

# FFGF-Net: A Forestry Fine-Grained Fusion Network for Tree Species Classification from UAV Imagery

Xuanzi Zhou<sup>a</sup>, Jingyi Liu<sup>a</sup>, Zhulin Chen<sup>b,c</sup>, Hongxin Yang<sup>d,e,\*</sup>, Sheng Xu<sup>a,\*</sup>

<sup>a</sup>College of Information Science and Technology & Artificial Intelligence, Nanjing Forestry University, Nanjing, 210037, Jiangsu, China

<sup>b</sup>Institute of Forest Resource Information Techniques, Chinese Academy of Forestry, Beijing, 100091, China

<sup>c</sup>State Forestry and Grassland Administration, Key Laboratory of Forest Management and Growth Modelling, Beijing, 100091, China

<sup>d</sup>Key Laboratory of Geographic Information Science (Ministry of Education), East China Normal University, Shanghai, 200241, China

<sup>e</sup>Key Laboratory of Spatial-temporal Big Data Analysis and Application of Natural Resources in Megacities, Ministry of Natural Resources, East China Normal University, Shanghai, 200241, China

## Abstract

Accurate classification of tree species is crucial for biodiversity conservation, forest resource management, and intelligent breeding. Traditional manual survey methods are inefficient and costly, while the AI-driven paradigm offers a new approach for intelligent forestry through computer vision technologies, enabling large-scale, high-precision tree species identification. However, fine-grained classification of tree species remains challenging due to high intra-class variation, low inter-class distinction, and complex environmental noise in real-world forest scenarios. To bridge the intra- and inter-class discrepancy, this study proposes a Forestry Fine-Grained Fusion Network (FFGF-Net), which integrates three key innovations: a CanopyPatchExtractor for local feature extraction, a CanopyTextureAnalyzer for multi-scale texture modeling, and a Cascaded Canopy Attention mechanism for noise suppression. For comprehensive evaluation, we introduce a newly collected UAV dataset, NJFUTreeData, and utilize established public datasets, including SZUTreeData and ETH dataset. Experimental results demonstrate that FFGF-Net achieves state-of-the-art performance, with Overall Accuracy (OA) of 75.00% on NJFUTreeData and 94.75% on the ETH dataset. Notably, our method consistently outperforms a wide range of existing model families, including representative CNN backbones (ConvNeXt-Tiny), lightweight architectures (e.g., EfficientNet-B0, MobileNet\_V3-Small), and advanced Vision Transformers (e.g., Swin-Tiny). The proposed network provides a robust technical foundation for automated forest resource surveys, ecological monitoring, and smart forestry applications.

**Keywords:** Computer vision, Forest resource monitoring, Fine-Grained visual classification, Tree species, UAV remote sensing

## 1. Introduction

Accurate tree species classification is crucial for biodiversity conservation (Sun et al., 2019; Poorter et al., 2014), forest inventory (Fang et al., 2020), and forest breeding (Mace et al., 2012). Traditional manual survey methods suffer from low efficiency and high costs (Zou et al., 2020). The AI-for-Science paradigm offers a new path for forestry intelligence by enabling large-scale, high-precision tree species identification through computer vision technology, thereby fostering interdisciplinary innovation between forestry and artificial intelligence (Zhang et al., 2023).

Fine-grained tree species classification in forestry faces dual challenges: environmental complexity and biodiversity (Yu et al., 2023). In real-world forest scenarios, species identification must overcome dynamic factors such as varying illumination and occlusion interference, while simultaneously distinguishing between morphologically highly similar, closely related species (e.g., *Pinus spp.*, *Quercus spp.*). This places higher

demands on the discriminative ability of classification models (Chen et al., 2023).

Current remote sensing-based tree species classification research primarily focuses on vegetation parameter inversion at the stand level or individual tree crown (contour) extraction. There is comparatively less involvement in fine-grained classification using high-resolution imagery at the crown scale. Current fine-grained classification in forestry primarily faces the following challenges: (1) Significant intra-class variance and small inter-class differences. (2) Discriminative features exist across multiple scales. (3) Complexity and noise in field environments (see Fig. 1).

To address these challenges, various deep learning approaches have been explored. In general visual classification, standard architectures like ResNet (He et al., 2016) and Vision Transformers (Dosovitskiy et al., 2020) extract global contextual features but often fail to adequately capture subtle local feature differences between sub-categories. Consequently, fine-grained visual classification (FGVC) has evolved to distinguish visually similar sub-categories (Zhao et al., 2017). In forestry remote sensing, initial works transferred pre-trained CNNs to tree crown imagery (Fricker et al., 2019) and evaluated

\*Corresponding authors:

E-mail addresses:

hxyang@geoai.ecnu.edu.cn (H. Yang), xusheng@njfu.edu.cn (S. Xu)

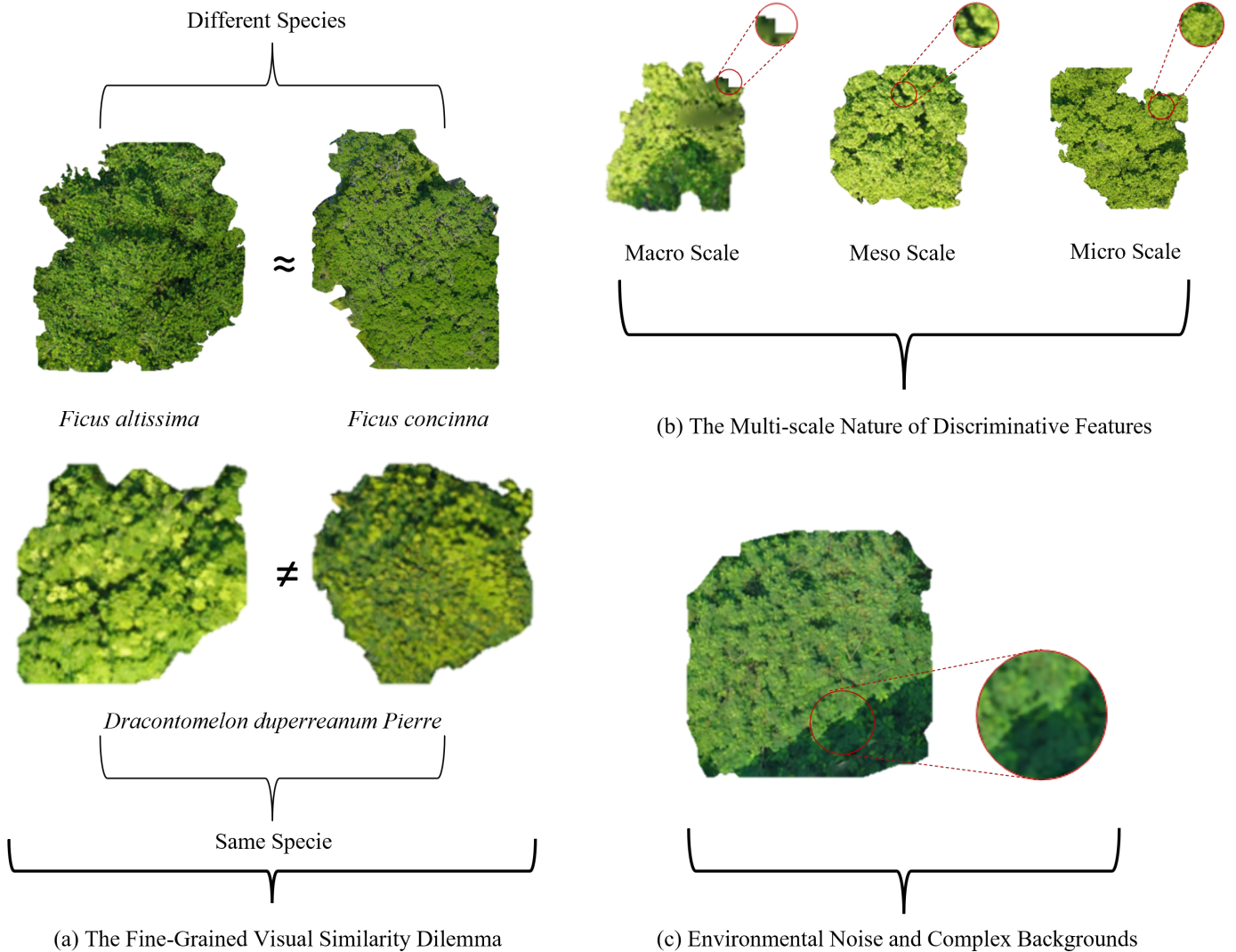


Figure 1: Challenges in fine-grained tree species classification

41 data augmentation strategies to improve adaptability (Neupane 59  
 42 et al., 2021). As understanding deepened, specialized archi- 60  
 43 tectures emerged, such as the Bilinear Squeeze-and-Excitation 61  
 44 Network (BiSENet) (He Z. & He D., 2020) and Transformer- 62  
 45 based Multi-scale Feature Fusion (Dong et al., 2025). Fur- 63  
 46 thermore, recent trends have explored multi-source data fusion, 64  
 47 combining LiDAR with optical imagery for urban tree invento- 65  
 48 ries (Ferreira et al., 2024; Liu et al., 2025).

49 Despite this significant progress, existing deep learning- 67  
 50 based forestry identification research still lacks specific im- 68  
 51 provements for crown-scale imagery captured under complex 69  
 52 field conditions. Most current methods are either limited to 70  
 53 coarse-grained categorization or exhibit limited discriminative 71  
 54 capability when differentiating visually similar, closely-related 72  
 55 species. They predominantly rely on generic feature extraction 73  
 56 mechanisms that struggle with severe background occlusion 74  
 57 and are constrained by their inability to comprehensively cap- 75  
 58 ture and model subtle yet discriminative features across multi-

76 ple scales.

In addition to overcoming the limitations of fine-grained vi-  
 77 sual similarity and texture, it is also necessary to consider the  
 78 role of 3D structural clues. Structural attributes, such as crown  
 79 geometry, canopy height variation, and vertical forest structure,  
 80 are often crucial for accurate tree species discrimination. Re-  
 81 cent large-scale datasets integrating structural forest attributes  
 82 further emphasize the importance of 3D information for im-  
 83 proving species discrimination (Ali et al., 2026; Calders et al.,  
 84 2022; Hollaus & Chen, 2022; Weiser et al., 2022). Furthermore,  
 85 recent studies have demonstrated that incorporating structural  
 86 features derived from LiDAR significantly improves vegetation  
 87 characterization and biomass estimation, highlighting the limi-  
 88 tations of purely image-based approaches (Ali et al., 2024; Wild  
 89 et al., 2026). Recognizing these structural dimensions provides  
 90 a more comprehensive perspective for forest inventory chal-  
 91 lenges. However, despite the undeniable value of 3D data, its  
 92 acquisition is often costly and complex to operate. Therefore,

77 maximizing the representativeness of cost-effective 2D optical  
78 images remains a highly practical and urgent research focus.

79 To bridge this specific gap and overcome the aforementioned  
80 challenges, this work introduces the Forestry Fine-Grained Fusion  
81 Network (FFGF-Net). The methodological novelty of our  
82 framework lies in its targeted design to address the shortcomings  
83 of existing models: (1) A CanopyPatchExtractor module that  
84 regularly segments canopy images into multiple local regions,  
85 mitigating large intra-class differences caused by growth  
86 environments while amplifying subtle local differences between  
87 species; (2) A CanopyTextureAnalyzer module, which utilizes  
88 parallel convolutional layers with different dilation rates to capture  
89 multi-scale features simultaneously, extracting texture information  
90 across different scales; (3) A Cascaded Canopy Attention mechanism  
91 that suppresses noise features representing background, shadows, and  
92 occlusions. On the NJFUTreeData dataset, the FFGF-Net achieved an  
93 accuracy of 75.00%, representing a noticeable improvement over existing  
94 state-of-the-art models and providing reliable technical support for  
95 forestry intelligence.  
96

## 97 2. Study Area and Data

98 This section overviews the three forestry remote sensing  
99 datasets from distinct geographical regions: SZUTreeData  
100 (SZU) (<http://szu-hsilab.com/szu-tree-dataset/>),  
101 NJFUTreeData (NJFU), and ETH dataset (ETH) ([https://form.ethz.ch/](https://form.ethz.ch/research/tree-ai-global-database/treeai-competition.html)  
102 [research/tree-ai-global-database/treeai-competition.html](https://form.ethz.ch/research/tree-ai-global-database/treeai-competition.html)), the  
103 locations of which are mapped in Fig. 2. These datasets encom-  
104 pass a variety of tree species and forest conditions, providing a  
105 comprehensive basis for validation.

106 SZUTreeData represents the world’s first multimodal remote  
107 sensing dataset designed for individual tree segmentation and  
108 species identification, featuring hyperspectral imagery, LiDAR  
109 data, and high-resolution optical imagery. Comprising tens of  
110 millions of annotated pixels across 25 common tree species, this  
111 dataset was collected and publicly shared by the Hyperspectral  
112 Remote Sensing Team at Shenzhen University using integrated  
113 UAV and ground surveys. It covers a total area of approxi-  
114 mately 1.5 square kilometers and includes detailed information  
115 on various land cover types such as roads, buildings, forest  
116 stands, grassland, and water bodies. SZUTreeData serves as  
117 a benchmark for precise species identification and multi-modal  
118 classification research. Fig. 3 (a1-a2) displays the true-color  
119 images of SZUTreeRGB, the individual tree annotation map,  
120 and the tree species category map derived from the SZUTree-  
121 Data (Jia et al., 2024; Li et al., 2024; Long et al., 2024; Zhang  
122 et al., 2024).

123 The ETH dataset, released by ETH Zurich, is a public dataset  
124 widely adopted in forestry remote sensing and computer vision  
125 research (Beloiu Schwenke et al., 2025). It integrates  
126 multi-source remote sensing data, including high-resolution op-  
127 tical imagery, LiDAR point clouds, and hyperspectral informa-  
128 tion, covering multiple tree species and stand types. Noted  
129 for its data diversity and annotation accuracy, the ETH dataset  
130 provides rich tree-level annotations such as crown contours,  
131 species labels, and height models, making it suitable for deep

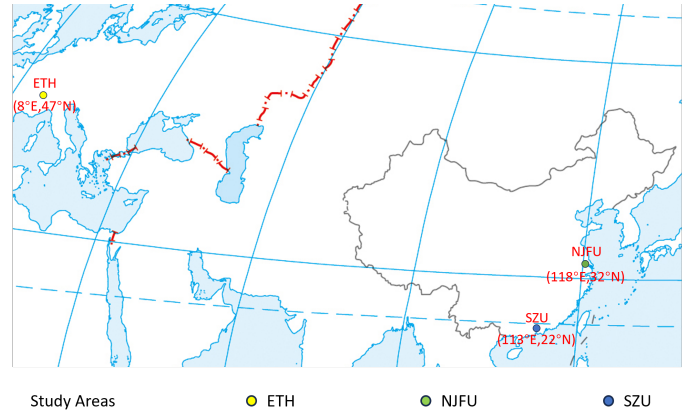


Figure 2: The geographical locations of study areas. The base map is sourced from the Standard Map Service website of the Ministry of Natural Resources of the People’s Republic of China (Approval Number: GS(2016)2945).

132 learning-based tasks including tree detection, segmentation,  
133 and classification. Its rigorous data acquisition and annotation  
134 protocols have established it as a benchmark dataset widely rec-  
135 ognized by the international academic and industrial communi-  
ties. Sample images from this dataset and their corresponding  
annotation schematics are exhibited in Fig. 3 (b1-b2).

Newly constructed specifically for this study, NJFUTreeData  
is a high-resolution UAV remote sensing dataset designed for  
forest resource investigation, focusing on the tree distribution  
area in Xuanwu District, Nanjing, China. Data acquisition was  
conducted using a DJI Matrice 3E series UAV equipped with  
an RTK module, flying at an altitude of 120 meters with a criss-  
cross flight plan. This rigorous acquisition strategy yielded an  
ultra-high spatial resolution of approximately 2 to 3 centime-  
ters per pixel, ensuring the precise capture of subtle structural  
and textural characteristics of tree canopies. Furthermore, the  
dataset exhibits remarkable taxonomic diversity, encompassing  
400 distinct plant samples distributed across 94 families  
and 222 genera. The dataset includes various data types, such  
as high-resolution orthoimagery and Digital Surface Models,  
providing detailed information on tree morphological structure,  
distribution density, and growth status. It is suitable for forestry  
remote sensing applications including individual tree segmen-  
tation, species identification, and biomass estimation, offering  
high-quality benchmark data for the precise monitoring and  
management of forest resources. In this paper, Fig. 3 (c1-c3)  
presents the overall orthomosaic of the dataset, sample cropped  
patches along with their corresponding annotation maps, and  
the final individual tree segmentation and species classification  
results.

## 3. Methodology

### 3.1. Problem Formulation

The task of automated tree species classification is funda-  
mentally a Fine-Grained Visual Categorization (FGVC) prob-  
lem. In contrast to generic image classification, which dis-  
tinguishes broadly divergent categories (e.g., cats vs. dogs),

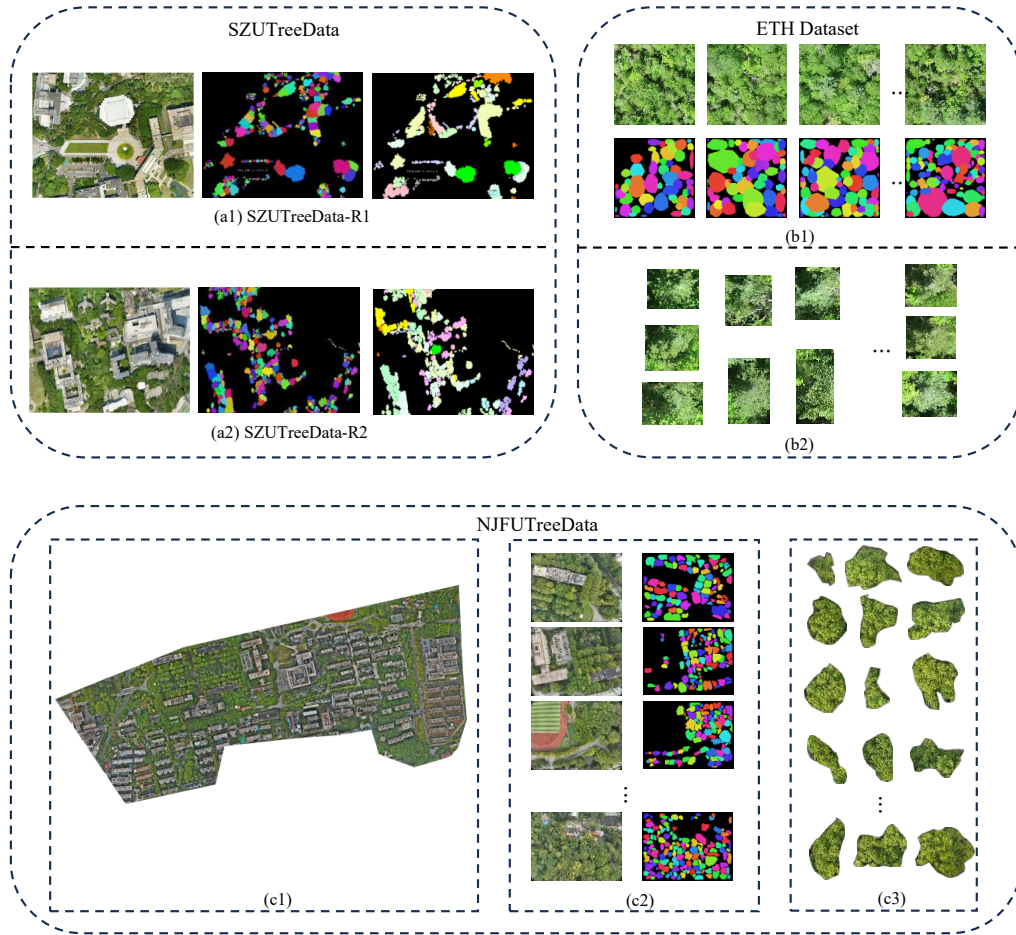


Figure 3: Three forestry remote sensing datasets used in this study. (a1) True-color images of SZUTreeRGB-R1 data, along with individual tree annotation map and tree species category map for SZUTreeDataR1 data; (a2) true-color images of SZUTreeRGB-R2 data, along with individual tree annotation map and tree species category map for SZUTreeDataR2 data. (b1) Individual tree annotation map for ETH dataset; (b2) the final map of individual tree segmentation. (c1) The overall orthomosaic of the NJFUTreeData dataset; (c2) individual tree annotation map for NJFUTreeData dataset; (c3) individual tree segmentation and crown contours.

168 FGVC requires discriminating between visually similar sub-167  
 169 categories under a common super-category. In the context of 168  
 170 forestry, this entails differentiating among tree species (all be-169  
 171 longing to the macro category “trees”) based on subtle and-190  
 172 localized discriminative features, such as leaf texture, crown-191  
 173 structure, or branching pattern. These distinguishing charac-192  
 174 teristics are often confined to specific regions of the canopy,193  
 175 while significant intra-class variations may arise due to factors-194  
 176 including growth environment, health status, and imaging con-195  
 177 ditions. Given a forestry image dataset  $D = \{(I_n, y_n)\}_{n=1}^N$ , where-196  
 178  $I_n \in \mathbb{R}^{H \times W \times 3}$  represents the input image and  $y_n \in \{1, 2, \dots, C\}$  is-197  
 179 the corresponding tree species label, Here,  $N$  denotes the total  
 180 number of samples in the dataset,  $H$  and  $W$  represent the spatial  
 181 height and width of the images, respectively, and  $C$  indicates the  
 182 total number of tree species categories. The goal of this study,198  
 183 is to construct a mapping function  $F : \mathbb{R}^{H \times W \times 3} \rightarrow \{1, 2, \dots, C\}$ ,199  
 184 such that the input image  $I$  is accurately mapped to the pre-200  
 185 dicted class  $\hat{y}$ . 201

186 As illustrated in Fig. 4, the FFGF-Net executes the mapping202

function  $F$  through a multi-stage pipeline. A Backbone CNN  
 ( $F_b$ ) performs initial feature learning, succeeded by Gated Fea-  
 ture Reduction ( $F_d$ ) for channel pruning. Multi-scale pattern  
 analysis and feature recalibration are then concurrently man-  
 aged by the CanopyTextureAnalyzer ( $F_t$ ) and Cascaded Canopy  
 Attention ( $F_a$ ), respectively. The CanopyPatchExtractor ( $F_p$ )  
 subsequently divides the enhanced features into localized to-  
 kens, with their mutual dependencies captured by Position-  
 aware Patch Attention ( $F_s$ ). The classification ( $F_c$ ) component  
 finalizes the process by integrating these tokens into a potent  
 representation for species categorization.

In summary, the tree species classification problem in  
 forestry demands that the model possesses the capability to  
 discover and utilize subtle, localized discriminative features  
 from highly similar categories, which is the core focus of fine-  
 grained visual recognition research.

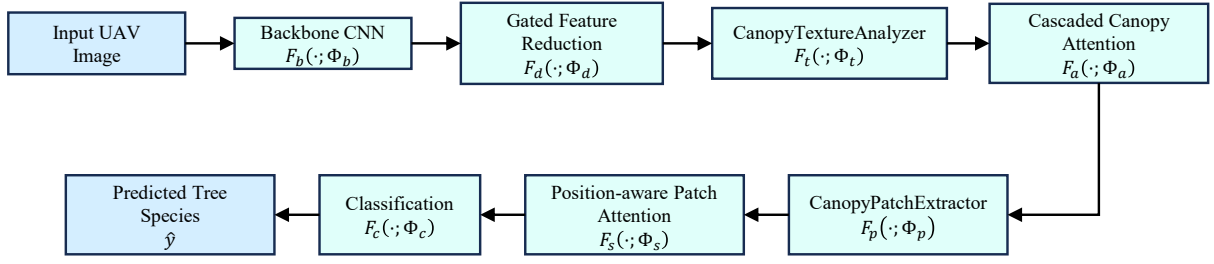


Figure 4: FFGF-Net architecture overview

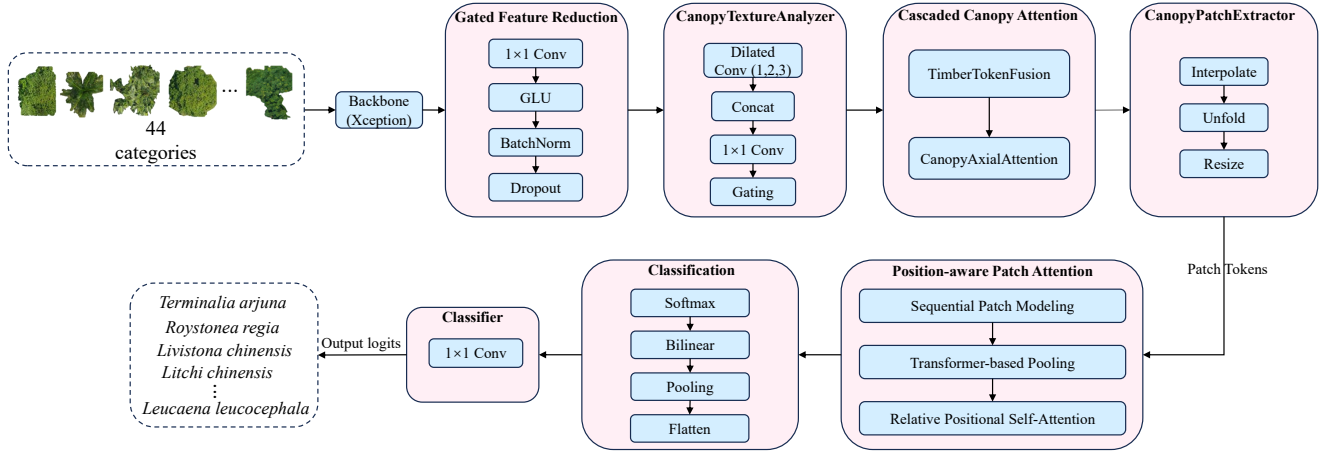


Figure 5: FFGF-Net: a local-global feature fusion network for fine-grained tree species classification from UAV imagery. The framework comprises a backbone CNN for feature extraction, a Gated Feature Reduction module for channel compression, a CanopyTextureAnalyzer for multi-scale texture modeling, a Cascaded Canopy Attention mechanism for discriminative region emphasis and noise suppression, a CanopyPatchExtractor for local token generation, a Position-aware Patch Attention block for spatial-semantic relation reasoning, and a Bilinear Attention Pooling module for feature aggregation and classification.

### 3.2. Gated Feature Reduction

High-dimensional features pose challenges in computational complexity and noise interference. Gated Feature Reduction (GFR) is inspired by GLU (Dauphin et al., 2017), strengthening the representation of foreground targets by compressing the feature dimensionality (as shown in Fig. 6). For an input feature  $X \in R^{H \times W \times 3}$ , GFR first maps it to a transitional channel space using a  $1 \times 1$  convolution:

$$U = W_p * X + b_p, U \in R^{h \times w \times c'} \quad (1)$$

Here,  $*$  denotes the convolution operation,  $W_p \in R^{1 \times 1 \times c \times 2c'}$  represents the convolutional kernel parameters,  $b_p$  is the bias term, and  $c'$  is the designated target number of channels. Subsequently, the tensor  $U$  is split evenly along the channel dimension into two parts  $U = [U_1, U_2]$ , where  $U_1, U_2 \in R^{h \times w \times c'}$ .

The core gating operation is:

$$\hat{U} = \sigma(U_1) \odot U_2 \quad (2)$$

Where  $\sigma(\cdot)$  is the Sigmoid function, which compresses the values of  $U_1$  between 0 and 1, forming a soft selection mask.  $\odot$  denotes element-wise multiplication. The features in  $U_2$  are

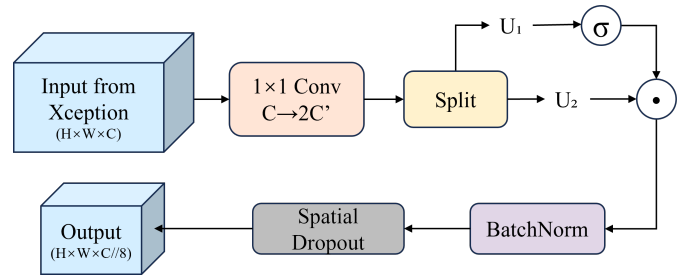


Figure 6: Gated Feature Reduction

reweighted by this mask, suppressing background responses irrelevant to classification and highlighting feature patterns relevant to tree species discrimination.

Finally, standardization and regularization are applied via Batch Normalization (BatchNorm) and Spatial Dropout to obtain the compressed feature output.

This module significantly reduces computational resource consumption while endowing the model with powerful feature selection capabilities.

### 229 3.3. CanopyTextureAnalyzer

230 Leaf texture is an important basis for distinguishing tree  
 231 species. To explicitly extract multi-scale texture patterns, we intro-  
 232 duce the CanopyTextureAnalyzer (CTA) module. This proces-  
 233 s can be seen in Fig. 7. The module takes the input fea-  
 234 ture  $\hat{X} \in \mathbb{R}^{h \times w \times c}$ , and extracts multi-scale contextual features  
 235 through three parallel dilated convolution paths:

$$236 F_1 = \text{ReLU}(W_{d1} * \hat{X}), \text{Dilation} = 1 \quad (3)$$

$$237 F_2 = \text{ReLU}(W_{d2} * \hat{X}), \text{Dilation} = 2 \quad (4)$$

$$238 F_3 = \text{ReLU}(W_{d3} * \hat{X}), \text{Dilation} = 3 \quad (5)$$

239 Here,  $W_{d1}, W_{d2}, W_{d3}$  are the convolution weights. Different di-  
 240 lation rates allow the kernels to have varying receptive fields  
 241 without increasing the number of parameters, thereby capturing  
 242 texture structures from fine to coarse. Subsequently, the multi-  
 scale features are concatenated and fused via convolution:

$$243 F_{\text{concat}} = [F_1, F_2, F_3], F_{\text{fused}} = W_f * F_{\text{concat}} \quad (6)$$

244 Here,  $[\cdot]$  denotes the concatenation operation along the channel  
 245 dimension, and  $W_f$  is a  $1 \times 1$  convolution kernel used to restore  
 the number of channels to  $c$ .

246 Simultaneously, an independent gating pathway utilizes  
 247 Global Average Pooling (GAP) to generate channel-wise statis-  
 248 tical summaries and produces modulation weights through a  
 249 two-layer fully connected network:

$$250 g = \sigma(W_2(\delta(W_1(\text{GAP}(\hat{X})))) \quad (7)$$

251 where  $\delta$  is the ReLU activation function,  $W_1$  and  $W_2$  are the  
 252 weights of the fully connected layers, and  $\sigma$  is the Sigmoid  
 253 function. The final output of the module is a weighted residual  
 sum of the input feature and the gated weighted texture features:

$$254 X_{\text{textured}} = \hat{X} + \alpha \cdot (g \odot F_{\text{fused}}) \quad (8)$$

255  $\alpha$  is a trainable scaling coefficient that controls the contribu-  
 256 tion strength of the texture features. This design injects en-  
 257 hanced multi-scale texture representations while preserving the  
 integrity of the original information.

### 279 3.4. Cascaded Canopy Attention

280 Similar to Xu & Wan (2024) and Zhang et al. (2024), we con-  
 281 struct a Cascaded Canopy Attention (CCA) mechanism composed  
 282 of two cascaded attention mechanisms to enable the model to  
 autonomously focus on discriminative regions.

283 1) *TimberTokenFusion*. This component achieves spatial and  
 284 channel-aware feature modulation through convolution and gat-  
 285 ing. The input features are first projected into Query (Q), Key  
 286 (K), and Value (V) triplets. Subsequently, Q and K are modu-  
 287 lated via spatial and channel operations respectively (see Fig. 8)  
 288 and fused via Depthwise Convolution (DWConv):

$$289 \text{TTF}(X) = \gamma \cdot \text{DWConv}(M_{\text{spa}}(Q) + M_{\text{cha}}(K)) \odot V + X \quad (9)$$

290  $M_{\text{spa}}$  and  $M_{\text{cha}}$  represent Spatial and Channel modulation.  $\gamma$   
 291 is a learnable scalar parameter, and  $\odot$  denotes element-wise mul-  
 292 tiplication. This operation allows the model to efficiently inte-  
 grate spatial context and recalibrate channel responses.

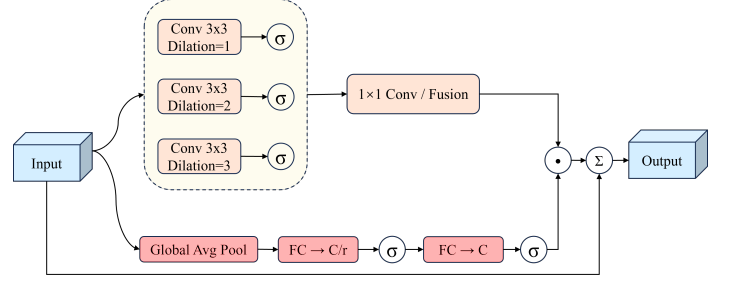


Figure 7: CanopyTextureAnalyzer

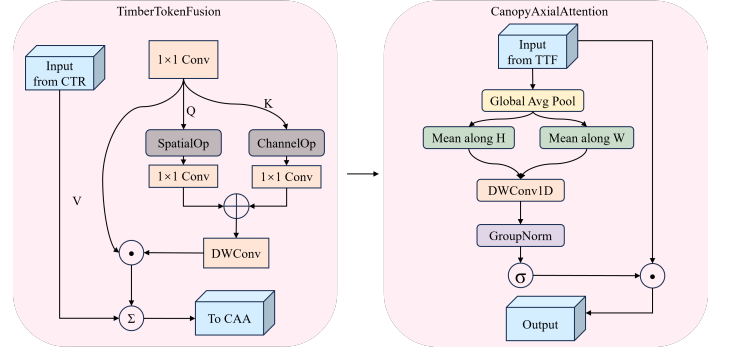


Figure 8: Cascaded Canopy Attention

2) *CanopyAxialAttention*. This component computes attention  
 weights along the image height and width directions separately  
 to establish long-range dependencies. First, Global Average  
 Pooling is applied to obtain global context:

$$273 g = \text{GAP}(X) \in \mathbb{R}^{c \times 1 \times 1} \quad (10)$$

Then, directional features are extracted with background sup-  
 274 pression:

$$275 x_h = \text{Mean}_w(X \odot g) \in \mathbb{R}^{c \times h}, \quad x_w = \text{Mean}_h(X \odot g) \in \mathbb{R}^{c \times w} \quad (11)$$

276 Shared-weight 1D depthwise convolution and Group Normal-  
 277 ization are applied:

$$278 a_h(x) = \sigma(\text{GroupNorm}(f_{\text{dw1d}}(x_h))) \quad (12)$$

$$279 a_w(x) = \sigma(\text{GroupNorm}(f_{\text{dw1d}}(x_w))) \quad (13)$$

The final output is the product of the attention maps from the  
 two directions:

$$280 \text{CAA}(X) = a_h(x) \cdot a_w(x) \cdot X \quad (14)$$

This cascaded design enables the model to sequentially filter  
 and reinforce features from different dimensions, thereby more  
 accurately locating key regions.

### 3.5. CanopyPatchExtractor

Inspired by the image patching strategy in Vision Transform-  
 281 ers (Dosovitskiy et al., 2020; Touvron et al., 2021), we design

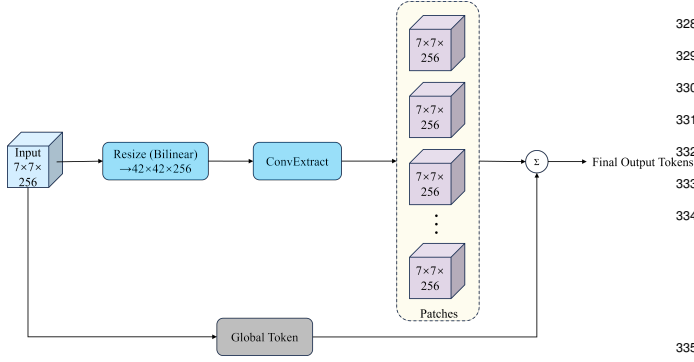


Figure 9: CanopyPatchExtractor

the CanopyPatchExtractor (CPE) module to structurally extract a series of overlapping local region tokens from the optimized feature map (see Fig. 9).

First, the feature map is upscaled to a preset resolution using an interpolation method to retain richer spatial details. A convolutional sliding-window operation with kernel size  $k$  then extracts multiple overlapping patches from the upscaled map, where the number of output channels corresponds to the desired patch count  $N$ . These patches are reorganized into a tensor and resized to a fixed spatial dimension  $p \times p$  via bilinear interpolation, yielding  $P \in \mathbb{R}^{B \times N \times d \times p \times p}$ .

To preserve global context, we augment the patch sequence with an additional global token derived from global average pooling of the original feature map, resulting in  $N+1$  tokens for subsequent processing. This design enables the model to capture both local discriminative details and holistic scene statistics, effectively mitigating environment-induced intra-class variance while enhancing fine-grained inter-species cues.

### 3.6. Feature Encoding and Classification

In the final stage of our methodology, region tokens are transformed into discriminative representations for fine-grained classification. The Position-aware Patch Attention module employs sequential modeling to capture global semantic dependencies, utilizes Transformer-based pooling for feature compression, and incorporates relative positional encoding to establish spatial relationships, forming a hierarchical processing mechanism that maintains structural awareness while enabling comprehensive contextual reasoning (Behera et al., 2021). The classification module adopts bilinear interaction modeling to aggregate features, fusing second-order statistics between semantic distributions and token characteristics to generate compact global representations, which are subsequently projected to produce final classification outputs.

### 3.7. Evaluation Metrics

To comprehensively and objectively evaluate the performance of the proposed model on the fine-grained tree species classification task, this study employs a set of widely accepted evaluation metrics. Overall Accuracy (OA) is adopted as the

primary and core metric for assessing model performance. Furthermore, to enable an in-depth analysis of the classification effectiveness and stability across different categories, supplementary metrics including Precision, Recall, and F1-Score are also utilized. OA directly measures the overall proportion of correctly classified samples out of the total predictions. The calculation formula is as follows:

$$OA = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \quad (15)$$

This metric reflects the overall classification correctness of the model and is the most intuitive and core performance measure for classification tasks. Precision is used to evaluate the accuracy of the model's positive predictions, i.e., the proportion of predicted positive samples that are truly positive. The calculation formula is:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (16)$$

where  $TP$  (True Positive) denotes the number of correctly predicted positive samples, and  $FP$  (False Positive) denotes the number of negative samples incorrectly predicted as positive. A high Precision indicates that the model is highly reliable when making positive predictions. Recall is used to assess the model's coverage ability for positive samples, i.e., the proportion of actual positive samples that are correctly predicted by the model. The calculation formula is:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (17)$$

where  $FN$  (False Negative) denotes the number of true positive samples incorrectly predicted as negative. A high Recall indicates that the model can effectively identify the majority of the target samples. F1-Score is the harmonic mean of Precision and Recall, used to comprehensively reflect the overall performance of the model. It is particularly important in fine-grained classification tasks where the sample distribution across classes may be imbalanced. The calculation formula is:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (18)$$

The overall performance values for the above metrics are obtained by calculating the Macro-average across all classes. This ensures an equal contribution from each category to the final metric, prevents dominant classes from skewing the evaluation results, and better aligns with the characteristics of fine-grained classification tasks.

### 3.8. Experiments Settings

We utilized a pre-trained Xception network as the backbone feature extractor. The input image resolution was uniformly adjusted to  $42 \times 42$  pixels. The number of feature channels was compressed to 256, and the number of regional patches was set to 32. Regarding training hyperparameters: the batch size was set to 8; the SGD optimizer (with a momentum of 0.99) and cross-entropy loss function were employed; the initial learning

371 rate was set to  $1e-4$  and dynamically decayed using a Lambda<sub>407</sub>  
 372 scheduler; the model was trained for a total of 200 epochs; and<sub>408</sub>  
 373 the random seed was fixed to ensure reproducibility. During<sub>409</sub>  
 374 training, the latest model was saved every 10 epochs, and the<sub>410</sub>  
 375 best model was automatically saved based on the highest vali-<sub>411</sub>  
 376 dation accuracy.

## 377 4. Results

### 378 4.1. Performance Comparison

379 To comprehensively evaluate the overall performance of<sub>418</sub>  
 380 FFGF-Net, we conducted extensive experiments on three chal-<sub>419</sub>  
 381 lenging public datasets for fine-grained tree species classifica-<sub>420</sub>  
 382 tion: SZUTreeData, NJFUTreeData, and ETH dataset. These<sub>421</sub>  
 383 datasets vary in scale, image provenance, and habitat complex-<sub>422</sub>  
 384 ity, constituting a comprehensive testing platform. We selected<sub>423</sub>  
 385 representative advanced methods as baseline comparisons, in-<sub>424</sub>  
 386 cluding generic deep network backbones (e.g., EfficientNet-<sub>425</sub>  
 387 B0, ConvNeXt-Tiny) and recently published state-of-the-art<sub>426</sub>  
 388 (SOTA) methods. As shown in Table 1, compared to current<sub>427</sub>  
 389 state-of-the-art methods, the FFGF-Net demonstrates consis-<sub>428</sub>  
 390 tent and significant performance improvements across all eval-<sub>429</sub>  
 391 uation metrics. This validates that our proposed architectural<sub>430</sub>  
 392 design, which combines the CanopyTextureAnalyzer and Cas-<sub>431</sub>  
 393 caded Canopy Attention mechanisms, is more effective than ex-<sub>432</sub>  
 394 isting techniques when handling the highly similar features be-<sub>433</sub>  
 395 tween tree species.

396 To provide a visual comparison of the performance trends<sub>435</sub>  
 397 across forestry datasets, we plot the overall accuracy of each<sub>436</sub>  
 398 method on NJFUTreeData, SZUTreeData, and ETH dataset<sub>437</sub>  
 399 in Fig. 10. The line graph clearly illustrates that FFGF-<sub>438</sub>  
 400 Net achieves the highest accuracy on SZUTreeData and ETH<sub>439</sub>  
 401 dataset, with a notable lead over other methods, especially on<sub>440</sub>  
 402 ETH dataset where it reaches 94.75%. This graphical represen-<sub>441</sub>  
 403 tation reinforces the superiority of FFGF-Net in forestry fine-<sub>442</sub>  
 404 grained classification.

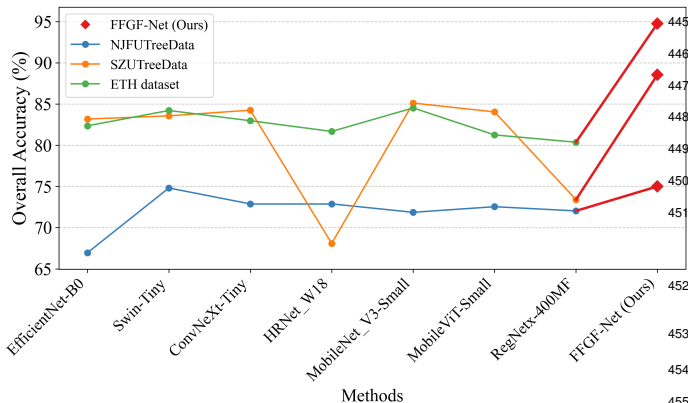


Figure 10: Performance comparison across forestry datasets

405 To evaluate the broad applicability of our method, FFGF-Net<sub>460</sub>  
 406 was evaluated against recent specialized models on standard<sub>461</sub>

412 fine-grained benchmarks including CUB-200-2011 (CUB),  
 413 FGVC-Aircraft (Aircraft), Stanford Cars (Cars), and Stanford  
 414 Dogs (Dogs). These models include Context-Aware Attentional  
 415 Pooling (CAP) (Behera et al., 2021), Counterfactual Attention  
 416 Learning (CAL) (Rao et al., 2021), Feature Fusion Vision  
 417 Transformer (FFVT) (Wang et al., 2021), Spatial Relation-  
 418 Aware Graph Neural Network (SR-GNN) (Bera et al., 2022)),  
 419 MetaFormer (Diao et al., 2022), and Inter and Intra Region  
 420 High-Order Feature Interaction (I2-HOFI) (Sikdar et al., 2025).  
 421 As shown in Table 2, FFGF-Net delivers favorable performance  
 422 in all evaluated domains, attaining state-of-the-art results on  
 423 CUB and maintaining robust accuracy across the other datasets.  
 424 It should be acknowledged, however, that despite its consis-  
 425 tent outcomes on general fine-grained benchmarks, the method  
 426 does not achieve top performance on certain generic tasks, sug-  
 427 gesting potential for further optimization in broader applica-  
 428 tions. These findings nonetheless validate that the integration  
 429 of multi-scale texture modeling with a cascaded attention archi-  
 430 tecture yields a generalizable and efficient strategy, which is  
 431 primarily designed for forestry contexts while retaining adapt-  
 432 ability to various fine-grained visual recognition challenges.

433 In addition to accuracy metrics, we further evaluated the  
 434 computational efficiency of the compared models. Although  
 435 generic methods provide powerful classification performance,  
 436 they often come with higher computational costs. In contrast,  
 437 our proposed method achieves highly competitive accuracy on  
 438 all datasets, while significantly reducing the required param-  
 439 eter count and computational cost. The significant reduction  
 440 in computational complexity means smaller memory usage and  
 441 shorter theoretical inference time. Consequently, our method  
 442 maintains a high balance between classification performance  
 443 and model efficiency, demonstrating its strong feasibility in  
 444 practical deployment.

445 Fig. 15 presents four line graphs comparing FFGF-Net with  
 446 other state-of-the-art methods on general fine-grained bench-  
 447 marks (CUB, Aircraft, Cars, Dogs). As illustrated, FFGF-Net  
 448 achieves leading performance on CUB while maintaining con-  
 449 sistent results across diverse tasks, further validating the ratio-  
 450 nality and practicality of our method’s design in forestry im-  
 451 age recognition. In this comparison, we adopted a strategy  
 452 of directly citing literature results to ensure that all compari-  
 453 son models were in their optimal state after being optimized by  
 454 their original authors. Although this results in missing metrics  
 455 not mentioned in some literature, it provides a broader macro  
 456 performance reference for our method.

## 457 4.2. Visualization Analysis

### 458 4.2.1. Confusion matrix analysis

459 Through an examination of the confusion matrix, the classifi-  
 460 cation behavior of the FFGF-Net is elucidated for eight charac-  
 461 teristic subtropical tree species in the NJFUTreeData collection  
 462 from East China (refer to Fig. 11, which displays a representa-  
 463 tive subset of the species for visual clarity). The network effec-  
 464 tively identifies species with distinctive morphological features,  
 465 as demonstrated by the strong diagonal distribution observed  
 466 in *Ginkgo biloba* and *Cinnamomum camphora*, indicating suc-

Table 1: FFGF-Net vs. SOTA in forestry fine-grained classification. The proposed Forestry Fine-Grained Fusion Network (FFGF-Net) is compared against various baseline models. Note: OA = Overall Accuracy. GFLOPs stands for Giga Floating-point Operations Per Second, Param. indicates the number of model parameters in millions, and IPS indicates Images Per Second. To provide a rigorously fair, head-to-head evaluation, all models in this table were locally re-implemented and evaluated by us under strictly identical training and testing conditions.

Methods	OA (%)			Precision (%)			Recall (%)			F1-score (%)			GFLOPs	Param. (M)	IPS
	NJFU	SZU	ETH	NJFU	SZU	ETH	NJFU	SZU	ETH	NJFU	SZU	ETH			
EfficientNet-B0 (2019)	66.95	83.17	82.35	59.61	73.67	76.68	57.31	68.59	73.30	57.86	71.73	74.08	0.4	5.3	742
Swin-Tiny (2021)	74.81	83.56	84.21	58.13	76.23	81.09	56.38	70.84	75.11	56.78	73.14	77.67	2.9	19.6	374
ConvNeXt-Tiny (2022)	72.88	84.25	82.97	54.35	77.29	75.64	50.82	72.80	74.58	52.26	74.83	74.90	4.5	28.6	451
HRNet-W18 (2019)	72.88	68.09	81.67	53.73	50.37	78.06	54.47	51.25	72.31	53.83	51.09	74.69	4.4	21.3	134
MobileNetV3-Small (2019)	71.86	85.11	84.52	56.52	63.32	82.26	54.40	60.68	77.42	54.59	61.34	79.63	0.1	1.5	1116
MobileViT-Small (2021)	72.54	84.04	81.25	53.64	70.81	77.83	52.68	68.26	75.19	52.65	69.53	76.87	1.4	5.6	570
RegNetx-400MF (2020)	72.03	73.40	80.36	40.26	58.14	71.47	49.68	54.99	70.64	44.37	56.53	70.52	4.0	22.0	338
FFGF-Net (Ours)	<b>75.00</b>	<b>88.54</b>	<b>94.75</b>	<b>67.05</b>	<b>81.25</b>	<b>89.24</b>	<b>62.25</b>	<b>74.08</b>	<b>83.61</b>	<b>64.30</b>	<b>77.26</b>	<b>85.86</b>	4.7	22.7	260

Table 2: FFGF-Net vs. SOTA in fine-grained visual classification. Note: To accurately reflect the optimal capabilities of these specialized architectures, baseline results are directly cited from their respective original publications. As underlying experimental conditions and hardware environments may vary across these independent studies, this table serves as a broader literature reference rather than a strictly controlled evaluation. Consequently, missing data points indicate specific metrics that were unavailable in the original papers.

Methods	CUB	Aircraft	Cars	Dogs	GFLOPs	Param. (M)
	OA (%)	OA (%)	OA (%)	OA (%)		
CAP (Behera et al. 2021)	91.8	95.7	94.1	96.1	10.2	34.2
generic CAL (Rao et al. 2021)	90.6	94.2	95.5	-	40.0	55.7
FFVT (Wang, Yu, and Gao 2021)	91.6	-	-	91.5	16.3	87.5
SR-GNN (Bera et al. 2022)	91.9	95.4	96.1	<b>97.3</b>	9.8	30.9
MetaFormer (Diao et al. 2022)	91.8	-	95.1	-	16.9	81.0
I2-HOFI (Sikdar et al. 2024)	91.6	<b>96.0</b>	<b>96.4</b>	-	14.2	<b>22.6</b>
FFGF-Net (Ours)	<b>92.4</b>	95.1	95.9	96.7	<b>4.7</b>	22.7

462 successful learning of their discriminative leaf patterns and canopy  
463 structures.

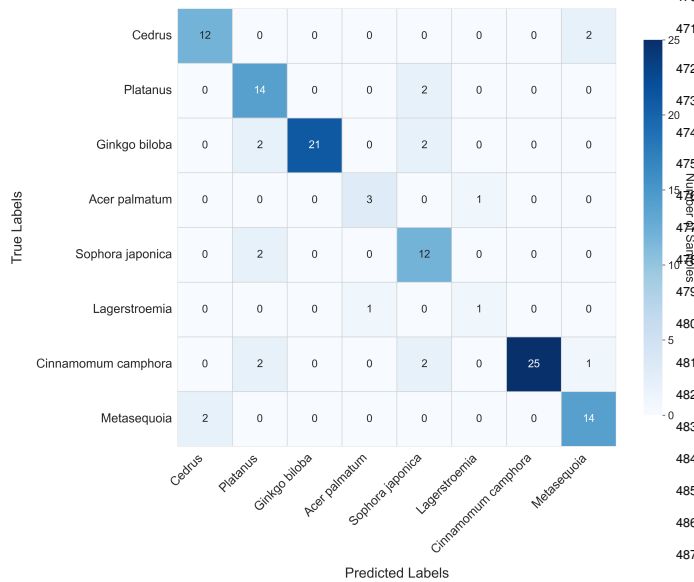


Figure 11: Confusion matrix of FFGF-Net tree species classification on the NJFUTreeData

464 The analysis identifies mutual misclassification between  
465 coniferous species *Cedrus* and *Metasequoia glyptostroboides*,  
466 highlighting the classification challenges posed by similar  
467 structural characteristics in needle-leaved trees within UAV im-

468 agery. Broadleaf species *Platanus* and *Sophora japonica* also  
469 exhibit reciprocal misclassification, revealing limitations in dis-  
470 tinguishing deciduous species with comparable canopy tex-  
471 tures.

472 Notably, the confusion between ornamental species *Acer*  
473 *palmatum* and *Lagerstroemia indica* is exacerbated by limited  
474 sample availability. *Cinnamomum camphora*, despite having  
475 the largest sample size, shows misclassification across multiple  
476 categories, suggesting either significant intra-species morpho-  
477 logical variation or feature overlap with other broadleaf species.  
478 These systematic error patterns provide clear directions for  
479 model optimization, particularly in enhancing the recognition  
480 of subtle features among morphologically similar tree species.

481 As shown in Fig. 12, the confusion matrix analysis was con-  
482 ducted on a selected, representative subset of tree species from  
483 the SZUTreeData. The results demonstrate that the model ex-  
484 hibits strong and consistent discriminative capabilities for most  
485 of the selected species. For instance, all test samples of *Firmi-*  
486 *ana simplex*, *Litchi chinensis*, *Araucaria cunninghamii*, *Livis-*  
487 *tona chinensis*, and *Cinnamomum camphora* tree were classif-  
488 iced correctly without any misidentification, underscoring the  
489 highly discriminative representations learned by the model for  
490 these categories.

491 However, the matrix also reveals subtle confusion between  
492 certain visually similar species pairs. Specifically, mutual mis-  
493 classification occurs between *Flcus spp.* and *Flcus microcarpa*,  
494 and one sample of *Flcus microcarpa* is incorrectly predicted as  
495 *Livistona chinensis*. Furthermore, one misclassification case is  
496 observed between *Melaleuca cajuputi* and *Cinnamomum cam-*

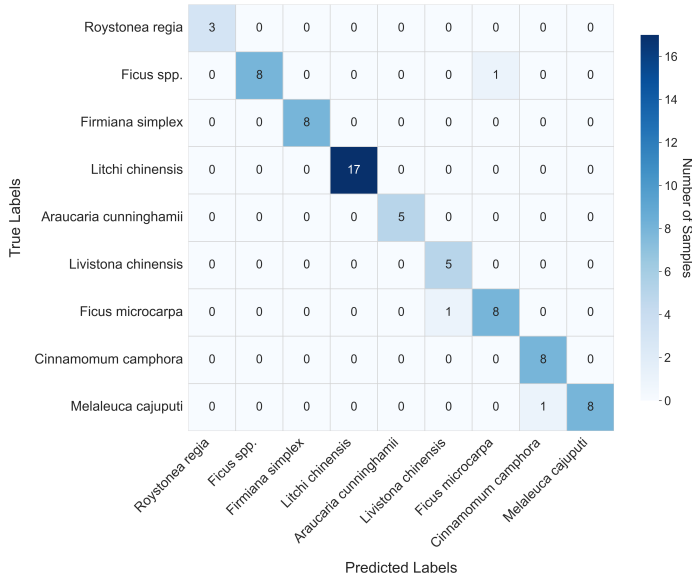


Figure 12: Confusion matrix of FFGF-Net tree species classification on the SZUTreeData



Figure 13: Confusion matrix of FFGF-Net tree species classification on the ETH dataset

phora. These errors likely stem from the intrinsic similarities in fine-grained features such as canopy texture and structure among these species, compounded by environmental noise, common in UAV imagery, such as occlusion and varying illumination conditions, which exacerbates the challenge of fine-grained distinction.

Although the analysis is confined to a representative subset, the results robustly indicate that the FFGF-Net possesses commendable classification consistency and robustness within complex forestry environments.

The confusion matrix analysis on the ETH dataset reveals the exceptional performance and specific challenges of the FFGF-Net in fine-grained classification of temperate tree species (see Fig. 13). Due to the large number of categories, only a representative subset of the tree species is visualized in the matrix to ensure readability. The network achieved a leading OA of 94.75% on the test set, demonstrating its powerful overall discriminative capability. The analysis indicates that the model exhibits near-perfect distinguishability for most species. Nearly all samples of *Betula papyrifera*, *Tsuga canadensis*, and *Picea abies* were classified correctly, with only minimal misidentifications. This confirms that the canopy features extracted by the FFGF-Net are highly robust for differentiating these common, yet morphologically diverse, species in temperate forests.

Nevertheless, the matrix captures rare misclassification cases, which are highly representative and concentrated on rare species with extremely limited samples. For instance, all 6 samples of *Pinus montezumae* were incorrectly predicted as *Abies alba*; conversely, one *Abies alba* sample was also misclassified as *Pinus montezumae*. This indicates the model’s difficulty in learning discriminative features for very uncommon species. This confusion likely stems from a “long-tail distribution” problem, where insufficient training samples for rare cate-

gories prevent the model from adequately learning their unique fine-grained characteristics.

In conclusion, the FFGF-Net achieves state-of-the-art classification performance on the ETH dataset, with particularly accurate and stable identification of dominant tree species. Simultaneously, the results clearly highlight the current model’s limitation in handling extremely imbalanced, rare species. This provides a clear direction for future work, suggesting the integration of targeted re-sampling strategies or few-shot learning techniques to enhance the recognition capability for rare species.

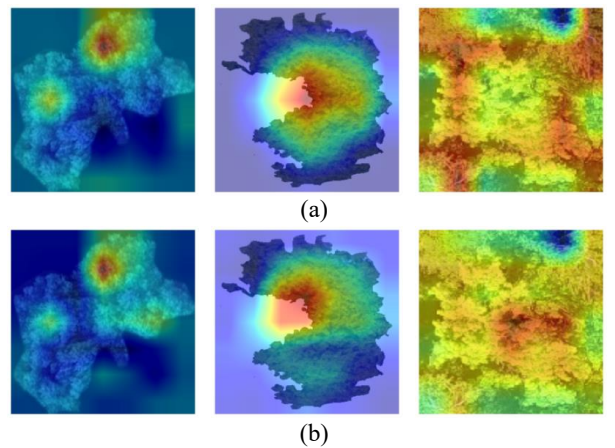


Figure 14: Cascaded Canopy Attention (CCA) Module Attention Visualization. (a) TTF Attention Responses. (b) CAA Attention Responses.

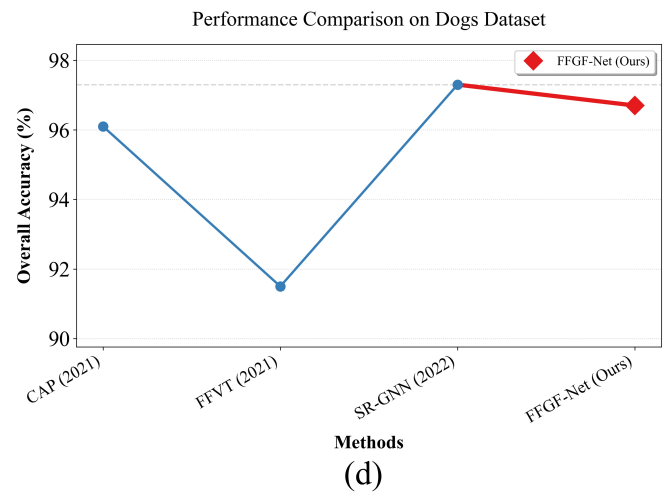
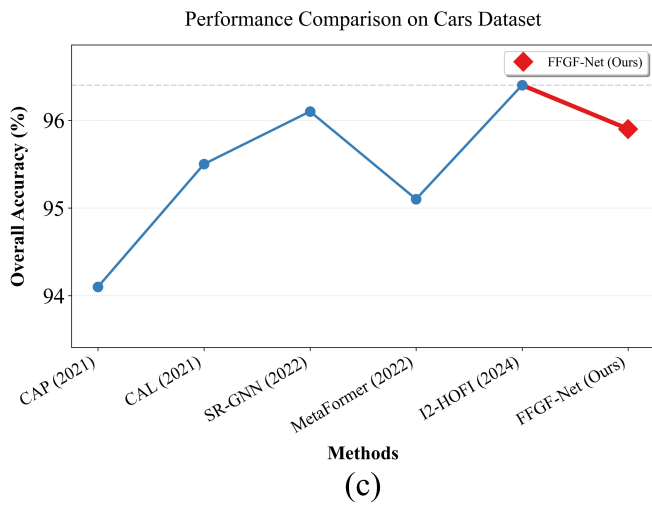
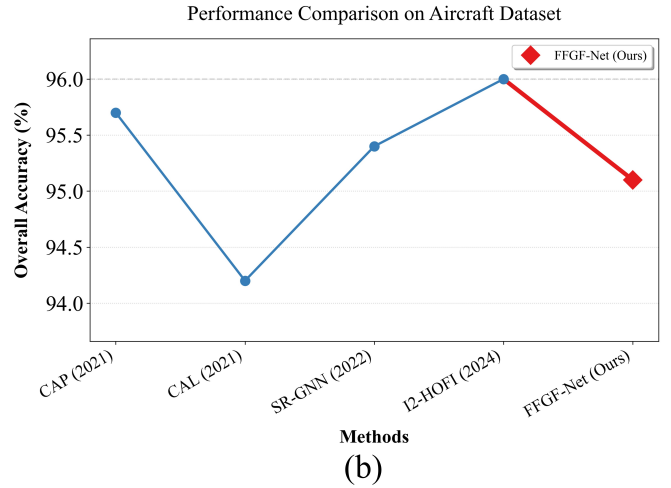
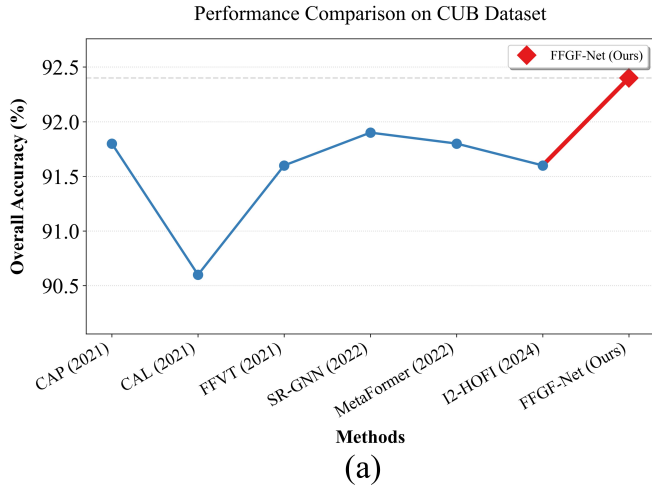


Figure 15: Performance comparison across general datasets. Note: Missing data points for certain baseline algorithms on specific datasets indicate that those performance metrics were not reported in their original publications. To ensure a fair and standardized evaluation, all baseline results presented here are directly cited from their respective original papers.

#### 4.2.2. Module response visualization

To empirically validate the capability of the Cascaded Canopy Attention (CCA) module in enhancing discriminative regions and suppressing noise, we performed a systematic visualization of the attention heatmaps generated by its two core sub-modules: TimberTokenFusion (TTF) and CanopyAxialAttention (CAA). The results are presented in Fig. 14.

An analysis of the spatial distribution characteristics within the attention heatmaps reveals that TTF and CAA exhibit complementary and cascading modulation effects. The TTF module, employing spatial-channel joint modulation, primarily amplifies local textural regions inside the canopy that demonstrate high discriminative power. Subsequently, the CAA module refines the global structural coherence of the attention distribution through axial attention computation along the height (H) and width (W) dimensions. By integrating TTF’s local discriminative enhancement with CAA’s global structural regularization, the CCA module establishes a progressive attention refinement

mechanism. This architecture provides essential support for fine-grained classification within complex forestry environments.

#### 4.2.3. Comparison of t-SNE visualization

To quantitatively analyze the contribution of each module in the FFGF-Net to the enhancement of feature discriminability, this study employed the t-SNE (t-distributed Stochastic Neighbor Embedding) algorithm to project high-dimensional features into a two-dimensional space. A comparison was conducted between the feature distributions of the “Xception Backbone Only” and the “Full FFGF-Net”, as illustrated in Fig. 16.

The visualization results clearly demonstrate that the features learned by the full FFGF-Net exhibit more distinct intra-class clustering and inter-class separation compared to those extracted using only the backbone network. Specifically, when solely relying on the backbone, feature points of different tree species show substantial overlap with blurred boundaries, indicating limited discriminative power. In contrast, after process-

ing with FFGF-Net, samples belonging to the same class form tighter clusters, while the margins between different categories are markedly enlarged. These findings empirically confirm the introduced mechanisms, including local texture modeling, attention modules, and feature fusion strategies, collectively contribute to enhancing the model’s ability to capture subtle inter-class distinctions. The FFGF-Net effectively improves the discriminative quality of the learned features, thereby providing a more reliable representation for high-accuracy fine-grained classification.

#### 4.2.4. Comparative performance visualization

As illustrated in Fig. 17, this study conducts comprehensive visual comparisons with three baseline methods on the SZUTreeData dataset, which comprises 25 tree species characterized by complex canopy structures and diverse morphological features. The visualization scheme employs a standardized color coding: ground truth labels are displayed in blue, correct predictions in green, and incorrect predictions in red. Representative samples were randomly selected for analysis, with each sample including both the original canopy image and an enlarged region of interest (ROI) that highlights texture details crucial for species discrimination. The experimental results demonstrate that FFGF-Net exhibits superior performance in recognizing species with distinctive morphological signatures and maintains high confidence scores for correct predictions, indicating robust feature representation learning. Furthermore, texture ROI analysis confirms that FFGF-Net effectively focuses on discriminative canopy regions by leveraging both global structural patterns and local texture details.

However, the analysis reveals persistent challenges in handling taxonomically complex scenarios. A notable limitation emerges in distinguishing congeneric species with high visual similarity, such as *Ficus virens* and *Ficus microcarpa* within the Moraceae family. These closely related species share numerous morphological traits inherited from common ancestors, causing models to preferentially learn these dominant shared characteristics while overlooking discriminative local features (e.g., leaf size, apex morphology) when training samples are insufficient for fine-grained differentiation. This phenomenon leads to “intra-genus confusion” in classification.

A more complex challenge arises from convergent evolution, where phylogenetically distant species develop similar morphological structures due to adaptation to analogous environments. For instance, *Terminalia arjuna* (Combretaceae) and *Cinnamomum camphora* (Lauraceae) both exhibit “full and regular” broad-ovate crown structures as a convergent adaptation to ecological niches such as light competition. During training, models tend to prioritize capturing these macroscopic, dominant “ecotype” features while struggling to focus on taxonomically significant subtle distinctions like bark texture and phyllotaxy, resulting in “inter-family confusion”.

The underlying mechanism for these misclassifications can be attributed to the FFGF-Net’s feature learning priorities. Many critical discriminative features belong to fine-grained or local characteristics. When the backbone network’s receptive field is either too large or too small, it may fail to effectively

capture these local high-resolution patterns, instead over-relying on macroscopic features such as overall crown shape. This fundamental limitation explains the systematic misclassification of different species that share similar macroscopic morphological structures.

These challenging cases highlight promising directions for future research in developing more sophisticated feature representation techniques that can better handle fine-grained botanical classification, particularly through enhanced local feature extraction and multi-scale representation learning.

#### 4.3. Ablation Study

This section presents an ablation study (see Table 3), evaluating the contributions of the CanopyTextureAnalyzer (CTA), Cascaded Canopy Attention (CCA), CanopyPatchExtractor (CPE), and Bilinear Attention Pooling (BAP) modules. The baseline model (utilizing only BAP) achieved the lowest accuracy. Incorporating CTA alone, which extracts multi-scale textures, yielded a marginal performance gain. The combination of CTA and CPE (without BAP) led to improved performance on the SZU dataset, highlighting CPE’s capability in discerning local discrepancies. The CCA+BAP combination further enhanced accuracy, demonstrating CCA’s efficacy in suppressing noise inherent in field environments. The full configuration of our model achieved superior performance on the SZUTreeData.

To rigorously verify the effectiveness of the proposed modules, we conducted robustness analysis on the ablation results. Specifically, all models were evaluated in three independent runs, each using a different random seed. The results demonstrate that the performance improvements yielded by our proposed components significantly exceed the standard deviations. This variability measurement confirms that the observed improvements in accuracy are robust and empirically significant, rather than just a result of random variance.

Table 3: FFGF-Net ablation study. To ensure statistical rigor, all models were evaluated over three independent runs with different random seeds. Results are reported as Mean  $\pm$  Std.

CTA	CCA	CPE	BAP	NJFU	SZU	ETH
			✓	70.42 $\pm$ 0.14	81.25 $\pm$ 0.11	85.36 $\pm$ 0.12
✓			✓	71.83 $\pm$ 0.10	82.29 $\pm$ 0.15	86.91 $\pm$ 0.13
✓		✓		72.54 $\pm$ 0.14	84.38 $\pm$ 0.11	88.25 $\pm$ 0.16
		✓	✓	73.16 $\pm$ 0.13	84.38 $\pm$ 0.13	89.47 $\pm$ 0.12
	✓		✓	73.89 $\pm$ 0.15	85.42 $\pm$ 0.15	91.32 $\pm$ 0.13
✓	✓	✓	✓	<b>75.00<math>\pm</math>0.14</b>	<b>88.54<math>\pm</math>0.17</b>	<b>94.75<math>\pm</math>0.14</b>

## 5. Discussion

### 5.1. Discussion

The proposed FFGF-Net was developed to overcome the inherent challenges of fine-grained tree species classification using UAV imagery, contributing a novel methodological approach to the broader field of Earth Observation (EO). Our framework shifts the paradigm towards domain-specific structural modeling. Experimental evaluations confirm that domain-specific feature extraction modules, namely the CanopyPatchExtractor and CanopyTextureAnalyzer, achieve substantial performance improvements over standard CNN and Vision

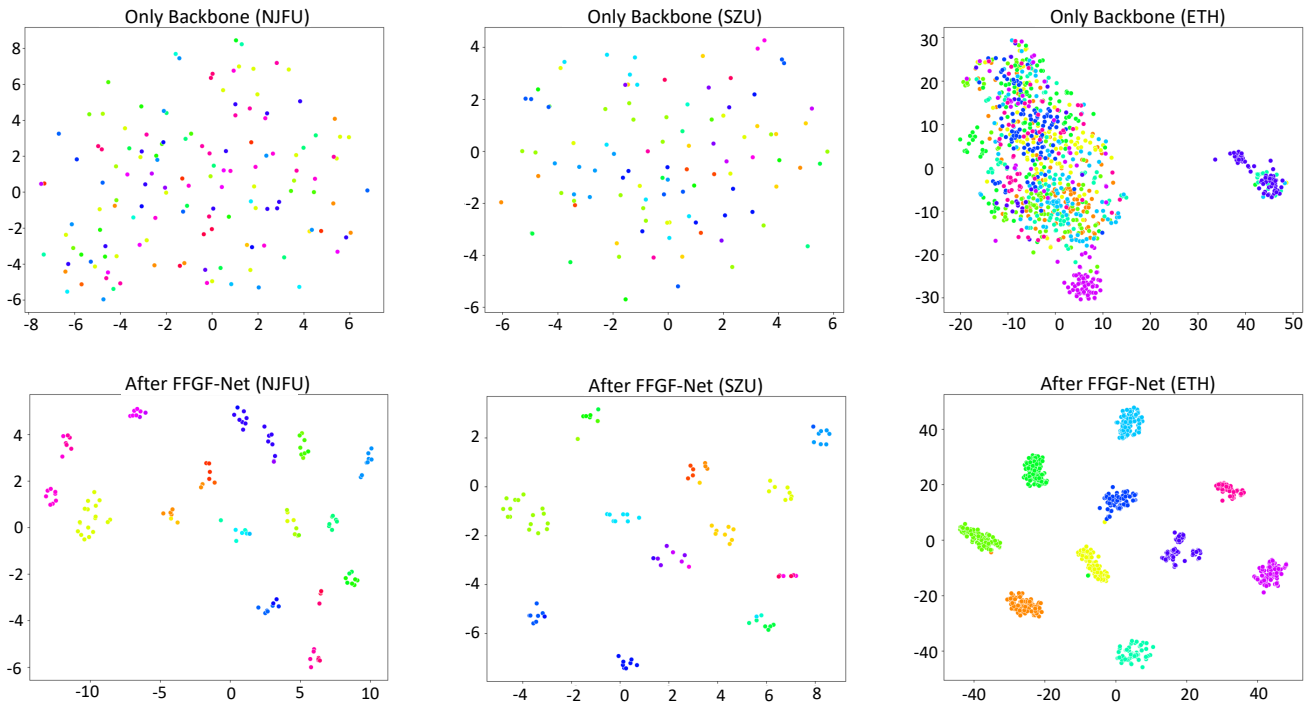


Figure 16: t-SNE visualization of the feature embedding space

677 Transformer baselines. By isolating subtle discriminative char-703  
 678 acteristics from complex canopy backgrounds and explicitly704  
 679 modeling multi-scale textures, this framework effectively dif-705  
 680 ferentiates morphologically similar and closely related species.706

681 Regarding operational deployment and existing technolo-707  
 682 gies, recent state-of-the-art forestry classification research in-708  
 683 creasingly relies on multi-modal fusion strategies (e.g., combin-709  
 684 ing LiDAR point clouds with optical imagery) to capture highly710  
 685 discriminative 3D canopy structures (Ferreira et al., 2024). As711  
 686 highlighted by Ali et al. (2024) and Wild et al. (2026), the in-712  
 687 tegration of such 3D structural features offers unparalleled ad-713  
 688 vantages in vegetation characterization and biomass estimation714  
 689 over purely image-based methods. Although these fusion meth-715  
 690 ods achieve exceptional accuracy in multi-layered canopies,716  
 691 they inherently suffer from prohibitive data acquisition costs,717  
 692 the need for bulky specialized UAV payloads, and highly com-718  
 693 plex, computationally intensive preprocessing pipelines. As a719  
 694 purely image-based framework, our proposed FFGF-Net pro-720  
 695 vides a highly accessible and cost-effective alternative for large-721  
 696 scale surveys. By directly extracting crucial spectral and tex-722  
 697 tural features from high-resolution UAV imagery, our method723  
 698 maximizes the potential of two-dimensional visual characteris-724  
 699 tics to document species-specific phenotypic traits. This image-725  
 700 centric solution maintains practical feasibility while offering an726  
 701 efficient technical pathway for routine forest resource monitor-727  
 702 ing (Sothe et al., 2019; Wallace et al., 2020).728

Beyond classification accuracy, the practical application value of FFGF-Net lies in its computational efficiency. Forestry inventory applications require rapid processing of large orthomosaic datasets on resource-limited devices. FFGF-Net maintains a highly favorable balance between performance and overhead, requiring only 4.7 GFLOPs and 22.7 million parameters. This lightweight architecture enables deployment on edge platforms like UAV onboard processors or mobile field tablets. An inference speed of 260 images per second delivers near real-time processing. This high-throughput efficiency overcomes the bottleneck of manual visual interpretation, positioning FFGF-Net as a deployable tool for automated large-area forest management.

Crucially, the extensibility of the FFGF-Net framework warrants emphasis. While currently optimized and validated on high-resolution UAV RGB datasets, the modular architecture holds significant potential for integrating emerging remote sensing data sources. For instance, coupling this texture-aware framework with hyperspectral imagery or sub-meter commercial satellite data could effectively address existing spectral data gaps, providing the critical physiological dimensions currently absent in standard RGB sensors. In addition, the core mechanism of localizing and analyzing fine-grained structural textures demonstrates high potential for broader transferability. Beyond forestry inventories, this EO-centric framework could be adapted to support wider applications, such as agricultural

● Ground Truth ● Correct ● Incorrect

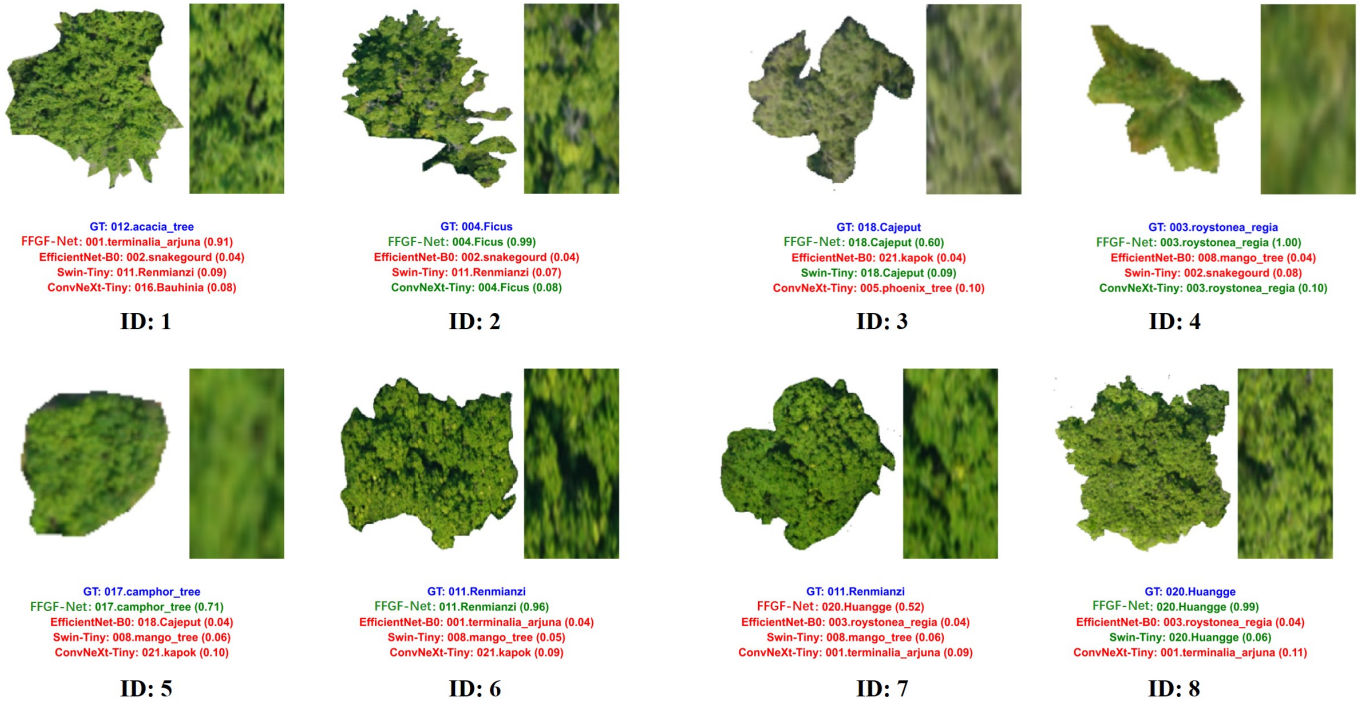


Figure 17: Comparative visualization of classification performance across different methods. In the visualization, blue text indicates the ground truth labels, while green and red texts represent correct and incorrect predictions, respectively.

yield estimation, urban green space phenotyping, and invasive plant species tracking across diverse geographical datasets.

## 5.2. Limitations

Despite the robust performance achieved by the FFGF-Net, several limitations remain, which also delineate directions for future research. Currently, as a purely image-based model, our framework’s primary drawback compared to multi-modal approaches is the absence of explicit three-dimensional spatial metrics. Consequently, in extremely dense forests with severe canopy overlap, the performance ceiling of FFGF-Net may still be constrained. Future research should explore lightweight multi-modal fusion strategies that combine our efficient 2D architecture with LiDAR point-cloud data to further enhance classification accuracy without entirely sacrificing deployment convenience.

Beyond structural constraints, a significant challenge is the prevalent "long-tail distribution" problem in forest environments, where numerous rare tree species are represented by extremely limited samples. Future research must focus on few-shot or even zero-shot learning paradigms to enhance the model’s capability for recognizing rare species.

Furthermore, the current model primarily relies on static imagery for identification and does not incorporate dynamic information related to tree canopy phenological changes. Developing temporal-aware models capable of integrating multi-temporal remote sensing data will be crucial for achieving more

stable classification and enabling early warning of pest and disease outbreaks.

Ultimately, the scope of the present study is restricted to taxonomic identification. Advancing toward comprehensive ecological monitoring necessitates extending the current architecture to assess physiological status. Developing inversion models to estimate biophysical parameters, directly from high-resolution UAV imagery represents a crucial trajectory for subsequent research.

## 6. Conclusion

This study focused on addressing the critical challenges of fine-grained tree species classification in forestry using computer vision techniques. Specifically, we aimed to overcome the inherent difficulties of crown-scale remote sensing imagery, including low feature discriminativity, high inter-class similarity, and complex noise in field environments.

To achieve this objective, we proposed a novel Forestry Fine-Grained Fusion Network (FFGF-Net) and constructed a new UAV-based dataset, NJFUTreeData. Our primary scientific contributions are twofold. First, the introduction of the NJFUTreeData, combined with public datasets such as SZUTreeData and ETH dataset, establishes a comprehensive multi-source evaluation framework, offering invaluable data assets for the research community. Second, the FFGF-Net effectively addresses the aforementioned challenges through its tailored ar-

chitecture. By integrating the CanopyPatchExtractor for local feature isolation, the CanopyTextureAnalyzer for multi-scale perception, and a Cascaded Canopy Attention mechanism for noise suppression, the framework significantly enhances the capacity to capture subtle discriminative features. Consequently, the model demonstrates robust and highly accurate performance in distinguishing visually similar tree species under challenging in-situ conditions.

The implications of this research are substantial from both theoretical and practical perspectives. Theoretically, the FFGF-Net offers a novel technical strategy for tackling common bottlenecks in fine-grained visual classification, shifting the paradigm towards domain-specific structural modeling. Practically, this purely image-based framework provides a reliable algorithmic foundation for high-precision, automated forest resource inventories and ecological monitoring. Ultimately, this research fosters the deep integration of computer vision technology within forestry applications, accelerating the digital transformation of intelligent forest management practices. Furthermore, a detailed analysis of real-world application scenarios, including the Hegyi competition index calculations for specific thinning strategies, is provided in the Supplementary Material.

## 7. Acknowledgments

This work was supported by the Key Research Projects of Yibin, Research and Integrated Demonstration of Key Technologies for Smart Bamboo Industry (YBZD2024-1), the Practice Innovation Training Program Projects for Jiangsu College Students under Grant 202510298056Z, the Biological Breeding-National Science and Technology Major Project (2023ZD0405605), the Shanghai Science and Technology Program (Grant No.25ZR1402133) and the China Postdoctoral Science Foundation (Grant No.GZC20250239).

## 8. Data Availability

To support the transparency of this research, a representative subset of the NJFUTreeData dataset and the core network implementation of the FFGF-Net are currently available at [https://github.com/zxz-xy/anonymous\\_code](https://github.com/zxz-xy/anonymous_code). The complete full-scale dataset, comprehensive training scripts, and pretrained models will be made fully publicly available upon the formal publication of this article.

## References

Sun, Y., Huang, J., Ao, Z., et al., 2019. Deep learning approaches for the mapping of tree species diversity in a tropical wetland using airborne LiDAR and high-spatial-resolution remote sensing images. *Forests* 10(11), 1047. <https://doi.org/10.3390/f10111047>

Poorter, L., McNeil, A., Hurtado, V. H., et al., 2014. Bark traits and life-history strategies of tropical dry- and moist forest

trees. *Funct. Ecol.* 28(1), 232-242. <https://doi.org/10.1111/1365-2435.12158>

Fang, F., McNeil, B.E., Warner, T.A., et al., 2020. Discriminating tree species at different taxonomic levels using multi-temporal WorldView-3 imagery in Washington DC, USA. *Remote Sens. Environ.* 246, 111811. <https://doi.org/10.1016/j.rse.2020.111811>

Mace, G.M., Norris, K., Fitter, A.H., 2012. Biodiversity and ecosystem services: A multilayered relationship. *Trends Ecol. Evol.* 27(1), 19-26. <https://doi.org/10.1016/j.tree.2011.08.006>

Zou, F., Li, A., Wang, Q., 2020. Comparative analysis of new forest resources survey methods. *J. Phys. Conf. Ser.* 1646(1), 012007. <https://doi.org/10.1088/1742-6596/1646/1/012007>

Zhang, X., Zhang, T., Wang, G., et al., 2023. Remote sensing object detection meets deep learning: A metareview of challenges and advances. *IEEE Geosci. Remote Sens. Mag.* 11(4), 8-44. <https://doi.org/10.1109/MGRS.2023.3312347>

Yu, F., Zhang, Q., Xiao, J., et al., 2023. Progress in the application of CNN-based image classification and recognition in whole crop growth cycles. *Remote Sens.* 15(12), 2988. <https://doi.org/10.3390/rs15122988>

Chen, M., Lin, M., Li, K., et al., 2023. CF-ViT: A general coarse-to-fine method for vision transformer. *Proc. AAAI Conf. Artif. Intell.* 37(6), 7042-7052. <https://doi.org/10.1609/aaai.v37i6.25860>

He, K., Zhang, X., Ren, S., et al., 2016. Deep residual learning for image recognition. *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pp. 770-778. [https://openaccess.thecvf.com/content\\_cvpr\\_2016/html/He\\_Deep\\_Residual\\_Learning\\_CVPR\\_2016\\_paper.html](https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html)

Dosovitskiy, A., Beyer, L., Kolesnikov, A., et al., 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. <https://doi.org/10.48550/arXiv.2010.11929>

Zhao, B., Feng, J., Wu, X., et al., 2017. A survey on deep learning-based fine-grained object classification and semantic segmentation. *Int. J. Autom. Comput.* 14(2), 119-135. <https://doi.org/10.1007/s11633-017-1053-3>

Fricker, G.A., Ventura, J.D., Wolf, J.A., et al., 2019. A convolutional neural network classifier identifies tree species in mixed-conifer forest from hyperspectral imagery. *Remote Sens.* 11(19), 2326. <https://doi.org/10.3390/rs11192326>

Neupane, B., Horanont, T., Aryal, J., 2021. Deep learning-based semantic segmentation of urban features in satellite images: A review and meta-analysis. *Remote Sens.* 13(4), 808. <https://doi.org/10.3390/rs13040808>

- 876 He, Z., He, D., 2020. Bilinear squeeze-and-excitation net-924  
877 work for fine-grained classification of tree species. IEEE925  
878 Geosci. Remote Sens. Lett. 18(7), 1139-1143. <https://doi.org/10.1109/LGRS.2020.2994952>926  
879 10.1109/LGRS.2020.2994952 927
- 880 Dong, Y., Ma, Z., Zi, J., et al., 2025. Multiscale feature fusion928  
881 and enhancement in a transformer for the fine-grained visual929  
882 classification of tree species. Ecol. Inform. 86, 103029. <https://doi.org/10.1016/j.ecoinf.2025.103029>930  
883 //doi.org/10.1016/j.ecoinf.2025.103029 931
- 884 Liu, Y., Chen, Y., Liu, Z., et al., 2025. A multi-feature fusion932  
885 network for tree species classification based on ground-based933  
886 LiDAR data. IEEE J. Sel. Top. Appl. Earth Obs. Remote934  
887 Sens. <https://doi.org/10.1109/JSTARS.2025.3527808>935  
888 936
- 889 Ferreira, M.P., dos Santos, D.R., Ferrari, F., et al., 2024. Im-937  
890 proving urban tree species classification by deep-learning938  
891 based fusion of digital aerial images and LiDAR. Urban939  
892 For. Urban Green. 94, 128240. <https://doi.org/10.1016/j.ufug.2024.128240>940  
893 941
- 894 Weiser, H., Schäfer, J., Winiwarter, L., Krašovec, N., Seitz, C.,942  
895 Schimka, M., et al., 2022. Terrestrial, UAV-borne, and air-943  
896 borne laser scanning point clouds of central European forest944  
897 plots, Germany, with extracted individual trees and manual945  
898 forest inventory measurements [dataset]. PANGAEA. <https://doi.org/10.1594/PANGAEA.942856>946
- 899 Hollaus, M., Chen, Y.-C., 2022. SilviLaser 2021 Benchmark947  
900 Dataset - Terrestrial Challenge [dataset]. TU Wien. <https://doi.org/10.48436/afdjq-ce434>948  
901 949
- 902 Ali, M., Biswas, A., Iglseider, A., Kumar, V., Kumar, S., Gupta,950  
903 S., et al., 2026. Terrestrial and airborne laser scanning dataset951  
904 of trees in the Shivalik Range, India with field measurements952  
905 and leaf-wood classifications. Sci. Data 13, 420. <https://doi.org/10.1038/s41597-026-06674-w>953  
906 954
- 907 Calders, K., Verbeeck, H., Burt, A., Origo, N., Nightingale,955  
908 J., Malhi, Y., et al., 2022. Terrestrial laser scanning data956  
909 Wytham Woods: individual trees and quantitative structure957  
910 models (QSMs) [dataset]. Zenodo. <https://doi.org/10.5281/zenodo.7307956>958  
911 959
- 912 Ali, M., Lohani, B., Hollaus, M., Pfeifer, N., 2024. Benchmark-960  
913 ing geometry-based leaf-filtering algorithms for tree volume961  
914 estimation using terrestrial LiDAR scanners. Remote Sens.962  
915 16(6), 1021. <https://doi.org/10.3390/rs16061021>963  
916 964
- 917 Wild, B., Özkan, T., Ali, M., Pöppel, F., Milenković, M.,965  
918 Hofhansl, F., et al., 2026. Evaluating RayCloudTools to es-966  
919 timate single-tree volume. Forestry, cpaf087. <https://doi.org/10.1093/forestry/cpaf087>967
- 920 Jia, S., Jiang, S., Zhang, S., et al., 2024. Graph-in-graph con-968  
921 volutional network for hyperspectral image classification.969  
922 IEEE Trans. Neural Netw. Learn. Syst. 35(1), 1157-1171.970  
923 <https://doi.org/10.1109/TNNLS.2022.3182715>971
- Li, N., Jiang, S., Xue, J., et al., 2024. Texture-aware self-attention model for hyperspectral tree species classification. IEEE Trans. Geosci. Remote Sens. 62, 1-15. Art no. 5502215. <https://doi.org/10.1109/TGRS.2023.3344787>
- Long, Y., Ye, S., Wang, L., et al., 2024. Scale pyramid graph network for hyperspectral individual tree segmentation. IEEE Trans. Geosci. Remote Sens. 62, 1-14. Art no. 5526014. <https://doi.org/10.1109/TGRS.2024.3439094>
- Zhang, S., Xu, M., Zhou, J., et al., 2022. Unsupervised spatial-spectral CNN-based feature learning for hyperspectral image classification. IEEE Trans. Geosci. Remote Sens. 60, 1-17. Art no. 5524617. <https://doi.org/10.1109/TGRS.2022.3153673>
- Belouï Schwenke, M., Xia, Z., Novoselova, I., et al., 2025. TreeAI Global Initiative - Advancing tree species identification from aerial images with deep learning [dataset]. Zenodo, Version TreeAI.V1.2. <https://doi.org/10.5281/zenodo.15351054>
- Dauphin, Y.N., Fan, A., Auli, M., et al., 2017. Language modeling with gated convolutional networks. Int. Conf. Mach. Learn., pp. 933-941. <https://proceedings.mlr.press/v70/dauphin17a.html>
- Xu, W., Wan, Y., 2024. ELA: Efficient local attention for deep convolutional neural networks. arXiv preprint arXiv:2403.01123. <https://arxiv.org/abs/2403.01123>
- Zhang, T., Li, L., Zhou, Y., et al., 2024. CAS-ViT: Convolutional additive self-attention vision transformers for efficient mobile applications. arXiv preprint arXiv:2408.03703. <https://arxiv.org/abs/2408.03703>
- Touvron, H., Cord, M., Douze, M., et al., 2021. Training data-efficient image transformers & distillation through attention. Int. Conf. Mach. Learn., pp. 10347-10357. <https://doi.org/10.48550/arXiv.2012.12877>
- Behera, A., Wharton, Z., Hewage, P.R.P.G., et al., 2021. Context-aware attentional pooling (CAP) for fine-grained visual classification. Proc. AAAI Conf. Artif. Intell. 35(2), 929-937. <https://doi.org/10.1609/aaai.v35i2.16176>
- Rao, Y., Chen, G., Lu, J., et al., 2021. Counterfactual attention learning for fine-grained visual categorization and re-identification. Proc. IEEE/CVF Int. Conf. Comput. Vis., pp. 1025-1034. <https://doi.org/10.1109/ICCV48922.2021.00109>
- Wang, J., Yu, X., Gao, Y., 2021. Feature fusion vision transformer for fine-grained visual categorization. arXiv preprint arXiv:2107.02341. <https://arxiv.org/abs/2107.02341>
- Bera, A., Wharton, Z., Liu, Y., et al., 2022. SR-GNN: Spatial relation-aware graph neural network for fine-grained image categorization. IEEE Trans. Image Process. 31, 6017-6031. <https://doi.org/10.1109/TIP.2022.3205215>

- 972 Diao, Q., Jiang, Y., Wen, B., et al., 2022. MetaFormer: A  
973 unified meta framework for fine-grained recognition. arXiv  
974 preprint arXiv:2203.02751. <https://arxiv.org/abs/2203.02751>
- 975 Sikdar, A., Liu, Y., Kedarisetty, S., et al., 2025. Interweav-  
976 ing insights: High-order feature interaction for fine-grained  
977 visual recognition. *Int. J. Comput. Vis.* 133(4), 1755-1779.  
978 <https://doi.org/10.1007/s11263-024-02260-y>
- 979 Sothe, C., Dalponte, M., Almeida, C. M., et al., 2019. Tree  
980 species classification in a highly diverse subtropical forest  
981 integrating UAV-based photogrammetric point cloud and hy-  
982 perspectral data. *Remote Sens.* 11(11), 1338. [https://doi.org/](https://doi.org/10.3390/rs11111338)  
983 [10.3390/rs11111338](https://doi.org/10.3390/rs11111338)
- 984 Wallace, L., Hally, B., Hillman, S., et al., 2020. Terrestrial  
985 image-based point clouds for mapping near-ground vegeta-  
986 tion structure: Potential and limitations. *Fire* 3(4), 59. [https:](https://doi.org/10.3390/fire3040059)  
987 [//doi.org/10.3390/fire3040059](https://doi.org/10.3390/fire3040059)