

Comparative Analysis of CNN and Transformer Architectures for Remote Sensing Scene Classification

Akhiyar Waladi Universitas Jambi
Jambi, Indonesia
akhiyar.waladi@unja.ac.id

Abstract—This study investigates how much network architecture actually matters for remote sensing scene classification by benchmarking eight deep learning models on two standard datasets: EuroSAT (10 classes, 27,000 Sentinel-2 images) and UC Merced (21 classes, 2,100 aerial images). The models span three design families: classical CNNs (ResNet-50, ResNet-101, DenseNet-121, EfficientNet-B0, EfficientNet-B3), vision transformers (ViT-B/16, Swin Transformer), and a modernized CNN (ConvNeXt-Tiny). Every model was trained with the same hyperparameters, the same augmentation pipeline, and the same ImageNet-pretrained initialization. ConvNeXt-Tiny reached the highest accuracy on EuroSAT (99.06%) and EfficientNet-B3 on UC Merced (99.76%), but the gap between the best and worst model was less than one percentage point on both datasets. McNemar’s test showed that most pairwise differences were not statistically significant. EfficientNet-B0, the smallest model at 4.0M parameters, reached 98.54% and 99.52%, which raises the question of whether these benchmarks can still meaningfully separate architectures. We argue that for standard scene classification tasks with transfer learning, the training recipe matters more than the specific architecture.

Index Terms—Scene classification, remote sensing, deep learning, convolutional neural networks, vision transformers, transfer learning, EuroSAT, UC Merced

I. INTRODUCTION

ASIGNING a single land-use label to a satellite or aerial image patch is one of the oldest problems in remote sensing, and one that has seen large accuracy gains since deep learning entered the field [1], [2]. The task matters because automated scene classification feeds into urban expansion tracking, environmental monitoring, disaster mapping, and national land cover inventories.

Before deep learning, the standard pipeline relied on hand-crafted features: color histograms, texture descriptors, and bag-of-visual-words representations [3]. Features were manually designed and fed to classifiers such as SVMs. While effective for simple cases, this approach was labor-intensive and did not scale well to large numbers of classes [4].

Convolutional neural networks fundamentally changed this paradigm. Networks such as VGGNet [5], ResNet [6], and DenseNet [7] learn features directly from pixels, and when pretrained on ImageNet [8] and fine-tuned on remote sensing data, they consistently outperform handcrafted approaches [9], [10]. Transfer learning proved especially useful because labeled satellite imagery is often scarce [11], [12].

More recently, transformers have arrived in computer vision. The idea, first proposed for language modeling [13], is to process an image as a sequence of patches and let self-attention learn which patches relate to which. Vision Transformers (ViTs) [14] and their hierarchical variants like the Swin Transformer [15] have matched or beaten CNNs on general benchmarks, and researchers quickly began testing them on remote sensing data [16], [17]. At the same time, ConvNeXt [18] showed that a purely convolutional network, when trained with modern recipes borrowed from transformers, can reach the same accuracy level.

This creates an uncomfortable situation for practitioners. There are now three competing families of architectures (classical CNNs, vision transformers, modernized CNNs), each with papers claiming superiority. But most published comparisons test only two or three models, or use different training protocols, or evaluate on a single dataset. It is hard to know whether reported differences come from the architecture itself or from differences in hyperparameters, augmentation, or training schedule.

We set out to remove these confounding factors. We took eight architectures from all three families, gave each one the same ImageNet-pretrained weights, applied the same augmentation and optimization, and measured accuracy on two datasets that cover different imaging modalities and class counts. We then applied McNemar’s test [19] to check whether any observed accuracy gaps were statistically real. We also recorded parameter counts and training times to assess efficiency.

II. RELATED WORK

A. CNN-Based Scene Classification

ResNet [6] introduced shortcut connections that let gradients flow directly through the network, and the 50- and 101-layer variants quickly became the default baselines in remote sensing. DenseNet [7] took a different route: every layer receives input from all preceding layers, which encourages feature reuse and keeps the parameter count low. The EfficientNet family [20] showed that scaling depth, width, and resolution together is more effective than scaling any one dimension alone. Nogueira et al. [9] provided early evidence that fine-tuning ImageNet-pretrained CNNs beats training from scratch for aerial scene recognition, a finding that has been replicated many times since [2], [21].

B. Transformer-Based Approaches

ViT [14] splits an image into fixed-size patches, embeds them, and feeds the sequence to a standard transformer encoder with self-attention. The appeal for remote sensing is that attention can capture relationships between distant image regions that local convolution filters would miss [17]. The main drawback is computational cost: self-attention scales quadratically with the number of patches. Swin Transformer [15] addresses this by computing attention inside local windows that shift across layers, producing a hierarchical feature pyramid at linear cost. Hong et al. [16] applied a transformer design specifically to hyperspectral data, showing that the architecture can also handle non-RGB inputs.

C. Modernized CNNs

ConvNeXt [18] is the result of a thought experiment: starting from a plain ResNet and, one design choice at a time, adopting ideas from transformers. Larger kernels (7×7), layer normalization, GELU activations, and an inverted bottleneck layout brought a standard convolution network to the same accuracy as Swin Transformer on ImageNet. This finding has direct implications for remote sensing: if architecture matters less than training recipe on ImageNet, the same pattern may hold for satellite and aerial imagery.

D. Benchmark Datasets

The UC Merced Land Use dataset [3] has 2,100 aerial images across 21 land-use classes at 0.3 m resolution, drawn from USGS National Map imagery. EuroSAT [22] provides 27,000 Sentinel-2 multispectral patches at 10 m resolution with 10 classes. Larger collections exist, including NWPU-RESISC45 [4] (31,500 images, 45 classes) and AID [23] (10,000 images, 30 classes), but EuroSAT and UC Merced remain popular because they are freely available and small enough to run full experiments quickly. We chose these two specifically because they differ in resolution, image source (satellite vs. aerial), and number of classes, giving us two complementary testbeds.

III. METHODOLOGY

Fig. 1 provides an overview of the research methodology. The pipeline consists of five phases: data acquisition, preprocessing, model training, performance evaluation, and comparative analysis.

A. Datasets

1) EuroSAT: We use EuroSAT [22], a collection of 27,000 Sentinel-2 satellite patches in 10 land-use categories. Each patch is 64×64 pixels at 10 m ground sampling distance. The classes range from natural covers (Forest, SeaLake, River) to agricultural types (AnnualCrop, PermanentCrop, Pasture) and built-up areas (Highway, Industrial, Residential). We use the RGB bands only, since all eight architectures expect three-channel input. Fig. 2 shows one sample per class.

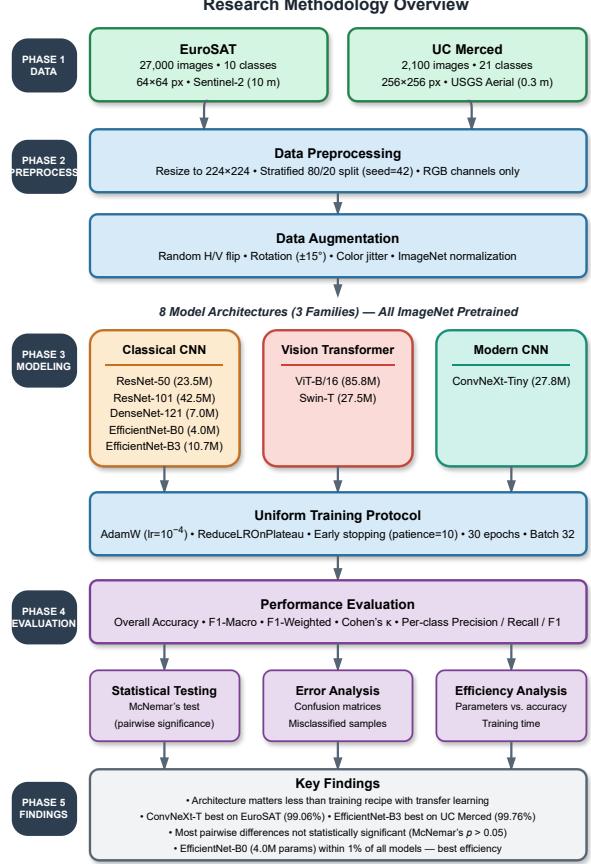


Fig. 1. Overview of the research methodology. Two benchmark datasets are preprocessed with a uniform pipeline and used to train eight architectures from three design families under identical conditions. The trained models are evaluated using classification metrics, statistical significance tests, error analysis, and computational efficiency measures.

2) *UC Merced Land Use:* The UC Merced dataset [3] has 2,100 aerial images at 0.3 m resolution, split evenly across 21 classes (100 images each, 256×256 pixels). The fine resolution means individual buildings, tennis courts, and storage tanks are clearly visible, but it also means that classes like denserresidential, mediumresidential, and sparseresidential differ only in the spacing between structures, which makes them easy to confuse. Samples from all 21 classes appear in Fig. 3.

3) *Data Partitioning:* We applied an 80/20 stratified random split with a fixed seed (42) for both datasets. This produces 21,600 training and 5,400 test images for EuroSAT, and 1,680 training and 420 test images for UC Merced.

B. Model Architectures

We selected eight architectures to cover three families. Table I lists them along with their parameter counts and publication years.

Classical CNNs. ResNet-50 and ResNet-101 are residual networks with 50 and 101 layers. DenseNet-121 connects each layer to every other in a feed-forward fashion, reusing features

EuroSAT Sentinel-2 Satellite Samples



UC Merced Aerial Image Samples

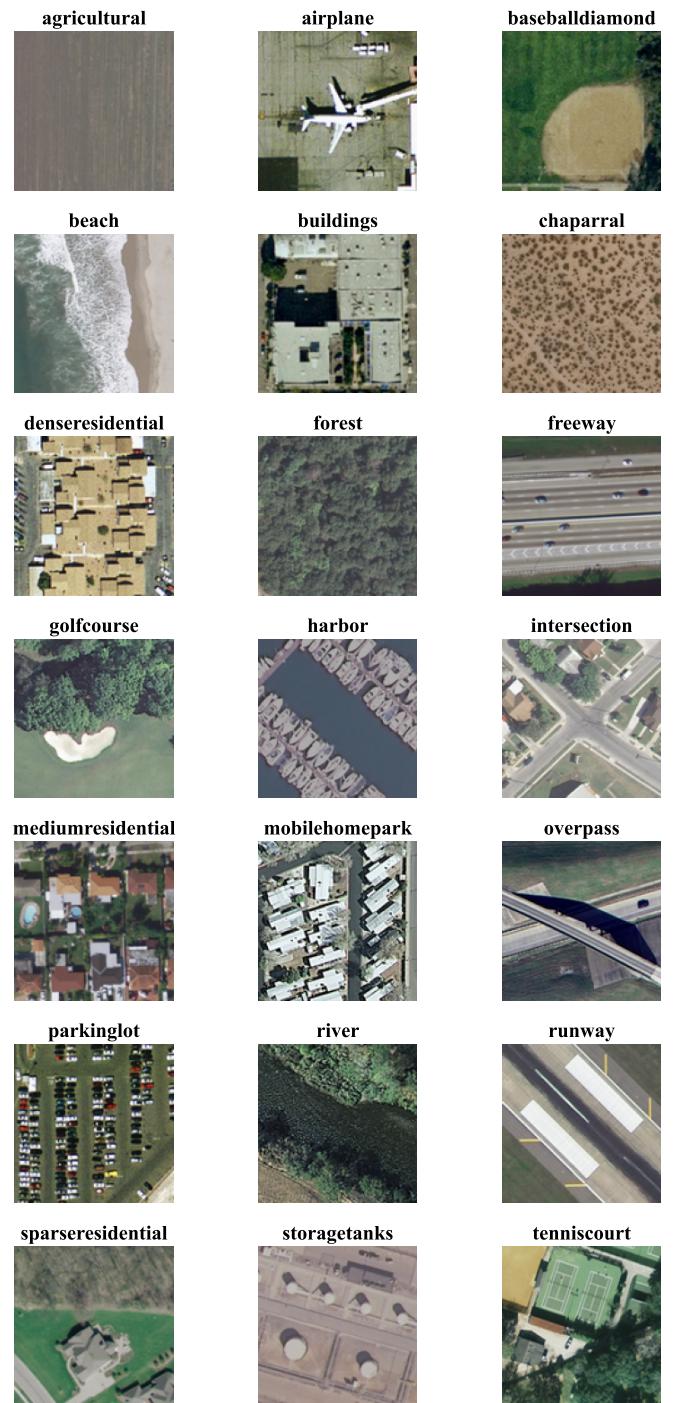


Fig. 2. Sample images from EuroSAT. One randomly selected Sentinel-2 patch (64×64 pixels, 10m resolution) is shown for each of the 10 classes.

Fig. 3. Sample images from UC Merced. One randomly selected aerial image (256×256 pixels, 0.3 m resolution) is shown for each of the 21 classes.

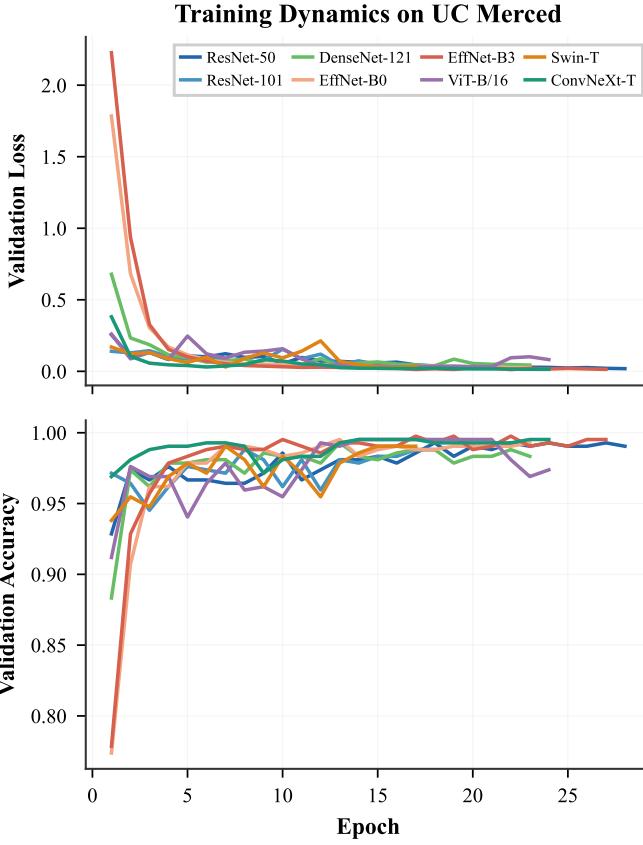


Fig. 7. Training dynamics on UC Merced. Top: validation loss. Bottom: validation accuracy. Convergence is faster due to the smaller dataset.

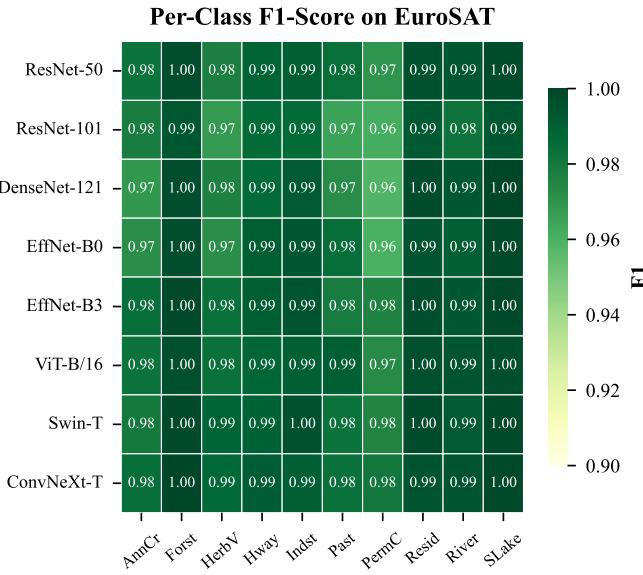


Fig. 8. Per-class F1-score heatmap on EuroSAT. Rows are models, columns are classes. Darker green is higher F1. PermanentCrop and River show the most variation.

2) *UC Merced*: On UC Merced (Fig. 9), the picture is even more uniform. Most cells in the heatmap are saturated at $F1 = 1.0$, meaning perfect classification. The exceptions are the residential classes: denseresidential, mediumresidential, sparseresidential, and to some extent buildings. These classes share similar visual content (rooftops, streets, trees), and the distinction between “dense” and “medium” residential is somewhat subjective even for a human interpreter.

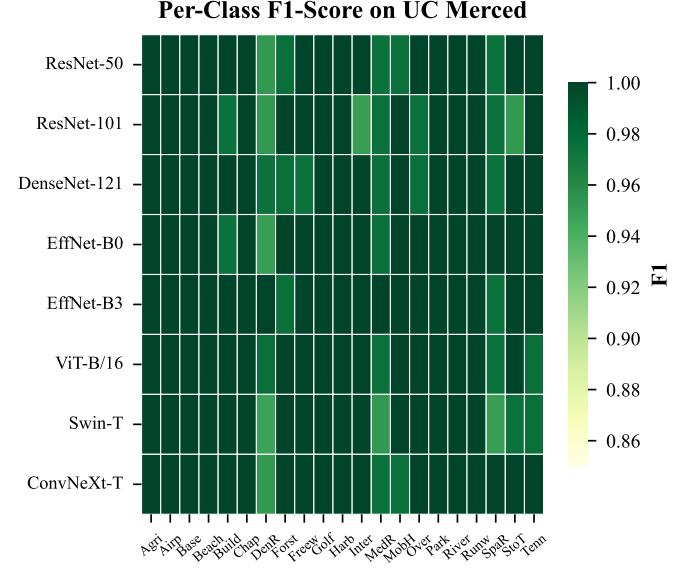


Fig. 9. Per-class F1-score heatmap on UC Merced (21 classes). Most cells are at $F1 = 1.0$ (dark green). Residential classes show the most variation.

D. Statistical Significance

Fig. 10 shows the McNemar p-value matrix for EuroSAT. Out of 28 pairwise comparisons, only a handful produce $p < 0.05$, mostly involving ResNet-101 versus the better-performing models. The four top models (ConvNeXt-Tiny, Swin-T, EfficientNet-B3, ViT-B/16) are not significantly different from each other: the green cells in the upper-left block of the matrix show $p > 0.05$ for every pair. This indicates that a different random test split could easily rearrange their ranking.

On UC Merced, even fewer pairs reach significance (not shown for brevity), which is expected given the smaller test set (420 images) and the tighter accuracy range. The McNemar results tell us that the accuracy ordering we observe is partly an artifact of which particular images ended up in the test set, not a reliable indicator of one architecture being strictly better than another.

E. Error Analysis

To identify the sources of remaining errors, Fig. 11 breaks down the predictions of ConvNeXt-Tiny on EuroSAT by class. Most classes have zero or near-zero misclassifications. The errors concentrate in Pasture and PermanentCrop: PermanentCrop gets confused with HerbaceousVegetation, and Pasture gets confused with AnnualCrop. Both patterns involve

Misclassified Examples on EuroSAT

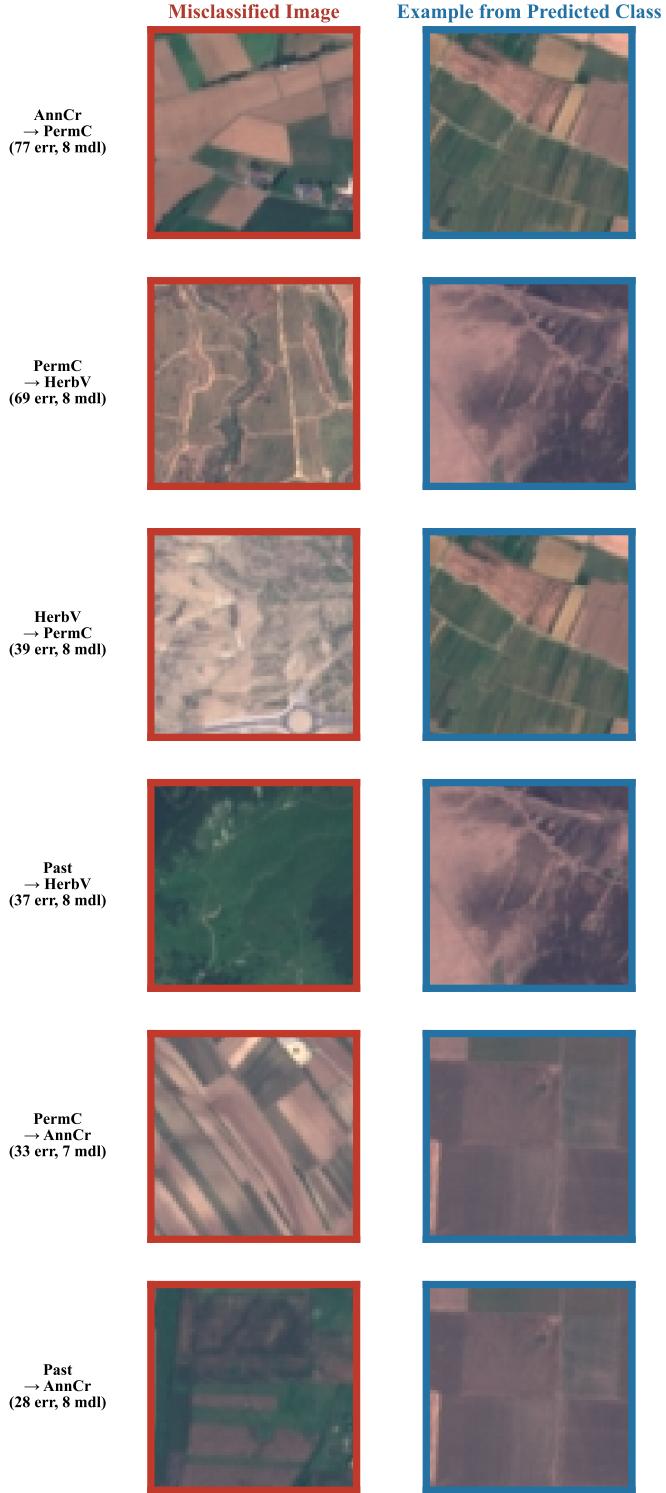


Fig. 12. Misclassified EuroSAT examples. Left (red border): misclassified test image. Right (blue border): correctly classified example from the predicted class. The visual similarity between each pair explains why models confuse them.

Misclassified Examples on UC Merced



Fig. 13. Misclassified UC Merced examples. Same layout as Fig. 12. Residential density classes are the main source of confusion.

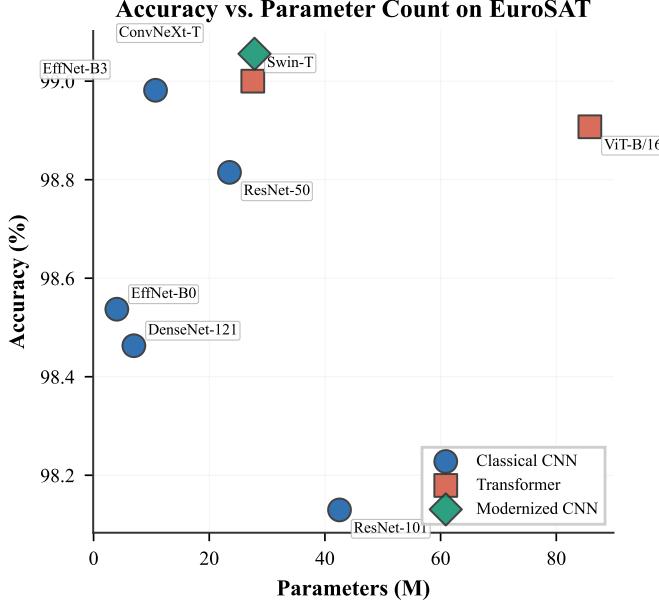


Fig. 14. Accuracy vs. parameter count on EuroSAT. Blue circles = classical CNNs, red squares = transformers, green diamonds = modernized CNN.

story. On EuroSAT, a modernized CNN (ConvNeXt-Tiny) finished first. On UC Merced, a classical CNN (EfficientNet-B3) finished first. The transformers performed well, but they did not dominate. This is consistent with the original ConvNeXt paper [18], which argued that the accuracy gap between CNNs and transformers closes once CNNs adopt modern training practices.

The likely explanation is the nature of the task. When images are resized to 224×224 pixels and each patch contains a single land-use type, local texture and color patterns may be sufficient for classification. Self-attention’s ability to relate distant patches may become more useful for tasks that require understanding spatial layout, such as detecting a harbor by noticing boats near a dock, rather than just recognizing a uniform texture.

B. The Depth Paradox

ResNet-101 consistently trailed ResNet-50 on both datasets (98.13% vs. 98.81% on EuroSAT, 98.81% vs. 99.29% on UC Merced), even though it has nearly twice as many parameters. This pattern, where a deeper pretrained network underperforms a shallower one after fine-tuning, has been reported before in transfer learning. The likely explanation is that 30 epochs of fine-tuning with early stopping is not enough to properly adapt all 101 layers. The deeper network starts from a good initialization but cannot move far from it before training stops, while the 50-layer version has a smaller parameter space that is easier to tune within the same budget.

C. Parameter Efficiency

EfficientNet-B0 deserves special attention. With 4.0M parameters ($18\times$ fewer than ViT-B/16), it reaches 98.54% on EuroSAT and 99.52% on UC Merced. For anyone deploying

a scene classification model on a mobile device, a drone, or an edge computing node, this is the clear choice. EfficientNet-B3, at 10.7M parameters, matches or beats every other model in our lineup while remaining small enough for practical deployment.

D. Dataset Saturation

When the worst model in a benchmark still scores above 98%, the benchmark has arguably reached its useful limit. The narrow accuracy range we observed (less than one percentage point on both datasets) means that the difference between the “best” and “worst” architecture is within noise for most practical applications. This echoes calls in the community for harder evaluation scenarios: larger-scale datasets, cross-domain generalization tests, few-shot settings, and tasks that go beyond single-label classification [1], [4].

E. Transfer Learning Dominance

The most striking pattern in our results is how little architecture matters once transfer learning is applied. The McNemar test (Fig. 10) shows that most accuracy differences are not statistically significant. Running the same experiment with a different random test split would likely produce a different model ranking. This implies that, at least for these benchmarks, the pretrained weights and the fine-tuning recipe are the main drivers of accuracy, not the architectural design itself. Further investigation is needed to determine whether this conclusion holds for more specialized tasks such as multi-label classification, change detection, or few-shot learning.

F. Limitations

We tested on two datasets only; results could differ on larger or more diverse collections. We used only the RGB bands of EuroSAT, even though the full multispectral data might benefit some architectures more than others. We fixed all hyperparameters across models, which is fair for comparison but may not produce the best possible result for each individual architecture. And we used a single train-test split without cross-validation, so the results depend on one particular partition of the data.

VI. CONCLUSION

We compared eight deep learning architectures across three design families on two remote sensing benchmarks and found that the differences between them are small. All models exceeded 98% accuracy with ImageNet pretraining and a uniform fine-tuning protocol. ConvNeXt-Tiny reached the top accuracy on EuroSAT (99.06%) and EfficientNet-B3 on UC Merced (99.76%), but McNemar’s test showed that most pairwise gaps were not statistically significant. The smallest model, EfficientNet-B0 (4.0M parameters), came within one percentage point of every other model on both datasets.

For practitioners choosing a model, we recommend EfficientNet-B3 as a general default: it is small (10.7M parameters), fast to train, and competitive everywhere. If model size is the main constraint, EfficientNet-B0 is a strong

