# HW4 experiment Report

Name:李崇楷  Student ID: 313834006

Code: https://github.com/kailee0422/RNN-Transformer/tree/main/HW4

## Introduction

In recent years, transformer-based architectures have shown remarkable success in various computer vision tasks. This assignment aims to compare two such architectures—Vision Transformer (ViT) and Swin Transformer (SWIN)—on the CIFAR-10 image classification dataset. By fine-tuning pre-trained models and visualizing their decision-making using Grad-CAM, we seek to understand not only their performance in terms of accuracy but also how each model interprets image features. Through this comparative analysis, we aim to highlight the strengths and limitations of ViT and SWIN for small-scale image classification tasks.

## Method

### Data Preprocessing

- Dataset: **CIFAR-10** (using PyTorch's built-in dataset).
- Preprocessing: Normalize images using ImageNet statistics (mean and standard deviation) and resize images from 32×32 to 224×224 to match input size required by pre-trained models.
- Data Splitting: Use the standard CIFAR-10 train/test split and further split the training set into **80% training** and **20% validation**.

### Data Preprocessing

- Models: Models are loaded from the **timm** library

  Vision Transformer: vit_base_patch32_224

  Swin Transformer: swin_base_patch4_window7_224
- Modification: Replace the final classification head to output 10 classes instead of the original 1000.
- Fine-tuning Options:

  Option 1: Fine-tune the entire model.

  Option 2: Freeze early layers and fine-tune only the classification head and later layers to reduce computational cost.

## Training

- Optimizer: AdamW
- Learning Rate Scheduler: Cosine Annealing
- Epoch: 20
- Monitoring: Track loss and accuracy on both training and validation sets and save checkpoints based on best validation accuracy.

## Evaluation

- Evaluate both models on the CIFAR-10 test set.

- Metrics:

  **Classification accuracy**

  **Top-1 error rate**

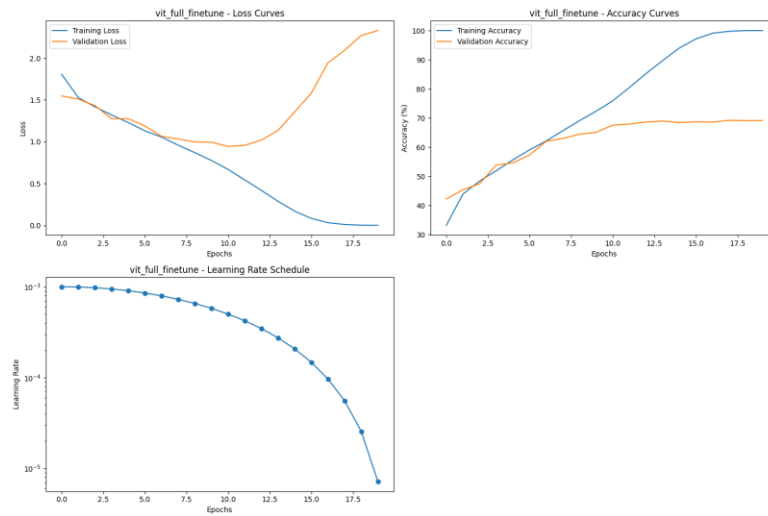  **Confusion matrix**

## Grad-CAM Visualization

- Input: One randomly selected image from the test set.
- Visualization for ViT: Initially targeted class token / last attention layer. Due to poor results, switched to **last transformer block's norm1 layer**.
- Visualization for Swin: Initially used feature maps from model.layers[3]. Also switched to **last transformer block's norm1 layer** due to better results.
- Output: Grad-CAM heatmaps showing model attention on the selected image.
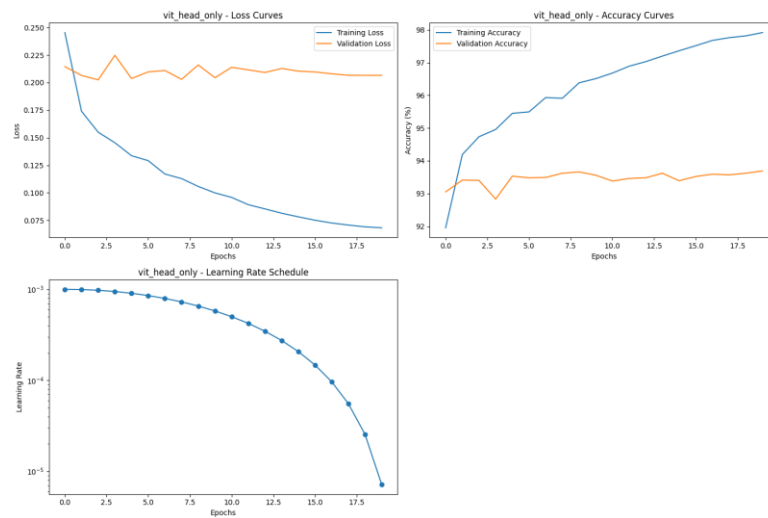
## Analysis and Comparison

- Performance Comparison:

  Compare accuracy and computational efficiency between ViT and Swin.
- Attention Analysis:

  Use Grad-CAM to understand differences in attention patterns.

  Discuss how each model interprets image features.
- Reflection:

  Evaluate strengths and limitations of Vision Transformer and Swin for the CIFAR-10 classification task.
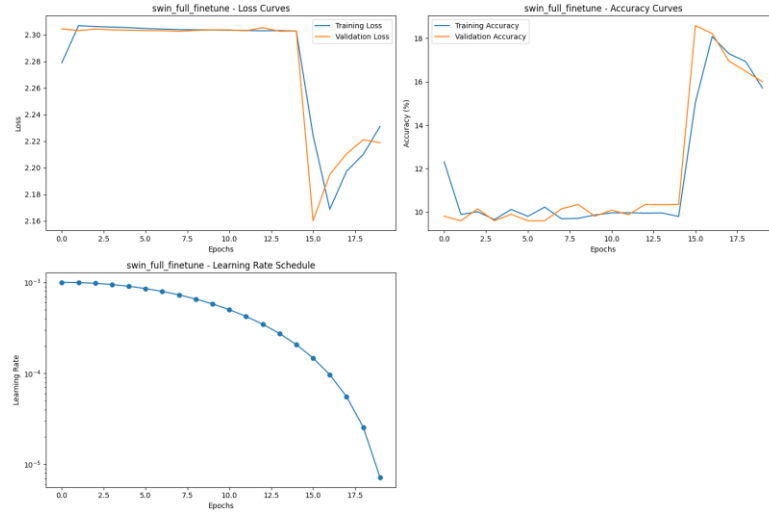
# Result

I conducted four experiments by training four different model configurations: fully fine-tuned Vision Transformer and Swin Transformer models, as well as Vision Transformer and Swin Transformer models where only the classification head and later layers were fine-tuned.
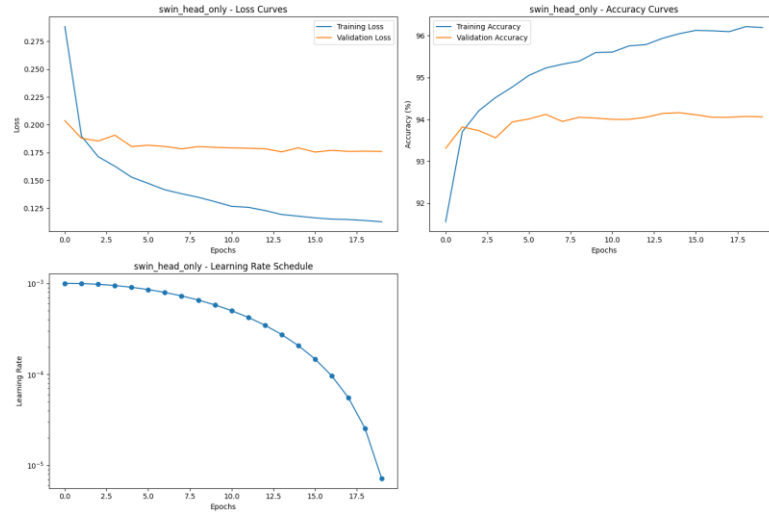


*Figure 1.* Training loss, accuracy and learning rate schedule over epochs for using fully fine-tuned Vision Transformer



*Figure 2.* Training loss, accuracy and learning rate schedule over epochs for using only head fine-tuned Vision Transformer.
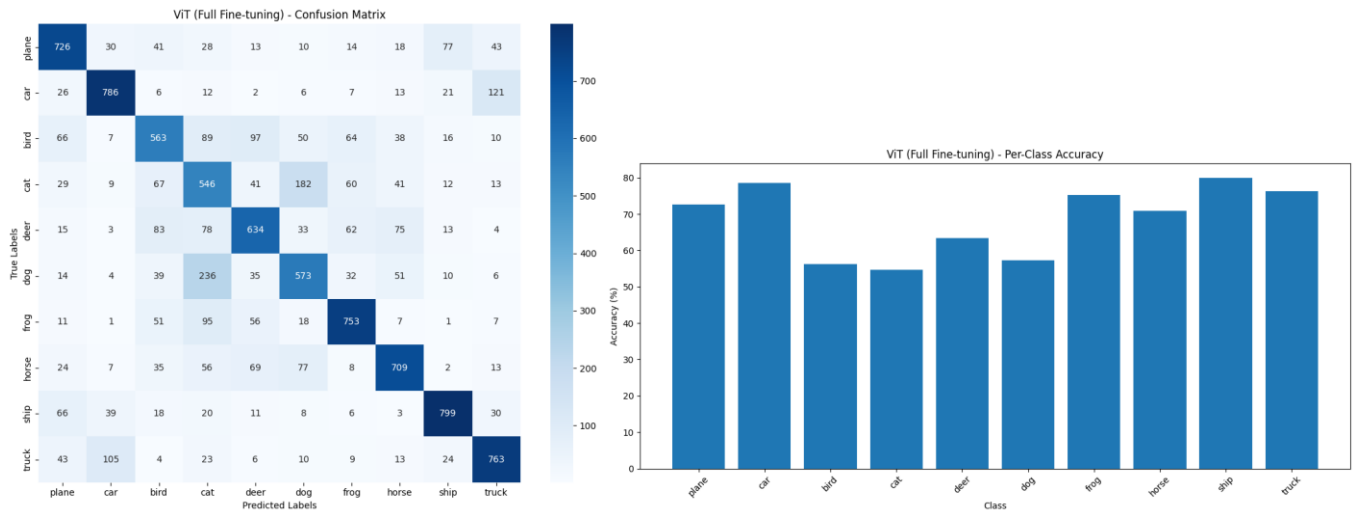
**Figure 3.** Training loss, accuracy and learning rate schedule over epochs for using fully fine-tuned Swin.
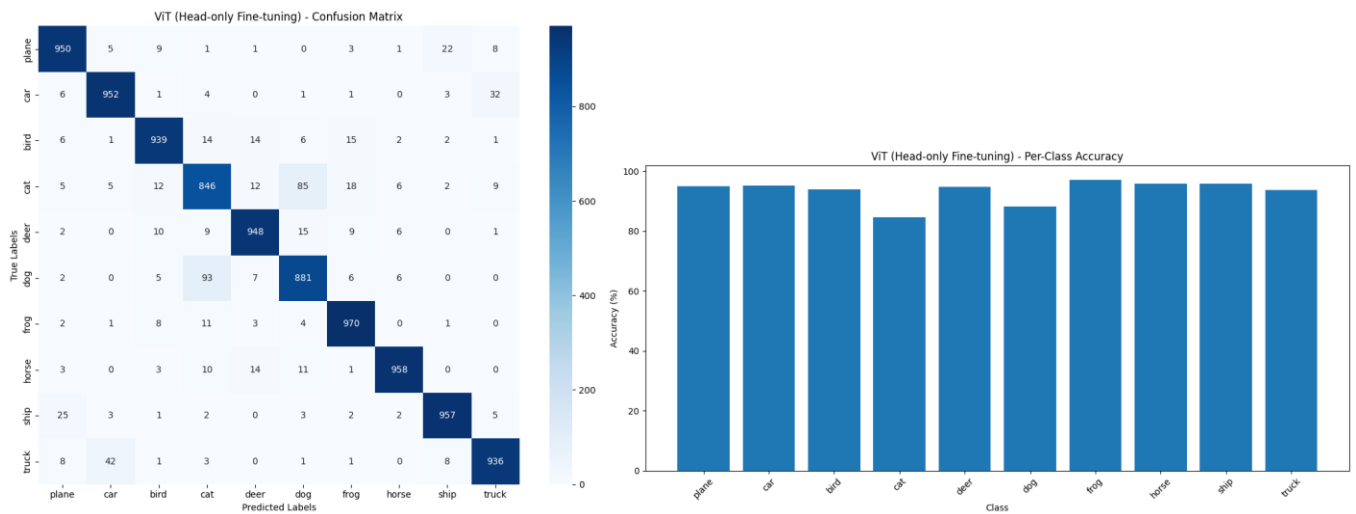


**Figure 4.** Training loss, accuracy and learning rate schedule over epochs for using only head fine-tuned Swin.

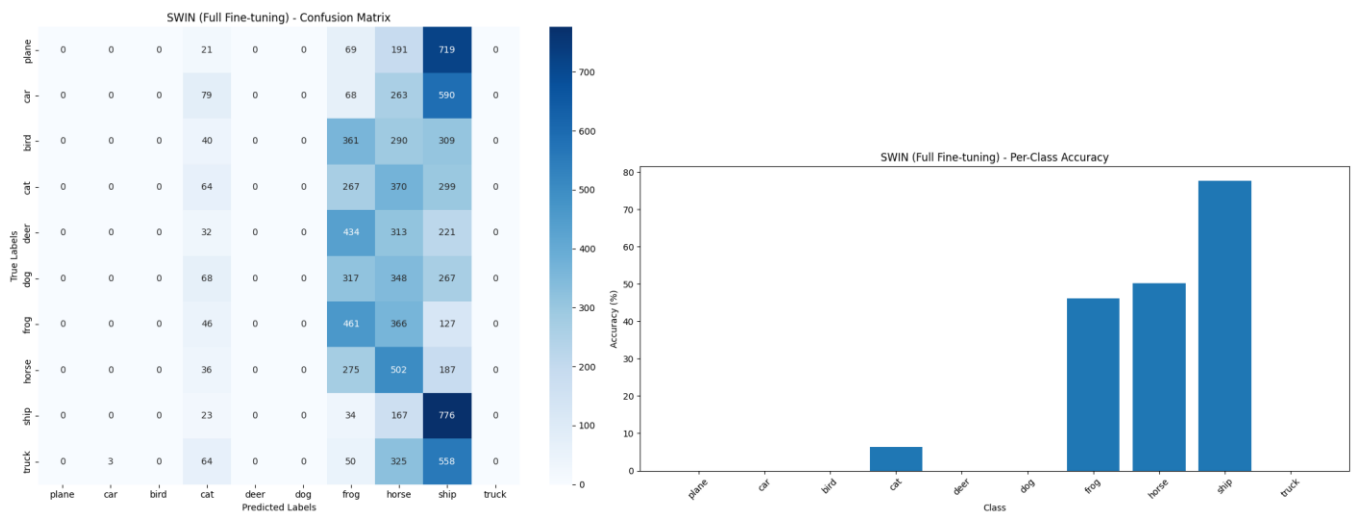|  | Accuracy↑ | Top-1 Error Rate↓ | Average inference time per batch↓ |
|---|---|---|---|
| ViT(full) | 68.52% | 31.48% | **3.18 ms** |
| ViT(head only) | 93.37% | 6.63% | 3.26 ms |
| Swin(full) | 18.03% | 81.97% | 13.37 ms |
| Swin(head only) | **94.21%** | **5.79%** | 13.74 ms |

**Table 1**. Accuracy, top-1 error rate and Average inference time per batch compared with four different models.
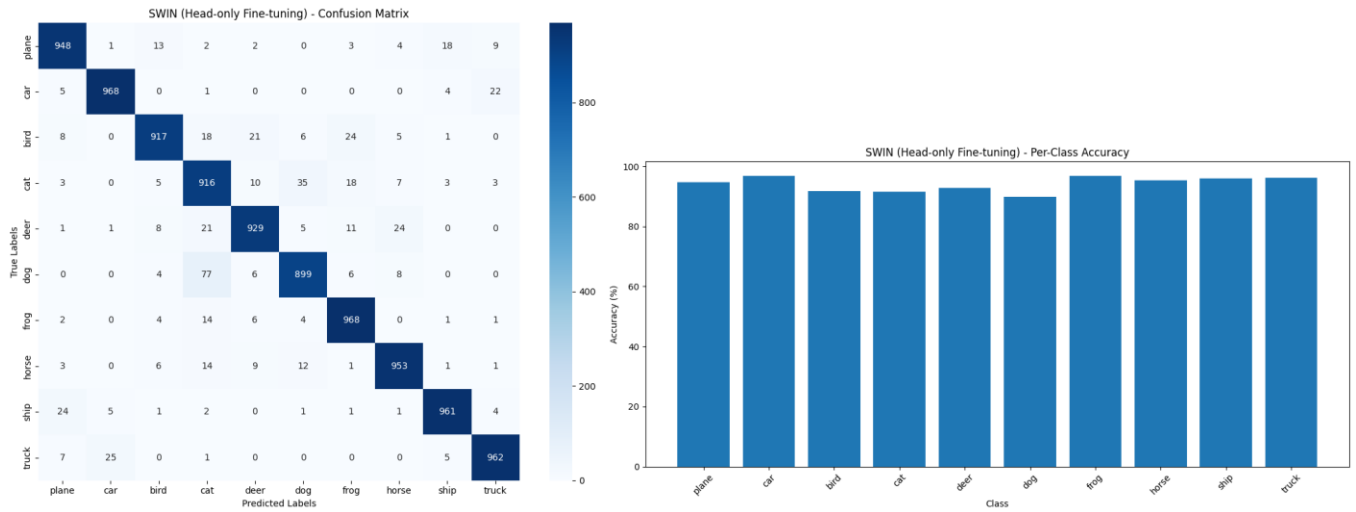
**Figure 5.** Confusion matrix and per class accuracy for using fully fine-tuned Vision Transformer
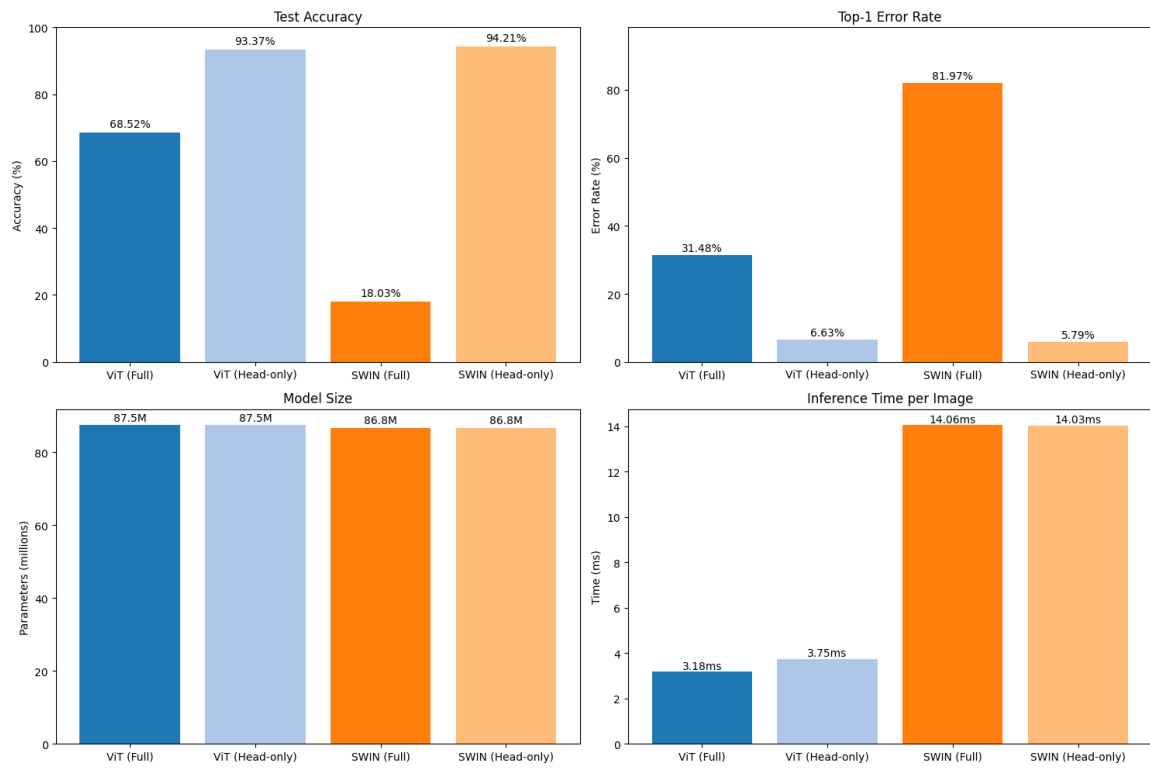


**Figure 6.** Confusion matrix and per class accuracy for using only head fine-tuned Vision Transformer.
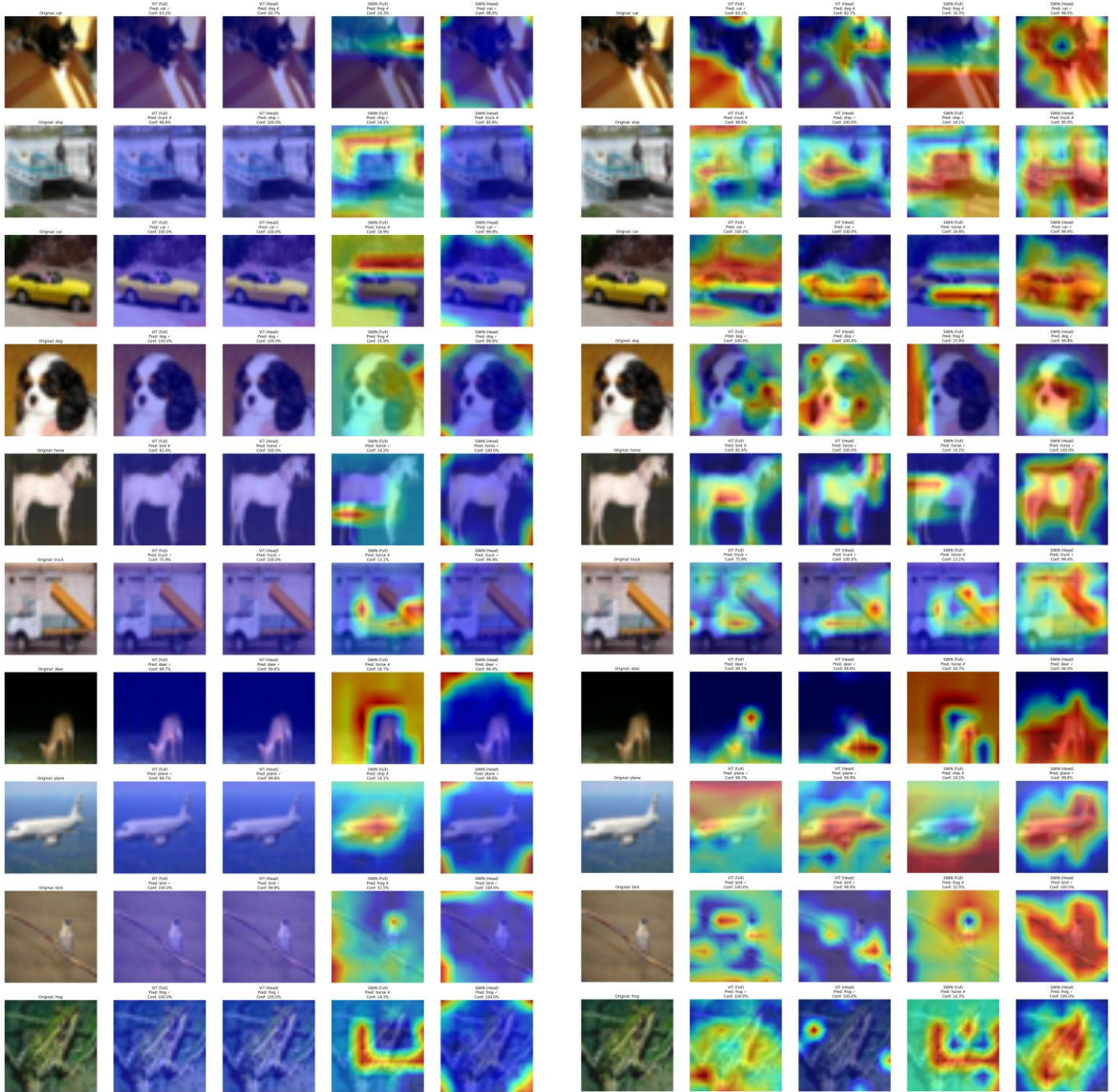


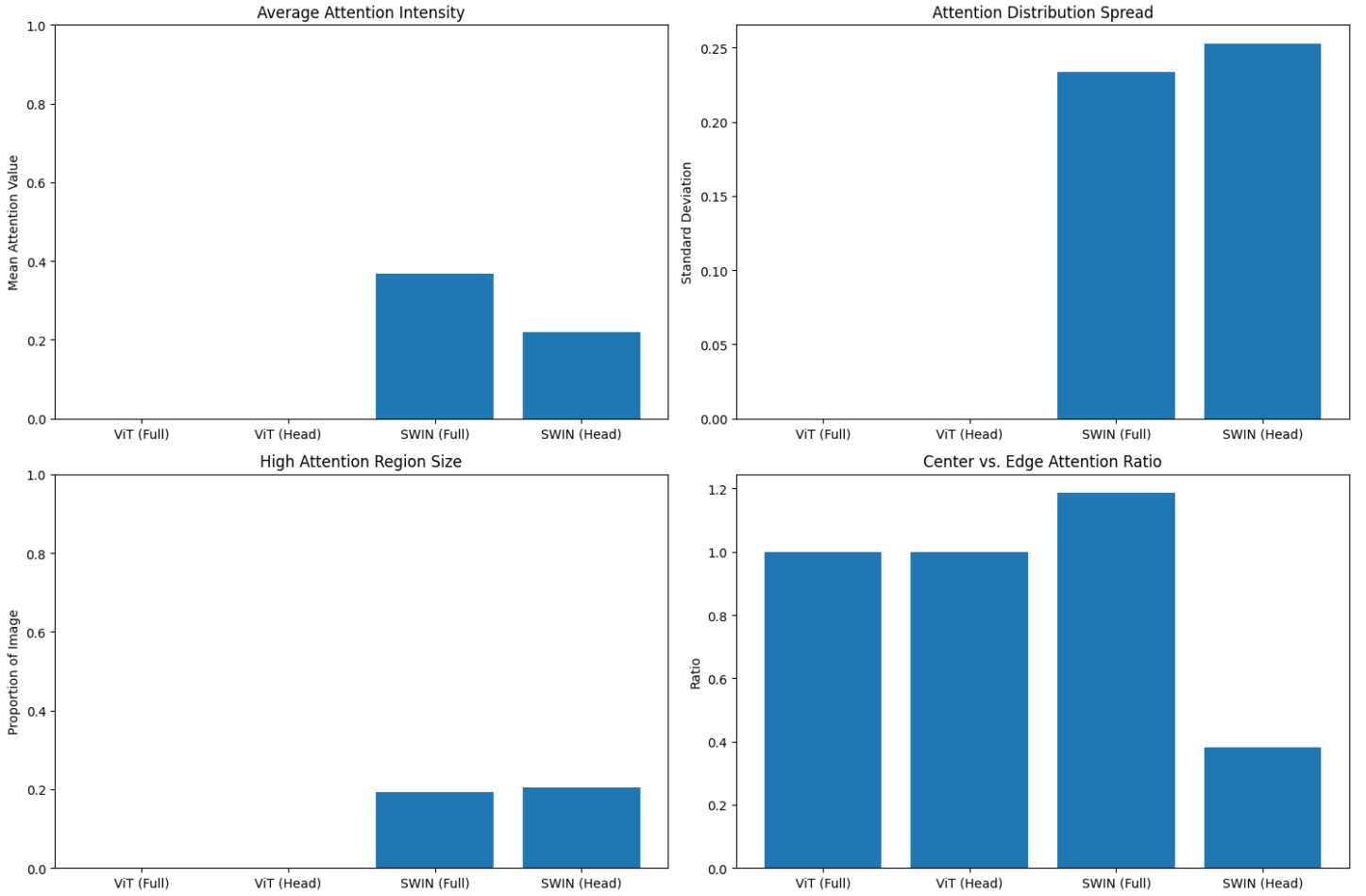**Figure 7.** Confusion matrix and per class accuracy using fully fine-tuned Swin.

***Figure8.*** Confusion matrix and per class accuracy using only head fine-tuned Swin.



***Figure9.*** Overall comparison with test accuracy, top1 error rate, model size and inference time per image.

*Figure10.* Grad Cam's comparison with baseline(left) and change layer with method(right).

***Figure11***. Attention pattern analysis compared with four different models.

| | Attn mean | Attn Std | High Attn Region Size | Center/Edge Ratio |
|---|---|---|---|---|
| ViT(full) | 39.34% | 25.99% | **20.41%** | 1.1743 |
| ViT(head only) | 23.54% | 26.64% | 19.39% | **1.9129** |
| Swin(full) | **46.95%** | 27.84% | **20.41%** | 1.0716 |
| Swin(head only) | 44.54% | **32.87%** | **20.41%** | **1.6750** |

***Table 2***. Attention mean, attention std, high attention region size and center/edge ratio compared with four different models.

## Conclusion

Among the four model configurations, partially fine-tuned models (head and later layers only) outperformed fully fine-tuned ones in both ViT and Swin. Swin with head-only fine-tuning achieved the highest accuracy (94.21%), while fully fine-tuned models suffered from lower performance, likely due to overfitting on the small CIFAR-10 dataset. ViT was more computationally efficient, showing faster inference time due to model selection. Grad-CAM visualizations confirmed that head-only fine-tuned models preserved better attention patterns. Overall, partial fine-tuning is more effective for adapting large pre-trained models to small datasets like CIFAR-10.