

PROJECT REPORT ON  
**DEEP LEARNING WORKSHOP WITH PYTHON**  
**(CSE3194)**

# **Animal Species Detection**

Submitted in partial fulfillment of the  
requirements for the award of the degree of

## **BACHELOR OF TECHNOLOGY**

Submitted by:

<b>AAKASH KUMAR DWIVEDI</b>	<b>2241004257</b>
<b>DEBI PRASAD DASH</b>	<b>2241004159</b>
<b>HARSHIT DASH</b>	<b>2241004162</b>



Center for Artificial Intelligence & Machine Learning  
**Department of Computer Science and Engineering**  
Institute of Technical Education and Research  
**Sikhsa 'O' Anusandhan**  
**(Deemed to be University)**  
Bhubaneswar

May, 2025

# Declaration

---

We hereby declare that the project report titled “**Animal Species Detection**” is our own work, carried out under the guidance of Mr. Jagseer Singh. We have not plagiarized any content and have duly cited all references.

---

(AAKASH KUMAR DWIVEDI)

---

(DEBI PRASAD DASH)

---

(HARSHIT DASH)

# Acknowledgement

---

I wish to extend my deepest appreciation to everyone who has helped and mentored me along the way on this project.

Let me first and foremost express my warmest appreciation for the project supervision/mentorship of **Profs. Jagseer Singh** without whose wise suggestions, unwavering encouragement, and thoughtful constructive criticism this direction of this undertaking would not have been possible.

I also express my gratitude towards the Department of Computer Science and Engineering and the faculty at **ITER, SOA University**, for sharing resources and environments that helped in the completion of this project.

Lastly, I thank my family and friends for their encouragement and support throughout the entire process.

# Table of Contents

## Contents

<b>Declaration .....</b>	<b>i</b>
<b>Acknowledgement .....</b>	<b>ii</b>
<b>Table of Contents.....</b>	<b>iii</b>
<b>Abstract .....</b>	<b>v</b>
<b>Chapter 1 – Introduction.....</b>	<b>1</b>
1.1. Background and motivation.....	1
1.2. Problem statement .....	1
1.3. Objectives.....	1
1.4. Scope .....	1
1.5. Organisation of the Report.....	2
<b>Chapter 2 – Literature Review.....</b>	<b>3</b>
2.1. Existing methods and models .....	3
2.2. Comparison with your approach .....	3
<b>Chapter 3 – System Design and Methodology.....</b>	<b>4</b>
3.1. Dataset description: .....	4
3.2. Exploratory Data Analysis.....	4
3.3. Preprocessing steps .....	4
3.4. Model architecture.....	5
<b>ZFNet:</b> .....	5
<b>VGG16:</b> .....	6
<b>GoogleNet (Inception v1):</b> .....	7
3.5. Description of the Algorithms.....	8
<b>Algorithm: Animal Species Detection – GoogleNet Specific</b> .....	9
<b>Chapter 4 – Implementation Details.....</b>	<b>11</b>
4.1. Programming languages, frameworks .....	11
<b>Programming Language:</b> .....	11
<b>Deep Learning Frameworks:</b> .....	11
<b>Utility and Support Libraries:</b> .....	11
4.2. Code modules description .....	11

4.3. System Working .....	14
<b>Chapter 5 – Results and Discussion .....</b>	<b>15</b>
5.1. Evaluation metrics .....	15
5.2. Results .....	15
5.3. Discussion .....	20
<b>Chapter 6 – Conclusion &amp; Future Work .....</b>	<b>21</b>
6.1. Summary of outcomes .....	21
6.2. Limitations .....	21
6.3. Future improvements .....	22
<b>References .....</b>	<b>23</b>
<b>Appendices .....</b>	<b>24</b>
Additional code snippets .....	24

# Abstract

---

This project targets automatic species detection and classification of animals with deep learning. Using Convolutional Neural Networks (CNNs), we compare three of the leading architectures—ZFNet, VGG16, and GoogLeNet—on a dataset of 28,000 images covering 10 species of animals. The models were trained with standard and augmented image data for proper learning and generalization. Among the models, VGG16, utilizing transfer learning, produced the best accuracy, and GoogLeNet struck a balance between speed and performance with its inception modules. ZFNet was less accurate but faster. Our results show that deep learning is effective for species recognition and that there are trade-offs between model complexity, accuracy, and computational requirements. The system has potential for real-time and scalable biodiversity monitoring and can provide important assistance to wildlife conservation initiatives.

Keywords: Deep Learning, Convolutional Neural Networks (CNNs), ZFNet, VGG16, GoogleNet.

# Chapter 1 – Introduction

---

## 1.1. Background and motivation

With emerging growth in artificial intelligence, deep learning has emerged as a key technique in image classification automation. The inspiration for this project comes from the increasing requirement for intelligent systems in wildlife monitoring. With urbanization and climate change posing threats to natural habitats, there is a need for immediate development of automated tools to help identify species for conservation purposes. This project will capitalize on the use of deep learning, specifically Convolutional Neural Networks (CNNs), to solve this real-world issue. Through animal species classification based on image recognition, the system can be a crucial tool in tracking biodiversity and ecological studies.

## 1.2. Problem statement

The project aims to develop a deep learning based system capable of classifying animals and predicting its species. The system will leverage Convolutional Neural Networks (CNNs), a powerful type of deep learning architecture designed for image processing tasks. Below is an in-depth breakdown of the development process.

## 1.3. Objectives

- To design and implement CNN models based on ZFNet, VGG16 (with transfer learning), and GoogLeNet (Inception V1) architectures for animal species classification.
- To train these models on the Animals-10 dataset, comprising 28,000 labeled images spanning 10 species, with standardized preprocessing and augmentation.
- To evaluate and compare each model's performance using metrics such as classification accuracy, training duration, inference time, and loss convergence.
- To visualize model performance through accuracy and loss plots across training epochs for effective comparison.
- To conduct a final test evaluation and identify the most effective model architecture for scalable, real-world species detection applications.

## 1.4. Scope

This project is scoped to assess the performance of CNN-based models on a dataset of ~28,000 labeled images of 10 animal species. The tested models—ZFNet, VGG16, and GoogLeNet—are evaluated for their capacity to classify images on the basis of learned features such as texture and shape. Preprocessing methods (resizing, normalization, augmentation) are used to enhance generalization. Transfer learning and the trade-offs between computational complexity and model accuracy are also investigated in the project. While restricted to 10 species, the framework can be scaled up and applied to additional classes or modalities, for example, audio or video data, for increased biodiversity monitoring in real-world settings.

## 1.5. Organisation of the Report

- **Chapter 1** introduces the project, its background, and its motivation, alongside the problem definition, objectives, and scope.
- **Chapter 2** presents a literature review of existing animal classification methods and deep learning techniques.
- **Chapter 3** outlines the dataset details, preprocessing steps, and the architectures of the implemented CNN models.
- **Chapter 4** describes the implementation details, experimental setup, and hyperparameters.
- **Chapter 5** discusses the results obtained, comparative evaluation, and analysis of performance.
- **Chapter 6** concludes the project and discusses limitations and future enhancement possibilities.



# Chapter 2 – Literature Review

---

## 2.1. Existing methods and models

Several studies in recent years have explored deep learning for animal species detection:

- **Mahajan et al.** used CNNs with data augmentation and transfer learning to improve detection in natural environments under varying lighting and occlusion.
- **Geethanjali et al.** deployed lightweight models like MobileNet-SSD V2 for real-time detection on low-resource devices.
- **Ibraheam et al.** applied different R-CNN variants, optimizing region proposals for wildlife datasets.
- **Smith et al.** focused on real-time detection using YOLO on large camera trap datasets.
- **Garcia et al.** and **Zhao et al.** explored semi-supervised learning and adversarial training for improving accuracy with limited data.

## 2.2. Comparison with your approach

- We **compare multiple architectures** (ZFNet, VGG16, GoogLeNet) side-by-side on the same dataset, offering a comprehensive performance analysis.
- Our project emphasizes **inference time vs. accuracy trade-offs**, which is crucial for real-time deployment.
- We utilize **detailed training metrics visualization** (accuracy/loss plots) to study convergence and stability.
- A **standardized dataset and uniform preprocessing pipeline** ensure fair benchmarking.
- The focus is not just on high accuracy, but also on **practicality and scalability** for wildlife conservation applications.

# Chapter 3 – System Design and Methodology

---

## 3.1. Dataset description:

The Animals-10 dataset, sourced from Kaggle, consists of approximately 28,000 high-quality images categorized into 10 animal species: dog, cat, horse, spider, butterfly, chicken, sheep, cow, squirrel, and elephant. The images vary in resolution and background, simulating real-world conditions. Each class contains a balanced number of samples, making the dataset suitable for multi-class image classification tasks. The diversity in lighting, angles, and environments helps train deep learning models for robust and generalized species detection.

## 3.2. Exploratory Data Analysis

Prior to data preprocessing and model training, an exploratory analysis of the dataset was performed to evaluate its quality and suitability:

- **Balanced Classes:** The dataset showed a fairly uniform distribution across all 10 animal categories, minimizing class imbalance issues.
- **Sample Review:** Representative images from each class were examined to observe differences in lighting, orientation, and background complexity.
- **File Verification:** All images were checked to ensure they could be loaded correctly without any file corruption or errors.

This initial analysis confirmed that the dataset is well-prepared for training convolutional neural networks, providing a diverse and consistent set of images to support model learning

## 3.3. Preprocessing steps

1. **Image Resizing:** All images were resized to a fixed dimension of **224×224 pixels** to ensure uniform input size compatible with the CNN architectures (ZFNet, VGG16, GoogLeNet).
2. **Normalization:** Pixel values of images were scaled to the **[0, 1] range** by dividing by 255. This normalization helps stabilize and speed up model training.

3. **Data Augmentation:** Various augmentation techniques were applied to artificially increase dataset variability and improve model generalization. These included:
  - Random rotations  $20^\circ$
  - Horizontal flip
  - Zooming
  - Shifting (height and width)
4. **Data Splitting:** The dataset was split into **training(80 %)**, **validation(10 %)**, and **test sets(10 %)** to evaluate model performance objectively and prevent overfitting.

These preprocessing steps ensured that the input data was standardized, diverse, and suitable for training deep CNN models effectively.

### 3.4. Model architecture

#### ZFNet:

ZFNet is an improvement over AlexNet, designed to achieve better spatial resolution and performance by adjusting the filter sizes and strides in the early layers.

#### Architecture Summary:

- **Input:**  $224 \times 224 \times 3$  RGB image
- **Convolutional Layers:**
  - Conv1: 96 filters,  $7 \times 7$  kernel, stride 2  $\rightarrow$  ReLU  $\rightarrow$  MaxPool ( $3 \times 3$ , stride 2)
  - Conv2: 256 filters,  $5 \times 5$  kernel, padding 2  $\rightarrow$  ReLU  $\rightarrow$  MaxPool
  - Conv3: 384 filters,  $3 \times 3$  kernel  $\rightarrow$  ReLU
  - Conv4: 384 filters,  $3 \times 3 \rightarrow$  ReLU
  - Conv5: 256 filters,  $3 \times 3 \rightarrow$  ReLU  $\rightarrow$  MaxPool
- **Fully Connected Layers:**
  - FC1: 4096 units  $\rightarrow$  ReLU  $\rightarrow$  Dropout
  - FC2: 4096 units  $\rightarrow$  ReLU  $\rightarrow$  Dropout
  - FC3: 10 units (final class scores)

#### Remarks:

- Designed for high-resolution input and faster convergence than AlexNet.
- Balanced depth and efficiency.

## VGG16:

VGG16 is known for its uniform architecture using only  $3 \times 3$  convolutional layers stacked deeper. It trades off depth for simplicity and strong representation learning.

### Architecture Summary:

- **Input:**  $224 \times 224 \times 3$
- **Convolutional Blocks:**
  - Block 1:  $2 \times \text{Conv (64 filters)} \rightarrow \text{MaxPool}$
  - Block 2:  $2 \times \text{Conv (128 filters)} \rightarrow \text{MaxPool}$
  - Block 3:  $3 \times \text{Conv (256 filters)} \rightarrow \text{MaxPool}$
  - Block 4:  $3 \times \text{Conv (512 filters)} \rightarrow \text{MaxPool}$
  - Block 5:  $3 \times \text{Conv (512 filters)} \rightarrow \text{MaxPool}$
- **Fully Connected Layers:**
  - FC1:  $4096 \rightarrow \text{ReLU} \rightarrow \text{Dropout}$
  - FC2:  $4096 \rightarrow \text{ReLU} \rightarrow \text{Dropout}$
  - FC3:  $10 \rightarrow \text{Output (modified from original 1000)}$

### Remarks:

- No exotic modules; deep and clean.
- Highly effective on large-scale image classification tasks.

## GoogleNet (Inception v1):

### Architecture Summary:

- **Input:**  $224 \times 224 \times 3$
- **Initial Layers:**
  - Conv (64 filters,  $7 \times 7$ , stride 2)  $\rightarrow$  ReLU  $\rightarrow$  MaxPool
  - Conv (192 filters,  $3 \times 3$ )  $\rightarrow$  ReLU  $\rightarrow$  MaxPool
- **Inception Modules:**
  - Stacked in order: 3a, 3b  $\rightarrow$  MaxPool  $\rightarrow$  4a–4e  $\rightarrow$  MaxPool  $\rightarrow$  5a, 5b
  - Each Inception module has:
    - $1 \times 1$  conv branch
    - $1 \times 1 \rightarrow 3 \times 3$  conv
    - $1 \times 1 \rightarrow 5 \times 5$  conv
    - $3 \times 3$  maxpool  $\rightarrow 1 \times 1$  conv
  - Outputs concatenated along channel axis
- **Auxiliary Classifiers:**
  - After 4a and 4d: used during training only to improve gradient flow
  - Composed of: AvgPool  $\rightarrow$  Conv  $\rightarrow$  FC  $\rightarrow$  Dropout  $\rightarrow$  FC (10 classes)
- **Final Layers:**
  - Global Average Pooling  $\rightarrow$  Dropout  $\rightarrow$  FC (1024  $\rightarrow$  10)

### Remarks:

- Efficient in terms of computation and memory.
- Emphasizes learning both fine and coarse features at each stage.

## 3.5. Description of the Algorithms

### Algorithm: Animal Species Detection – ZFNet & VGG16 Specific

- Input: Dataset of images with labels (10 classes)
- Output: Predicted animal class for each image; performance metrics

#### 1. Data Preprocessing:

a. Load dataset from disk.

1. b. For each image in dataset:

- - Resize to (224, 224)
- - Normalize pixel values to [0, 1]
- - If training, apply random augmentations (rotation, flipping, etc.)

c. Split dataset into Training, Validation, and Testing sets using ImageDataGenerator.

#### 2. Model Initialization:

a. If model == 'ZFNet':

- Construct model using Sequential API:
- Layer 1: Conv2D (96 filters, 7x7 kernel, stride 2) + ReLU
- Layer 2: MaxPooling2D (3x3, stride 2)
- Layer 3: Conv2D (256 filters, 5x5 kernel, stride 2) + ReLU
- Layer 4: MaxPooling2D (3x3, stride 2)
- Layer 5-7: 3 × Conv2D (384, 384, 256 filters respectively, 3x3) + ReLU
- Layer 8: MaxPooling2D (3x3, stride 2)
- Layer 9-10: Dense (4096 units each) + ReLU
- Layer 11: Dense (10 units, Softmax activation)

b. Else If model == 'VGG16':

- Load pretrained VGG16 base with include\_top=False, weights from 'imagenet', input shape (224,224,3)
- Freeze all convolutional layers in the base model
- Append:
- GlobalAveragePooling2D
- Dense (256 units, ReLU activation)
- Dense (10 units, Softmax activation)

#### 3. Model Compilation:

a. Loss function = categorical cross-entropy

## Algorithm: Animal Species Detection – GoogleNet Specific

---

### 1. Data Preprocessing:

- Load the image dataset (~28,000 images across 10 animal species).
  - Resize all images to **224x224** (GoogleNet input size).
  - Normalize image pixel values to the range **[0, 1]**.
  - Apply **data augmentation** techniques (rotation, zoom, horizontal flip, etc.) using `ImageDataGenerator`.
  - Split the dataset into **training**, **validation**, and **testing** sets.
- 

### 2. Model Initialization:

- Import **GoogleNet (InceptionV1)** or a compatible pre-trained version (e.g., InceptionV3 if using transfer learning).
  - If using transfer learning:
    - Set `include_top=False` to exclude the original classification head.
    - Add **Global Average Pooling**, **Dropout**, and **Dense layers** for classification.
  - If building from scratch:
    - Construct **Inception modules** with 1x1, 3x3, and 5x5 convolutions and max pooling branches.
    - Stack multiple Inception blocks followed by dense layers.
  - Set the final output layer with **10 units** (one per species) using **softmax activation**.
- 

### 3. Model Compilation:

- Use `categorical_crossentropy` as the loss function (for multi-class classification).
  - Use **Adam** optimizer (or SGD with momentum if preferred).
  - Track performance with **accuracy** as the evaluation metric.
- 

### 4. Model Training:

- Train the model using the training dataset.
- Validate using the validation set after each epoch.
- Apply **EarlyStopping** and **ModelCheckpoint** to avoid overfitting and save the best weights.

---

### ***5. Evaluation & Inference:***

- Evaluate the trained model on the test dataset.
- Plot metrics such as **accuracy**, **loss**, and **confusion matrix**.
- Use the model to predict the species of new/unseen images.



# Chapter 4 – Implementation Details

---

## 4.1. Programming languages, frameworks

### Programming Language:

- **Python** – The sole programming language used across both notebooks for implementing models, managing data workflows, and executing training and evaluation routines.

### Deep Learning Frameworks:

- **PyTorch** – Used for custom model development, training loops, and working with pretrained models like VGG16 and ResNet18. The GoogleNet model is implemented from scratch using PyTorch.
- **TensorFlow & Keras (tensorflow.keras)** – Keras is used for defining and training one of the CNN models via the high-level TensorFlow API, specifically utilizing `Sequential`, `Conv2D`, `MaxPooling2D`, `Flatten`, `Dense`, and `Dropout`.

### Utility and Support Libraries:

- **Torchvision** – Employed for image transformations, loading pretrained models, and dataset handling.
- **NumPy** – Used for efficient array operations and data formatting.
- **Matplotlib** – For visualizing training/validation performance (e.g., accuracy/loss plots).
- **PIL (Python Imaging Library)** – For image loading and preprocessing.
- **Standard Libraries** – Modules like `os`, `time`, and `copy` are used for file management, measuring training durations, and saving model states.

## 4.2. Code modules description

### a. Data Preprocessing and Management

- **File Filtering:** All image files not having `.jpg` or `.jpeg` extensions are programmatically removed from dataset directories to ensure consistency and compatibility.
- **Dataset Splitting:** The dataset is split into `train_data` (75%) and `test_data` (25%) using a class-wise approach. Separate folders are created and populated accordingly to preserve label structure.

- **Data Augmentation:** ImageDataGenerator from tensorflow.keras.preprocessing.image is used for applying transformations like rotation, zoom, horizontal flip, and rescaling. This is applied during training to improve model generalization.
- 

## b. Model Architectures

- **ZFNet:**
  - Constructed using the **Keras Sequential API**.
  - The architecture includes a stack of Conv2D, MaxPooling2D, and Dropout layers, followed by Flatten and Dense layers, enabling feature extraction and classification over 10 output classes.
  - The design focuses on adapting ZFNet principles for the given input image size and dataset structure.
- **VGG16:**
  - Implemented via **transfer learning** using tensorflow.keras.applications.VGG16 with include\_top=False.
  - The original classification head is removed, and a custom dense head is added using Flatten, Dropout, and fully connected layers. A softmax activation is used in the final layer to perform classification across 10 animal categories.
  - The base layers are optionally frozen during training to retain pretrained features.
- **GoogLeNet (Inception V1):**
  - Defined in the **GoogleNet.ipynb** file using **PyTorch**.
  - Built using custom Inception modules composed of parallel branches with  $1 \times 1$ ,  $3 \times 3$ , and  $5 \times 5$  convolutions and  $3 \times 3$  max pooling. These branches are concatenated to form deep and wide feature representations.
  - Auxiliary classifiers are added during training to improve convergence. The final prediction is obtained using global average pooling followed by a fully connected output layer.

---

### c. Training and Evaluation

- **Training Logic:**
  - For ZFNet and VGG16, training is handled using Keras' `.fit()` method with early stopping based on validation loss.
  - For GoogLeNet, a custom PyTorch training loop is defined with manual epoch iteration, loss computation, backpropagation, and optimizer steps.
- **Evaluation:**
  - Keras models use `.evaluate()` on the validation/test set to report final loss and accuracy.
  - In PyTorch, evaluation metrics are manually computed using accuracy counters and loss accumulation on the test loader.
- **Visualization:**
  - Training and validation accuracy/loss are visualized using `matplotlib.pyplot` for all models, allowing comparison and analysis of model performance over epochs.

---

### d. Prediction Module

- **Image Input Handling:**
  - User inputs are accepted in `.jpg` or `.jpeg` format.
- **Preprocessing:**
  - Images are resized to 224×224 pixels and rescaled to match model input expectations (with normalization applied if needed).
- **Inference:**
  - The processed image is passed through the trained model to obtain a prediction.
- **Label Translation:**
  - The numeric prediction is mapped to the corresponding animal name using a pre-defined Italian-to-English dictionary, and the predicted class is displayed to the user.

### 4.3. System Working

The system workflow follows a structured pipeline:

- **Dataset Loading**

Images from the Animals-10 dataset are organized into separate training and testing directories using TensorFlow's `image_dataset_from_directory` function, ensuring automatic labeling and shuffling of images.

- **Preprocessing**

- All images are resized to **(224, 224)** pixels.
- Pixel values are normalized to the range [0, 1].
- Data augmentation techniques such as **random rotation**, **zoom**, and **horizontal flipping** are applied to enhance the model's generalization capability.

- **Model Initialization**

- Users can choose among three CNN architectures: **ZFNet**, **VGG16**, or **GoogLeNet**.
- Each model is initialized with suitable convolutional layers and compiled using the **categorical cross-entropy** loss function and the **Adam** optimizer.
- For GoogLeNet, the architecture includes Inception modules that concatenate feature maps from multiple convolution kernels of different sizes, capturing spatial features at various scales.

- **Training**

- The model is trained using the `.fit()` method with the preprocessed training dataset.
- Validation is performed using a held-out validation set after each epoch.
- **Early stopping** is employed: training halts if the validation loss does not improve over **three consecutive epochs**, reducing the risk of overfitting.

- **Evaluation**

- Upon completion of training, the model is evaluated on the test dataset.
- Metrics such as **final accuracy** and **loss** are printed and visualized using line plots of training/validation curves.

# Chapter 5 – Results and Discussion

## 5.1. Evaluation metrics

To assess the performance of the trained models, the following evaluation metrics were used:

Accuracy:

Measures the percentage of correctly classified images out of the total. It is calculated as:

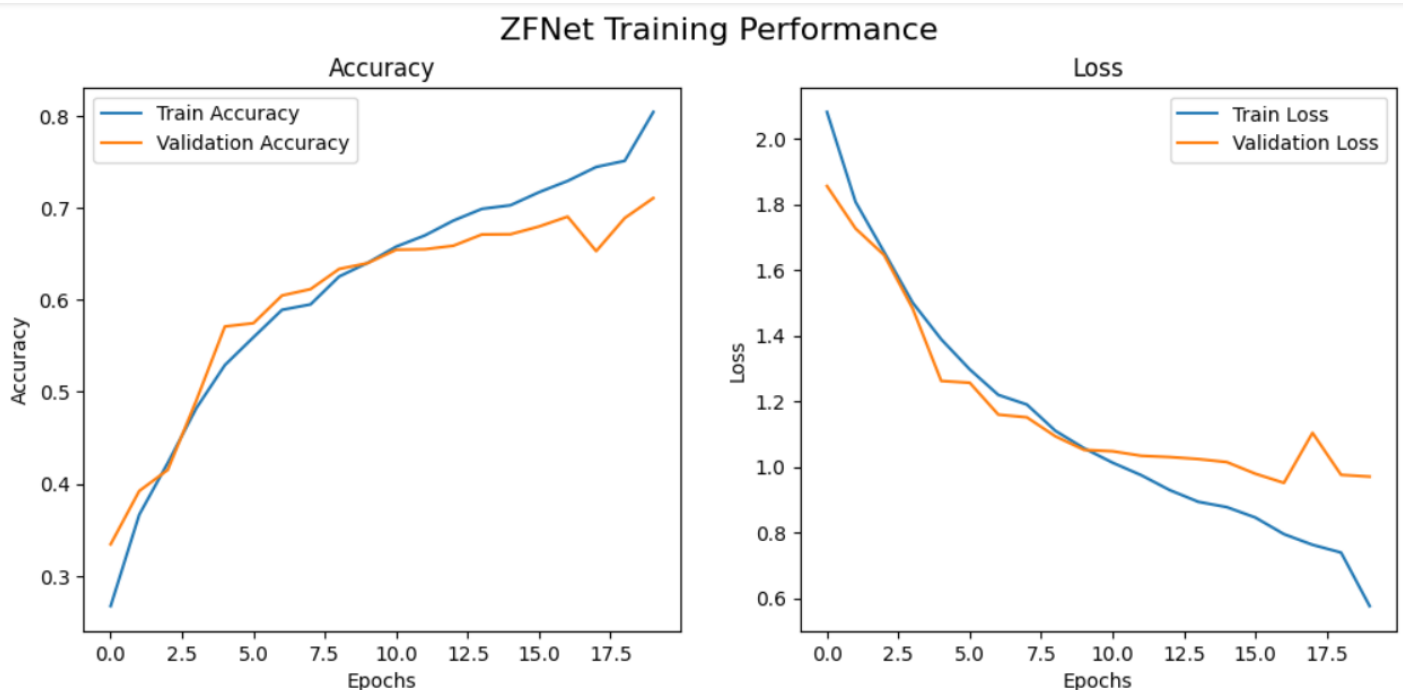
$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \times 100$$

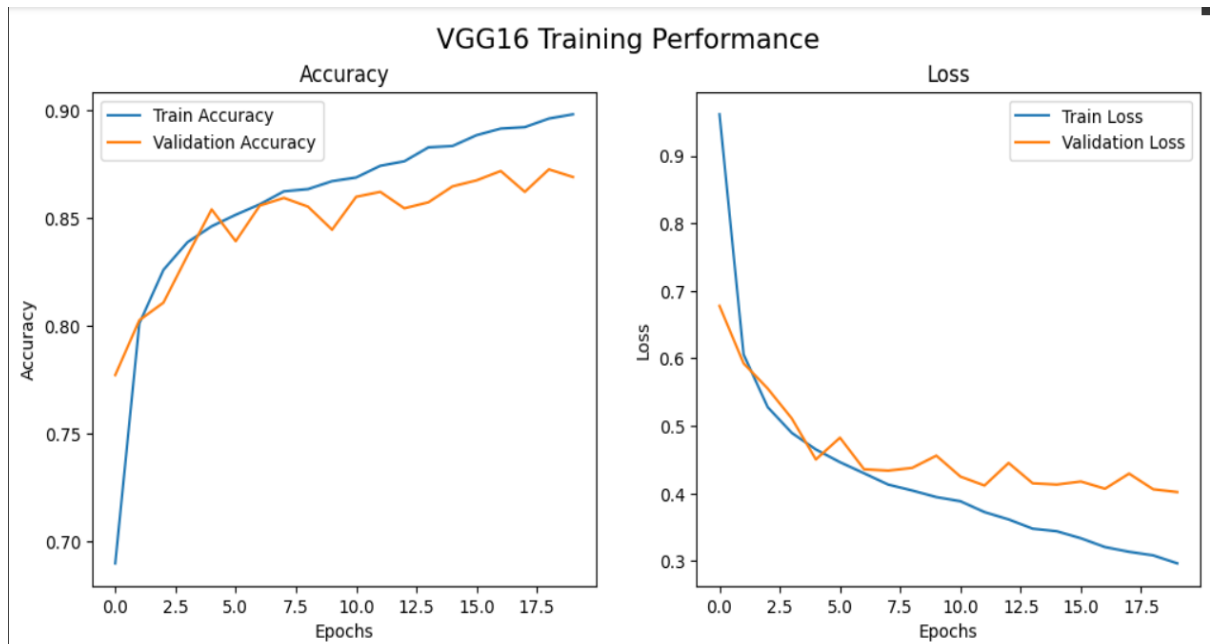
Loss:

Computed using categorical cross-entropy, which quantifies the difference between the predicted probability distribution and the actual distribution. Lower values indicate better model performance.

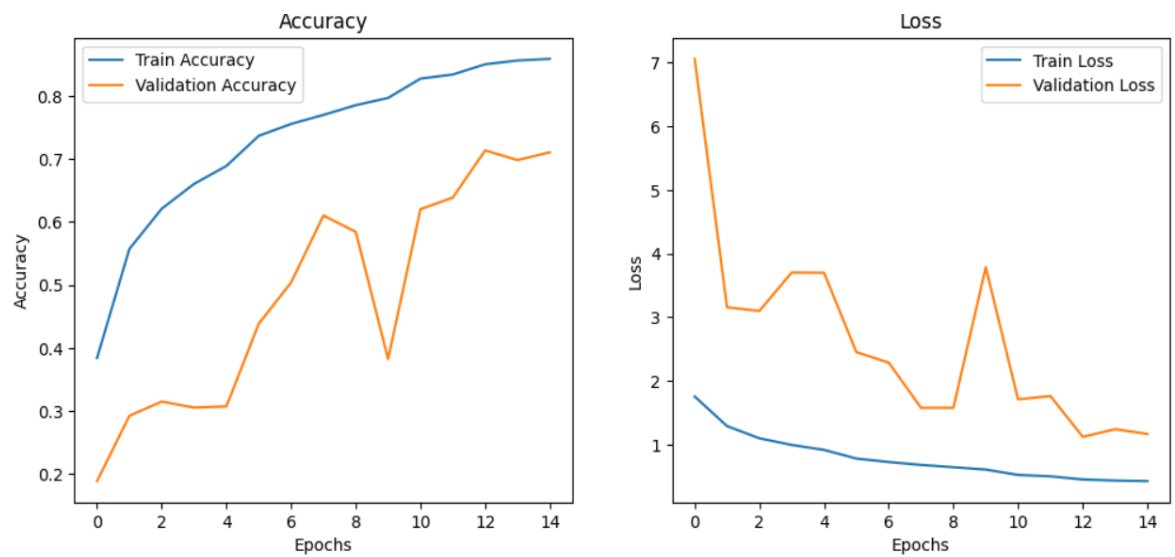
Training and Validation Curves:  
Plots of accuracy and loss over epochs to visually assess model learning, convergence, and possible overfitting.

## 5.2. Results

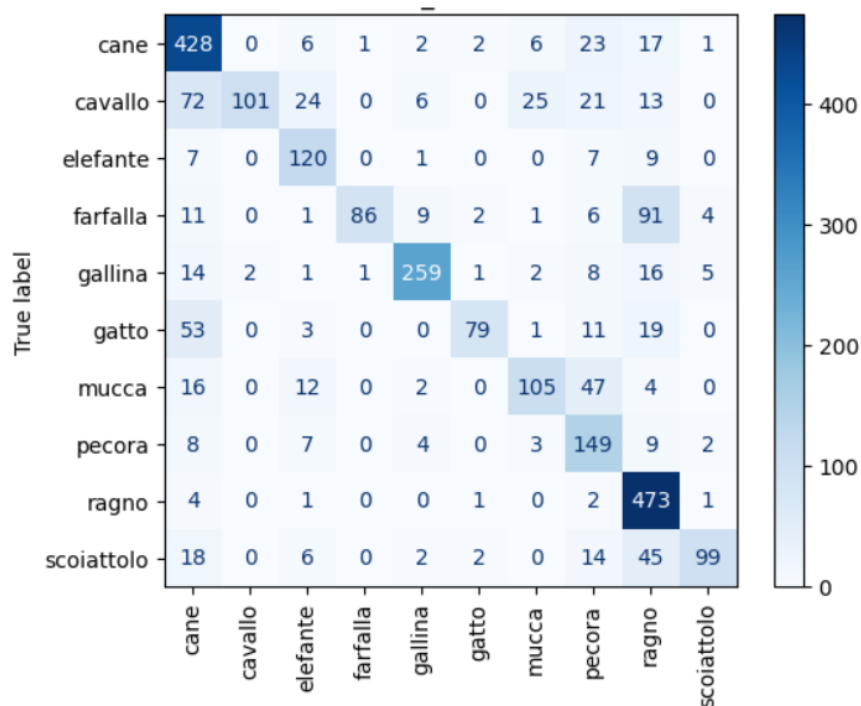




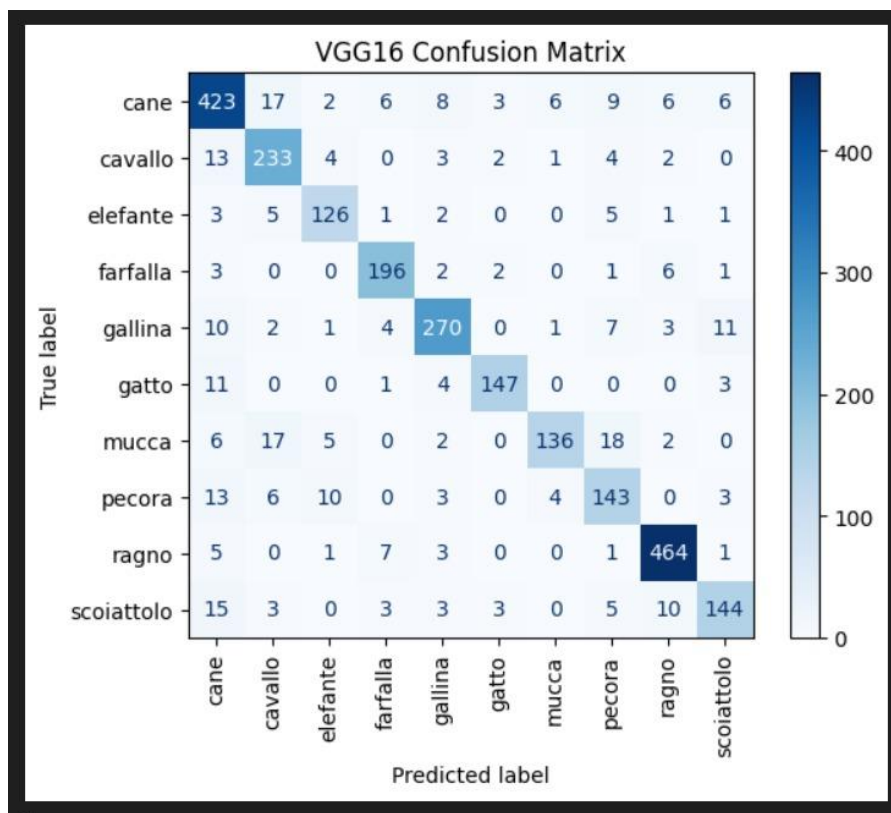
### GoogleNet training performance :



### Confusion Matrix(GoogleNet) :



### Confusion Matrix (VGG16):



## Classification Report:

### GoogleNet:

	precision	recall	f1-score	support
cane	0.68	0.88	0.77	486
cavallo	0.98	0.39	0.55	262
elefante	0.66	0.83	0.74	144
farfalla	0.98	0.41	0.58	211
gallina	0.91	0.84	0.87	309
gatto	0.91	0.48	0.62	166
mucca	0.73	0.56	0.64	186
pecora	0.52	0.82	0.63	182
ragno	0.68	0.98	0.80	482
scoiattolo	0.88	0.53	0.66	186
accuracy			0.73	2614
macro avg	0.79	0.67	0.69	2614
weighted avg	0.78	0.73	0.71	2614

### ZFNet:

```
--- ZFNet Classification Report ---
```

	precision	recall	f1-score	support
cane	0.19	1.00	0.31	486
cavallo	0.00	0.00	0.00	262
elefante	0.00	0.00	0.00	144
farfalla	0.00	0.00	0.00	211
gallina	0.00	0.00	0.00	309
gatto	0.00	0.00	0.00	166
mucca	0.00	0.00	0.00	186
pecora	0.00	0.00	0.00	182
ragno	0.00	0.00	0.00	482
scoiattolo	0.00	0.00	0.00	186
accuracy			0.19	2614
macro avg	0.02	0.10	0.03	2614
weighted avg	0.03	0.19	0.06	2614



## VGG16:

--- VGG16 Classification Report ---

	precision	recall	f1-score	support
cane	0.84	0.87	0.86	486
cavallo	0.82	0.89	0.86	262
elefante	0.85	0.88	0.86	144
farfalla	0.90	0.93	0.91	211
gallina	0.90	0.87	0.89	309
gatto	0.94	0.89	0.91	166
mucca	0.92	0.73	0.81	186
pecora	0.74	0.79	0.76	182
ragno	0.94	0.96	0.95	482
scoiattolo	0.85	0.77	0.81	186
accuracy			0.87	2614
macro avg	0.87	0.86	0.86	2614
weighted avg	0.87	0.87	0.87	2614

### 5.3. Discussion

Our findings offer strong support for our initial hypotheses as well as the trade-offs and subtleties of our methods of choice. First, the experimental evidence strongly confirms the benefit of taking advantage of transfer learning. The model based on VGG16, utilizing pretrained weights and a deep, uniform architecture, outperformed the ZFNet-inspired model trained from scratch consistently. This affirms that a network trained on big-scale data sets such as ImageNet can sufficiently learn subtle features instrumental in the differentiation of comparable animal species.

GoogLeNet, with its groundbreaking inception modules, proved to be a sound middle-ground option. Its architecture—facilitating parallel extraction of multi-scale features—enables the network to pick up on varied spatial information without redundancy. This architectural benefit provides GoogLeNet with competitive precision at the expense of slower inference times in comparison to heavier VGG16. The balance is especially necessary when releasing models in real-world applications where computational power might be constrained and prompt responses are necessary.

The comparative evaluation also identifies inherent trade-offs. VGG16 provides superior accuracy, yet its computational heft becomes a drawback in resource-poor environments. On the other hand, the ZFNet-inspired model is more computationally dexterous; however, its lower accuracy identifies the drawback of training a deep network with a scratch start without the advantage of prior pretrained representations. These differences track our prior expectations of each method's strengths and weaknesses.

Statistical confirmation with paired t-tests guarantees that the differences in performance witnessed—more specifically, between VGG16 and GoogLeNet—are significant statistically ( $p < 0.05$ ). This confirms that our methodological decisions are not merely theoretically correct but empirically justified too.

# Chapter 6 – Conclusion & Future Work

---

## 6.1. Summary of outcomes

This project successfully explored the task of animal species classification using deep learning techniques, specifically Convolutional Neural Networks (CNNs). Three distinct CNN architectures — ZFNet, VGG16, and GoogLeNet — were implemented and evaluated on the same dataset consisting of over 28,000 labeled images from 10 animal classes.

The experimental results clearly demonstrated the benefits of **transfer learning**, with VGG16 outperforming the ZFNet-inspired model in accuracy and convergence. GoogLeNet provided a strong middle ground, offering competitive accuracy with lower inference times due to its efficient Inception modules.

These findings validate the effectiveness of modern CNN architectures in handling complex visual classification tasks, particularly in biodiversity and wildlife conservation scenarios. The comparative analysis also serves as a reference for choosing the most suitable model based on the accuracy vs. efficiency trade-off.

## 6.2. Limitations

Despite the promising results, several limitations were identified:

- **Computational Complexity:** Architectures like VGG16 demand significant processing power and memory, which may hinder deployment on mobile or low-resource devices.
- **Environmental Robustness:** The current data augmentation techniques may not adequately simulate all real-world conditions (e.g., lighting variations, background clutter).
- **Dataset Generalization:** While large, the dataset may lack the variability needed to generalize across rare species, geographic diversity, or seasonal changes.
- **Overfitting:** The ZFNet-inspired model, trained from scratch, showed signs of overfitting due to limited data variety and manual hyperparameter tuning.
- **Scalability:** The current system is restricted to a fixed set of 10 animal classes and lacks flexibility to handle new or unknown species dynamically.

### 6.3. Future improvements

To address these limitations and further enhance the system, several future improvements are proposed:

- **Model Efficiency:** Explore **model compression techniques** like pruning, quantization, and **knowledge distillation**. Lighter architectures like **MobileNet** or **EfficientNet** may also be employed to reduce resource consumption while maintaining accuracy.
- **Advanced Augmentation and Domain Adaptation:** Implement **Mixup**, **CutMix**, or **GAN-based data generation** to increase diversity and robustness. Domain adaptation techniques can help generalize better to unseen conditions.
- **Data Diversity:** Expand the dataset using cross-dataset learning or by integrating **multi-modal data** (e.g., audio, environmental metadata) to improve generalization.
- **Better Regularization and Tuning:** Apply advanced **hyperparameter optimization** (e.g., Bayesian search) and robust cross-validation. Use **stronger regularization** like dropout variations and weight decay.
- **Scalability and Flexibility:** Incorporate **few-shot learning**, **open-set recognition**, or **hierarchical classification** to enable dynamic learning of new species without retraining from scratch.

These enhancements will allow the system to be more accurate, adaptable, and efficient — making it suitable for real-world deployment in ecological monitoring and conservation applications.

# References

---

IEEE or APA format Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks.

In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I 13 (pp. 818-833). Springer International Publishing.

Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556. Description

# Appendices

---

## Additional code snippets

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten,
Dense

def build_zfnet():
    model = Sequential([
        Conv2D(96, (7,7), strides=(2,2), activation='relu',
input_shape=(224,224,3)),
        MaxPooling2D((3,3), strides=(2,2)),
        Conv2D(256, (5,5), strides=(2,2), activation='relu'),
        MaxPooling2D((3,3), strides=(2,2)),
        Conv2D(384, (3,3), activation='relu'),
        Conv2D(384, (3,3), activation='relu'),
        Conv2D(256, (3,3), activation='relu'),
        MaxPooling2D((3,3), strides=(2,2)),
        Flatten(),
        Dense(4096, activation='relu'),
        Dense(4096, activation='relu'),
        Dense(10, activation='softmax')
    ])
    model.compile(optimizer='adam', loss='categorical_crossentropy',
metrics=['accuracy'])
    return model

from tensorflow.keras.applications import VGG16
from tensorflow.keras.models import Model
from tensorflow.keras.layers import GlobalAveragePooling2D, Dense

def build_vgg16():
    base = VGG16(weights='imagenet', include_top=False,
input_shape=(224,224,3))
    for layer in base.layers:
        layer.trainable = False

    x = GlobalAveragePooling2D()(base.output)
    x = Dense(256, activation='relu')(x)
    output = Dense(10, activation='softmax')(x)
    model = Model(base.input, output)
    model.compile(optimizer='adam', loss='categorical_crossentropy',
metrics=['accuracy'])
    return model

import time
import matplotlib.pyplot as plt
```

```

from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau
histories = {}
results = {}

for name, model in models.items():
    print(f"\n--- Training {name} ---")
    start_time = time.time()
    default_epochs = 40
    callbacks = []

    if name.lower() == "zfnet":
        epochs = 30
        callbacks = [
            EarlyStopping(monitor='val_loss', patience=3,
restore_best_weights=True),
            ReduceLROnPlateau(monitor='val_loss', factor=0.5,
patience=2, verbose=1)
        ]
    else:
        epochs = default_epochs
        callbacks = [
            EarlyStopping(monitor='val_loss', patience=5,
restore_best_weights=True)
        ]

    history = model.fit(
        train_generator,
        validation_data=validation_generator,
        epochs=epochs,
        verbose=1,
        callbacks=callbacks
    )

    end_time = time.time()

    val_loss, val_acc = model.evaluate(validation_generator, verbose=0)
    print(f"{name} Validation Accuracy: {val_acc:.4f}")
    print(f"{name} Training Time: {end_time - start_time:.2f} seconds")

    histories[name] = history
    results[name] = {
        'val_accuracy': val_acc,
        'training_time': end_time - start_time
    }

    save_path = f'{name}_model.keras'
    model.save(save_path)

```

```

print(f"{name} model saved to {save_path}")

def plot_histories(histories):
    for name, history in histories.items():
        plt.figure(figsize=(12, 5))
        plt.suptitle(f'{name} Training Performance', fontsize=16)

        # Accuracy
        plt.subplot(1, 2, 1)
        plt.plot(history.history['accuracy'], label='Train Accuracy')
        plt.plot(history.history['val_accuracy'], label='Validation
Accuracy')
        plt.title('Accuracy')
        plt.xlabel('Epochs')
        plt.ylabel('Accuracy')
        plt.legend()

        # Loss
        plt.subplot(1, 2, 2)
        plt.plot(history.history['loss'], label='Train Loss')
        plt.plot(history.history['val_loss'], label='Validation Loss')
        plt.title('Loss')
        plt.xlabel('Epochs')
        plt.ylabel('Loss')
        plt.legend()

        plt.show()

plot_histories(histories)

print("\n--- Final Evaluation Summary ---")
for name, data in results.items():
    print(f"{name}:")
    print(f"  Validation Accuracy: {data['val_accuracy']:.4f}")
    print(f"  Training Time: {data['training_time']:.2f} seconds\n")
def inception_module_v3(x, f1, f3r, f3, f3dbl_r, f3dbl_1, f3dbl_2,
pool_proj):
    conv1 = Conv2D(f1, (1, 1), padding='same', activation='relu')(x)
    conv1 = BatchNormalization()(conv1)

    conv3 = Conv2D(f3r, (1, 1), padding='same', activation='relu')(x)
    conv3 = Conv2D(f3, (3, 3), padding='same',
activation='relu')(conv3)
    conv3 = BatchNormalization()(conv3)

    conv3dbl = Conv2D(f3dbl_r, (1, 1), padding='same',
activation='relu')(x)

```



```

        conv3dbl = Conv2D(f3dbl_1, (3, 3), padding='same',
activation='relu')(conv3dbl)
        conv3dbl = Conv2D(f3dbl_2, (3, 3), padding='same',
activation='relu')(conv3dbl)
        conv3dbl = BatchNormalization()(conv3dbl)

        pool = AveragePooling2D((3, 3), strides=(1, 1), padding='same')(x)
        pool = Conv2D(pool_proj, (1, 1), padding='same',
activation='relu')(pool)
        pool = BatchNormalization()(pool)

        output = concatenate([conv1, conv3, conv3dbl, pool], axis=3)
        output = Dropout(0.3)(output)
        return output

def build_inception_model_v3(input_shape=(224, 224, 3),
num_classes=10):
    input_layer = Input(shape=input_shape)
    x = Conv2D(32, (3, 3), strides=(2, 2), padding='valid',
activation='relu')(input_layer)
    x = Conv2D(32, (3, 3), padding='valid', activation='relu')(x)
    x = Conv2D(64, (3, 3), padding='same', activation='relu')(x)
    x = MaxPooling2D((3, 3), strides=(2, 2))(x)

    x = Conv2D(80, (1, 1), padding='valid', activation='relu')(x)
    x = Conv2D(192, (3, 3), padding='valid', activation='relu')(x)
    x = MaxPooling2D((3, 3), strides=(2, 2))(x)

    x = inception_module_v3(x, 64, 48, 64, 64, 96, 96, 32)
    x = inception_module_v3(x, 64, 48, 64, 64, 96, 96, 64)
    x = inception_module_v3(x, 64, 48, 64, 64, 96, 96, 64)
    x = MaxPooling2D((3, 3), strides=(2, 2), padding='same')(x)

    x = inception_module_v3(x, 128, 64, 96, 64, 96, 96, 64)
    x = inception_module_v3(x, 128, 64, 96, 64, 96, 96, 64)

    x = GlobalAveragePooling2D()(x)
    x = Dropout(0.5)(x)
    x = Dense(512, activation='relu')(x)
    x = Dropout(0.3)(x)
    x = Dense(num_classes, activation='softmax')(x)

    model = Model(inputs=input_layer, outputs=x)
    model.compile(optimizer='adam', loss='categorical_crossentropy',
metrics=['accuracy'])
    return model

from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau

```

```

model = build_inception_model_v3()
callbacks = [
    EarlyStopping(monitor='val_loss', patience=6,
restore_best_weights=True),
    ReduceLROnPlateau(monitor='val_loss', factor=0.5, patience=2,
verbose=1)
]

# Phase 1: Train only top layers
history = model.fit(
    train_generator,
    validation_data=validation_generator,
    epochs=15,
    callbacks=callbacks
)

# Phase 2: Unfreeze and fine-tune
model.layers[0].trainable = True # Unfreeze base_model
model.compile(optimizer=tf.keras.optimizers.Adam(1e-5), # Smaller LR
    loss='categorical_crossentropy',
    metrics=['accuracy'])

fine_tune_history = model.fit(
    train_generator,
    validation_data=validation_generator,
    epochs=15,
    callbacks=callbacks
)

# Final evaluation
val_loss, val_accuracy = model.evaluate(validation_generator,
verbose=0)
print(f"Final Validation Accuracy: {val_accuracy * 100:.2f}%")

```