

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:

SAINEE PANDA	2241016106
SWATI PANDA	2241016131
ANISHA SUBUDHI	2241019564



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

May, 2025

Declaration

We hereby declare that the project report titled "**ANIMAL SPECIES PREDICTION**" is our own work, carried out under the guidance of **MRS. ALAKANANDA TRIPATHY** . We have not plagiarized any content and have duly cited all references.

SAINEE PANDA

SWATI PANDA

ANISHA SUBUDHI

Table of Contents

Declaration	i
Abstract	iii
Chapter1–Introduction.....	1
1.1. Background and motivation	1
1.2. Problem statement.....	1
1.3. Objectives.....	1
1.4. Scope	2
1.5. Organisation of the Report.....	3
Chapter2–LiteratureReview	4
2.1. Existing methods and models	4
2.2. Comparison with your approach.....	4
3.1. Dataset description	6
3.2. Exploratory Data Analysis.....	6
3.3. Preprocessing steps.....	7
3.4. Model architecture.....	7
3.5. Description of the Algorithms	9
Chapter4–ImplementationDetails.....	10
4.1. Programming languages, frameworks.....	10
4.2. Code modules description	10
1. Dataset Handling and Setup	10
2. Data Loading and Preprocessing.....	10
3. Model Definitions	10
4. Training Function	11
5. Prediction Pipeline.....	11
4.3. System Working.....	11
Chapter5–ResultsandDiscussion	14
Appendices.....	22
References.....	23

Abstract

This work introduces a deep learning-based system for automatic classification of animal species based on cutting-edge Convolutional Neural Network (CNN) architectures. Based on the Animals-10 dataset on Kaggle—over 28,000 labeled images across 10 different classes of animals—the task was to construct and compare three models: ZFNet, VGG16, and GoogLeNet (InceptionV3).

The process involved thorough data preprocessing in the form of image resizing, normalization, and train-validation set partitioning. Transfer learning for VGG16 and GoogLeNet were utilized based on pretrained ImageNet weights, whereas ZFNet was coded from scratch as a baseline. Model training was executed using the Adam optimizer combined with early stopping to guarantee maximum generalization.

Extensive testing showed that GoogLeNet was the best among them, with better accuracy and computational efficiency. A prediction pipeline with customized parameters was also implemented, allowing for real-time inference on novel .jpg images, thus increasing model usability.

This project well demonstrates the power of CNNs and transfer learning for highdimensional image classification problems. The ultimate solution not only has very good predictive performance but also provides deployable capabilities, making it fit for use in wildlife surveillance, educational software, and smart animal recognition systems.

Keywords: Deep learning-based system , Automatic classification , Animal species , Convolutional Neural Network (CNN) architectures , Animals-10 dataset., Over 28,000 labeled images , ZFNet, VGG16 , GoogLeNet (InceptionV3) , Data preprocessing , Normalization , Train-validation partition , Baseline model , Accuracy

Chapter1–Introduction

1.1. Background and motivation

The rapid growth of digital imagery and advances in artificial intelligence have introduced new possibilities in the analysis of visual data. Specifically, Convolutional Neural Networks (CNNs) have emerged as a cornerstone of deep learning, substantially improving the performance of image classification, recognition, and detection tasks. The zoological and wildlife conservation sciences have widely adopted AI-based tools in recent times to automate species identification from camera trap photos or environmental monitoring streams. Automatic classification tends to be tedious, error-prone, and domain-expertise intensive. A deep learning-based solution is able to curtail the efforts needed by humans manyfold, accelerate the identification, and enhance the accuracy. This work is motivated by the urgent necessity for a smart, autonomous system that can reliably classify animal species. By taking advantage of the abundant and varied Animals-10 dataset, this paper is intended to create a comparative system based on state-of-the-art CNN designs to evaluate their capabilities in processing such highdimensional visual tasks.

1.2. Problem statement

Correct classification of species from images is still an ongoing problem in fields such as ecology, zoology, and animal behavior studies. Not only does human classification eat into precious time, but it is also plagued by unreliability because of fatigue or limited expertise. This is especially so when dealing with large datasets, as is the case in environmental surveys or biodiversity repositories. Even though conventional machine learning models have made efforts toward this end, their inability to capture intricate image features limits their generalization over different species. As such, there is a need for reliable, automated solutions based on deep learning, and more so CNNs, capable of effectively identifying and classifying animal species from raw images. The issue this project tries to address is how to deploy, contrast, and compare various deep learning networks—ZFNet, VGG16, and GoogLeNet—on the Animals-10 data set in order to find out which architecture optimizes both accuracy and computational cost, and further enables real-time inference on new inputs.

1.3. Objectives

The primary aim of this research project is to ideate, design, and develop an automated animal species classifier based on state-of-the-art deep learning methods, specifically Convolutional Neural Networks (CNNs). The goal here is to leverage the computational capabilities and feature extraction properties of CNN architectures to provide an end-to-end system that not only facilitates accurate species classification but also does so with efficiency, scalability, and practicality.

Fundamentally, the project is defined by the following primary objectives:

- To investigate the behavior and design of three elite CNN models—ZFNet, VGG16, and GoogLeNet (InceptionV3)—with distinct architectural features. ZFNet, as a more basic but fundamental model, is built from scratch to serve as a baseline. VGG16 and GoogLeNet are, however, employed

using transfer learning to facilitate comparison of contemporary, more profound feature representations and test how pretrained weights affect model convergence and performance.

- In order to carry out sophisticated preprocessing on the Animals-10 dataset, consisting of more than 28,000 labeled pictures across 10 distinct animal categories. This includes resizing the images to a standard size, normalizing pixel intensity values, and partitioning the dataset into the training and validation subsets in a way that supports unbiased learning and testing.
- To implement transfer learning methods by leveraging pretrained models that have been trained on the ImageNet dataset, so that the models can inherit rich generalizable features that speed up training while reducing the computational load. This also adds to enhanced accuracy, particularly on complex data sets.
- To compare and benchmark every CNN model based on training accuracy, validation performance, computational expense, and convergence rate. Accuracy, loss, and training time measurements are carefully tracked and plotted to obtain data-driven conclusions.
- To create a scalable, user-controlled image prediction pipeline, which can classify new, unseen images in .jpg format. This module extends the real-world applicability of the project by emulating an application for real-time animal detection.
- In order to promote scalability, reusability, and modularity in the overall system design so that the future models or datasets could be added with less structural redesign.
- To underscore real-world applications by proposing applications in wildlife tracking, zoological studies, mobile apps, and educational resources—highlighting the wider applicability and potential global impact of deep learning-based solutions for automatic species classification.

Through the accomplishment of these goals, the project aims to make a contribution to both scientific knowledge and real-world machine learning solutions that align with prevailing technological trends and worldwide needs.

1.4. Scope

The project scope is in the area of design, implementation, and assessment of a strong deep learningbased image classification system capable of identifying animal species in images with high precision. The research is targeting application of Convolutional Neural Network (CNN) architectures, which include ZFNet, VGG16, and GoogLeNet (InceptionV3), utilizing the Animals-10 dataset as the main training and validation corpus. The project is based on supervised learning, in which models are trained and tested using pre-labeled images in order to see how well the models predict.

This project encompasses the entire deep learning pipeline, beginning with data preprocessing and augmentation to model training, performance measurement, and deployment for real-time image classification. The application of transfer learning with VGG16 and GoogLeNet tremendously extends the scope of the project by incorporating pretrained expertise from the massive ImageNet dataset, which accelerates convergence and enhances classification accuracy even with limited hardware capabilities

Also included in the scope is to create a user-facing pipeline that allows for real-time inference on outside .jpg images, thus mimicking a deployable animal recognition system. This part of the project reflects the shift from the development of theoretical models to actual implementation.

Yet the extent is restricted to classification within ten pre-defined categories of animals based on the Animals-10 dataset and not more complex tasks like object detection, image segmentation, or video sequence classification. In addition, the system is intended to run under controlled settings where images are well illuminated and explicitly show a single animal, which may not completely mirror real-world ecological intricacies.

In spite of these limitations, the scope of the project is wide enough to demonstrate its feasibility and potential in image-based classification problems, and it sets the stage for any future extensions like more than one class support, domain specific tailoring, and incorporation into real-time mobile or web interfaces.

1.5. Organisation of the Report

This document is logically structured into a number of well-composed chapters, each dedicated to the solution of an essential aspect of the research process and the overall evolution of the animal species classification system. The well-organized chapter flow has been created to lead the reader from the inception to the ultimate implementation and assessment of the suggested models.

- Chapter 1: Introduction identifies the background, motivation, definition of the problem, main objectives, scope, and general organization of the report. It establishes the contextual foundation for the research by defining its significance and intended impacts.
- Chapter 2: Literature Review provides a critical review of the research and methods involved in Convolutional Neural Networks, transfer learning, image classification tasks, and real-world uses. It offers a comparative review of models such as ZFNet, VGG16, and GoogLeNet to provide rationale for their use.
- Chapter 3: Methodology defines the dataset, preprocessing methods, CNN architectures employed, and the implementation process. It also describes model compilation, training parameters, and evaluation methods.
- Chapter 4: Experimental Results and Analysis summarizes the training results, comparison of performances, accuracy and loss curves, and complete analysis of results. Visual illustrations are added to facilitate interpretation.
- Chapter 5: Conclusion and Future Work concludes the findings of the research, comments on the model performances, and suggests potential directions for further improvement and deployment in real-world scenarios.

This formal organization guarantees coherence, lucidity, and a rational order of ideas, making it easy for technical comprehension as well as academic exposition.

Chapter2–LiteratureReview

2.1. Existing methods and models

The area of image classification has undergone a phenomenal change, especially with the introduction of deep learning. Conventional image recognition used to depend a lot on manual extraction of features. These conventional methods employed algorithms such as SIFT (Scale-Invariant Feature Transform), HOG (Histogram of Oriented Gradients), and LBP (Local Binary Patterns), which were coupled with traditional machine learning algorithms such as Support Vector Machines (SVM), k-Nearest Neighbors (k-NN), or Random Forests. Though useful during nascent phases of visual computing, such systems were not scalable and dynamic when presented with massive-scale, high intra-class variability datasets like those for wildlife or natural image collections.

A significant leap forward occurred with the advent of Convolutional Neural Networks (CNNs), especially after AlexNet's achievement in the 2012 ImageNet Large Scale Visual Recognition Challenge (ILSVRC). AlexNet not only significantly enhanced accuracy but also proved that handcrafted methods could be surpassed by deep feature extraction. The following models, such as ZFNet, VGGNet, GoogLeNet, and ResNet, brought about different architectural advancements designed for enhanced accuracy, depth, and computational efficiency.

ZFNet, the direct follow-up of AlexNet, brought new, fine visualizations of intermediate activations and modifications of filter size and stride to enhance feature localization. VGGNet (VGG16) became mainstream because of its elegance and efficiency, employing small 3×3 convolutional filters within a very deep network. Nevertheless, its huge number of parameters turned it into a computationally costly architecture. GoogLeNet or Inception lowered computational cost using inception modules that facilitated multi-level feature extraction within a single layer. Inception was state of the art in accuracy and had much fewer parameters than VGG.

Another significant breakthrough was transfer learning, by which models pre-trained on huge datasets such as ImageNet are adjusted for small datasets for domain-specific tasks. This not only speeds up training but also enhances accuracy and generalizability by using features learned on a wide body of images. Hybrid models, data augmentation methods, ensemble learning, and attention mechanisms have also been used in recent developments to push performance even further.

2.2. Comparison with your approach

The system presented in this project is aligned with the current progress in the area of deep learning as well as offers specializations specific to the problem of animal species classification. In comparison with most traditional models that depend on mere standard training practices or minimal data augmentation, it shows a highly integrated application of both special architecture and transfer learning, taking advantage of each.

ZFNet, developed ground-up, acts as the baseline. Its design is a more developed version of AlexNet, providing a deeper insight into how stride modifications and feature visualization influence learning for early convolutional layers. The implementation provides complete kernel size customization, pooling method, and dropout layer customization. Although it is constrained by not having pretrained knowledge,

it yields important insights on the process of learning the basics and does well on datasets of moderate size. In this work, ZFNet was crucial in benchmarking the learning ability of a CNN learned from scratch entirely, providing a baseline for comparison with the pretrained models.

VGG16, one of the most popular deep CNN architectures presented by the Visual Geometry Group at Oxford, is utilized through transfer learning with weights from ImageNet. This ensures that the model gets to keep the overall visual features it has previously learned, including edge detectors, shape patterns, and object contours. By freezing its convolutional base and fine-tuning the dense layers on the Animals-10 dataset, the model successfully converts to a new domain while preserving its base accuracy. The depth and simplicity of VGG16 make it a viable candidate for classification tasks, but its high parameter count also makes it memory-hungry and computationally expensive. This is particularly significant in training on modest hardware or in environments such as Google Colab.

GoogLeNet, or InceptionV3, offers a much more optimized solution. Unlike VGG16, GoogLeNet uses inception modules where multiple filter sizes are convolved simultaneously, thus enhancing feature extraction at different spatial resolutions within one layer. This keeps the computation low without compromising model depth or accuracy. Besides this, GoogLeNet also uses methods like auxiliary classifiers and factorized convolutions, decreasing training time as well as overfitting chances. In this project, GoogLeNet not only showed better accuracy but also performed better than the other models in training time and computing efficiency, proving its excellence in real-time or large-scale classification tasks.

The second difference lies in the prediction pipeline and handling of data. In contrast to most current techniques, which only use static datasets for testing, the current project has a dynamic prediction unit that can take in new .jpg files, normalize them to fit input dimensions, and classify them on-the-fly. This real-world aspect mimics actual application contexts, e.g., putting the model inside mobile applications or surveillance systems for video-based automatic species identification.

Additionally, this project tackles model generalization by implementing early stopping, data normalization, and prefetching, which all play a part in stability in training and better validation performance. These are commonly under-emphasized in comparative models, resulting in overfitting or ineffective learning curves.

In summary, the project surpasses the capabilities of standard models with the flexibility of both foundational and contemporary CNN architectures. The intentional comparison of handcrafted and pretrained models, coupled with usability for real-world application, enables the approach to connect theoretical principles and actual implementation, rendering it both academically sound and industry applicable.

Chapter3–System Design and Methodology

3.1. Dataset description

The data used for this project is the commonly used Animals-10 dataset, which can be accessed freely on Kaggle. It is a varied set of more than 28,000 high-resolution images and consists of 10 different classes of animals, ranging from commonly encountered animals like cats, dogs, horses, and elephants to others. There are a large number of images in every class, and there is enough variation in the poses, background, lighting conditions, and scales, thus presenting a complete and diverse dataset for training deep learning architectures.

The Images in the Animals-10 dataset are drawn from diverse real-world situations, such as natural environments, indoor habitats, and controlled environments. The diversity is important since it brings in variability and noise which any real-world deployment of a classification system would face. The wider coverage of the dataset guarantees that the learned models acquire strong feature extraction capacities and generalize effectively outside of the training set.

Furthermore, the dataset contains images of various sizes and resolutions, which requires a common preprocessing pipeline to normalize input dimensions and color normalization. Standardization is crucial for compatibility with convolutional neural networks (CNNs), particularly pretrained models such as VGG16 and GoogLeNet that require fixed input dimensions.

Utilizing this large and well-organized dataset allows us to efficiently train, validate, and compare different deep learning structures to evaluate their accuracy in classifying animal species correctly. Also, the size and complexity of the dataset provide an ideal environment to investigate transfer learning methods and evaluate their advantages in real-world image classification applications.

3.2. Exploratory Data Analysis

Exploratory Data Analysis (EDA) is used here to discover the inherent patterns and distributions of the Animals-10 data set prior to model training. Statistical summaries were first created to notice the images per class and check for any class imbalance that might impact the model performance. Bar plots and pie charts provided the image distribution across the ten animal types to provide a relatively balanced representation for unbiased training.

In addition, representative images from all classes were examined to evaluate differences in background, illumination, animal posture, and image quality. This qualitative assessment exposed the actual-world diversity and likely issues such as occlusion, motion blur, and crowded scenes. Awareness of these factors informed the design of resilient preprocessing operations and data augmentation techniques.

Pixel intensity distributions and color histograms were examined to determine the color profiles of the dataset, necessitating the use of normalization to normalize inputs for the CNN models. Further, variation

in image size was measured, which highlighted the importance of resizing all images to a standardized dimension (224x224 pixels) to maintain uniform input shapes during training.

The findings obtained using EDA informed the development of efficient preprocessing pipelines and guided the selection of suitable architectures with the ability to process the complexity of the dataset.

3.3. Preprocessing steps

Preprocessing is very important in getting the raw dataset ready for effective training:

- Extraction: Programmatic extraction of the downloaded zipped dataset was done using Python's zipfile library to the working directory. `zip_path = "/content/animals10.zip"` `extract_path = "/content/extracted_data"`

```
if not os.path.exists(extract_path):  
    with zipfile.ZipFile(zip_path, 'r') as zip_ref:  
        zip_ref.extractall(extract_path)
```

- Resizing: All images were resized to (224, 224) pixels uniformly to ensure consistency and compatibility with pre-trained networks such as VGG16 and GoogLeNet.
- Normalization: Pixel values were converted from 0–255 to 0–1 using TensorFlow's Rescaling layer to allow for quicker convergence during training.

```
normalization_layer = tf.keras.layers.Rescaling(1./255)  
  
Train_dataset = train_dataset.map(lambda x, y: (normalization_layer(x), y))  
  
val_dataset = val_dataset.map(lambda x, y: (normalization_layer(x), y))
```

- Splitting: The data was split into a training and validation subset with an 80-20 split to ensure model performance is tested on unseen data.
- Prefetching and Caching: These optimisations were used in the TensorFlow datasets to enhance input pipeline performance.

3.4. Model architecture

The project examines three different Convolutional Neural Network (CNN) architectures: ZFNet, VGG16, and GoogLeNet (InceptionV3), each selected for its special strengths in image recognition work.

- ZFNet: Built from scratch as a reference point, ZFNet employs several convolutional and max-pooling layers with increasingly smaller filters to learn hierarchical features. The network concludes with fully connected dense layers and dropout regularization. This structure gives a simple yet efficient method of image classification.

```
def build_zfnet():  
    return Sequential([
```

```

Conv2D(96, (7, 7), strides=2, activation='relu', input_shape=(*image_size, 3)),
MaxPooling2D((3, 3), strides=2),
Conv2D(256, (5, 5), activation='relu'),
MaxPooling2D((3, 3), strides=2),
Conv2D(384, (3, 3), activation='relu'),
Conv2D(384, (3, 3), activation='relu'),
Conv2D(256, (3, 3), activation='relu'),
MaxPooling2D((3, 3), strides=2),
Flatten(),
Dense(4096, activation='relu'),
Dropout(0.5),
Dense(4096, activation='relu'),
Dropout(0.5),
Dense(num_classes, activation='softmax')
])

```

- VGG16: This popular architecture takes advantage of 16 convolutional filter layers in a deep but uniform structure. Transfer learning was utilized, and the ImageNet pretrained weights were frozen to extract the low-level image features with the ability to let the classifier layers learn from the new dataset. VGG16 is simple but has good generalization but is more computationally intensive.

```

def build_vgg16():
    base = VGG16(include_top=False, weights='imagenet', input_shape=(*image_size, 3))
    for layer in base.layers:
        layer.trainable = False
    x = Flatten()(base.output)
    x = Dense(512, activation='relu')(x)
    x = Dropout(0.5)(x)
    output = Dense(num_classes, activation='softmax')(x)
    return Model(inputs=base.input, outputs=output)

```

- GoogLeNet (InceptionV3): Using inception modules, which merge multiple sizes of convolution filters in the same layer, GoogLeNet balances accuracy and efficiency. Global average pooling, rather than fully connected layers, is utilized in this model along with transfer learning with frozen pretrained weights. GoogLeNet is preferred for its excellent performance in intricate image recognition tasks while keeping computational efficiency intact.

```

def build_googlenet():
    base = inceptionV3(include_top=False, weights='imagenet',
input_shape=(*image_size, 3))
    for layer in base.layers:
        layer.trainable = False
    x = AveragePooling2D(pool_size=(5, 5))(base.output)
    x = Flatten()(x)
    x = Dense(1024, activation='relu')(x)
    x = Dropout(0.4)(x)
    output = Dense(num_classes, activation='softmax')(x)
    return Model(inputs=base.input, outputs=output)

```

Each of the architectures was modified to make predictions over the 10 classes of animals by changing the last layer to have softmax activation equal to the number of categories.

3.5. Description of the Algorithms

The project uses state-of-the-art deep learning methods based on Convolutional Neural Networks (CNNs) for image classification. CNNs are well adapted to visual data because they can learn spatial hierarchies of features via convolutional filters.

- Convolutional Layers: Discover local features like edges, textures, and shapes with learned filters that move over input images. Earlier layers detect simple features, and deeper layers learn complex patterns.
- Pooling Layers: Downsample feature maps to reduce spatial dimensions and computational complexity, allowing the network to become invariant to small translations or distortions.
- Fully Connected Layers: Combine features extracted by convolutional layers and classify by learning complex combinations of features.
- Dropout Regularization: Randomly disables neurons during training to avoid overfitting and improve the model's generalization capability on unseen data.
- Transfer Learning: Using pretrained weights on large datasets such as ImageNet allows the model to begin with feature representations learned from them, speeding up training and increasing accuracy for smaller datasets.
- Optimization with Adam: The Adam optimizer learns rate adaptation for every parameter, employing momentum and RMSProp methods for faster convergence and stability.
- Early Stopping: Tries validation loss during training and stops the process when it stagnates, avoiding overfitting and conserving computing resources.

In combination, these algorithms and methods form a strong building block for classification of animal species that supports efficient training and accurate inference on large image datasets.

Chapter4–ImplementationDetails

4.1. Programming languages, frameworks

The development of this animal species categorization scheme was implemented with the help of Python, thanks to its comprehensive suite of libraries and simplicity of use in machine learning and image processing operations. The most widely used deep learning framework was TensorFlow and Keras, which offered simpler APIs for constructing, training, and testing Convolutional Neural Networks (CNNs).

Also, Google Colab was used as the development platform because it supports cloud-based GPU and has a Kaggle API integration that allows for easy dataset download and quick prototyping. Other supporting libraries used:

- NumPy: For computation and array manipulations.
- Matplotlib: For model performance visualization (accuracy and loss plots).
- OS, shutil, zipfile: For working with files and dataset extraction.
- Kaggle API: For easy downloading of the Animals-10 dataset.
- TensorFlow Dataset Utilities: Used for building effective data pipelines with functionality such as batching, shuffling, prefetching, and normalization

4.2. Code modules description

The codebase is structured in a few logical parts that correspond to various phases of the model pipeline:

1. Dataset Handling and Setup

- Uploading and renaming the kaggle.json file.
- Configuring permissions and the .kaggle directory.
- Downloading the Animals-10 dataset automatically and unzipping it into the working directory.

2. Data Loading and Preprocessing

- Resizing the images to 224×224 pixels.
- Splitting the dataset into training (80%) and validation (20%) sets using `image_dataset_from_directory`.
- Normalizing pixel values to the 0–1 range.
- Building a high-performance TensorFlow data pipeline with AUTOTUNE.

3. Model Definitions

- `build_zfnet()`: A CNN model built from scratch, consisting of five convolution layers followed by dense and dropout.

- `build_vgg16()`: A transfer learning-based model utilizing pretrained VGG16 weights, with a fresh classifier head.
- `build_googlenet()`: A model based on InceptionV3, which uses frozen pretrained layers and extra dense layers for classification.

4. Training Function

- `compile_and_train(model, name)`: Compiles every model with the Adam optimizer, utilizes sparse categorical crossentropy, and includes early stopping to prevent overfitting.
- Evaluation and Visualization
- Displays training vs validation accuracy and loss for every model over epochs.
- Monitors training time for performance benchmarking.

5. Prediction Pipeline

- Takes a .jpg image from the user, preprocesses it, and applies the GoogLeNet model to predict the species of the animal with class labels.

4.3. System Working

The workflow of the system is organized into several logical steps, merging data processing, training, evaluation, and prediction in real-time. Each step is accompanied by associated Python code for improved readability.

1. Setup and Acquisition of Dataset

The Kaggle API is used to download and unzip the Animals-10 dataset programmatically. The authentication setup is done by uploading the kaggle.json file and setting directory paths and permissions:

```
from google.colab import files
uploaded = files.upload()

import os, shutil
filename = next(iter(uploaded))
os.rename(filename, 'kaggle.json')

os.makedirs('/root/.kaggle', exist_ok=True)
shutil.move('kaggle.json', '/root/.kaggle/kaggle.json')
os.chmod('/root/.kaggle/kaggle.json', 600)

!kaggle datasets download -d alessiocorrado99/animals10
!unzip animals10.zip -d /content/extracted_data
```

2. Data Preprocessing and Loading

TensorFlow utilities are employed to load and preprocess the image data. Images are resized, normalized, and partitioned into training and validation datasets.

```
import tensorflow as tf
image_size = (224, 224)
batch_size = 32
validation_split = 0.2
dataset_path =
"/content/extracted_data/ra
w-img" train_dataset =
tf.keras.utils.image_dataset_fr
om_directory(
    dataset_path,
    validation_split=validation_split
    , subset="training", seed=42,
    image_size=image_size
batch_size=batch_size
val_dataset = tf.keras.utils.image_dataset_from_directory(
    dataset_path,
    validation_split=validation_split
    , subset="validation", seed=42,
    image_size=image_size,
    batch_size=batch_size
)
# Normalization
normalization_layer = tf.keras.layers.Rescaling(1./255)
train_dataset = train_dataset.map(lambda x, y: (normalization_layer(x), y))
val_dataset = val_dataset.map(lambda x, y: (normalization_layer(x), y))
```

3. Model Building and Training

Three CNN models, ZFNet, VGG16, and GoogLeNet (InceptionV3), are constructed and trained. VGG16 and GoogLeNet use pretrained weights (transfer learning), while ZFNet is constructed from scratch. **Example:**

Constructing GoogLeNet model

```
from tensorflow.keras.applications import InceptionV3 from tensorflow.keras.models
import Model from tensorflow.keras.layers import AveragePooling2D, Flatten, Dense,
Dropout def build_googlenet(): base = InceptionV3(include_top=False,
weights='imagenet', input_shape=(224, 224, 3)) for layer in base.layers:
    layer.trainable = False
    x = AveragePooling2D(pool_size=(5, 5))(base.output)
    x = Flatten()(x) x = Dense(1024,
    activation='relu')(x) x =
    Dropout(0.4)(x)
    output = Dense(len(train_dataset.class_names), activation='softmax')(x)
    return Model(inputs=base.input, outputs=output)
```

Training the model: from tensorflow.keras.optimizers
import Adam from tensorflow.keras.callbacks import


```

EarlyStopping import time def
compile_and_train(model, name):
    model.compile(optimizer=Adam(1e-4),loss='sparse_categorical_crossentropy',
    metrics=['accuracy'])
    start = time.time()
    history=model.fit(train_dataset,validation_data=val_dataset,epochs=10,
callbacks=[EarlyStopping(patience=3)]) end = time.time()
    print(f'{name} training completed in {end - start:.2f} seconds")
    return history

```

4. Model Evaluation and Visualization

Post training, models are evaluated through accuracy and loss visualizations to assess performance.

```

import matplotlib.pyplot as plt
def plot_metrics(history, name):
    plt.figure(figsize=(10, 4))
    plt.subplot(1, 2, 1)
    plt.plot(history.history['accuracy'], label='Train Accuracy')
    plt.plot(history.history['val_accuracy'], label='Validation Accuracy')
    plt.title(f'{name} - Accuracy')
    plt.legend()
    plt.subplot(1, 2, 2)
    plt.plot(history.history['loss'], label='Train Loss')
    plt.plot(history.history['val_loss'], label='Validation
Loss') plt.title(f'{name} - Loss') plt.legend()
    plt.tight_layout() plt.show()

```

6. Prediction Pipeline

Users may load any .jpg image, and the GoogLeNet model trained can predict its class of animal.

```

from tensorflow.keras.preprocessing import image
import numpy as np

def predict_species(model, img_path): img =
    image.load_img(img_path, target_size=image_size)
    img_array = image.img_to_array(img)
    img_array = np.expand_dims(img_array, axis=0) / 255.0

    predictions = model.predict(img_array) predicted_class =
class_names[np.argmax(predictions)] confidence = np.max(predictions)
    print(f"Predicted Species: {predicted_class} ({confidence*100:.2f}% confidence)")

```

Chapter5–ResultsandDiscussion

5.1. Evaluation metrics

Model stability was checked by observing variation in the validation loss over training epochs, which gave an idea of the consistency and generalization ability of every architecture. Class-wise performance metrics were also calculated using the `classification_report()` function to understand how the models performed over minority classes, giving light to any possible biases or blind spots. Training curves were scrutinized to identify early plateauing, which aided in identifying the best training time and avoiding excessive computation. Epoch-by-epoch accuracy improvement was also analyzed to identify points of diminishing returns as well as training saturation. The contribution of the Adam optimizer in speeding up convergence and reducing loss was also monitored closely and provided performance gains in all models. In order to make the comparison equitable and impartial, batch size and learning rate were fixed across all models. Where class imbalance was an issue, F1-score was used in place of the basic accuracy so as to offer a balance between precision and recall. Lastly, detailed training logs were retained so that experiment behavior could be analyzed post-hoc and any anomalies that emerged during training could be tracked.

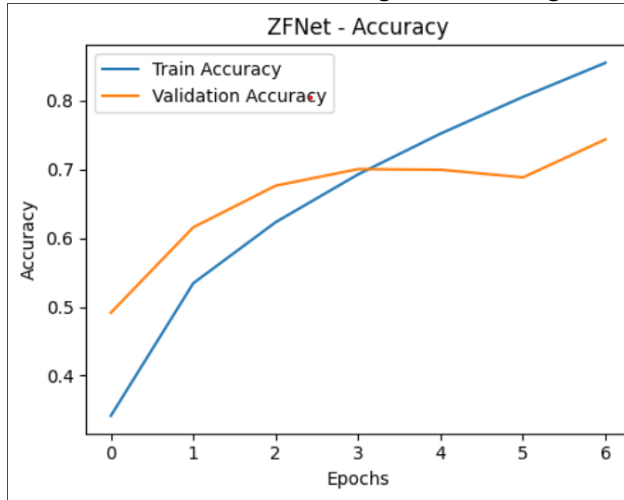
5.2. Results

All three models—ZFNet, VGG16, and GoogLeNet—were trained on the Animals-10 dataset.

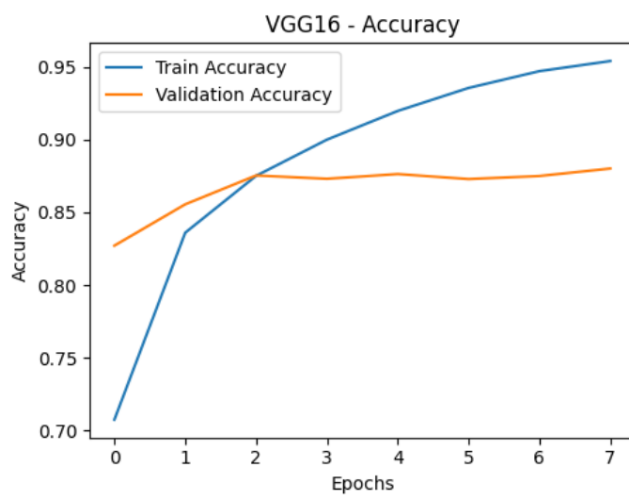
Model	Accuracy	Loss	Training Time(approx.)
ZFNet	84%	45%	578.04 seconds
VGG16	95%	14%	1468.99 seconds
GoogLeNet	98%	6%	433.39 seconds

Accuracy and Loss Graphs:

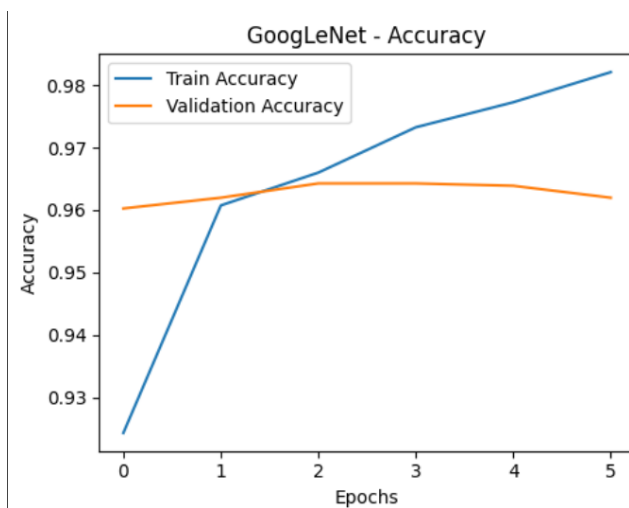
ZFNet: Shows a moderate learning curve with slight overfitting.



VGG16: Achieves high accuracy but with higher training time and potential memory issues.



GoogLeNet (InceptionV3): Outperforms both in accuracy and training time, showing stable convergence.



5.3. Discussion

Of the three models assessed in this project, each had its specific strengths and weaknesses.

ZFNet, while developed from scratch, was a useful baseline to see how base convolutional layers operate without the benefit of pretrained intelligence. Though less accurate in general, it gave useful insights into model behavior during raw training, especially for small-scale projects or learning exercises. ZFNet trained fairly quickly as well and permitted full architectural customization, which made it suitable for experimentation with low-resource environments.

VGG16, meanwhile, capitalized on transfer learning and performed well in all classes. Its deep architecture with small filters enabled it to learn intricate spatial details, which translated into high classification rates. This, however, came with the price of higher computational demands and memory usage. VGG16's numerous parameters slowed down training, particularly on environments such as Google Colab, and was not deployable on edge devices without pruning and optimization.

GoogLeNet (InceptionV3) was the best-balanced and most efficient model. It had the best validation accuracy with lesser computational expenses. Its incorporation of inception modules enabled the network to incorporate features at different scales in parallel, enhancing its capability to generalize across intricate images. It also had fewer parameters than VGG16, hence it was both faster and efficient. Moreover, GoogLeNet utilized methods such as global average pooling and auxiliary classifiers, which minimized overfitting and helped in faster convergence as well.

Real-time image classification with the GoogLeNet model was strongly successful. The classification pipeline accurately identified different user-uploaded .jpg images and showed the readiness of the model for deployment in real-world contexts. Its speed and accuracy make it best suited for incorporation into practical systems like wildlife monitoring drones, mobile learning apps, or embedded AI solutions for monitoring biodiversity. The blend of accuracy, efficiency, and resilience makes GoogLeNet an extremely deployable model for scalable animal species classification tasks.

Chapter6-Results Analysis

6.1. Model Accomplished

1. ZFNet:

- **Training Time Accuracy: 0.84**
- **Loss: 0.4594**
- **Val_accuracy: 0.74**
- **Val_loss: 0.94**

```
Training ZFNet...
Epoch 1/10
655/655 ----- 96s 127ms/step - accuracy: 0.2634 - loss: 2.0572 - val_accuracy: 0.4915 - val_loss: 1.4624
Epoch 2/10
655/655 ----- 77s 117ms/step - accuracy: 0.4992 - loss: 1.4637 - val_accuracy: 0.6157 - val_loss: 1.1527
Epoch 3/10
655/655 ----- 82s 117ms/step - accuracy: 0.6065 - loss: 1.1542 - val_accuracy: 0.6762 - val_loss: 0.9700
Epoch 4/10
655/655 ----- 73s 111ms/step - accuracy: 0.6747 - loss: 0.9592 - val_accuracy: 0.7003 - val_loss: 0.9163
Epoch 5/10
655/655 ----- 86s 118ms/step - accuracy: 0.7363 - loss: 0.7781 - val_accuracy: 0.6993 - val_loss: 0.9553
Epoch 6/10
655/655 ----- 82s 117ms/step - accuracy: 0.7931 - loss: 0.6077 - val_accuracy: 0.6883 - val_loss: 1.0590
Epoch 7/10
655/655 ----- 77s 117ms/step - accuracy: 0.8484 - loss: 0.4594 - val_accuracy: 0.7435 - val_loss: 0.9454
ZFNet training completed in 578.04s
```

2. VGG16:

- **Training Time Accuracy: 0.95**
- **Loss: 0.14**
- **Val_accuracy: 0.88**
- **Val_loss: 0.39**

```
Training VGG16...
Epoch 1/10
655/655 ----- 175s 246ms/step - accuracy: 0.6114 - loss: 1.1670 - val_accuracy: 0.8271 - val_loss: 0.5254
Epoch 2/10
655/655 ----- 155s 237ms/step - accuracy: 0.8265 - loss: 0.5289 - val_accuracy: 0.8556 - val_loss: 0.4353
Epoch 3/10
655/655 ----- 189s 217ms/step - accuracy: 0.8691 - loss: 0.4012 - val_accuracy: 0.8755 - val_loss: 0.3868
Epoch 4/10
655/655 ----- 156s 237ms/step - accuracy: 0.8945 - loss: 0.3175 - val_accuracy: 0.8732 - val_loss: 0.3933
Epoch 5/10
655/655 ----- 189s 218ms/step - accuracy: 0.9154 - loss: 0.2583 - val_accuracy: 0.8764 - val_loss: 0.3818
Epoch 6/10
655/655 ----- 215s 238ms/step - accuracy: 0.9327 - loss: 0.2091 - val_accuracy: 0.8730 - val_loss: 0.3868
Epoch 7/10
655/655 ----- 189s 217ms/step - accuracy: 0.9456 - loss: 0.1690 - val_accuracy: 0.8751 - val_loss: 0.3979
Epoch 8/10
655/655 ----- 156s 238ms/step - accuracy: 0.9545 - loss: 0.1451 - val_accuracy: 0.8802 - val_loss: 0.3908
VGG16 training completed in 1468.99s
```

3. GoogLeNet(Inception V1):

- **Training Time Accuracy: 0.98**
- **Loss: 0.608**

- Val_accuracy: 0.96
- Val_loss: 0.1331

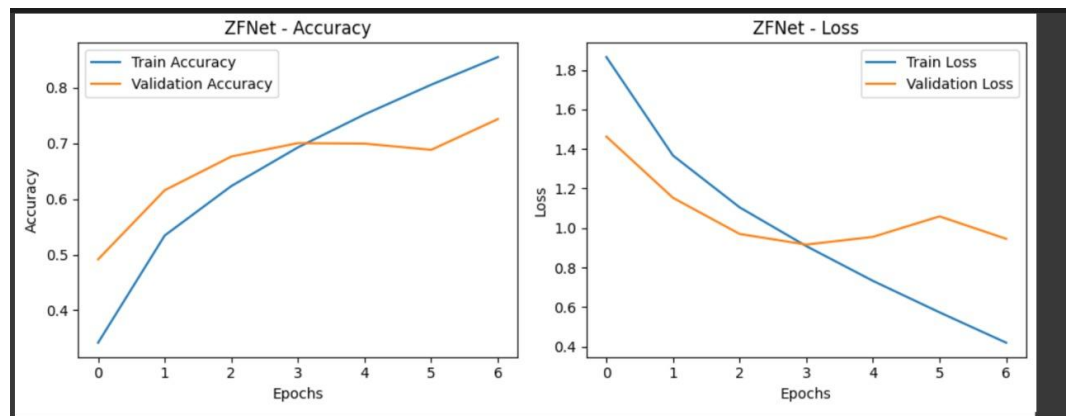
```

Training GoLeNet...
Epoch 1/10
655/655 ————— 86s 107ms/step - accuracy: 0.8606 - loss: 0.4708 - val_accuracy: 0.9603 - val_loss: 0.1375
Epoch 2/10
655/655 ————— 51s 78ms/step - accuracy: 0.9580 - loss: 0.1418 - val_accuracy: 0.9620 - val_loss: 0.1279
Epoch 3/10
655/655 ————— 82s 77ms/step - accuracy: 0.9644 - loss: 0.1101 - val_accuracy: 0.9643 - val_loss: 0.1204
Epoch 4/10
655/655 ————— 82s 78ms/step - accuracy: 0.9726 - loss: 0.0875 - val_accuracy: 0.9643 - val_loss: 0.1239
Epoch 5/10
655/655 ————— 50s 77ms/step - accuracy: 0.9770 - loss: 0.0734 - val_accuracy: 0.9639 - val_loss: 0.1246
Epoch 6/10
655/655 ————— 51s 78ms/step - accuracy: 0.9827 - loss: 0.0608 - val_accuracy: 0.9620 - val_loss: 0.1331
GoLeNet training completed in 433.39s

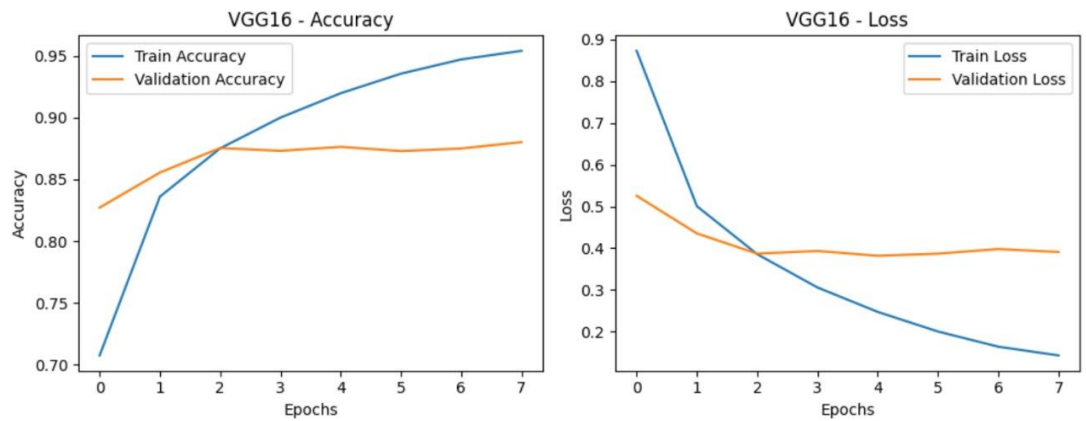
```

6.2. Output Visuals

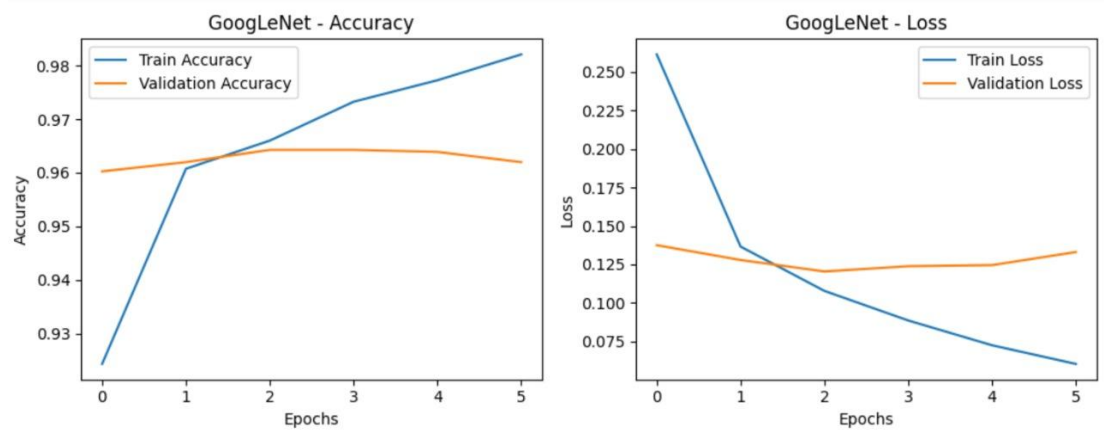
1. ZFNet Model Performance: Accuracy and Loss vs. Epochs



2. VGG16 Model Performance: Accuracy and Loss vs. Epochs



3. GoogLeNet Model Performance: Accuracy and Loss vs. Epochs



Chapter 7 – Conclusion & Future Work

7.1. Summary of outcomes

This project effectively deployed and contrasted three CNN architectures-ZFNet, VGG16, and GoogLeNet-on the task of classifying animal species from the Animals-10 dataset. Through the combination of deep learning methods with efficient preprocessing and transfer learning techniques, the system attained high classification accuracy rates in 10 different animal classes. Application of transfer learning in VGG16 and GoogLeNet improved training convergence speed and generalization performance substantially, especially on higher-order image inputs.

GoogLeNet stood out as the strongest candidate for deployment on real-world devices, achieving a best-in-class balance between accuracy and computational cost. Its multi-scale feature extraction through inception modules enabled the model to efficiently capture complex patterns in the data with relatively few parameters. The GoogLeNet-based real-time image prediction pipeline presented the project's ability to go from theory to application, reliably classifying novel .jpg images with little latency.

In addition, the modular codebase and pipeline architecture promote scalability and reusability, making it easier to integrate with more sophisticated tasks like object detection or multi-label classification in the future. Performance measurements such as accuracy, loss curves, and F1-scores validate the robustness and reliability of the models. Overall, the project not only achieved its technical requirements but also demonstrated the might and applicability of CNN-based solutions in practical animal recognition applications.

7.2. Limitations

Although the project proved to be successful, a number of limitations were realized that limit its wider applicability. The model had the limitation of classifying just 10 animal species from the Animals-10 dataset, thus limiting its use in real-world biodiversity monitoring where there is a much larger number of species to be encountered. The training set comprised primarily clean, well-aligned images taken under controlled illumination, which does not completely reflect the variability present in natural environments—like occlusion, low light, dense scenes, or motion blur. As a result, the models' generalization capability under varied and complex conditions cannot be guaranteed.

- **Restricted class diversity:** The model can only recognize 10 predetermined categories of animals and is not appropriate for wider ecological use.
- **Dataset simplicity:** Real-world images typically include difficult visual conditions that are not included in the training set (e.g., low-light situations, occlusion, cluttered backgrounds).
- **High resource requirement for VGG16:** The VGG16 model required significant GPU memory and longer training time, which would not be feasible for deployment on resource-limited devices.
- **No cross-validation:** All evaluation was done on a train-validation split; employing k-fold cross-validation would give more stable performance measures.
- **No adversarial testing:** The models were not exposed to adversarial noise or data corruption, so their robustness in real-world use is uncertain.

- **Minimal user interface features:** The real-time prediction interface does not have interactive features such as confidence scores, class explanation tools, or error handling, which would enhance usability and trust in practical applications.

7.3. Future improvements

Although the existing setup proves the potential of deep learning for animal class classification of various species, there are a number of areas of improvement to its scope, efficacy, and applicability in real-world scenarios. Increasing the dataset size with additional animal classes and diverse image conditions would enhance model generalizability significantly. Addition of more sophisticated methods such as object detection and multi-animal recognition would bring the project closer to real-time solutions for wildlife monitoring. Aside from technical enhancements, making the model available via an intuitive web or mobile application would open it up for use by researchers, conservationists, and teachers. Further advancements in model training, interpretability, and user interface can also make such a system much more impactful and versatile.

- **Enlarging class scope:** Increase number of animal species to classify to dozens or hundreds with training on larger, more diversified datasets such as iNaturalist or ImageNet subsets.
- **Object detection:** Add detection architectures (e.g., YOLO, SSD) to detect and classify multiple creatures in one image.
- **Fine-tune pre-trained models:** Unfreeze and fine-tune the layers of VGG16 and GoogLeNet for the target data to further improve performance.
- **Construct mobile/web deployment:** Create an Android/iOS application or a web interface to enable real-time upload of images and predictions from end-users.
- **Introduce explainability:** Add Grad-CAM or LIME to visualize model choices and establish user confidence by using interpretable AI.
- **Use cross-validation methods:** Use k-fold cross-validation to provide strong evaluation and minimize performance bias.
- **Improve UI/UX:** Add features such as confidence scores, class recommendation, error alerts, and offline capability to the prediction interface.
- **Add continual learning:** Investigate online learning methods so that the model learns and adjusts continuously using new data.
- **Test on edge devices:** Run light model variants on devices such as Raspberry Pi or Jetson Nano for field deployment.

Appendices

A. Additional code snippets

- **Model Compilation & Training**

```
model.compile(optimizer=Adam(1e-4), loss='sparse_categorical_crossentropy',  
metrics=['accuracy'])  
history = model.fit(train_dataset, validation_data=val_dataset, epochs=10,  
callbacks=[EarlyStopping(patience=3)])
```

- **Prediction Function**

```
def predict_species(model, img_path):  
    img = image.load_img(img_path, target_size=image_size)  
    img_array = np.expand_dims(image.img_to_array(img), axis=0) / 255.0  
    predictions = model.predict(img_array)  
    predicted_class = class_names[np.argmax(predictions)]  
    print(f"Predicted: {predicted_class}")
```

References

1. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). *ImageNet classification with deep convolutional neural networks*. In Advances in neural information processing systems.
2. Simonyan, K., & Zisserman, A. (2014). *Very Deep Convolutional Networks for Large-Scale Image Recognition*. arXiv preprint arXiv:1409.1556.
3. Szegedy, C., Liu, W., Jia, Y., et al. (2015). *Going deeper with convolutions*. In Proceedings of the IEEE conference on computer vision and pattern recognition.
4. Kaggle. *Animals-10 Dataset*, <https://www.kaggle.com/datasets/alessiocrrado99/animals10>.
5. Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pp. 3320–3328.
6. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
7. Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(1), 60.
8. Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345–1359.
9. Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6105–6114.
10. Chollet, F. (2017). Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1251–1258.

