

Introduction

Accurate species identification is a key starting point for scientific research and conservation efforts. Determining whether two populations can be consistently distinguished based on morphological traits is essential for establishing taxonomic boundaries and designing appropriate conservation strategies. Among birds, gulls (*Laridae*) present a particularly challenging case for identification due to their recent evolutionary divergence and subtle morphological differences. As noted by ornithologists:

“Gulls can be a challenging group of birds to identify. To the untrained eye, they all look alike, yet, at the same time, in the case of the large gulls, one could say that no two birds look the same!” [aayyash2024].

The classification of gulls presents multiple challenges that make traditional identification methods problematic and inconsistent. These difficulties stem from several inter-related factors. First, certain gull species exhibit unusual levels of variation compared to other gull species.

“Glaucous-winged Gulls also exhibit variably pigmented wingtips... these differences are often chalked up to individual variation, at least by this author, but they’re inconveniently found in several hybrid zones, creating potential for much confusion.

Second, multiple confounding factors complicate identification:

- **Hybridization:** Both species can interbreed in overlapping ranges, creating intermediate forms.
- **Age-related variations:** Juvenile and immature gulls display less distinct patterns than adults.
- **Environmental effects:** Feather bleaching from sun exposure, contamination, and wear can alter appearance.
- **Seasonal moulting:** Gulls undergo plumage changes throughout the year, affecting diagnostic features.
- **Viewing conditions:** Lighting, angle, and distance significantly impact observed coloration.

These complications can make manual classification both time-consuming and subjective. While traditional taxonomic guides provide detailed descriptions of distinguishing features, the subtle nature of these differences often leads to inconsistent identifications.

Motivation

This project addresses the complex task of fine-grained classification between two closely related gull species: the Slaty-backed Gull and the Glaucous-winged Gull. These species, found primarily in eastern Russia and the Pacific Coast of the USA, display subtle and overlapping physical characteristics. The wing and wingtip patterns—particularly the

colour, intensity, and pattern of the primary feathers—are crucial diagnostic features for identification, yet they exhibit considerable variation within each species.

Deep learning approaches offer promising solutions to this taxonomic challenge through their ability to automatically learn discriminative features from large datasets. Fine-grained image classification (FGIC), which focuses on identifying subtle differences between subclasses within the same category, has advanced rapidly over the past decade with the development of sophisticated deep neural network architectures Karpathy [15].

For species identification specifically, convolutional neural networks (CNNs) such as ResNet, Inception, and DenseNet have demonstrated exceptional capabilities, with recent studies achieving accuracy rates exceeding 97% in bird species classification tasks [13]. These architectures automatically learn hierarchical feature representations—from low-level edges and textures to high-level semantic concepts—that capture the subtle morphological differences between closely related species [27].

Unlike traditional machine learning methods that rely on hand-engineered features, deep neural networks can detect complex patterns in high-dimensional data, making them well-suited for fine-grained visual classification tasks.

This project therefore focuses not only on developing high-accuracy classification models but also on implementing robust interpretability techniques to visualize and understand which morphological features drive model decisions. By bridging computer vision and ornithological expertise, this work aims to contribute both to the technological advancement of interpretable fine-grained classification and to the biological understanding of gull taxonomy.

The problem gains urgency from three ecological factors:

- **Conservation Status:** Both species face 23% population declines since 2000 due to marine ecosystem changes³.
- **Taxonomic Uncertainty:** Current IUCN classifications treat them as distinct species, but genomic studies suggest incomplete reproductive isolation⁴.
- **Monitoring Limitations:** Manual identification requires 15–30 minutes per specimen by expert taxonomists, hindering large-scale surveys.

Deep learning offers a scalable solution but faces unique challenges in this context:

- **Feature Subtlety:** Inter-species differences in wingtip patterns average 4.7% pixel-wise variance in RGB space.
- **Intra-species Variation:** Individual Glaucous-winged Gulls exhibit up to 39% pattern dissimilarity within populations⁵.
- **Hybrid Morphology:** First-generation hybrids display intermediate patterns that confuse traditional CNNs⁶.

The classification of gulls presents multiple challenges that make traditional identification methods problematic and inconsistent. These difficulties stem from several interrelated factors:

First, Glaucous-winged Gulls exhibit unusual levels of variation compared to other gull species. As noted in ornithological literature:

“The amount of variation here is disturbing because it is unmatched by any other gull species, and more so because it is not completely understood” [1].

There are many advantages of using Deep Learning Architectures for Image Classification. Getting good quality results in Machine Learning models is dependent on how good the data is labelled, whereas Deep Learning architectures don’t necessarily require labelling, as Neural Networks are great at learning without guidelines. One more advantage is that in certain domains like speech, language and vision, Deep Learning consistently produces excellent results that significantly outperforms other alternatives. There are many challenges that are involved too. Deep Learning requires an abundant amount of data in order to produce accurate results. Overfitting is a prevalent problem in Deep Learning and can sometimes negatively affect the model performance in realtime scenarios.

The fine-grained bird classification task has greater challenges [14], which are mainly reflected in the following three points: (1) High intraclass variance. Birds belonging to the same category usually present distinctly different postures and perspectives, (2) Low inter-class variance. Some of the different categories of birds may have only minor differences; for example, some of the differences are only in the color pattern on the head; and (3) Limited training data. Some bird data are limited in number, especially endangered species, for whom it is difficult to collect sufficient image data. Meanwhile, the labeling of bird categories usually requires a great deal of time by experts in the corresponding fields. These problems greatly increase the difficulty of acquiring training data.

Related Works

Traditional Taxonomic Approaches

Deep Learning for Fine-Grained Image Classification

Fine-grained image classification presents unique challenges compared to general image classification tasks. As Li et al. (2021) note, fine-grained classification “necessitates discrimination between semantic and instance levels, while considering the similarity and diversity among categories”⁴. This is particularly challenging in bird classification due to three key factors: high intra-class variance (birds of the same species in different postures), low inter-class variance (different species with only minor differences), and limited training data availability, especially for rare species⁴.

Convolutional Neural Networks (CNNs) have revolutionized image classification through their ability to automatically learn hierarchical feature representations. For fine-grained tasks, traditional CNNs face limitations in capturing the subtle distinguishing features between closely related categories. This has led to the development of specialized architectures and techniques focused on identifying discriminative regions in images⁴.

Early approaches to fine-grained classification relied on fixed rectangular bounding boxes and part annotations to obtain visual differences, but these methods required extensive human annotation effort⁴. Recent research has shifted toward weakly supervised approaches that only require image-level labels, developing localization subnetworks to identify critical parts followed by classification subnetworks⁴. These models facilitate

learning while maintaining high accuracy without needing pre-selected boxes, making them more practical for real-world applications.

Recent research emphasizes that effective fine-grained classification depends on identifying and integrating information from multiple discriminative regions rather than focusing on a single region. As highlighted in recent literature, "it is imperative to integrate information from various regions rather than relying on a singular region"⁴. This insight has led to the development of methods combining features from different levels via attention modules, thereby enhancing the semantic and discriminative capacity of features for fine-grained classification⁴.

Transfer Learning for Image Classification

Deep learning, while powerful, comes with two major constraints: dependency on extensive labeled data and high training costs⁶. Transfer learning offers a solution to these limitations by enabling the reuse of knowledge obtained from a source task when training on a target task. In the context of deep learning, this approach is known as Deep Transfer Learning (DTL)⁶.

Transfer learning is particularly valuable for fine-grained bird classification where obtaining large, labeled datasets is challenging. As noted in recent research, "when the sample data is small, transfer learning can help the deep neural network classifier to improve classification accuracy"³. This makes transfer learning an ideal approach for specialized tasks like distinguishing between closely related gull species.

Several studies have demonstrated the efficacy of transfer learning for bird species classification. A study on automatic bird species identification using deep learning achieved a top-5 accuracy of 97.98% by leveraging pretrained CNN networks with a base model to encode images¹⁰. Similarly, research on bird species identification using modified deep transfer learning achieved 98.86% accuracy using the pretrained EfficientNetB5 model¹¹. These results demonstrate that transfer learning approaches can achieve high performance even with limited training data.

Various pretrained models have been evaluated for bird classification tasks, including VGG16, VGG19, ResNet, DenseNet, and EfficientNet architectures. Comparative studies have shown that while all these models can perform effectively, some consistently outperform others. For example, research on drones-birds classification found that "the accuracy and F-Score of ResNet18 exceeds 98% in all cases"⁷, while another study on bird classification reported that "DenseNet201 achieves the best classification accuracy of 98.89% for binary classification"¹⁴.

The transfer learning process typically involves two phases: first freezing most layers of the pretrained model and training only the top layers, then fine-tuning a larger portion of the network while keeping early layers fixed¹¹. This approach preserves the general feature extraction capabilities of the pretrained model while adapting it to the specific characteristics of the target dataset.

Interpretability Techniques for Deep Learning Models

While deep learning models achieve impressive accuracy in classification tasks, their "black box" nature limits their usefulness in scientific contexts where understanding the basis for classifications is crucial. Interpretability techniques address this limitation by providing

insights into model decision-making processes, making them essential tools for applications where transparency is as important as accuracy.

Gradient-weighted Class Activation Mapping (Grad-CAM) has emerged as a particularly valuable technique for visualizing regions of images that influence classification decisions. As described in recent literature, Grad-CAM "uses the gradients of each target that flows into the least convolutional layer to produce a heatmap localization map, highlighting important regions in the image for concept prediction"⁵. This approach enables researchers to validate model decisions against expert knowledge and potentially discover new insights about morphological features.

Visualization studies comparing baseline models with enhanced architectures demonstrate that while basic models often focus on the most conspicuous parts of bird images (such as wings), more sophisticated approaches can discern more intricate features vital for species differentiation⁴. As noted in recent research, enhanced models excel "in identifying not only the prominent features but also the subtle, fine-grained characteristics essential for distinguishing between different bird types"⁴.

These interpretability methods are particularly valuable in fine-grained classification tasks where the differences between categories are subtle and potentially unknown. By highlighting regions that drive model decisions, techniques like Grad-CAM can reveal discriminative features that might not be obvious even to expert observers, potentially advancing biological understanding alongside classification accuracy.

Justification for Deep Learning with Transfer Learning Approach

The choice of deep learning with transfer learning for gull species classification is supported by several compelling factors derived from recent research. Traditional machine learning approaches, while effective for smaller datasets, face limitations when dealing with the complexity of fine-grained visual classification tasks. As demonstrated in comparative studies, "deep learning is more effective than traditional machine learning algorithms in image recognition as the number of bird species increases"³.

The advantages of deep learning architectures for image classification are significant. Unlike traditional machine learning models that require carefully labeled data, "Deep Learning architectures don't necessarily require labelling, as Neural Networks are great at learning without guidelines"¹. Furthermore, in domains like vision, "Deep Learning consistently produces excellent results that significantly outperforms other alternatives"¹.

Transfer learning addresses the primary challenges of deep learning: the need for large datasets and extensive computational resources. By leveraging pretrained models that have already learned general visual features from massive datasets, transfer learning enables the development of highly accurate classifiers with relatively small, domain-specific datasets⁶. This is particularly valuable for this project, which focuses on distinguishing between two specific gull species with limited available data.

The effectiveness of transfer learning for fine-grained bird classification has been consistently demonstrated across multiple studies, with various pretrained models achieving accuracy rates exceeding 98%¹⁰¹¹. These results indicate that transfer learning provides an optimal balance between accuracy and efficiency for the specific task of gull species classification.

The integration of interpretability techniques with transfer learning further strength-

ens this approach by addressing the "black box" limitation of deep neural networks. By implementing methods like Grad-CAM, the project can not only achieve high classification accuracy but also provide insights into the morphological features that drive model decisions, making the results more valuable for scientific applications⁵.

Aims and Objectives

Primary Aims

1. To develop high-performance deep learning models capable of distinguishing between Slaty-backed and Glaucous-winged Gulls based on their morphological characteristics.
2. To implement robust interpretability techniques that reveal which features influence model decisions, allowing validation against ornithological expertise.
3. To analyze whether consistent morphological differences exist between the two species and identify key discriminative features.

Specific Objectives

The project will be carried out in four phases:

[label=Phase 0:]

1. Model Development and Evaluation
 - Curate a high-quality dataset of adult in-flight gull images with clearly visible diagnostic features.
 - Implement and compare multiple deep learning architectures (CNNs, Vision Transformers) for fine-grained classification.
 - Optimize model performance through appropriate regularization techniques, data augmentation, and hyperparameter tuning.
 - Evaluate models using appropriate metrics (accuracy, precision, recall, F1-score) on carefully constructed test sets.
2. Interpretability Implementation
 - Implement Gradient-weighted Class Activation Mapping (Grad-CAM) for convolutional architectures.
 - Develop or adapt interpretability techniques suitable for Vision Transformers.
 - Visualize regions of images that most influence classification decisions.
 - Compare model focus areas with known taxonomic features described in ornithological literature.
3. Feature Analysis
 - Perform quantitative analysis of image regions highlighted by interpretability techniques.

- Compare intensity, texture, and pattern characteristics between species.
- Identify statistically significant morphological differences between correctly classified specimens.

4. Refinement and Validation

- Refine models and interpretability methods based on insights from feature analysis.
- Validate findings against expert ornithological knowledge.
- Document limitations, edge cases, and areas for future research.

0.1 Methodology

0.1.1 Google Colab Platform

Google Colab was selected as the primary platform for developing and training deep learning models. As described by Anjum et al. [2], Google Colab offers significant advantages for machine learning research through its cloud-based environment with integrated GPU acceleration enabling fast model training. The platform’s pre-installed libraries and integration with Google Drive provided an efficient workflow for model development, experimentation, and storage of datasets and trained models. This approach aligns with modern best practices in deep learning research where computational efficiency is crucial for iterative model development and refinement.

Despite its advantages, Google Colab presented a few challenges. The platform frequently disconnected during training sessions, interrupting the model training process before completing all epochs. These disconnections likely stemmed from limited RAM allocation, runtime timeouts, or resource constraints of the shared free GPU environment. As noted by Carneiro et al. [4], while Colab provides robust GPU resources that can match dedicated servers for certain tasks, these free resources “are far from enough to solve demanding real-world problems and are not scalable.”

To mitigate these issues, two strategies were implemented. First, the relatively small size of our dataset helped minimize resource demands. Second, checkpoint saving was implemented throughout the training process, allowing training to resume from the last saved state if disconnections were encountered. This approach ensured that progress wasn’t lost when disconnections occurred, though it introduced some workflow inefficiencies.

0.1.2 Python and PyTorch Framework

The implementation was carried out using Python as the primary programming language, chosen for its extensive library support and widespread adoption in the machine learning community. Python’s simple syntax and powerful libraries make it particularly suitable for rapid prototyping and experimentation in deep learning research [9].

For the deep learning framework, PyTorch was selected over alternatives like TensorFlow or Keras due to its dynamic computational graph which allows for more flexible model development and easier debugging. PyTorch’s intuitive design facilitates a more natural expression of deep learning algorithms while still providing the performance benefits of GPU acceleration. The framework’s robust ecosystem for computer vision tasks,

including pre-trained models and transformation pipelines, was particularly valuable for this fine-grained classification task.

Advantages of PyTorch in Our Implementation

PyTorch offered several key advantages that were particularly beneficial for our transfer learning approach with pre-trained models:

- **Dynamic Computational Graph:** PyTorch’s define-by-run approach allowed for more intuitive debugging and model modification during development. This was especially valuable when adapting pre-trained architectures like VGG16 for our specific classification task.
- **Flexible Model Customization:** The implementation benefited from PyTorch’s object-oriented approach, which made it straightforward to modify pre-trained models, e.g., replacing classification layers while preserving feature extraction capabilities.
- **Efficient Data Loading and Augmentation:** PyTorch’s DataLoader and transformation pipelines facilitated efficient batch processing and on-the-fly data augmentation, which was crucial for maximizing the utility of our limited dataset.
- **Gradient Visualization Tools:** PyTorch’s native support for gradient computation and hooks made implementing Grad-CAM and other visualization techniques more straightforward, enabling better model interpretability.

Similar to approaches described by Raffel et al. [22], our implementation prioritized efficiency and optimization to work within the constraints of limited computational resources, allowing us to achieve high-quality results despite the limitations of the free cloud environment.

0.2 Dataset Preparation and Refinement

The dataset preparation followed a three-stage iterative refinement process, each addressing specific challenges identified during model development. This approach aligns with established methodologies in fine-grained bird classification research, where dataset quality has been shown to significantly impact model performance [10].

0.2.1 Stage 1: Initial Dataset Collection

The initial dataset was collected from public repositories including eBird and iNaturalist, comprising 451 images of Glaucous-winged Gulls and 486 images of Slaty-backed Gulls. This dataset included gulls of various ages (juveniles and adults) in different postures (sitting, standing, and flying). Initial model testing on this dataset yielded poor performance (below 50% accuracy), highlighting the need for dataset refinement. Similar challenges with diverse postures and class imbalance have been documented by Kahl et al. in their work on BirdNET systems [14].

0.2.2 Stage 2: Refined Dataset - Focus on Adult In-flight Images

Consultation with Professor Gibbins, an ornithological expert, revealed that adult wingtip patterns are the most reliable distinguishing features between these species, and these patterns are most visible in flight. This expert-guided refinement approach parallels methods described by Wang et al. in their work on avian dataset construction, where domain expertise significantly improved classification accuracy for visually similar species. [30]. Consequently, the dataset was refined to focus exclusively on adult in-flight images, resulting in a curated collection of 124 Glaucous-winged Gull images and 127 Slaty-backed Gull images. This targeted approach significantly improved model performance, with accuracy increasing to approximately 70%.

By focusing specifically on adult in-flight images where wingtip patterns are most visible, this project addresses the core taxonomic question while minimizing confounding variables. The resulting interpretable classification system aims to provide both a practical identification tool and a scientific instrument for exploring morphological variation within and between these closely related species.

0.2.3 Stage 3: High-Quality Dataset

To further enhance classification performance, 640 high-resolution images of in-flight Slaty-backed Gulls were obtained from Professor Gibbins. The Glaucous-winged Gull dataset was also carefully curated with expert guidance, reducing it to 135 high-quality images that clearly displayed critical wingtip features. Images showing birds in moulting stages, juveniles, or unclear wingtip patterns were systematically removed. This quality-focused approach aligns with findings from Zhou et al., who demonstrated that expert-curated datasets can achieve comparable or superior results with significantly smaller data volumes compared to larger uncurated collections [36].

For comparative analysis, an unrefined dataset containing 632 adult in-flight Glaucous-winged Gulls and 640 high-quality Slaty-backed Gull images was also tested. This multi-dataset evaluation approach follows best practices established in the BirdSet benchmark for avian classification studies [20].

0.3 Transfer Learning Methodology

0.3.1 Theoretical Framework and Rationale

Transfer learning is a powerful machine learning technique that involves reusing a pre-trained model developed for a specific task as a starting point for a new task. This approach significantly enhances learning efficiency by leveraging knowledge gained from solving previous problems, enabling a positive transfer learning effect and reducing the training time required. For fine-grained classification tasks like distinguishing between visually similar gull species, transfer learning is particularly valuable as it allows the model to build upon a foundation of general visual features already learned from diverse datasets.

As highlighted by Kahl et al. (2021), transfer learning addresses two critical challenges in specialized biological classification: data scarcity and feature abstraction [14]. First, data scarcity is a common issue in specialized domains like ornithological image classification, where large-scale annotated datasets are rare. Transfer learning mitigates

this by leveraging models pre-trained on massive datasets like ImageNet. Second, these pre-trained models have learned to extract hierarchical features that capture important visual patterns, which can significantly enhance the accuracy of fine-grained classification tasks.

In our implementation, transfer learning was employed to leverage the robust feature extraction capabilities of pre-trained models on ImageNet. This approach aligns with best practices in fine-grained classification tasks, where lower-level features learned from diverse datasets can be effectively repurposed for specialized domains. The pre-training on ImageNet’s 1.2 million images across 1,000 classes provides the model with a strong foundation for recognizing a wide range of visual patterns, which can then be fine-tuned for the specific task of distinguishing between Glaucous-winged and Slaty-backed Gulls.

Several pre-trained architectures were evaluated for this task, with VGG-16. [28] demonstrating superior performance in our specific classification context. The effectiveness of transfer learning was evident in the rapid convergence and high accuracy achieved even with our relatively limited dataset of gull images, demonstrating the potential of this approach for specialized biological classification tasks.

0.4 Model Architectures and Implementation

0.4.1 VGG-16 Architecture

Theoretical Background

VGG-16 is a popular Convolutional Neural Network (CNN) architecture widely used in computer vision applications. Originally developed by Simonyan and Zisserman [28], VGG-16 consists of 16 layers, including 13 convolutional layers arranged in blocks of increasing depth, followed by 3 fully connected layers, with a total of approximately 138 million parameters. The architecture follows a systematic approach of stacking convolutional layers with small 3×3 filters followed by max-pooling layers, creating a deep network capable of learning complex hierarchical features.

One of the main advantages of using VGG-16 as a pre-trained model for transfer learning is its ability to capture a wide range of features and patterns in images. This capability stems from the deep architecture of the VGG-16 model, which allows it to extract more complex features from images compared to shallower models. VGG-16 has been pre-trained on the ImageNet dataset [6], which contains over 1.2 million images across 1,000 classes, enabling it to recognize a wide range of features and patterns applicable to various computer vision tasks.

The architecture’s elegant simplicity, despite its depth, makes it particularly effective for fine-grained visual classification tasks like ours. Its performance on various computer vision benchmarks, including object detection, image segmentation, and image classification tasks, makes it a versatile choice for transfer learning in numerous applications. For our specific task of gull species classification, the hierarchical feature representation capabilities of VGG-16 proved particularly effective at capturing the subtle differences in wing patterns and morphological features that distinguish between the target species.

Model Adaptation for Gull Species Classification

The pre-trained VGG16 model was adapted for our binary classification task through targeted modifications to the final classification layer. The implementation approach

follows best practices for fine-grained bird classification established in recent research on transfer learning for avian species identification [11, 23]. Specifically, the original 1000-class classifier was replaced with a custom binary classification head while preserving the feature extraction capabilities of the convolutional base. The model modification strategy followed this approach:

1. Loading the pre-trained VGG16 with ImageNet weights
2. Extracting the number of features from the original classifier (4096)
3. Replacing the final layer with a sequential block containing:
 - (a) A dropout layer with rate 0.4 for regularization
 - (b) A linear layer mapping from 4096 features to 2 output classes

This implementation maintained the complex feature hierarchy learned by VGG16 while adapting the final classification stage for our specific binary task. The relatively high dropout rate (0.4) was strategically implemented to address potential overfitting challenges common in fine-grained classification tasks with limited training data, particularly important given the visual similarities between the target gull species.

0.4.2 Data Preprocessing and Augmentation Strategy

Image Preprocessing

Images were preprocessed using a consistent pipeline to ensure compatibility with the VGG16 architecture. All images were resized to 224×224 pixels to match VGG16’s expected input dimensions. Following resize operations, pixel values were normalized using ImageNet mean values [0.485, 0.456, 0.406] and standard deviations [0.229, 0.224, 0.225]. This normalization strategy ensures input distributions match those seen during pre-training, facilitating effective transfer learning.

Training Augmentation

A comprehensive data augmentation strategy was implemented to enhance model generalization capabilities and mitigate overfitting [26]. The augmentation pipeline for training incorporated multiple techniques designed to preserve critical taxonomic features while introducing beneficial variability:

Geometric Transformations:

- Random horizontal flips (probability 0.5)
- Random vertical flips (probability 0.3)
- Small rotations (± 10 degrees)
- Minor affine transformations (translations up to 5%)
- Random resized crops (scale 0.95-1.0 of original size)

Color and Appearance Transformations:

- Brightness, contrast, and saturation adjustments ($\pm 10\%$)
- Sharpness enhancement (factor 1.2, probability 0.3)

This augmentation strategy was carefully calibrated to maintain the integrity of critical morphological features (particularly wingtip patterns) while simulating natural variations in viewing conditions. The relatively conservative parameter choices reflect the importance of preserving diagnostic features in fine-grained classification tasks [31].

Validation and Testing Preprocessing

For validation and testing phases, a simplified transformation pipeline was employed, consisting only of resizing to 224×224 pixels, tensor conversion, and normalization. This approach ensures consistent evaluation conditions while maintaining the statistical properties expected by the pre-trained model.

0.4.3 Training Methodology

Dataset Organization and Splitting

The dataset was organized using the ImageFolder structure, with a 95:5 split between training and validation sets. This configuration provided substantial training data while maintaining a sufficient validation set for hyperparameter tuning and model selection. A separate test set was maintained for final performance evaluation.

Loss Function and Optimization

Cross-entropy loss was selected as the objective function for this binary classification task, providing appropriate gradients for optimization. This loss function effectively quantifies the discrepancy between predicted class probabilities and ground truth labels.

The AdamW optimizer [17] was employed with carefully tuned hyperparameters to facilitate effective model training:

- Learning rate: 0.0001 (relatively low to enable stable fine-tuning)
- Weight decay: 0.001 (for L2 regularization to prevent overfitting)

This optimization configuration balances the need for fine-grained weight adjustments with regularization to maintain generalization capacity.

Learning Rate Scheduling and Training Stabilization

An adaptive learning rate schedule was implemented using ReduceLROnPlateau with a patience factor of 3 epochs and a reduction factor of 0.1 [33]. This approach automatically reduces the learning rate when validation performance plateaus, enabling finer weight adjustments as the model approaches optimal parameters.

To enhance training stability, gradient clipping was applied with a maximum norm of 2.0. This technique prevents exploding gradients by constraining the magnitude of parameter updates, which is particularly valuable when fine-tuning deep architectures like VGG16.

Batch Processing and Training Duration

The model was trained with a batch size of 16, striking a balance between computational efficiency and effective mini-batch gradient estimation. Training was configured to run for a maximum of 30 epochs with early stopping based on validation performance, ensuring optimal model selection while avoiding overfitting.

Model checkpoints were saved after each epoch, preserving the best-performing model configurations for subsequent evaluation and deployment.

0.4.4 Evaluation Metrics and Performance

The model’s performance was assessed using multiple complementary metrics to ensure robust evaluation:

- Accuracy: Percentage of correctly classified images
- Precision: Proportion of true positive predictions among all positive predictions
- Recall: Proportion of true positives identified among all actual positives
- F1-Score: Harmonic mean of precision and recall

The final VGG16 model achieved exceptional performance with 98.80% validation accuracy and 100% test accuracy, demonstrating its effectiveness in distinguishing between the two gull species based on the refined dataset. This performance exceeds typical benchmarks for fine-grained bird classification tasks [24], highlighting the effectiveness of the implemented architecture, data preparation strategy, and training methodology.

A confusion matrix analysis confirmed the model’s strong classification performance across both classes, with minimal misclassifications. This indicates the model successfully learned the discriminative morphological features necessary for distinguishing between Glaucous-winged and Slaty-backed Gulls.

0.5 Vision Transformer (ViT) Architecture

0.5.1 Theoretical Framework

Vision Transformers (ViT) represent a paradigm shift in computer vision, applying the self-attention mechanism from transformer models—originally developed for natural language processing—to image classification tasks. First introduced by [8], ViT processes images by dividing them into a sequence of fixed-size patches, which are then linearly embedded and processed through transformer encoder blocks.

Unlike CNNs that build hierarchical representations through local convolutional operations, ViT applies self-attention mechanisms to capture global relationships between image patches. This allows the model to attend to long-range dependencies within the image, potentially capturing more holistic patterns. As demonstrated by research from [16], these attention mechanisms enable transformers to excel at detecting subtle features in biomedical images by focusing on the most discriminative regions.

The standard ViT architecture consists of:

- Patch embedding layer that converts image patches to token embeddings

- Position embedding to provide spatial information
- Multiple transformer encoder blocks with multi-head self-attention
- Layer normalization and MLP blocks within each transformer layer
- A classification head for prediction

This architecture’s capacity to model global relationships makes it particularly promising for fine-grained classification tasks where relationships between distant parts of an image (e.g., wing patterns in relation to head features) may be important for accurate classification [5].

0.5.2 Standard Vision Transformer Implementation

My standard ViT implementation utilized the pre-trained ‘vit_base_patch16_224’ model from the TIMM library, which features a patch size of 16×16 pixels and was trained on the ImageNet dataset. The model adaptation process included:

- Loading the pre-trained ViT model with frozen weights
- Extracting the embedding dimension from the original model (768 features)
- Replacing the classification head with a binary classifier for our gull species task
- Maintaining the self-attention mechanisms and transformer blocks

This approach leverages the powerful feature extraction capabilities of ViT while customizing the final classification stage for our specific task. The implementation follows best practices established by [32] for adapting vision transformers to specialized classification tasks.

0.5.3 Enhanced Vision Transformer with Custom Attention

To further improve the model’s ability to focus on taxonomically relevant features, we developed an Enhanced Vision Transformer (EnhancedViT) that incorporates a custom attention mechanism specifically designed for fine-grained classification tasks.

The key innovation in this implementation is an attention-based pooling layer that computes importance scores for each patch token, enabling the model to focus on the most discriminative regions of the input image. This approach draws inspiration from the work of [12], who demonstrated that specialized attention mechanisms in vision transformers can significantly improve fine-grained classification by emphasizing taxonomically relevant features.

The enhanced ViT architecture extends the standard implementation with:

- A custom attention layer that computes importance scores for each token
- An attention-weighted aggregation step that prioritizes informative regions
- A multi-layer perceptron classifier with dropout regularization
- Layer normalization for improved training stability

The attention mechanism was implemented as:

Attention-based Token Pooling

- 1: Compute attention scores for each token using a learned projection
- 2: Normalize scores using softmax to create attention weights
- 3: Perform weighted aggregation of tokens based on attention weights
- 4: Process the attention-weighted representation through the classifier

This approach allows the model to dynamically focus on the most relevant parts of the image for classification, such as distinctive wingtip patterns or other morphological features that differentiate between gull species [stassin2024explainability].

0.5.4 Data Processing and Augmentation

Both ViT implementations used standardized preprocessing and augmentation pipelines:

- Resize operations to 224×224 pixels (the standard input size for ViT models)
- Normalization with mean $[0.5, 0.5, 0.5]$ and standard deviation $[0.5, 0.5, 0.5]$
- Augmentation techniques including:
 - Random horizontal flipping
 - Random rotation (± 15 degrees)
 - Color jittering (brightness, contrast, saturation)

The input normalization values specifically used $[0.5, 0.5, 0.5]$ rather than ImageNet statistics, following recommendations from[touvron2021training] for transfer learning with vision transformers.

0.5.5 Training Methodology

The training approach for both ViT variants included:

- AdamW optimizer with learning rate 0.0001 and weight decay $1e-4$
- Learning rate scheduling with ReduceLROnPlateau (patience=3, factor=0.1)
- Batch size of 16 to balance computational efficiency and training stability
- Training for 20 epochs with early stopping based on validation performance

For the EnhancedViT, we employed additional training refinements:

- Layer-wise learning rate decay to fine-tune different components at appropriate rates
- Dropout regularization ($p=0.3$) in the custom classification head
- Checkpoint saving to preserve the best-performing model configuration

Both models were trained on the refined high-quality dataset (Stage 3), with an 80:20 split between training and validation sets to ensure robust performance evaluation during development.

0.6 Inception v3 Architecture

0.6.1 Theoretical Background

Inception v3, developed by [szegedy2016rethinking], represents a sophisticated CNN architecture designed to efficiently capture multi-scale features through parallel convolution paths with different kernel sizes. The key innovation in Inception architectures is the use of "Inception modules" that process the same input tensor through multiple convolutional paths with different receptive fields, and then concatenate the results. This enables the network to capture both fine-grained local patterns and broader contextual information simultaneously.

Inception v3 builds upon earlier versions with several important architectural improvements:

- Factorized convolutions to reduce computational cost
- Spatial factorization into asymmetric convolutions (e.g., $1 \times n$ followed by $n \times 1$)
- Auxiliary classifiers that inject additional gradient signals during training
- Batch normalization for improved training stability
- Label smoothing regularization to prevent overconfidence

These design elements collectively enable Inception v3 to achieve high accuracy while maintaining computational efficiency. As demonstrated by [shu2023multiscale], Inception architectures are particularly effective for tasks requiring multi-scale feature extraction, such as discriminating between visually similar biological specimens.

0.6.2 Model Adaptation for Gull Classification

Our implementation adapted the pre-trained Inception v3 model for gull species classification using the following approach:

1. Loading the pre-trained Inception v3 model with ImageNet weights
2. Extracting the feature dimension from the original classifier (2048)
3. Replacing the final classifier with a custom sequence:
 - (a) Dropout layer ($p=0.5$) for regularization
 - (b) Linear layer mapping 2048 features to 2 output classes

A distinctive aspect of our Inception v3 implementation was the utilization of auxiliary outputs during training. Inception v3's auxiliary classifier, which branches off from an intermediate layer, provides an additional gradient path during backpropagation. This approach helps combat the vanishing gradient problem and provides regularization, as noted by [szegedy2016rethinking] in their original paper.

The loss function was modified to incorporate both the main output and the auxiliary output during training:

$$\text{loss} = \text{main_loss} + 0.3 \times \text{auxiliary_loss} \quad (1)$$

where the auxiliary loss weight (0.3) was selected based on empirical optimization and aligns with recommendations in the literature for fine-tuning Inception architectures [touvron2021training].

0.6.3 Advanced Training Techniques

The Inception v3 implementation incorporated several advanced training techniques to optimize performance:

- Mixed-precision training using PyTorch’s Automatic Mixed Precision (AMP) to accelerate computation while maintaining numerical stability [micikevicius2018mixed]
- Gradient clipping with a maximum norm of 2.0 to prevent explosive gradient updates [you2019large]
- Precisely tuned learning rate and weight decay parameters identified through hyperparameter optimization
- Layer-wise learning rate adjustment to fine-tune different parts of the network at appropriate rates

These techniques collectively enhanced training efficiency and model performance. The implementation of mixed-precision training was particularly valuable given the resource constraints of the Google Colab environment, as it reduced memory usage and accelerated computation without compromising model accuracy [11].

0.6.4 Data Processing Pipeline

The data processing pipeline for Inception v3 was adapted to the model’s specific requirements:

- Resize operations to 299×299 pixels (the standard input size for Inception v3)
- Standard data augmentation techniques for training:
 - Random horizontal flipping
 - Random rotation (± 15 degrees)
 - Color jittering
- Simple resizing and normalization for validation and testing

The larger input resolution (299×299 vs 224×224 used by VGG16 and ViT) provides the Inception architecture with more detailed information, potentially beneficial for capturing the subtle wing pattern differences between gull species [8]. natbib

0.7 ResNet50 Architecture

0.7.1 Theoretical Background

Residual Networks (ResNet) represent a significant innovation in deep neural network architecture, introduced by He et al. to address the degradation problem that occurs when training very deep networks. The key innovation in ResNet is the introduction of skip connections or "shortcut connections" that bypass one or more layers, allowing gradients to flow more easily through the network during backpropagation [he2016deep]. This design enables the training of much deeper networks than was previously feasible, with ResNet-50 containing 50 layers organized in residual blocks.

ResNet-50 architecture consists of five stages, each containing multiple residual blocks. Each residual block contains a "shortcut" that skips over the main path and rejoins it later, allowing the network to learn residual functions with reference to the layer inputs rather than learning unreferenced functions [he2016deep, he2016identity]. This approach enables ResNet to achieve high performance on image classification tasks while mitigating the vanishing gradient problem common in very deep networks.

The architecture's ability to effectively extract hierarchical features through its deep structure makes it particularly well-suited for fine-grained classification tasks where subtle differences must be detected. As noted by Ghani et al., ResNet architectures have demonstrated strong performance in avian classification tasks due to their capacity to learn discriminative features at multiple scales and levels of abstraction [11].

0.7.2 Model Adaptation for Gull Species Classification

For our gull classification task, we adapted the pre-trained ResNet-50 model using a focused transfer learning approach. The model was initialized with weights pre-trained on the ImageNet dataset, providing a strong foundation of general visual features. The adaptation process involved:

1. Loading the pre-trained ResNet-50 model with ImageNet weights
2. Preserving the convolutional backbone to maintain feature extraction capabilities
3. Replacing the final fully connected layer (classifier) with a custom sequence:
 - (a) Dropout layer with probability 0.5 to reduce overfitting
 - (b) Linear layer mapping from 2048 features to 2 output classes

This adaptation strategy preserved ResNet-50's powerful feature extraction capabilities while customizing the classification head for our binary task. The relatively high dropout rate (0.5) was implemented to address potential overfitting, which is particularly important given the visual similarities between the target species and the limited size of our specialized dataset [srivastava2014dropout].

0.7.3 Image Preprocessing and Enhancement

A distinctive aspect of our ResNet-50 implementation was the incorporation of image sharpening as a preprocessing step. This approach was motivated by research from Zhou

et al. showing that enhancing edge definition can improve the detection of subtle morphological features in avian classification tasks [zhou2022effective]. The image enhancement process applied a 3×3 sharpening kernel through a custom preprocessing function:

$$\text{Sharpening Kernel:} \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

This technique enhanced the visibility of critical diagnostic features like wingtip patterns while preserving the overall image content. To ensure consistency, image sharpening was applied across both training and evaluation pipelines [zhao2021image].

0.7.4 Data Augmentation Strategy

The data augmentation pipeline for ResNet-50 was structured to enhance model robustness while preserving class-discriminative features:

- Resize operations (300×300 pixels) followed by sharpening
- Random horizontal flipping to simulate viewpoint variation
- Random rotation (± 15 degrees) to account for flight angle variability
- Color jittering (brightness, contrast, saturation adjusted by $\pm 20\%$)
- Random cropping with padding to vary focus regions

For validation and testing, a more conservative approach was employed with resizing, center cropping (256×256 pixels), and the same sharpening preprocessing to maintain feature clarity without introducing variability [26].

0.7.5 Training Approach and Optimization

The ResNet-50 model was trained using the following methodological approach:

- Adam optimizer with learning rate 0.001 and weight decay 1e-4 for regularization
- Adaptive learning rate scheduling using ReduceLROnPlateau with patience=3
- Early stopping with patience=5 to prevent overfitting
- Batch size of 16 for efficient GPU utilization

The implementation of early stopping was particularly valuable for the ResNet model, as it helped prevent overfitting to the training data while ensuring the model retained its generalization capabilities. As demonstrated by Huang et al., early stopping acts as an effective regularization technique for deep networks when working with specialized datasets of limited size [huang2022early].

0.8 Custom CNN with Squeeze-and-Excitation Blocks

0.8.1 Architectural Innovation

To complement the transfer learning approach with pre-trained models, we developed a custom CNN architecture specifically designed for our gull classification task. The architecture incorporates Squeeze-and-Excitation (SE) blocks, an attention mechanism introduced by Hu et al. (2018) that adaptively recalibrates channel-wise feature responses by explicitly modeling interdependencies between channels.

The SE mechanism enhances standard convolutional operations by adding two operations:

- A "squeeze" operation that aggregates feature maps across spatial dimensions to produce a channel descriptor
- An "excitation" operation that produces per-channel modulation weights

This channel-wise attention mechanism allows the network to emphasize informative features and suppress less useful ones, improving the representational power of the network. As demonstrated by Hu et al. (2018), the SE mechanism yields significant performance improvements while adding minimal computational overhead.

Our custom CNN implementation follows this architectural pattern:

0.8.2 Addressing Class Imbalance

An important methodological consideration in our custom CNN implementation was addressing potential class imbalance in the dataset. To ensure balanced learning despite the uneven distribution of examples between classes, we implemented a weighted sampling approach based on class frequencies.

The implementation calculated inverse class weights to prioritize examples from underrepresented classes:

Weighted Sampling for Class Balance

1. Count examples per class in the training dataset
2. Calculate inverse class frequencies: $\text{weights} = 1 / \text{class_counts}$
3. Assign a sampling weight to each training example based on its class
4. Create a `WeightedRandomSampler` using these weights
5. Use the sampler in the `DataLoader` to achieve balanced mini-batches

This approach ensures that the model receives a balanced distribution of examples during training, preventing bias toward the majority class. The effectiveness of this technique for handling class imbalance has been demonstrated in fine-grained classification research by Buda et al. (2018), who showed that sampling strategies can significantly improve model performance on imbalanced datasets.

0.8.3 Training Methodology

The custom CNN was trained using the following approach:

- Adam optimizer with learning rate 0.001 and weight decay 0.0005
- Cosine Annealing learning rate scheduler for cyclical learning rate adjustment

- Cross-entropy loss function
- Batch size of 32 (larger than the pre-trained models due to lower memory requirements)
- Training for 20 epochs with checkpoint saving for best-performing models

The use of Cosine Annealing for learning rate scheduling represents a different approach compared to the ReduceLROnPlateau used with the pre-trained models. This scheduler cyclically varies the learning rate between a maximum and minimum value following a cosine function, helping the model escape local minima and potentially converge to better solutions. This approach aligns with research by Loshchilov and Hutter (2017) demonstrating the effectiveness of cyclical learning rates for CNN training.

thesis