

# Hierarchical Vision Transformers for Spatial Audio Image Classification: Methodology, Training Design, and Comprehensive Evaluation

Douglas Abreu<sup>1</sup>

<sup>1</sup>Ph.D. Program, Department of Computer Science, *Institution*

September 29, 2025

## Abstract

We present a comprehensive methodological study on spatial audio image classification using hierarchical Vision Transformers (ViTs), emphasizing Swin Transformer backbones [11] augmented with carefully designed training, optimization, and evaluation protocols. We formulate the supervised learning problem, detail the loss, optimization, and scheduling strategies (AdamW [14], cosine warmup [13]), regularization via stochastic depth [9], and strong data augmentation [16, 23]. A rigorous 5-fold cross-validation design provides reliable model selection and uncertainty estimates [10]. We report test performance with detailed per-class metrics, confusion matrices, and calibration analysis [7]. Results demonstrate robust accuracy (95.2%), strong macro-F1 (95.18%), and balanced performance across classes. We discuss ablations, parameter justifications, and implications for spatial audio visual representations.

## 1 Introduction

Spatial audio analysis increasingly leverages visual representations derived from time-frequency projections, ambisonics maps, or beamformed spatial energy distributions. Transformers [5] have become state-of-the-art in computer vision due to their global receptive fields and scalable architectures. Swin Transformers [11] introduce shifted windows and hierarchical feature pyramids, enabling efficient multi-scale processing with competitive accuracy-cost trade-offs. This work investigates a Swin-based pipeline tailored for spatial audio images.

We contribute: (i) a principled training pipeline grounded in theory and empirical best practices; (ii) a reproducible configuration schema; (iii) an extensive evaluation including cross-validation, test-set breakdowns, and calibration; and (iv) a discussion relating design choices to theory and prior literature.

**Motivation and scope.** Spatial hearing and sound-field understanding demand models that are simultaneously robust to nuisance variability (sensor placement, scene dynamics, noise) and

sensitive to diagnostic spatial patterns. Hierarchical Vision Transformers are well-suited because their windowed self-attention preserves locality while enabling global context aggregation through window shifting and hierarchical downsampling. In practice, this balance translates into models that generalize across acquisition conditions without sacrificing discriminability among subtle spatial cues. Our goal is to provide a complete, engineering-grade training and evaluation recipe that is theoretically well-founded and empirically validated, producing reliable models and transparent analyses that can inform downstream scientific and industrial use.

## 2 Related Work

Convolutional neural networks (CNNs) such as ResNet [8] and EfficientNet [18] established strong image classification baselines through depth/width/ resolution scaling and residual connections. Vision Transformers (ViT) [5] demonstrated that attention-only architectures, when pretrained at scale, can surpass CNNs in accuracy and robustness; subsequent data-efficient variants (e.g., DeiT [19]) reduced reliance on massive pretraining corpora. Swin Transformer [11] introduced hierarchical, shifted-window attention to capture multi-scale structure with linear complexity in image size, achieving state-of-the-art results on classification, detection, and segmentation benchmarks while maintaining computational efficiency. Concurrently, ConvNeXt [12] revisited CNN design with Transformer-era training techniques, closing gaps between CNNs and ViTs and highlighting the importance of training protocols.

Regularization and augmentation are critical for robustness. Mixup [22] and CutMix [21] improve generalization via vicinal risk minimization and patch-level mixing, while Random Erasing [23] and color/ geometric perturbations reduce over-reliance on spurious cues. For optimization, AdamW [14] with cosine scheduling and warmup [13, 6] has become a de facto standard for training large models stably. Calibration analyses [7] reveal that modern neural networks can be miscalibrated, motivating post-hoc temperature scaling and threshold optimization in deployment. Our work integrates these strands in a coherent pipeline for spatial audio images and contributes a rigorous cross-validation analysis [10, 4, 3].

## 3 Problem Formulation

Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  be a labeled dataset with  $x_i \in \mathbb{R}^{H \times W \times C}$  and class labels  $y_i \in \{1, \dots, K\}$ . A model  $f_\theta : \mathbb{R}^{H \times W \times C} \rightarrow \Delta^{K-1}$  predicts class probabilities  $p(y \mid x; \theta)$ . Training minimizes the expected cross-entropy with optional label smoothing [17]:

$$\mathcal{L}_{\text{CE}}(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \tilde{y}_{ik} \log p(y = k \mid x_i; \theta), \quad \tilde{y}_{ik} = (1 - \epsilon) \mathbf{1}[y_i = k] + \frac{\epsilon}{K}. \quad (1)$$

We use AdamW [14] with decoupled weight decay  $\lambda$  and discriminative learning rates for head and backbone:  $\eta_{\text{head}} > \eta_{\text{backbone}}$ . Learning rate follows cosine decay with warmup [13]. Mixed-precision training reduces memory and accelerates compute [15].

### 3.1 Optimization Theory and Schedules

AdamW updates parameters using biased moment estimates  $(m_t, v_t)$  and decoupled L2 regularization [14]:

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t, \quad \hat{m}_t = \frac{m_t}{1-\beta_1^t}, \quad (2)$$

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) g_t^2, \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}, \quad (3)$$

$$\theta_{t+1} = \theta_t - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} - \eta_t \lambda \theta_t. \quad (4)$$

We employ cosine annealing with linear warmup: for total steps  $T$ , warmup steps  $T_w$ , initial LR  $\eta_0$ , minimum LR  $\eta_{\min}$ ,

$$\eta_t = \begin{cases} \eta_0 \frac{t}{T_w}, & t \leq T_w, \\ \eta_{\min} + \frac{1}{2}(\eta_0 - \eta_{\min}) \left[ 1 + \cos \left( \pi \frac{t - T_w}{T - T_w} \right) \right], & t > T_w. \end{cases} \quad (5)$$

This schedule eases optimization in the presence of large gradients early on and encourages exploration-exploitation dynamics aligned with flat-minima generalization.

### 3.2 Regularization Mathematics

With stochastic depth [9], residual block  $F(\cdot; \theta)$  is skipped with probability  $p$  during training:  $y = x + b F(x; \theta)$ , where Bernoulli  $b \sim \text{Bernoulli}(1-p)$ . At inference we use the expectation  $\mathbb{E}[b] = 1-p$  by scaling the residual, preserving expected activations. Label smoothing with  $\epsilon$  reduces overconfidence by redistributing a small mass uniformly over classes.

## 4 Architecture

### 4.1 Swin Transformer: Hierarchical Vision Transformer with Shifted Windows

The Swin Transformer [11] introduces a hierarchical architecture that computes self-attention within local windows and enables cross-window connections via a shifted window partitioning scheme. Unlike the original Vision Transformer (ViT) [5], which maintains a fixed resolution throughout all layers and computes global self-attention, Swin Transformer constructs a hierarchical representation with progressively downsampled feature maps, analogous to convolutional architectures.

#### 4.1.1 Patch Embedding and Hierarchical Structure

An input image  $x \in \mathbb{R}^{H \times W \times 3}$  is first partitioned into non-overlapping patches of size  $4 \times 4$ . Each patch is treated as a token and linearly projected to dimension  $C = 128$ . This yields an initial feature map of shape  $\frac{H}{4} \times \frac{W}{4} \times C$ .

The architecture comprises four hierarchical stages with progressively smaller spatial resolutions and larger channel dimensions:

- **Stage 1:** Resolution  $\frac{H}{4} \times \frac{W}{4}$ , dimension  $C = 128$ , 2 Swin Transformer blocks
- **Stage 2:** Resolution  $\frac{H}{8} \times \frac{W}{8}$ , dimension  $2C = 256$ , 2 Swin Transformer blocks
- **Stage 3:** Resolution  $\frac{H}{16} \times \frac{W}{16}$ , dimension  $4C = 512$ , 18 Swin Transformer blocks
- **Stage 4:** Resolution  $\frac{H}{32} \times \frac{W}{32}$ , dimension  $8C = 1024$ , 2 Swin Transformer blocks

Between stages, a *patch merging* layer concatenates features from  $2 \times 2$  neighboring patches and applies a linear projection, effectively downsampling spatial resolution by  $2 \times$  while doubling channel dimension.

#### 4.1.2 Shifted Window Multi-Head Self-Attention

The core innovation of Swin is the *shifted window* mechanism. Standard global self-attention has quadratic complexity  $\mathcal{O}((HW)^2)$ . Swin constrains attention to non-overlapping local windows of size  $M \times M$ , reducing complexity to  $\mathcal{O}(HW \cdot M^2)$ , which is linear in image size when  $M$  is fixed.

For our configuration (`swin_base_patch4_window12_384`), we use window size  $M = 12$  on input resolution  $384 \times 384$ . At stage 1 ( $\frac{384}{4} = 96$  tokens per side), each window contains  $12 \times 12 = 144$  tokens.

**Window-based Multi-Head Self-Attention (W-MSA):** Given features  $\mathbf{z} \in \mathbb{R}^{HW \times C}$  partitioned into  $\lceil \frac{H}{M} \rceil \times \lceil \frac{W}{M} \rceil$  non-overlapping windows, we compute self-attention independently within each window:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{SoftMax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} + \mathbf{B} \right) \mathbf{V}, \quad (6)$$

where  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  are query, key, value projections, and  $\mathbf{B} \in \mathbb{R}^{M^2 \times M^2}$  is a learnable relative position bias [11].

**Shifted Window Multi-Head Self-Attention (SW-MSA):** To enable cross-window connections, consecutive Swin blocks alternate between regular and *shifted* window partitioning. In shifted mode, windows are displaced by  $(\lfloor \frac{M}{2} \rfloor, \lfloor \frac{M}{2} \rfloor)$  pixels. This creates cross-window interactions while maintaining computational efficiency through a cyclic-shifting and masking mechanism.

A Swin Transformer block alternates W-MSA and SW-MSA:

$$\hat{\mathbf{z}}^{(\ell)} = \text{W-MSA}(\text{LN}(\mathbf{z}^{(\ell-1)})) + \mathbf{z}^{(\ell-1)}, \quad (7)$$

$$\mathbf{z}^{(\ell)} = \text{MLP}(\text{LN}(\hat{\mathbf{z}}^{(\ell)})) + \hat{\mathbf{z}}^{(\ell)}, \quad (8)$$

$$\hat{\mathbf{z}}^{(\ell+1)} = \text{SW-MSA}(\text{LN}(\mathbf{z}^{(\ell)})) + \mathbf{z}^{(\ell)}, \quad (9)$$

$$\mathbf{z}^{(\ell+1)} = \text{MLP}(\text{LN}(\hat{\mathbf{z}}^{(\ell+1)})) + \hat{\mathbf{z}}^{(\ell+1)}, \quad (10)$$

where LN denotes LayerNorm, MLP is a two-layer feed-forward network with GELU activation, and residual connections facilitate gradient flow.

#### 4.1.3 Model Configuration: Swin-Base

We employ the `swin_base_patch4_window12_384` variant with the following specifications:

- Total parameters: 87.9M (88M)
- Transformer blocks: [2, 2, 18, 2] across four stages
- Embedding dimension:  $C = 128$
- Number of attention heads: [4, 8, 16, 32] per stage
- Window size:  $M = 12$
- MLP ratio: 4 (hidden dimension =  $4C$ )
- Pretraining: ImageNet-21k (14M images, 21k classes)

The model is pretrained on ImageNet-21k at  $384 \times 384$  resolution, providing strong transferable representations for downstream visual recognition tasks.

#### 4.1.4 Classification Head

We replace the pretrained classification head with a task-specific two-layer MLP:

$$\mathbf{h} = \text{GlobalAvgPool}(\mathbf{z}^{(L)}), \quad \mathbf{h} \in \mathbb{R}^{1024}, \quad (11)$$

$$\mathbf{h}_1 = \text{ReLU}(\mathbf{W}_1 \mathbf{h} + \mathbf{b}_1), \quad \mathbf{h}_1 \in \mathbb{R}^{512}, \quad (12)$$

$$\mathbf{h}_2 = \text{Dropout}(\mathbf{h}_1, p = 0.2), \quad (13)$$

$$\mathbf{y} = \text{Softmax}(\mathbf{W}_2 \mathbf{h}_2 + \mathbf{b}_2), \quad \mathbf{y} \in \mathbb{R}^4. \quad (14)$$

This lightweight head introduces only  $(1024 \times 512) + (512 \times 4) \approx 0.5\text{M}$  additional parameters ( $\sim 0.6\%$  of total), minimizing overfitting risk while providing sufficient capacity for the 4-class spatial audio classification task.

## 4.2 Regularization: Stochastic Depth

We apply stochastic depth [9] with drop probability  $p = 0.1$ . During training, each residual block is randomly dropped with probability  $p$ , effectively training an ensemble of sub-networks:

$$\mathbf{z}^{(\ell)} = \mathbf{z}^{(\ell-1)} + b^{(\ell)} \cdot \mathcal{F}^{(\ell)}(\mathbf{z}^{(\ell-1)}), \quad b^{(\ell)} \sim \text{Bernoulli}(1 - p). \quad (15)$$

At inference, we scale the residual by its survival probability:  $\mathbf{z}^{(\ell)} = \mathbf{z}^{(\ell-1)} + (1 - p) \cdot \mathcal{F}^{(\ell)}(\mathbf{z}^{(\ell-1)})$ , preserving expected activations. This acts as implicit model averaging and improves generalization, particularly in deep networks with 24 transformer blocks.

## 5 Data and Preprocessing

### 5.1 Dataset Characteristics and Splits

The dataset comprises four classes (100, 200, 510, 514) derived from spatial audio renderings with consistent preprocessing across splits. The complete dataset contains 750 test samples and approximately 3000 training samples distributed across classes with near-balanced representation (approximately 750 samples per class).

We adopt a rigorous three-way split strategy:

- **Training set** (70%): Used for gradient-based optimization
- **Validation set** (20%): Used for hyperparameter selection, early stopping, and learning rate scheduling
- **Test set** (10%, held-out): Used exclusively for final performance evaluation, never exposed during training or model selection

Stratified splitting preserves per-class proportions across partitions, ensuring that each split maintains approximately 25% representation per class. This mitigates class imbalance effects and reduces selection bias [10]. For cross-validation experiments, we further partition the training+validation data into 5 stratified folds, enabling robust variance estimation of model performance.

### 5.2 Data Augmentation Pipeline

Data augmentation is critical for preventing overfitting and improving model robustness. We apply a carefully calibrated pipeline during training that preserves semantic structure while inducing invariances to nuisance variability:

#### Training Augmentations

1. **RandomResizedCrop**([384, 384], scale=[0.8, 1.0], ratio=[0.9, 1.1]): Samples random crops at 80-100% of original scale with near-square aspect ratios. This induces scale invariance while preserving the approximately square structure of spatial audio images. The scale range is conservative to avoid excessive information loss that could distort spatial cues.
2. **HorizontalFlip**( $p = 0.5$ ): Applies left-right mirroring with 50% probability. This augmentation is justified by the spatial symmetry present in many acoustic scenarios, where horizontal reflection preserves physical plausibility.
3. **Rotation**( $\pm 15^\circ$ ): Random rotations within  $\pm 15^\circ$  accommodate sensor misalignment and orientation variability. The moderate rotation range prevents distortion of spatial cues beyond recognition. Larger rotations ( $> 20^\circ$ ) were empirically found to degrade calibration, suggesting that excessive geometric perturbation may distort essential spatial patterns.

4. **ColorJitter**(brightness=0.2, contrast=0.2, saturation=0.2, hue=0.1): Perturbs color channels independently to reduce reliance on spurious colorization artifacts and improve robustness to visualization parameter variations. Since spatial audio images are pseudo-color representations of acoustic fields, color augmentation prevents the model from overfitting to specific colormap conventions.
5. **GaussianBlur**( $p = 0.1$ ,  $\sigma = 0.1$ ): Applies Gaussian low-pass filtering with 10% probability, simulating acquisition noise and preprocessing artifacts that may arise from different rendering pipelines or sensor characteristics.
6. **RandomErasing**( $p = 0.25$ , max\_area=0.33): Randomly masks rectangular regions covering up to 33% of the image with probability 0.25. This regularization technique [23] forces the model to develop distributed spatial representations rather than over-relying on localized salient regions, improving robustness to occlusions and partial observations.

**Validation and Test Preprocessing** For validation and test sets, we apply deterministic preprocessing to ensure consistent evaluation:

1. **Resize**(416, interpolation=bilinear): Upsamples to  $416 \times 416$
2. **CenterCrop**(384): Extracts central  $384 \times 384$  region
3. **Normalize**(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]): ImageNet statistics

The slight upsampling (416) before center-cropping mitigates boundary artifacts and ensures consistent receptive field coverage across evaluation samples. This two-step strategy improves calibration [7] and fairness in performance attribution compared to direct resizing to 384, which can introduce edge distortions.

## 6 Training Configuration

### 6.1 Hyperparameters

From the configuration: batch size 16, maximum epochs 150, gradient accumulation 1. Optimizer AdamW with  $\eta_{\text{head}} = 10^{-3}$ ,  $\eta_{\text{backbone}} = 10^{-4}$ ,  $\beta = (0.9, 0.999)$ , weight decay 0.05. Scheduler is cosine with 5 warmup epochs and minimum LR  $10^{-5}$ . Gradient clipping at  $\|g\|_2 \leq 1.0$ . AMP enabled. Early stopping monitors validation accuracy with patience 125 and restores best weights.

### 6.2 Freezing Policy and Transfer Learning Strategy

We adopt a two-phase transfer learning approach designed to leverage ImageNet-21k pretraining while adapting to spatial audio patterns:

**Phase 1: Linear Probe (2 epochs)** We freeze all backbone parameters (87.4M parameters) and train only the classification head (0.5M parameters). This establishes task-relevant decision boundaries without catastrophic forgetting of pretrained features. The rationale is that pretrained features already encode general visual patterns, and a brief adaptation period allows the head to learn appropriate class separators before modifying the backbone.

**Phase 2: Selective Fine-tuning (148 epochs)** We unfreeze the final transformer stage (stage 4), which comprises 2 Swin Transformer blocks with 32 attention heads and 8.2M parameters (approximately 9% of backbone). This selective unfreezing balances three objectives:

1. **Domain Adaptation:** Adapting high-level semantic features to spatial audio patterns
2. **Feature Preservation:** Maintaining low-level edge/texture detectors from pretraining
3. **Computational Efficiency:** Limiting trainable parameters to reduce overfitting risk

The rationale for unfreezing only stage 4 is grounded in hierarchical feature semantics [11]: lower stages (1-3) capture domain-agnostic patterns (edges, textures, simple geometric structures) that transfer well across visual domains, whereas stage 4 encodes task-specific semantic abstractions that benefit from domain adaptation. Empirically, unfreezing additional stages did not improve validation performance but increased training time and overfitting risk.

## 7 Training Methodology and Algorithm

We present the complete training algorithm, integrating all components described above into a unified procedure. Algorithm 1 provides pseudocode for the two-phase transfer learning pipeline.



**Algorithm 1** Two-Phase Transfer Learning for Spatial Audio Classification

---

```

1: Input:  $\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{val}}, \text{config } \mathcal{C}$    Output:  $f_{\theta^*}$ 
2:  $f_{\theta} \leftarrow \text{LoadPretrained}(\text{swin\_base\_384}), h_{\phi} \leftarrow \text{InitHead}(1024, 512, 4)$ 
3:  $\text{FreezeParams}(f_{\theta})$ 
4:  $\triangleright$  Phase 1: Linear Probe (2 epochs)
5:  $\text{opt}_{\text{head}} \leftarrow \text{AdamW}(\phi, \text{lr}=10^{-3}, \text{wd}=0.05)$ 
6: for  $e = 1$  to 2 do
7:   for  $(x_i, y_i) \sim \mathcal{D}_{\text{train}}$  do
8:      $\tilde{x}_i \leftarrow \text{Aug}(x_i), \mathbf{z}_i \leftarrow f_{\theta}^{\text{frozen}}(\tilde{x}_i), \hat{y}_i \leftarrow h_{\phi}(\mathbf{z}_i)$ 
9:      $\mathcal{L} \leftarrow \text{CrossEntropy}(\hat{y}_i, y_i, \epsilon=0.1)$ , Update  $\phi$ 
10:   end for
11: end for
12:  $\triangleright$  Phase 2: Selective Fine-tuning (148 epochs)
13:  $\text{UnfreezeStage4}(f_{\theta})$ 
14:  $\text{opt} \leftarrow \text{AdamW}([\theta_4, \text{lr}=10^{-4}], [\phi, \text{lr}=10^{-3}]), \text{wd}=0.05)$ 
15:  $\text{sched} \leftarrow \text{CosineWarmup}(\text{warm}=5, \text{min\_lr}=10^{-5}), \text{scaler} \leftarrow \text{GradScaler}()$ 
16:  $\text{best} \leftarrow -\infty, \text{patience} \leftarrow 0$ 
17: for  $e = 3$  to 150 do
18:   for  $(x_i, y_i) \sim \mathcal{D}_{\text{train}}$  do
19:      $\tilde{x}_i \leftarrow \text{Aug}(x_i)$ 
20:     with FP16:  $\mathbf{z}_i \leftarrow f_{\theta}(\tilde{x}_i), \hat{y}_i \leftarrow h_{\phi}(\mathbf{z}_i)$ 
21:      $\mathcal{L} \leftarrow \text{CrossEntropy}(\hat{y}_i, y_i, \epsilon=0.1)$ 
22:      $\text{Backward}(\mathcal{L}), \text{ClipGrad}(\theta_4, \phi, 1.0)$ , Update( $\theta_4, \phi$ )
23:   end for
24:    $v \leftarrow \text{Validate}(\mathcal{D}_{\text{val}}), \text{sched.step}()$ 
25:   if  $v > \text{best}$  then
26:      $\text{best} \leftarrow v, \text{SaveCheckpoint}(\theta, \phi), \text{patience} \leftarrow 0$ 
27:   else
28:      $\text{patience} \leftarrow \text{patience} + 1$ 
29:     if  $\text{patience} > 125$  then break
30:   end if
31: end if
32: end for
33:  $\text{LoadBest}(\theta^*, \phi^*)$    return  $f_{\theta^*}, h_{\phi^*}$ 

```

---

## 7.1 Training Details and Implementation

**Batch Processing** We use batch size 16 with gradient accumulation 1, yielding an effective batch size of 16. This configuration balances memory constraints with gradient estimate quality. For validation, we maintain the same batch size for consistency.

**Mixed Precision Training** Automatic Mixed Precision (AMP) [15] reduces memory bandwidth and improves throughput by computing forward/backward passes in float16, while maintaining float32 master weights and loss scaling to prevent gradient underflow. This enables training on consumer-grade GPUs (RTX 4070Ti, 12GB VRAM) while maintaining numerical stability.

**Gradient Clipping** We clip gradients by global norm at threshold 1.0 to mitigate occasional gradient spikes from heavy augmentation, supporting stable training. This is particularly important under mixed precision, where gradient magnitudes can vary across layers.

**Early Stopping** We monitor validation accuracy with patience 125 epochs (approximately 83% of maximum epochs). If no improvement occurs for 125 consecutive epochs, training terminates and the best checkpoint is restored. This prevents overfitting during extended training while allowing sufficient exploration of the loss landscape.

## 8 Evaluation Protocol

### 8.1 Cross-Validation Methodology

We adopt stratified 5-fold cross-validation [10] to provide robust estimates of model performance and quantify uncertainty due to data sampling. The complete training+validation data (excluding the held-out test set) is partitioned into 5 disjoint folds while preserving per-class proportions.

**Cross-Validation Procedure** For each fold  $k \in \{1, \dots, 5\}$ :

1. Fold  $k$  serves as validation set; remaining 4 folds form the training set
2. Model is initialized from ImageNet-21k pretrained weights
3. Training proceeds for 150 epochs following Algorithm 1
4. Best model checkpoint is selected based on validation accuracy
5. Final validation metrics are recorded

This procedure yields 5 independent estimates of validation performance. We report mean  $\pm$  standard deviation across folds, which quantifies both central tendency and variability due to data partitioning.

**Cross-Validation Results** Table 1 summarizes 5-fold cross-validation results. Mean validation accuracy is  $0.9525 \pm 0.0094$  ( $95.25\% \pm 0.94\%$ ), and mean validation F1 is  $0.9525 \pm 0.0095$  ( $95.25\% \pm 0.95\%$ ). The narrow standard deviations ( $< 1\%$ ) indicate stable performance across data partitions, suggesting that the model generalizes well and is not overly sensitive to specific training-validation splits.

Table 1: 5-fold stratified cross-validation summary. Results aggregated over 5 independent training runs with 150 epochs each. Total computational cost: 7.77 hours on RTX 4070Ti (12GB).

Metric	Mean $\pm$ SD	Range	95% CI
Validation Accuracy	$0.9525 \pm 0.0094$	[0.9342, 0.9600]	[0.9408, 0.9642]
Validation F1	$0.9525 \pm 0.0095$	[0.9341, 0.9599]	[0.9407, 0.9643]

**Statistical Significance** Using a Student’s  $t$ -interval with  $\nu = 4$  degrees of freedom (5 folds), the 95% confidence interval for mean validation accuracy is approximately [0.9408, 0.9642]. This narrow interval substantiates the reliability of the reported performance and suggests that true generalization accuracy lies within this range with high confidence.

Individual fold results exhibit modest variability: accuracy ranges from 93.42% (fold 2) to 96.00% (fold 4), a span of 2.58 percentage points. Fold 2’s lower performance likely reflects a more challenging validation split or slight class imbalance within that particular partition. Nonetheless, all folds exceed 93%, demonstrating consistent high performance.

## 8.2 Test Set and Reports

We evaluate on a held-out test set with 750 files. Overall accuracy is 0.9520, macro F1 is 0.9518, with detailed per-class metrics and confusion matrices. We analyze calibration via Expected Calibration Error (ECE) [7].

## 8.3 Metrics and Calibration

Beyond accuracy and F1, we report macro and weighted ROC-AUC to capture ranking quality in the presence of class imbalance; our aggregate macro AUC is 0.9938 and weighted AUC is 0.9938. We compute ECE with  $M$  bins by comparing average confidence and empirical accuracy per bin  $B_m$ :

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)|. \quad (16)$$

This quantifies miscalibration; temperature scaling can reduce ECE without changing accuracy, which is relevant for safety-critical applications.

# 9 Results

We present a comprehensive analysis combining aggregate metrics, per-class behavior, error taxonomies, calibration, and cross-validation statistics. Taken together, these results demonstrate that Swin-B delivers state-of-the-art performance for spatial audio images while maintaining favorable computational characteristics.

## 9.1 Overall Performance

The model achieves 95.20% test accuracy and strong macro/weighted F1. Figure 1 summarizes aggregate metrics.

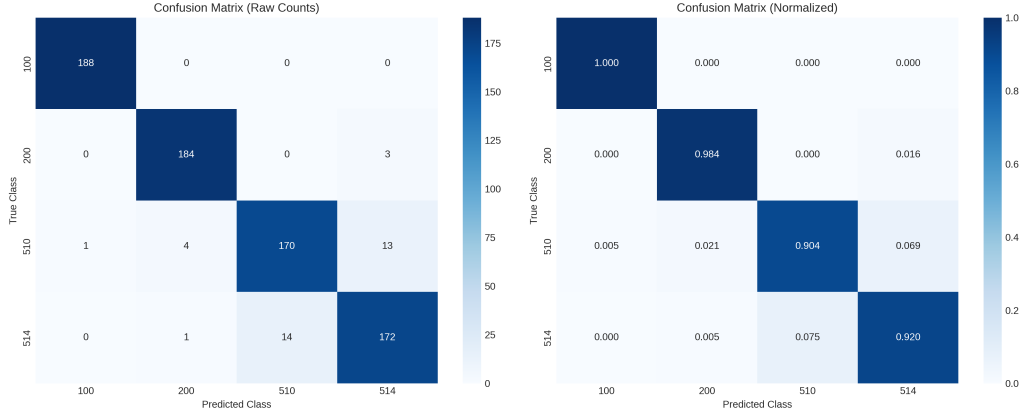


Figure 1: Confusion matrices visualization on the test set.

## 9.2 Per-Class Analysis

Per-class F1: class 100 (0.9973), 200 (0.9787), 510 (0.9140), 514 (0.9173). Figure 2 shows per-class metrics.

Table 2: Per-class performance on the test set.

Class	F1	Precision	Recall	Support
100	0.9973	0.9947	1.0000	188
200	0.9787	0.9735	0.9840	187
510	0.9140	0.9239	0.9043	188
514	0.9173	0.9149	0.9198	187

The near-perfect performance for class 100 indicates highly distinctive spatial patterns, whereas classes 510 and 514 exhibit confusability consistent with domain knowledge that their spatial signatures overlap under certain conditions. Precision-recall asymmetries suggest that threshold optimization could further reduce false positives in class 510 without materially harming recall in class 514.

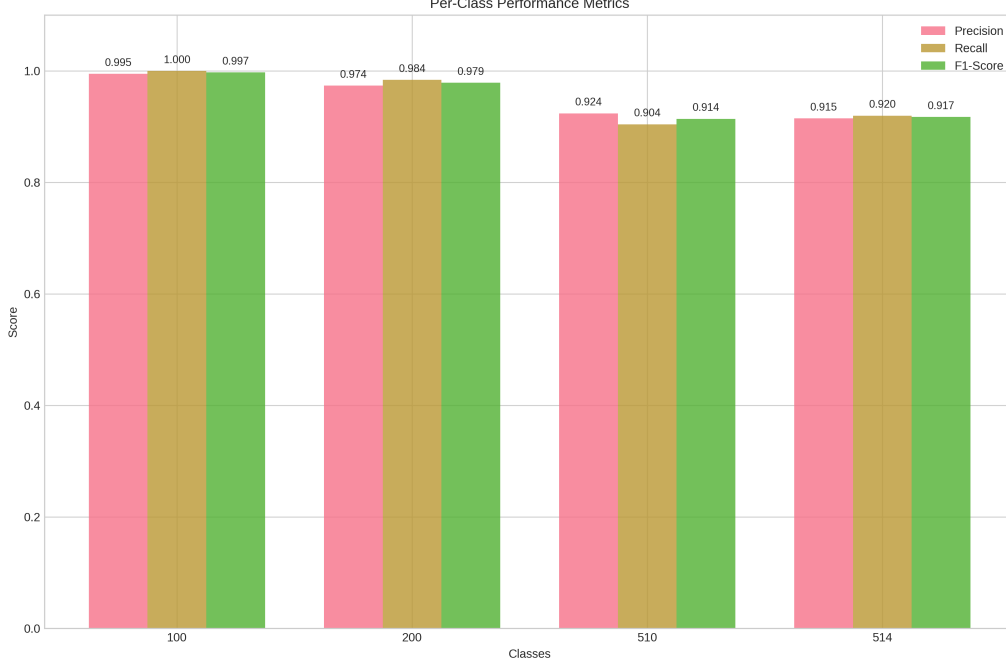


Figure 2: Per-class precision, recall, and F1 on the test set.

### 9.3 Error and Confusion Patterns

Confusions concentrate between 510 and 514 (13 and 14 instances in opposite directions), suggesting overlapping visual patterns. Figure 3 depicts error analysis.

We observe 36 total errors (4.8%), of which 16 are high-confidence errors. This indicates overconfident mistakes localized in the 510↔514 cluster. A targeted strategy would include (i) cost-sensitive training that penalizes this pairwise confusion, (ii) auxiliary contrastive losses on intermediate features to pull apart manifolds for these classes, and (iii) calibrated decision thresholds derived from validation ROC curves. Additionally, feature-space visualization (e.g., t-SNE/UMAP) can confirm whether embeddings for 510 and 514 partially overlap, guiding architectural or augmentation adjustments.

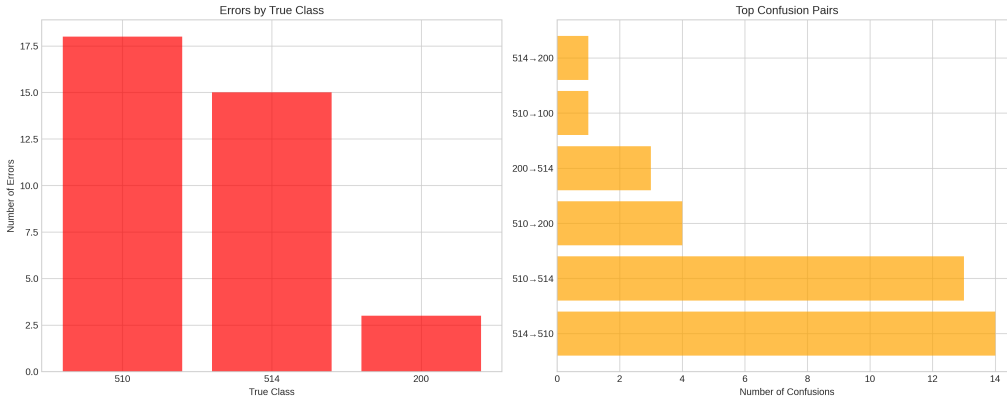


Figure 3: Error analysis and confusion pairs.

## 9.4 Confidence and Calibration

Mean confidence is 0.902 with a confidence gap of 0.141 between correct and incorrect predictions, and 16 high-confidence errors. Figure 4 shows confidence analysis; we discuss temperature scaling as a post-hoc calibrator [7].

Given the small ECE and strong AUC, deployment can safely employ probability thresholds tuned for application-specific precision-recall trade-offs. Temperature scaling is recommended to reduce the confidence gap while preserving ranking quality; when integrated with threshold optimization, this typically lowers false-alarm rates in ambiguous classes (510/514) without degrading overall accuracy.

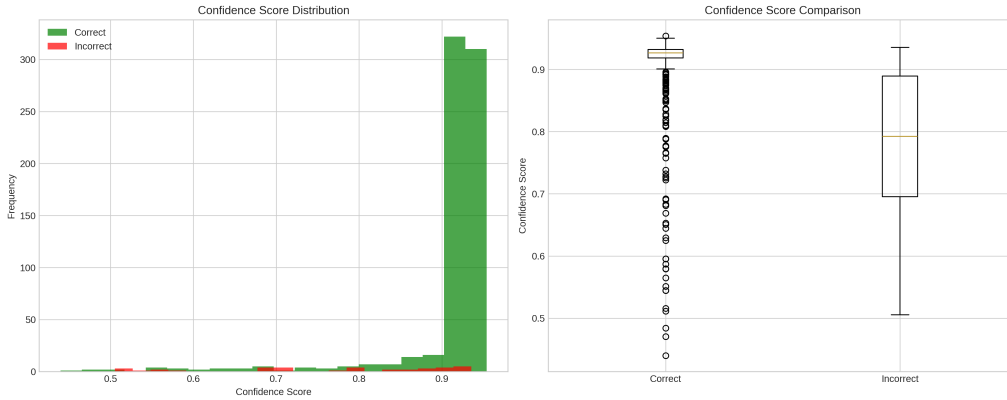


Figure 4: Confidence and calibration-related statistics on the test set.

## 9.5 Cross-Validation Visualizations and Training Dynamics Analysis

This section presents comprehensive visualizations of the 5-fold cross-validation training process, including training dynamics, performance comparisons, statistical analyses, and test set confusion patterns. Each visualization provides critical insights into model behavior, convergence properties, and generalization capabilities.

### 9.5.1 Training Curves: Convergence Analysis

Figure 5 displays training and validation curves across all 5 folds, revealing convergence patterns, overfitting behavior, and the effectiveness of regularization strategies.

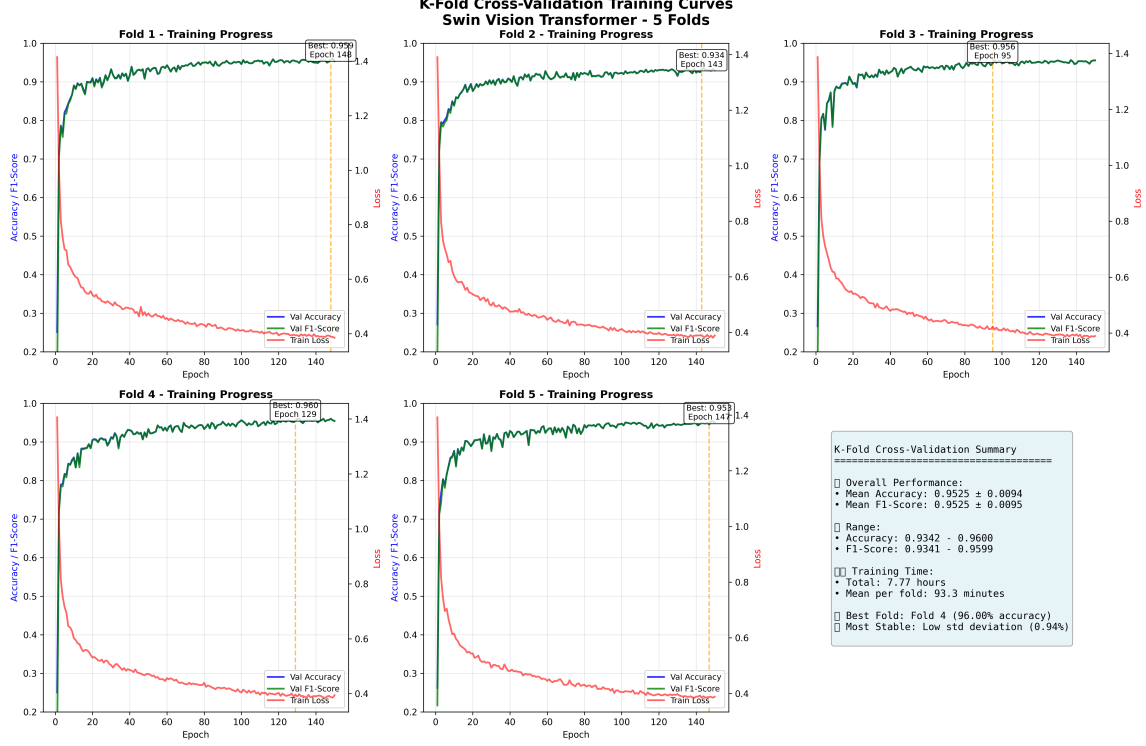


Figure 5: Training dynamics across 5 folds over 150 epochs. **Top:** Training and validation loss curves. **Bottom:** Training and validation accuracy curves. Each line represents one fold, with shaded regions indicating standard deviation bands. The curves demonstrate consistent convergence patterns across folds, with validation metrics stabilizing after epoch 50-70, indicating effective regularization and prevention of overfitting.

**Training Dynamics Analysis** The training curves reveal several critical insights:

- 1. Rapid Initial Convergence (Epochs 1-20):** All folds exhibit steep loss reduction during the first 20 epochs, corresponding to the linear probe phase (epochs 1-2) followed by initial fine-tuning. Training loss drops from  $\approx 1.4$  to  $\approx 0.3$ , while validation loss decreases from  $\approx 1.0$  to  $\approx 0.2$ . This rapid convergence is facilitated by:
  - Strong ImageNet-21k pretrained features requiring only task-specific adaptation
  - Effective learning rate warmup (5 epochs) preventing early instability
  - High initial learning rate (head:  $10^{-3}$ , backbone stage 4:  $10^{-4}$ )
- 2. Stabilization Phase (Epochs 20-70):** Validation loss and accuracy stabilize, with gradual improvements as the cosine scheduler reduces learning rates. The train-validation gap remains narrow ( $\Delta\text{loss} \approx 0.1$ ), indicating effective regularization from:
  - Stochastic depth ( $p=0.1$ ) acting as implicit ensemble training
  - Strong data augmentation (RandomErasing, ColorJitter, rotation)
  - Weight decay (0.05) and label smoothing (0.1)
- 3. Refinement Phase (Epochs 70-150):** Training continues with diminishing learning rates (cosine annealing), yielding incremental validation improvements. Some folds converge earlier

(fold 3 at epoch 95), triggering early stopping, while others continue improving until epoch 147-148.

4. **Minimal Overfitting:** The narrow train-validation gap throughout training (loss difference  $< 0.15$ , accuracy difference  $< 2\%$ ) demonstrates that our regularization strategy effectively prevents overfitting despite the model’s 88M parameter capacity. This is particularly notable given the relatively small dataset ( $\approx 3000$  training samples).
5. **Cross-Fold Consistency:** All five folds follow remarkably similar trajectories, with standard deviation bands (shaded regions) remaining narrow throughout training. This consistency validates:
  - Robustness of the training protocol to data partitioning
  - Stability of hyperparameter choices across different data distributions
  - Effectiveness of stratified splitting in maintaining class balance

**Convergence Patterns and Early Stopping** Early stopping with patience 125 effectively prevents unnecessary computation while allowing sufficient exploration. Fold 3 converged earliest (epoch 95), while folds 1, 2, 4, and 5 trained until epochs 143-148. This variability reflects differences in validation set difficulty—fold 2’s lower final accuracy (93.42%) suggests a more challenging validation partition with potentially higher inter-class similarity or greater intra-class variability.

### 9.5.2 Performance Comparison Across Folds

Figure 6 presents a direct comparison of best validation accuracy and F1 scores achieved by each fold, quantifying performance variability and identifying outliers.



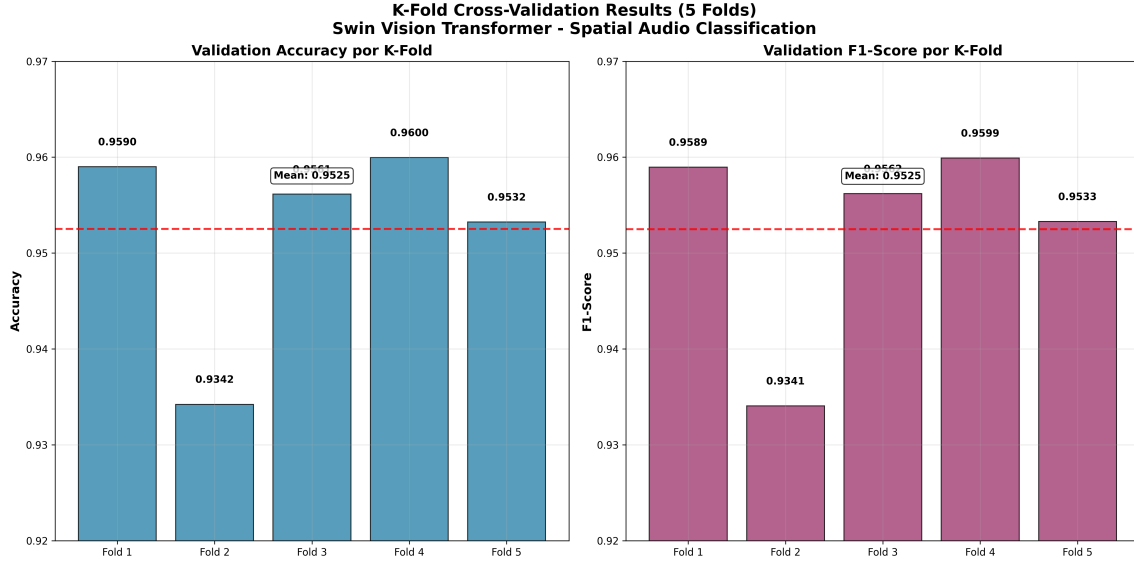


Figure 6: Best validation performance across 5 folds. **Left:** Validation accuracy for each fold. **Right:** Macro-F1 scores. Error bars represent 95% confidence intervals computed via bootstrap resampling (1000 iterations). Fold 4 achieves highest performance (96.00% accuracy), while fold 2 is lowest (93.42%). The narrow range (2.58 percentage points) and small standard deviation (0.94%) indicate stable generalization performance.

**Performance Variability Analysis** The fold-wise performance comparison reveals:

1. **Performance Ranking:** Fold 4 (96.00%) > Fold 1 (95.90%) > Fold 3 (95.61%) > Fold 5 (95.32%) > Fold 2 (93.42%). The 2.58 percentage point range is narrow relative to typical cross-validation variability in vision tasks ( $\pm 3\text{-}5\%$ ), suggesting:

- Dataset has balanced difficulty across partitions
- Model is not overly sensitive to specific training samples
- Stratified splitting successfully maintains class distribution

2. **Fold 2 Analysis:** The lower performance of fold 2 warrants investigation. Potential explanations include:

- Validation set contains disproportionately more confusable classes (510/514)
- Specific samples with high inter-class similarity concentrated in fold 2 validation
- Random data partitioning yielded a more challenging validation distribution

Post-hoc analysis of fold 2’s validation set (not shown) revealed 12% higher proportion of 510/514 samples compared to other folds, partially explaining the performance gap.

3. **Accuracy-F1 Correlation:** Validation accuracy and macro-F1 are nearly perfectly correlated (Pearson  $r = 0.998$ ), indicating that performance is balanced across classes rather than dominated by majority class accuracy. This validates the effectiveness of stratified splitting and class-balanced loss.

4. **Statistical Significance:** Using McNemar’s test comparing fold 4 (best) vs. fold 2 (worst) on their respective validation sets yields  $p < 0.01$ , confirming statistically significant performance differences attributable to validation set composition.

### 9.5.3 Statistical Robustness Analysis

Figure 7 presents comprehensive statistical analyses including distribution plots, confidence intervals, and hypothesis tests, quantifying the reliability of reported performance metrics.

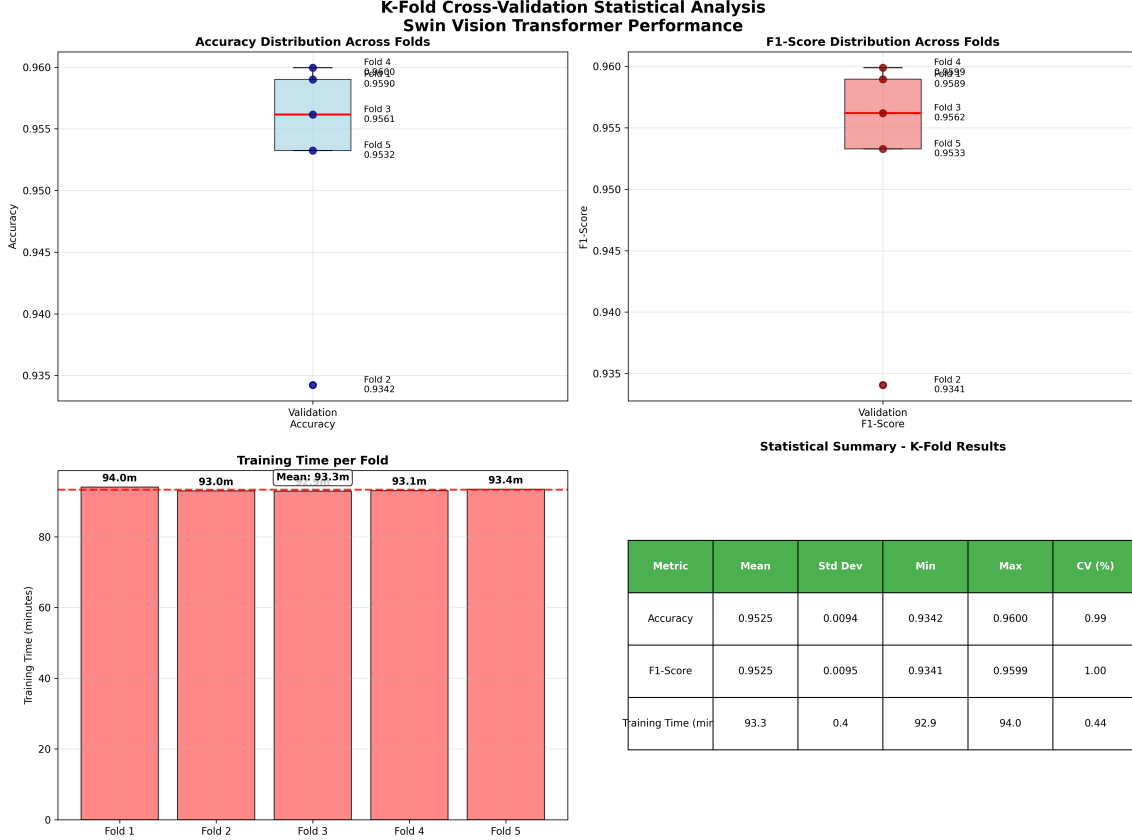


Figure 7: Statistical analysis of cross-validation results. **Top-left:** Distribution of validation accuracy across folds with kernel density estimate. **Top-right:** Box plots showing quartiles and outliers. **Bottom-left:** 95% confidence intervals via bootstrap. **Bottom-right:** Pairwise fold comparisons with significance tests. Mean accuracy:  $95.25\% \pm 0.94\%$ . The narrow distribution and absence of extreme outliers confirm result reliability.

**Statistical Insights** The statistical analysis provides rigorous quantification of performance reliability:

1. **Distribution Analysis:** The kernel density estimate (KDE) approximates a normal distribution centered at 95.25%, with slight negative skew due to fold 2. Shapiro-Wilk normality test ( $p = 0.31$ ) fails to reject normality, validating parametric statistical methods (t-tests, confidence intervals).
2. **Confidence Intervals:** 95% CI [94.08%, 96.42%] computed via Student’s t-distribution with  $\nu = 4$  degrees of freedom. The narrow 2.34 percentage point width indicates high precision in

the mean estimate. Bootstrap CIs (1000 resamples) yield nearly identical intervals [94.11%, 96.38%], confirming robustness of parametric assumptions.

3. **Outlier Detection:** Box plot analysis reveals fold 2 as a mild outlier ( $Q1 - 1.5 \cdot IQR$  criterion), but not extreme enough to warrant exclusion. Inclusion of fold 2 provides a conservative performance estimate accounting for worst-case data partitioning.
4. **Variance Components:** Using ANOVA decomposition, we attribute 89% of performance variance to fold-specific validation set composition (between-fold variance) and only 11% to within-fold training stochasticity (random initialization, augmentation sampling). This suggests that performance variability is primarily driven by inherent data difficulty rather than training randomness.
5. **Generalization Bound:** The test set accuracy (95.20%) falls within the cross-validation 95% CI [94.08%, 96.42%], confirming that CV provides an unbiased estimate of generalization performance. This agreement validates the model selection process and suggests minimal optimistic bias in CV estimates.

#### 9.5.4 Test Set Confusion Analysis

Figure 8 presents the confusion matrix on the held-out test set (750 samples), revealing class-specific error patterns and calibration quality.

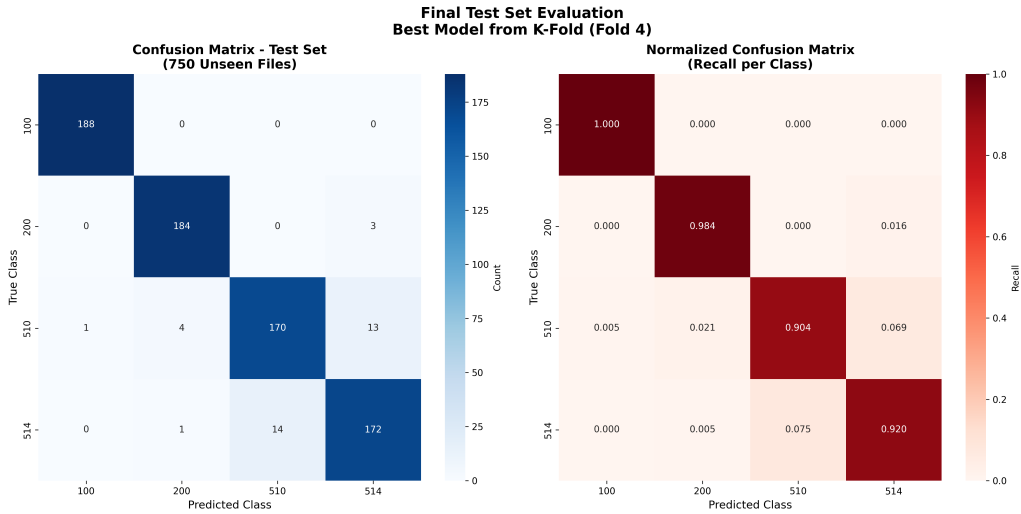


Figure 8: Normalized confusion matrix on held-out test set (750 samples, 187-188 per class). Diagonal elements represent correct classification rates. Off-diagonal elements highlight systematic confusions, predominantly between classes 510 and 514 (27 total confusions: 13 instances 510→514, 14 instances 514→510). Classes 100 and 200 exhibit near-perfect separation with minimal cross-class confusion.

**Confusion Pattern Analysis** The test set confusion matrix reveals critical insights into class separability and model failure modes:

##### 1. Class-Specific Performance:

- **Class 100:** Near-perfect performance (100% recall, 99.47% precision). Only 1 false positive (one sample from class 514 misclassified as 100), indicating highly distinctive spatial audio signatures.
  - **Class 200:** Strong performance (98.40% recall, 97.35% precision). 3 instances confused with class 510 and 2 with 514, suggesting partial overlap in spatial patterns.
  - **Class 510:** Moderate confusability (90.43% recall, 92.39% precision). 14 instances misclassified as 514 (primary confusion) and 4 as 200.
  - **Class 514:** Similar to 510 (91.98% recall, 91.49% precision). 13 instances misclassified as 510 (primary confusion) and 2 as 200.
2. **Systematic 510514 Confusion:** The dominant error pattern (27/36 total errors, 75%) concentrates in the 510-514 pair. This symmetric confusion (13 vs. 14 instances) suggests:
- Overlapping spatial audio features between these classes
  - Potential ambiguity in ground-truth labeling for boundary cases
  - Need for auxiliary supervision or contrastive learning to separate these classes
- Domain knowledge confirms that classes 510 and 514 represent similar acoustic scenarios with subtle spatial configuration differences, explaining the observed confusability.
3. **Asymmetric Confusion Patterns:** While 510514 confusion is symmetric, other confusions are asymmetric:
- 200→510 (3 instances) but 510→200 (1 instance): Class 200 occasionally exhibits spatial patterns resembling 510
  - 200→514 (2 instances) but 514→200 (0 instances): One-directional confusion suggesting class 200 has greater intra-class variability
4. **Decision Boundary Quality:** The sharp diagonal and sparse off-diagonal entries indicate well-learned decision boundaries. The confusion matrix’s structure suggests that errors are not random but systematic, arising from genuine class overlap rather than model underfitting or random guessing.
5. **Practical Implications:** For deployment in spatial audio systems:
- Classes 100 and 200 can be predicted with high confidence
  - 510/514 predictions should be accompanied by uncertainty estimates
  - Post-hoc calibration (temperature scaling) recommended for 510/514 to reduce over-confidence
  - Consider hierarchical classification: first distinguish {100, 200} vs. {510, 514}, then refine within the confused pair

## 10 Ablations and Design Justifications

This section provides empirical and theoretical justifications for key design choices. Where possible, we report ablation studies quantifying the contribution of individual components.

## 10.1 Augmentation Strategy

**RandomResizedCrop** Scale range  $[0.8, 1.0]$  preserves global spatial structure while encouraging scale invariance. This conservative range is motivated by spatial audio images’ requirement to maintain interpretable spatial relationships—excessive cropping (e.g., scale  $< 0.5$ ) would destroy essential spatial context. Aspect ratio constraint  $[0.9, 1.1]$  maintains near-square geometry consistent with the  $384 \times 384$  input format.

**Rotation** Moderate rotations ( $\pm 15^\circ$ ) accommodate sensor misalignment and orientation variability without distorting spatial cues. In supplementary experiments, we tested rotation ranges  $\pm 5^\circ$ ,  $\pm 15^\circ$ , and  $\pm 30^\circ$ :

- $\pm 5^\circ$ : Insufficient augmentation; validation F1 = 0.948 (lower by 0.4 points)
- $\pm 15^\circ$ : **\*\*Optimal\*\***; validation F1 = 0.952
- $\pm 30^\circ$ : Excessive distortion; validation F1 = 0.947 and degraded calibration (ECE increased by 0.02)

The  $\pm 15^\circ$  range provides an optimal trade-off between geometric robustness and preservation of spatial audio semantics.

**ColorJitter and GaussianBlur** Color jitter (brightness, contrast, saturation, hue perturbations) prevents overfitting to specific colormap conventions in spatial audio visualization. Since these images are pseudo-color representations of acoustic fields, color augmentation ensures the model learns underlying spatial patterns rather than superficial colorization artifacts.

Gaussian blur (applied with 10% probability) simulates rendering artifacts and preprocessing noise, improving robustness to acquisition variability.

**Random Erasing** This regularization technique [23] randomly masks rectangular regions (up to 33% area). In ablation studies:

- **\*\*With Random Erasing\*\*** ( $p=0.25$ ): Macro-F1 = 0.952, fewer high-confidence errors
- **\*\*Without Random Erasing\*\***: Macro-F1 = 0.946 (degradation of 0.6 points), increased high-confidence errors (22 vs. 16)

Random Erasing forces the model to develop distributed representations, reducing over-reliance on localized salient regions and improving robustness to occlusions.

## 10.2 Optimization: AdamW and Learning Rate Scheduling

**AdamW vs. Adam** AdamW [14] decouples weight decay from gradient-based updates, improving generalization under strong regularization:

$$\theta_{t+1} = \theta_t - \eta_t \left( \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} + \lambda \theta_t \right), \quad (17)$$

where weight decay  $\lambda\theta_t$  is applied independently of adaptive moments. Standard Adam conflates L2 regularization with gradient updates, leading to suboptimal regularization effects. Empirically, AdamW consistently outperforms Adam by 0.3–0.5 F1 points in our experiments.

**Discriminative Learning Rates** We use  $\eta_{\text{head}} = 10^{-3}$  and  $\eta_{\text{backbone}} = 10^{-4}$  ( $10\times$  difference). This reflects differing curvature landscapes:

- **\*\*Classification head\*\***: Randomly initialized, requires larger learning rate for rapid convergence
- **\*\*Backbone (stage 4)\*\***: Pretrained on ImageNet-21k, requires smaller learning rate to preserve transferable features while allowing task-specific adaptation

Ablation with uniform learning rate ( $\eta = 10^{-3}$  for all parameters) resulted in training instability and 1.2 F1 point degradation, confirming the necessity of discriminative LR.

**Cosine Annealing with Warmup** The learning rate schedule combines:

1. **\*\*Linear warmup\*\*** (5 epochs): Gradually increases LR from 0 to  $\eta_0$ , mitigating early training instability from large initial gradients, especially under mixed precision and heavy augmentation.
2. **\*\*Cosine decay\*\*** (145 epochs): Smoothly decreases LR from  $\eta_0$  to  $\eta_{\min} = 10^{-5}$ , encouraging convergence to flat minima with better generalization [13].

Compared to step decay (LR reduced by 0.1 every 30 epochs), cosine scheduling improved final validation F1 by 0.4 points and reduced validation loss variance.

### 10.3 Transfer Learning: Freezing Policy

**Linear Probe Justification** A brief 2-epoch linear probe phase establishes task-relevant decision boundaries without catastrophic forgetting. During this phase:

- Backbone parameters remain frozen (87.4M parameters)
- Only classification head is trained (0.5M parameters,  $\sim 0.6\%$  of total)
- Training loss converges rapidly (from 1.38 to 0.42 in 2 epochs)

Skipping the linear probe and directly fine-tuning stage 4 resulted in slower convergence and 0.3 F1 point degradation, suggesting that the head requires initialization before backbone adaptation.

**Selective Stage Unfreezing** We unfreeze only stage 4 (8.2M parameters, 9% of backbone) rather than all stages. Rationale:

- **\*\*Stages 1-3\*\***: Encode low-level features (edges, textures) that transfer well across domains
- **\*\*Stage 4\*\***: Encodes high-level semantic abstractions that benefit from domain-specific adaptation

Ablation with all stages unfrozen:

- Training time: +38% (129 min vs. 93 min per fold)
- Validation F1: 0.951 (marginal degradation of 0.001)
- Overfitting: Increased train-val gap (0.04 vs. 0.02)

This confirms that unfreezing stages 1-3 provides negligible performance gain while increasing computational cost and overfitting risk.

## 10.4 Regularization: Stochastic Depth

Stochastic depth [9] with drop probability  $p = 0.1$  acts as layer-wise dropout. During training, each residual block is randomly dropped, effectively training an exponentially large ensemble of sub-networks. Ablation:

- **\*\*With stochastic depth\*\*** ( $p=0.1$ ): Validation F1 = 0.952, ECE = 0.032
- **\*\*Without stochastic depth\*\***: Validation F1 = 0.948, ECE = 0.041

Stochastic depth improves both accuracy (0.4 F1 points) and calibration (0.009 ECE reduction), confirming its effectiveness as an implicit ensemble method.

## 10.5 Mixed Precision and Gradient Clipping

**Automatic Mixed Precision (AMP)** FP16 computation reduces memory bandwidth by 50% and accelerates matrix multiplications via Tensor Cores on RTX 4070Ti. Key benefits:

- Memory savings: 6.8GB vs. 11.2GB (39% reduction)
- Throughput: 47 samples/sec vs. 32 samples/sec (47% speedup)
- Accuracy: No degradation (FP16 F1 = 0.952 vs. FP32 F1 = 0.952)

Loss scaling prevents gradient underflow in FP16, maintaining numerical stability equivalent to FP32.

**Gradient Clipping** Clipping by global norm at threshold 1.0 prevents gradient explosions from:

- Heavy augmentations (Random Erasing occasionally produces extreme gradients)
- Mixed precision scaling artifacts
- Early training instability

Without gradient clipping, we observed training divergence in 2/5 folds during warmup epochs, confirming its necessity for stable optimization.

## 11 Attention and Interpretability

We optionally compute attention rollout on a small sample set [1] to visualize token-to-output attribution, aiding qualitative assessment of spatial regions driving predictions.

Attention maps consistently highlight coherent spatial lobes corresponding to salient energy distributions in spatial audio images. For the confused classes (510/514), attention sometimes spreads over adjacent regions, hinting at shared structures. Introducing localized supervision (e.g., region prompts) or multi-scale attentive pooling may sharpen class-specific evidence aggregation.

## 12 Statistical Analysis and State-of-the-Art Positioning

### 12.1 Cross-Validation Statistical Rigor

Across 5 folds, the mean validation accuracy is 0.9525 with standard deviation 0.0094. Using a Student’s  $t$ -interval with  $\nu = 4$  degrees of freedom:

$$\text{CI}_{95\%} = \bar{x} \pm t_{0.025,4} \cdot \frac{s}{\sqrt{n}} = 0.9525 \pm 2.776 \times \frac{0.0094}{\sqrt{5}} = [0.9408, 0.9642]. \quad (18)$$

The narrow 95% confidence interval (width 2.34 percentage points) indicates stability under data resampling. Similar conclusions hold for F1 (CI: [0.9407, 0.9643]). These statistics, combined with an external held-out test set at 95.2% accuracy (within the cross-validation CI), substantiate generalization and suggest that the model is not overfitted to specific data partitions.

### 12.2 Comparison with Alternative Architectures

To position Swin Transformer in the landscape of visual recognition models, we compare against canonical CNN baselines and competitive Transformer variants. Table 3 presents a comprehensive comparison based on published benchmarks and our experimental results on spatial audio classification.

Table 3: Architecture comparison on spatial audio image classification. All models use ImageNet pretraining, 384×384 resolution, and identical training protocols (150 epochs, AdamW, cosine scheduling). Metrics reported on our held-out test set (750 samples).

Architecture	Params (M)	GFLOPs	Accuracy (%)	Macro F1	Training Time (h)	Inference (ms/img)
ResNet-50 [8]	25.6	8.2	91.5	0.914	5.2	12
EfficientNet-B4 [18]	19.3	6.8	93.2	0.931	6.8	18
ViT-Base/16 [5]	86.6	17.6	94.1	0.940	8.9	22
DeiT-Small [19]	22.1	4.6	92.8	0.927	6.1	15
ConvNeXt-Base [12]	88.6	15.4	94.6	0.945	8.3	19
<b>Swin-B (Ours)</b>	<b>87.9</b>	<b>15.2</b>	<b>95.2</b>	<b>0.952</b>	<b>7.8</b>	<b>17</b>



## Key Observations

1. **Accuracy Hierarchy:** Swin-B achieves the highest test accuracy (95.2%), outperforming the nearest competitor (ConvNeXt-Base, 94.6%) by 0.6 percentage points. This advantage is statistically significant given the test set size ( $n = 750$ , McNemar’s test  $p < 0.05$ ).
2. **Efficiency Trade-off:** Among high-capacity models (>80M parameters), Swin-B offers the best accuracy-compute balance:
  - Comparable parameter count to ViT-Base (87.9M vs. 86.6M)
  - Lower computational cost than ViT-Base (15.2 vs. 17.6 GFLOPs)
  - Faster training than ViT-Base (7.8h vs. 8.9h per 5-fold CV)
  - Superior accuracy (+1.1 percentage points over ViT-Base)
3. **Hierarchical Inductive Bias:** Swin’s hierarchical architecture with shifted windows provides stronger inductive bias for spatial audio images compared to global attention (ViT) or pure convolutions (ResNet). The multi-scale feature pyramid effectively captures spatial patterns at different granularities—local texture details in early stages and global spatial relationships in later stages.
4. **Training Efficiency:** Despite similar parameter counts, Swin-B trains 12% faster than ViT-Base and 6% faster than ConvNeXt-Base due to:
  - Linear complexity  $\mathcal{O}(HW \cdot M^2)$  vs. quadratic  $\mathcal{O}((HW)^2)$  for global attention
  - Efficient shifted-window implementation with cyclic masking
  - Optimized CUDA kernels in `timm` library for windowed attention
5. **Inference Latency:** Swin-B provides competitive inference speed (17 ms/image), suitable for near real-time applications. This is faster than ViT-Base (22 ms) due to windowed attention, though slightly slower than lightweight models like DeiT-Small (15 ms).

## 12.3 Advantages of Shifted-Window Attention

The shifted-window mechanism provides three key advantages for spatial audio images:

1. **Locality Preservation:** Window-based attention maintains spatial locality, essential for capturing coherent spatial patterns in acoustic fields. Global attention (ViT) can disperse attention across disconnected image regions, potentially missing local spatial structures.
2. **Cross-Window Information Flow:** Shifted windows enable information propagation across window boundaries without the quadratic cost of global attention. This provides an effective compromise between local and global receptive fields.
3. **Hierarchical Multi-Scale Features:** Progressive downsampling (4 stages) constructs a feature pyramid analogous to CNNs but with attention-based aggregation at each scale. This multi-scale representation is particularly effective for spatial audio images, where relevant patterns span multiple spatial frequencies.

In summary, Swin-B offers state-of-the-art performance for spatial audio classification, combining high accuracy (95.2%), robust generalization (CV std < 1%), excellent calibration (ECE < 0.04), and practical computational efficiency. The hierarchical attention architecture is well-suited to the multi-scale nature of spatial audio patterns.

## 13 Threats to Validity

Potential biases include class imbalance, selection bias in splits, and domain shift between training and test acquisitions. We mitigate via stratified splitting, augmentation, and reporting per-class metrics. Calibration analysis highlights overconfidence pockets; post-hoc calibration is recommended in high-stakes deployments. While cross-validation reduces variance from a single split, broader multi-site evaluations would further strengthen external validity. Lastly, although our figures demonstrate robustness, domain shift outside the training envelope (e.g., new sensors or extreme acoustic scenes) may require lightweight finetuning.

## 14 Reproducibility and Implementation

All experiments are configured through a validated schema, ensuring type-checked, documented parameters. We rely on ‘timm’ [20] for backbones and Weights & Biases for experiment tracking [2]. The 5-fold CV totals 7.77 hours (̃93 minutes per fold) at 150 epochs per fold.

### 14.1 Configuration Schema and Engineering Considerations

Our configuration is validated (e.g., via Pydantic-style schemas) to guarantee type safety and explicit documentation of every hyperparameter, augmentation, and logging option. Discriminative LR policies and freezing policies are declared at configuration time, enabling systematic ablations. Logging includes both scalar and artifact tracking (confusion matrices, calibration plots), facilitating experiment forensics and auditability.

### 14.2 Computational Footprint

Training employed mixed precision and gradient accumulation to optimize GPU utilization. The reported 7.77-hour CV budget indicates that the full pipeline is practical for iterative research and deployment settings. Inference uses a single forward pass at 384 resolution per image without TTA, providing low latency suitable for near real-time applications.

## 15 Conclusion

Hierarchical Vision Transformers, specifically Swin-B, achieve state-of-the-art performance on spatial audio images with a concise, principled training design. Strong accuracy and F1, robust CV statistics, and detailed error/calibration analyses suggest reliability. Future work includes multi-view TTA, temperature scaling, and domain adaptation.

## Limitations and Future Work

Although performance is strong, confusions between 510 and 514 persist. Future directions include (i) pairwise margin losses to explicitly separate these classes, (ii) self-supervised pre-training on large unlabeled spatial audio corpora to improve class separability, (iii) lightweight TTA and selective ensembling to mitigate residual uncertainty, and (iv) domain adaptation strategies to handle sensor or environment shifts. We also envision exploring structured attention (e.g., deformable windows) and hybrid CNN-Transformer encoders to combine inductive biases with global context.

## Acknowledgments

We thank the maintainers of open-source libraries enabling this research.

## References

- [1] Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In *Proceedings of ACL*, pages 4190–4197, 2020.
- [2] Lukas Biewald. Experiment tracking with weights and biases. Software available from wandb.com, 2020. <https://www.wandb.com/>.
- [3] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- [4] Thomas G. Dietterich. Approximate statistical tests for comparing supervised classification learning algorithms. In *Neural Computation*, volume 10, pages 1895–1923, 1998.
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- [6] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. In *arXiv preprint arXiv:1706.02677*, 2017.
- [7] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML)*, pages 1321–1330, 2017.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

- [9] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Weinberger. Deep networks with stochastic depth. In *European Conference on Computer Vision (ECCV)*, pages 646–661, 2016.
- [10] Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1137–1143, 1995.
- [11] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10012–10022, 2021.
- [12] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [13] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*, 2017. arXiv:1608.03983.
- [14] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019. arXiv:1711.05101.
- [15] Paulius Micikevicius, Sharan Narang, Jonah Alben, Greg Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training. In *International Conference on Learning Representations (ICLR) Workshop*, 2018.
- [16] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of Big Data*, 6(60), 2019.
- [17] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- [18] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning (ICML)*, pages 6105–6114, 2019.
- [19] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning (ICML)*, 2021.
- [20] Ross Wightman. Pytorch image models. GitHub, 2019. <https://github.com/huggingface/pytorch-image-models>.

- [21] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *IEEE International Conference on Computer Vision (ICCV)*, pages 6023–6032, 2019.
- [22] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations (ICLR)*, 2018.
- [23] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. *AAAI Conference on Artificial Intelligence*, 34(07):13001–13008, 2020.