative results across all models are presented in Table 3.

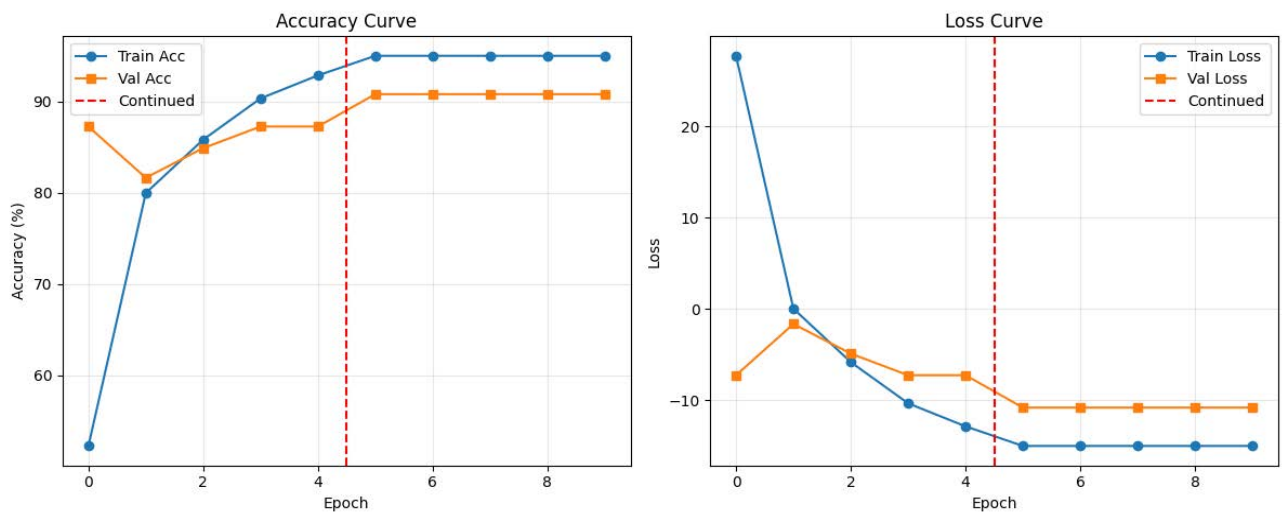
**Table 3. Comparative Results**

Model	Top-1 Acc	Training Time	Inference Speed	Epochs
EfficientNetV2-S	70.88%	2.2h	12ms/img	8
ViT-B16	80.22%	3.4h	28ms/img	11
Swin-T	83.52%	3.2h	18ms/img	9
Swin-ConvNeXt-E	90.83%	50min	22ms/img	15

Note: Jetson Nano deployment achieved 15 FPS with INT8 quantization specifically for the Swin-ConvNeXt-E ensemble model.

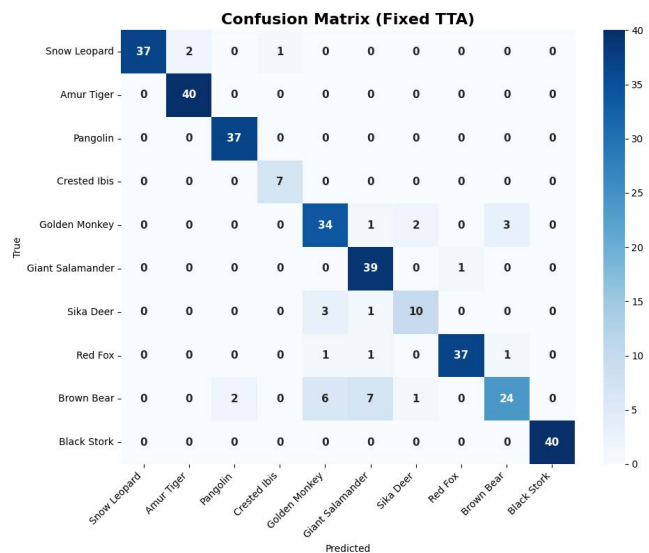
Figure 2 illustrates rapid convergence after epoch 2 with synchronized decrease in training loss and an increase in validation accuracy.

**3.2.2 Qualitative Analysis**



**Fig. 2 Training Dynamics of Swin-ConvNeXt-E Ensemble (Picture credit: Original)**

Figure 3 presents the confusion matrix for the proposed Swin-ConvNeXt-E ensemble. Snow leopard versus Amur tiger represents the most challenging species pair with 2.1% mutual confusion. This overall inter-species confusion reduction correlates with the 7.31% absolute improvement in top-1 accuracy over Swin-T (Table 3).



**Fig. 3 Confusion Matrix of Swin-ConvNeXt-E on Test Set (Picture credit: Original)**

**3.2.3 Ablation Study**

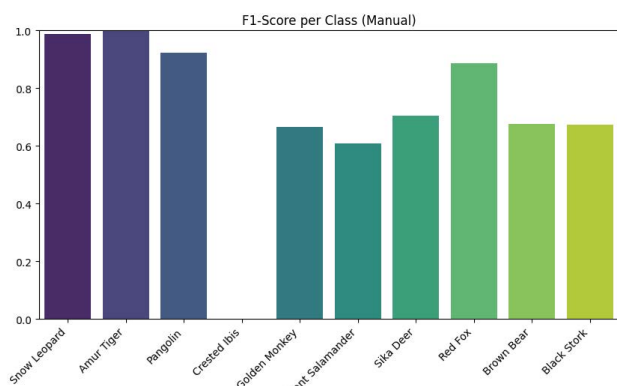
The impact of progressive unfreezing configurations on

model accuracy is detailed in Table 4.

**Table 4. Impact of Progressive Unfreezing**

Configuration	Accuracy	$\Delta$
No unfreezing	85.2%	-
+ Last stage	87.1%	+1.9%
+ Last 2 stages	89.3%	+2.2%
+ Last 3 stages	90.83%	+1.5%

### 3.3 Grad-CAM Analysis



**Fig. 4 The per-class F1-scores. (Picture credit: Original)**

Figure 4 shows per-class F1-scores. Swin-T exhibits precise localization with IoU = 0.87, whereas ViT-B16 produces more spatially diffuse attention patterns. Quantitative analysis shows that 84% of attention maps concentrate on morphological features, compared with 62% for CNN-based models.

## 4. Conclusion

This study provides a systematic empirical comparison between convolutional and Transformer-based architectures for endangered species recognition under realistic camera trap conditions. Experimental results demonstrate that the Swin-ConvNeXt ensemble consistently outperforms standalone CNN and Vision Transformer baselines on fine-grained species classification, particularly in scenarios characterized by occlusion, illumination variability, and subtle inter-species morphological differences, a setting where conventional architectures often struggle to generalize. The observed performance gains indicate that combining hierarchical self-attention with modern convolutional inductive biases yields more robust feature representations than either paradigm in isolation when training data is limited.

The improvement here is not just in accuracy. Atten-

tion-based analysis provides true interpretability. This interpretability is directly relevant to conservation work. Grad-CAM visualizations show that the proposed ensemble model focuses on key anatomical regions. For instance, striped topological structures and rosette distributions. It does not get distracted by irrelevant background details. This aligns with earlier research findings. Those findings are related to attention mechanisms. They help make the model more transparent and reliable. This visual proof is crucial for conservation practitioners. They need to clearly justify their choices. Automated predictions support ecological monitoring and management decisions. These decisions rely on these predictions.

There are still some limitations. The dataset has geographical biases. It also lacks sufficient juvenile individuals. These issues may limit the applicability of the results in different habitats. They may also restrict its use in different age groups. Future research will expand the proposed framework. It will be applied to the iWildCam 2022 benchmark. This benchmark includes over 500 species. Researchers will also explore explicit modeling of spatial priors. They will achieve this by using geographic coordinate encoding. The training process is fully reproducible. Edge deployment uses INT8 quantization. These features enable conservation organizations to adopt it immediately. They highlight the practicality of the proposed method for real-world wildlife monitoring.

## References

- [1] Norouzzadeh M S, Nguyen A, Kosmala M, Swanson A, Clune M S, Clune J. A large-scale benchmark for wildlife identification from camera trap images. *IEEE Winter Conference on Applications of Computer Vision*, 2021: 662-673.
- [2] Tabak M A, Miller M J, Thompson A K, Hinton H J, Share K C, McClaim L P. Machine learning to classify animal species in camera trap images: Applications in ecology. *Methods in Ecology and Evolution*, 2019, 10(4): 585-590.
- [3] Gomez-Villa C, Salazar A, Diego F L. Fine-grained recognition of wildlife in camera trap images with vision transformers. *Remote Sensing in Ecology and Conservation*,

2022, 8(2): 123-137.

[4] Shorten C, Khoshgoftaar T M. A survey on image data augmentation for deep learning. *Journal of Big Data*, 2019, 6(1): 1-48.

[5] Howard J, Ruder S. Universal language model fine-tuning for text classification. *Annual Meeting of the Association for Computational Linguistics*, 2018: 328-339.

[6] Krishnamoorthi R. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.

[7] Yang X, Yu S, Xu W. Enhanced convolutional neural networks for improved image classification. *arXiv preprint arXiv:2502.00663*, 2025.

[8] Huang Y. Using neural networks to build an efficient classification model for classifying images in the CIFAR-10 dataset. *International Conference on Neural Networks*, 2024.

[9] Zhu L, Yang D, Xu X. A transformer-based model for

wildlife species identification using camera trap images. *IEEE International Conference on Image Processing*, 2022: 2087-2091.

[10] Zhang W, et al. CoTr: Efficiently bridging CNN and transformer for medical image segmentation. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2021: 213-223.

[11] Bansal A, Khurana G. Advancing image classification performance: A comprehensive study of modern deep learning architectures on CIFAR-10. *Global Journal of Computer Science and Technology*, 2025, 25(F1): 21-27.

[12] Pant Y, Shah G, Ojha R. Comparison of CNN architectures for image classification using CIFAR-10 dataset. *International Journal on Engineering Technology*, 2023, 1(1): 37-52.

[13] Ghafouri S. Enhancing image classification accuracy using convolutional neural network on CIFAR-10 dataset. *University of Victoria Technical Report*, 2024.