

Deep Pre-trained Time Series Features for Tree Species Classification in the Dutch Forest Inventory

Takayuki Ishikawa^a, Carmelo Bonannella^b, Bas J. W. Lerink^c, Marc Rußwurm^a

^aWageningen University, Droevendaalsesteeg 3, Wageningen, 6708 PB, The Netherlands

^bOpenGeoHub Foundation, Waldeck Pyrmontlaan 14, Doorwerth, 6865 HK, The Netherlands

^cWageningen Environmental Research, P.O. Box 47, Wageningen, 6700 AA, The Netherlands

Abstract

National Forest Inventory (NFI)s serve as the primary source of forest information, providing crucial tree species distribution data for carbon storage estimation and biodiversity assessments. However, maintaining these inventories requires labor-intensive on-site campaigns by forestry experts to identify and document tree species. Remote sensing approaches, particularly when combined with machine learning, offer opportunities to update NFIs more frequently and at larger scales. While the use of Satellite Image Time Series (SITS) has proven effective for distinguishing tree species through seasonal canopy reflectance patterns, current approaches rely primarily on Random Forest (RF) classifiers with hand-designed features and phenology-based metrics. Using deep features from an available pre-trained remote sensing foundation models offers a complementary strategy. These pre-trained models leverage unannotated global data and are meant to be used for general-purpose applications and can then be efficiently fine-tuned with smaller labeled datasets for specific classification tasks. This work systematically investigates *how deep features improve tree species classification accuracy in the Netherlands with few annotated data* (1,462 pure species plots). We evaluate this question using three tree species datasets of varying difficulty: First the unbalanced Dutch NFI data with 13 species classes, second, unbalanced with 7 species classes, and, third, a balanced NFI dataset with 7 classes used by a RF algorithm used by recent related work. Data-wise, we extracted time-series data from Sentinel-1 (S1), Sentinel-2 (S2) and European Centre for Medium-Range Weather Forecasts Re Analysis v5 (ERA5) satellites data (January-December 2020) and Shuttle Radar Topography Mission (SRTM) data using Google Earth Engine. Our results demonstrate that fine-tuning a publicly available remote sensing time series foundation model outperforms the current state-of-the-art in NFI classification in the Netherlands by a large margin of up to 10% across all datasets. This demonstrates that classic hand-defined harmonic features are too simple for this task and highlights the potential of using deep AI features for data-limited application like NFI classification. By leveraging openly available satellite data and pre-trained models, this approach significantly improves classification accuracy compared to traditional methods and can effectively complement existing forest inventory processes.

Keywords: National Forest Inventory, deep learning, pre-trained remote sensing foundation model, time series classification

1. Introduction

Forests play a significant role in mitigating climate change, adopting disaster prevention strategies, and providing ecosystem services, including sequestering carbon dioxide (CO₂), providing wood materials, and serving as a source of biodiversity (Tomppo et al., 2010; FAO, 2020; Francini et al., 2024). Tree species diversity improves productivity of materials and resistance to natural disturbance (Jactel et al., 2017). Therefore, monitoring and tracking records of forests including spatial distribution of tree species are necessary for a sustainable forest management. National Forest Inventory (NFI)s are the primary source of information for various purposes such as sustainable forest management, industry investment planning, biodiversity monitoring, and Greenhouse Gas accounting (Tomppo et al., 2010; Bonannella, 2024). The Paris Agreement also requires

the submission of a national carbon inventory, including carbon removals and reductions from forest lands (UNFCCC, 2015). Spatial tree species distribution information plays an important role in NFIs for various applications such as carbon storage estimation, forest management, and biodiversity assessments (Hermosilla et al., 2022; Blickensdörfer et al., 2024). Additionally, detailed tree species information is essential for national reports to the Forest Resource Assessment (FRA) of the Food and Agriculture Organization (FAO) and Forest Europe (Schelhaas et al., 2014). One of the key challenges in managing forest inventory is frequency and scale. Traditional inventory methods are based on sample-based field measurements conducted every 5 to 10 years (Tomppo et al., 2010). Current climate change and land-use changes due to economic growth cannot be captured in a timely manner using these traditional methods (Bonannella, 2024). Remote sensing helps improve and update NFIs (Hermosilla et al., 2022; Francini et al., 2024). Satellite sensors capture information about Earth every few days to several weeks, allowing us to monitor Earth's condition more frequently and

Email address: takayuki.ishikawa@wur.nl, wildflowers315@gmail.com (Takayuki Ishikawa)

extensively, even in remote areas. However, the global coverage and complex data volumes in remote sensing make human visual interpretation impractical for large-scale analysis. Machine learning techniques have achieved significant improvements in Earth observations, including tree species classification, particularly when using multiple timestamps of images over large areas (Blickensdörfer et al., 2024; Francini et al., 2024; Hermosilla et al., 2022). The RF algorithm is one of the most popular machine learning algorithms for tree species classification due to its robustness, interpretability, and ability to handle high-dimensional data (Breiman, 2001). However, RF models require well-designed input features engineering, and this selection of appropriate features is crucial for model performance (Heaton, 2016). These choices depend on domain knowledge and target area characteristics such as climate and tree species variety, and often fail to include all necessary features (Ahlsweide et al., 2023). In the context of NFIs, several state-of-the-art machine learning models utilizing RF have emerged (Hermosilla et al., 2022; Blickensdörfer et al., 2024; Francini et al., 2024)). These current methods for tree species classification rely on country-specific knowledge for input features, parameter settings, and high-quality data. However, scaling these methods to other countries at a national level presents challenges due to cost constraints and data availability limitations. Furthermore, deep learning models, which generally require large datasets to enhance performance, have remained largely unexplored for national-level tree species classification, primarily due to the limited size of available training datasets. These constraints collectively hinder the transferability of existing methods across different regions or countries. Deep learning models such as transformer architectures (Vaswani et al., 2023) have been recently introduced for forest monitoring, including tree species classification. This adoption is driven by increasing interest in multimodal and time-series data fusion in Remote Sensing (RS), enabled by the availability of big data and advancements in deep learning models (Li et al., 2022). Deep learning models can capture complex patterns in input data and create deep features that can be used for downstream tasks with classifiers and regressors including RF (Basu et al., 2015). While recent studies have achieved success in regional-scale tree species classification using high-quality labeled data, significant gaps remain in large-scale classifications (Fassnacht et al., 2016) for NFIs due to limited labeled data availability and high computational cost for training. Freely available pre-trained models, trained on large unlabeled datasets containing millions of pixels or images, have emerged as powerful tools for various downstream tasks. These models can achieve comparable or superior accuracy to traditional state-of-the-art machine learning approaches through fine-tuning without computationally expensive pre-training (Bommasani et al., 2022). Self-supervised learning, where models are trained without labels, has gained particular attention in RS applications (Wang et al., 2022). Using self-supervised learning as a model backbone with fine-tuning on limited labeled data has demonstrated significant accuracy improvements (Yu et al., 2022), particularly in time-series analysis tasks, while requiring less inductive bias (Dosovitskiy et al., 2021). The rapid growth of both labeled

and unlabeled datasets for RS (Gorelick et al., 2017; Ahlsweide et al., 2023) has enabled the development of various pre-trained models for tasks including tree classification (Lu et al., 2024). However, research comparing performance between fine-tuned pre-trained models and traditional approaches for tree species classification in NFIs remains limited.

This paper evaluates the effectiveness of deep features from freely available fine-tuned pre-trained models compared to a recently published random forest approach by Francini et al. (2024), who set a new state-of-the-art for Tree Species Identification in the Netherlands through a RF fitted on median reflectances and harmonic features.

The paper is structured as follows: section 2 details the three datasets from the Dutch Forest Inventory used in this study. Section 3 details the deep features we fine-tuned together with the harmonic features, as proposed by prior work.

2. Data

We use three datasets from the Dutch National Forest Inventory (NFI) at different complexity levels that we will use to compare the proposed fine-tuned deep time series features with the RF by Francini et al. (2024) fitted on harmonic time series features, as detailed later in section 3.

Complex & imbalanced (COMB) 13 tree species classes, imbalanced distribution, 1,462 time series samples.

Simplified & imbalanced (SIMB) 7 tree species classes, imbalanced class distribution, 1,479 time series samples.

Simplified & balanced (SIBA) 7 species classes, a balanced class distribution, by Francini et al. (2024). In total 13,790 samples and contains additional samples augmented by photointerpretation.

Table 1 shows a detailed list of the tree species across all three datasets.

Table 1: Species count and grouping in NFI and Francini datasets.

Dominant Species	COMB	Aggregated	SIMB	SIBA
<i>Pinus sylvestris</i>	513	Pinus	603	1,970
<i>Other Pinus</i>	89			
<i>Larix</i> spp	56	Larix	56	1,970
<i>Quercus robur petraea</i>	255	Quercus	288	1,970
<i>Other Quercus</i>	33			
<i>Fagus</i> spp	58	Beech	58	1,970
<i>Populus</i> spp	72	Populus	72	1,970
<i>Alnus</i> spp	30	Other Broadleaves	242	1,970
<i>Betula</i> spp	58			
<i>Fraxinus</i> spp	40			
<i>Other broadleaved</i>	102			
<i>Pseudotsuga menziesii</i>	90	DarkConifer	160	1,970
<i>Picea</i> spp	66			
total	1,462		1,479	13,790

2.1. Study area

This study focuses on forests in the Netherlands, which were described in the annual Greenhouse Gas accounting report (Arets et al., 2023). According to the latest National Forest Inventory 7, forests covered 363,801 ha in 2021, corresponding to 11% of the land use in the Netherlands (Schelhaas et al., 2022) and comprising approximately 36.4 million $10 \times 10\text{m}$ pixels, as shown in fig. 1d.

2.2. National Forest Inventory

Label Data Origin Ground truth data were collected manually through field measurements at 3,062 plots for the Dutch National Forest Inventory 6 between 2012 and 2013. Each plot contained a circular area with a variable radius (5 to 20 m) to ensure inclusion of at least 20 trees (Schelhaas et al., 2022). These plots correspond to 1 to 16 pixels at a $10 \times 10\text{m}$ resolution of Sentinel-1 and Sentinel-2. Due to privacy considerations, the precise coordinates of the plots were obtained under a confidentiality agreement with the Dutch government, with the requirement that the data be discarded after project completion. Prior to data disposal, we utilized the plot center coordinates from the total 3,062 pixel-level data points to extract satellite data at $10 \times 10\text{m}$ pixel resolution from Google Earth Engine (GEE) (Gorelick et al., 2017).

Data Quality Filtering. Following the identical protocol to Francini et al. (2024), we selected plots where a single dominant species represented more than 80% of the Basal Area (BA), which is the cross-sectional area of trees at breast height. This process removed the *Castanea spp* dominant species class that had no samples at this threshold.

Class Aggregation The original NFI data contains 19 dominant tree species classes. However, we noticed that six dominant species classes have less than 10 samples and we decided to remove these classes to not extensively rely in our evaluation on underrepresented tree species, similar to Kang et al. (2017). This resulted in the dataset **Complex & Imbalanced (COMB)** with 1,462 samples with 13 classes for the dominant species classification task. For the **Simplified & Imbalanced (SIMB)**, we further aggregated the classes into the 7 tree species identical to Francini et al. (2024). This resulted in 1479 data samples.

Finally, we also acquired the identical dataset used by Francini et al. (2024) which contains 13,790 data points evenly distributed across the aggregated classes (1,970 points per class). This dataset originated from the same NFI data but was augmented with additional labeled data points through visual interpretation of satellite imagery. We call **Simplified & Balanced (SIBA)** compared to aforementioned ones, as it contains 7 tree species classes and an even number of samples per class.

2.3. Remote Sensing Time Series

Since mono-temporal data cannot adequately capture seasonal leaf phenology for tree species classification, we utilized time-series data spanning January through December 2020. This temporal range aligns with both the Pretrained Remote Sensing Transformer (Presto) model requirements (Tseng et al.,

2024) and previous Dutch tree species classification research (Francini et al., 2024).

Input Data Modalities S2, S1, ERA5, SRTM. This study employs a comprehensive multi-source remote sensing (RS) approach that integrates satellite imagery and environmental data to characterize forest vegetation across space and time. The model relies on several key data sources: Sentinel-1 SAR GRD data with VH and VV polarizations at a 6-day revisit frequency; Sentinel-2 multispectral Level 1C data with 5–10 day revisit frequency, including 10 spectral bands (B2–B8, B8A, B11, B12) and NDVI computed during preprocessing; ERA5 climate data providing monthly measurements of 2m air temperature and total precipitation; and static SRTM terrain data comprising elevation and slope.

Remote Sensing Data Preparation For S2 data, cloud masking was performed using the S2 cloud probability dataset (Pasquarella et al., 2023), excluding pixels with cloud cover probability exceeding 65% (Francini et al., 2024). Geographic coordinates were utilized both for satellite data extraction from the GEE archive and as input features for the model. Raw time-series data was downloaded as individual CSV files for each plot location across all datasets (NFI, Francini, and non-labeled data), ensuring direct correspondence between remote sensing observations and ground truth measurements. We calculated monthly median values for each band in S1 and S2 except monthly ERA5 data. After getting monthly median values for S1 and S2, Normalized Difference Vegetation Index (NDVI) was calculated from S2 bands to capture vegetation dynamics.

3. Methods

In this study, we fundamentally compare results between a classic machine learning pipeline using a RF fitted on harmonic features and seasonal reflectances by Francini et al. (2024) with pre-trained deep-learning features obtained by fine-tuning the Presto time series foundation model proposed by Tseng et al. (2024). Hence, we outline the methodological details in both complimentary approaches in the next two sections before detailing the experimental setup in section 4.

3.1. Harmonic and Seasonal Features

We reproduced closely the classification approach of Francini et al. (2024), who recently published a new state-of-the-art classification approach for tree species mapping in the Netherlands. Their approach employ a RF Classifier (Breiman, 2001) with 500 trees that is fitted on 209 features containing *Seasonal Medoid Statistics* and *Harmonic Features* that we describe in the next two paragraphs. While Francini et al. (2024) used only features extracted from Sentinel-2 time series, we found experimentally (detailed later in section 5.1 and table 2) that including features from Sentinel-1 further slightly improves the results. Since we will also use Sentinel-1 to extract deep features (section 3.2), we also include Sentinel-1 features to ensure a fair comparison.

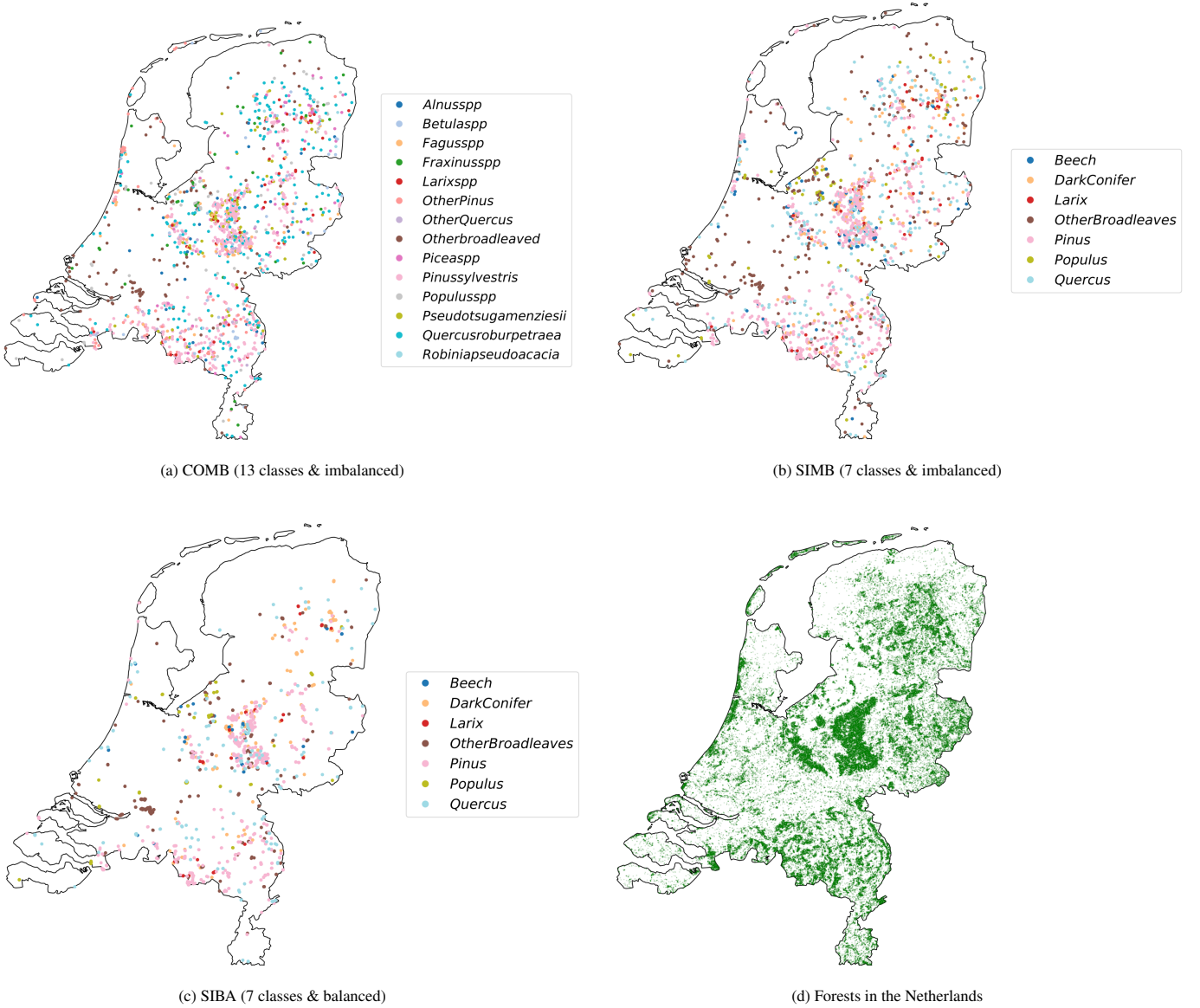


Figure 1: Overview of the study area and the three datasets used in this study. The simplified, class-balanced dataset is identical to related work by Francini et al. (2024).

76 Seasonal Medoid Statistics. In total, 12 bands from S1 and S2 are selected and additional 7 indices are calculated from S2 including NDVI, Normalized Burn Ratio (NBR), Enhanced Vegetation Index (EVI), and Tasseled Cap transformations: Brightness (TCB), Wetness (TCW), Greenness (TCG), Angle (TCA). This results in 19 bands, where for each band the median band value in each of the 4 seasons, winter (January, February, December), spring (March-May), summer (June-August), and autumn (September-November), the medoids values are extracted. These features describe the data through season-wise statistics, but do not capture the dynamic change in seasonality itself.

133 Harmonic Features. To capture seasonal changes, sine and cosine harmonics

$$P_t = \beta_0 + \beta_1 t + \beta_2 \cos(2\pi\omega t) + \beta_3 \sin(2\pi\omega t) \quad (1)$$

are fitted to the underlying monthly signal of each band according. The seven harmonic parameters $\beta_0, \beta_1, \beta_2, \beta_3, A, \phi$ for each of the 19 band are used as classification features alongside the residual RMSE: β_0 is the constant, β_1 is the time coefficient, and β_2 and β_3 are the frequency sine and cosine coefficients, respectively. To fit these four coefficients and to select P_t , the pixel p harmonic values at time t , we used a least squares regression to fit Eq. 1. $A = \sqrt{\beta_2^2 + \beta_3^2}$ is the amplitude of the harmonic curve on the y-axis, $\phi = \arctan(\beta_3/\beta_2)$ is the phase of the curve on the x-axis to the origin and $\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (P_t - X_t)^2}$ is root mean square error between P_t and the actual pixel values X_t . The frequency $\omega = 1$ (in years) and indicates one cycle per unit of time, initial guess values before fitting four coefficients were set as $\beta_0 = 0.1, \beta_1 = 0.1, \beta_2 = 0.4, \beta_3 = 0.4$.

This approach was proposed by Francini et al. (2024) based

on Sentinel-2 features. We added of Sentinel-1 features, as we found them to be slightly beneficial in accuracy (see section 5.1).

3.2. Deep Pre-trained Time Series Features

Overview of Presto deep features. In contrast to harmonic features that fit mathematical curves to individual bands, we evaluate the effectiveness of deep features obtained by the pre-trained time series foundation model, which is Lightweight, Pre-trained Transformers for Remote Sensing Timeseries (Presto) developed and trained by Tseng et al. (2024) for general-purpose remote sensing time series classification and regression. This model is based on transformer architecture and was originally pretrained on 21.5 million pixel time-series with 12-month contiguous intervals and static topographic and coordinates data. Each month’s composite satellite data included S1, S2 and its NDVI, ERA5, Dynamic World (DW), extracted between 2020-01-01 and 2021-12-31. The static SRTM and coordinates data are added to the time-series data as input of the Presto encoder. The model was originally pretrained through *Masked Auto Encoding (MAE)* (He et al., 2022) reconstruct masked sequences across different satellite modalities. The advantages of this model include the ability to handle freely available multi-source and multi-temporal data even when some sources (e.g., ERA5) or temporal data (e.g., S2 values in November) are missing, and computational efficiency when processing large areas. We obtained 128 dimensional deep features from the encoder output.

For this study, we used the pre-trained weights to initialize the Presto Encoder.

Multi Layer Perceptron (MLP) classifier was attached for fine-tuning of Presto encoder to the downstream task of tree species classification, which is described in section 4. We used the fine-tuned encoder to extract deep features from the test set for the comparison between *Harmonic and Seasonal Features* and *Presto deep features*.

As shown in fig. 2, we directly input 12-monthly time series ($T = 12$ timesteps) for (i) 2 Sentinel-1 bands (S1), (ii) 10 Sentinel-2 bands (S2) and (iii) NDVI (NDVI), (iv) 2 ERA5 features (ERA5), (v) Dynamic world tensor with all tree classes value (equal to value of 1. DW). To every pixel-time series we appended two static-in-time products: (i) elevation and slope from the SRTM digital elevation model (Farr et al., 2007) (TG) and (ii) location coordinates (Loc) of each pixel as input tokens to fine-tuned Presto encoder. Location coordinates (Loc) is 3D static in time Cartesian coordinates computed from the latitude and longitude of the pixel’s geographical location.

One training sample \mathbf{x} , comprising a pixel-timeseries $\mathbf{t} \in \mathbb{R}^{T \times 15}$, land cover classes $\mathbf{v} \in \mathbb{V}^{T \times 1}$, and static variables $\mathbf{s} \in \mathbb{R}^{1 \times 5}$, is summarized as follows:

$$\mathbf{x} = \left[\{t_i^{S1}; t_i^{S2}; t_i^{ERA5}; t_i^{NDVI}; t_i^{DW} \mid i = 1, \dots, 12\}; s^{TG}; s^{Loc} \right] \quad (2)$$

The pixel-timeseries \mathbf{x} were transformed into a sequence of tokens, each represented by an embedding \mathbf{e} , suitable for processing by the Presto transformer. For each timestep $0 \leq i < T$,

input variables were split into channel groups C based on their sensor or data source as follows;

- \mathbf{e}_{S1} : S1 VV and VH bands.
- \mathbf{e}_{S2-rgb} : S2 RGB(B2, B3, B4) bands.
- \mathbf{e}_{S2-re} : S2 Red Edge (B5, B6, B7) bands.
- $\mathbf{e}_{S2-nir10}$: S2 NIR (B8) 10m band.
- $\mathbf{e}_{S2-nir20}$: S2 NIR (B8A) 20m band.
- $\mathbf{e}_{S2-swir}$: S2 SWIR(B11, B12) bands.
- \mathbf{e}_{NDVI} : Computed from S2 B4 and B8.
- \mathbf{e}_{ERA5} : ERA5 Precipitation and 2m temperature.
- \mathbf{e}_{TG} : Elevation and slope from SRTM DEM.
- \mathbf{e}_{Loc} : Cartesian coordinates from latitude/longitude.

Each real-valued group was projected into a common latent space of dimension d_e using group-specific learned linear projections h_c , such that $e_i^{S1} = h_{S1}(t_i^{S1})$. Categorical variables like Dynamic World classes were embedded via indexing into a learnable embedding matrix. To convey the (i) location, (ii) timestamp, and (iii) channel group of each token, specific encodings were added to the embeddings \mathbf{e} . The final embedding of dimension d_e is a concatenation of:

- **Positional Encoding:** Sinusoidal encoding as in Vaswani et al. (2023).
- **Month Encoding:** To represent seasonal similarity, each month (0 to 11) is encoded as:

$$p_{\text{month}, 2i} = \sin\left(\frac{2\pi \cdot \text{month}}{12}\right)$$

$$p_{\text{month}, 2i+1} = \cos\left(\frac{2\pi \cdot \text{month}}{12}\right)$$

For static variables, positional and month encodings were set to zero.

- **Channel Group Encoding:** A learnable encoding was added for each channel group $c \in C = \{S1, S2 \text{ RGB}, \dots, ERA5, TG, Loc\}$, to capture source-specific information.

The final transformer input matrix $\mathbf{E} \in \mathbb{R}^{(T \cdot |C_{\text{dynamic}}| + |C_{\text{static}}|) \times d_e}$ is composed of:

- **Dynamic Variables:** For each timestep $i < T$ and channel group $c \in C$:

$$\mathbf{e}_i^c = h_c(t_i^c) + [p_{\text{channel}}^c; p_{\text{sin}}(i); p_{\text{month}}(i)]$$

- **Topographical Data:**

$$\mathbf{e}_{TG} = h_{TG}(s_{TG}) + [p_{\text{channel}}^{TG}; 0; 0]$$

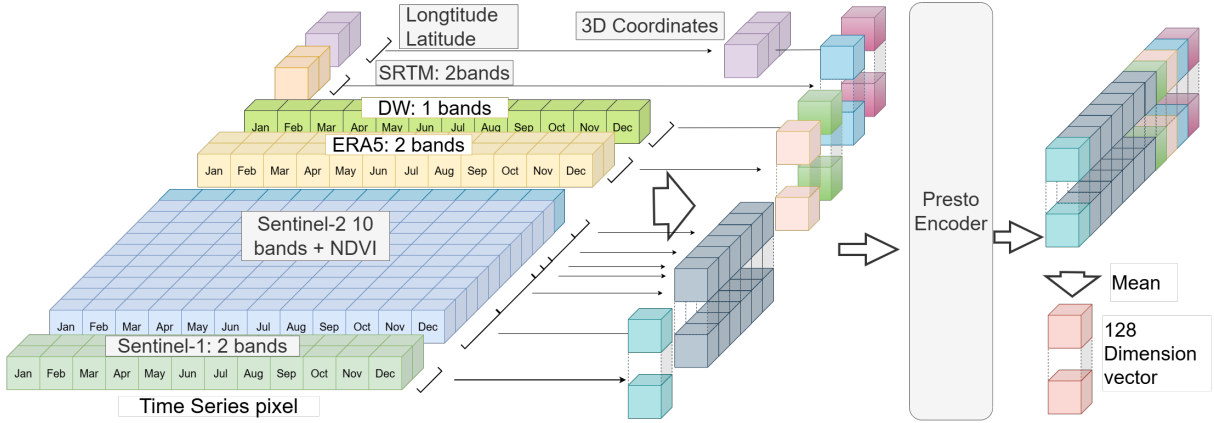


Figure 2: Presto encoder feature extraction for downstream tasks. Figure modified from (Tseng et al., 2024).

- **Coordinates:**

$$\mathbf{e}_{\text{Loc}} = h_{\text{Loc}}(s_{\text{Loc}})$$

Data Normalization (for deep features). The data normalization follows the original *Presto* setup (Tseng et al., 2024), with adjustments based on the observed value ranges for each source. Sentinel-1 VH and VV bands, originally ranging from -31 to 17 dB, are normalized by shifting $+25$ and dividing by 25 , resulting in a range of -0.24 to 1.68 . Sentinel-2 spectral bands, with raw values from 10 to $15,769$, are scaled by dividing by $10,000$, yielding normalized values between 0.00 and 1.58 . ERA5 temperature data, originally 278 to 295 K (equivalent to 6 to 23°C), are shifted by -272.15 and divided by 35 , producing a range of 0.17 to 0.66 . ERA5 precipitation, ranging from 0.008 to 0.208 m, is normalized by dividing by 0.03 , resulting in values from 0.27 to 6.93 . SRTM elevation data, spanning -27 to 331 m, are divided by $2,000$ to give a range of -0.01 to 0.17 , while slope values from 0.0 to 39.3° are scaled by dividing by 50 , yielding normalized values between 0.00 and 0.79 .

4. Experimental Setup

This study compares classic machine learning with hand-crafted feature engineering against deep learning approaches for tree species classification. Our experimental design follows a systematic progression to evaluate different methodological approaches and their combinations.

4.1. Comparative evaluation framework

Our experimental design implements four distinct comparison scenarios to systematically evaluate feature representations and classifier architectures:

1. **Sensor contribution analysis:** S1 vs S2 vs S1+S2 using harmonic features with RF
2. **Feature representation comparison:** Harmonic features vs deep features using RF
3. **Classifier architecture evaluation:** RF vs MLP using deep features

This systematic approach allows us to isolate the contributions of different methodological components while building toward the optimal configuration for operational NFI applications, extending the work of Francini et al. (2024) with comprehensive deep learning comparisons.

4.2. Model fine-tuning

The fine-tuning approach leverages the pre-trained Presto encoder’s learned representations while adapting the model to the specific characteristics of Dutch forest species. This transfer learning strategy balances computational efficiency with task-specific optimization.

For fine-tuning the pre-trained models, the encoder block from the pre-trained model was reused and trained with newly attached Multi Layer Perceptron (MLP) classifier for the downstream task. This fine-tuning model is trained on the $1,479$ and $13,790$ samples for 7 aggregated groups classification task, and the $1,462$ samples for the 13 dominant species classification task. Pre-trained encoder was set as trainable in this process.

Based on the current state-of-the-art classifier of MLP on tree species classification (Mouret et al., 2024), this classifier architecture consists of a 3-layer MLP with no dropout (rate: 0.0). The layer configuration comprises 1024 nodes in the first layer, 512 nodes in the second layer, and 256 nodes for the last layer, with batch normalization and ReLU activation applied throughout. The training parameters are configured as follows: a learning rate of 0.0001 is used with a weight decay of 0.00746 . The model is trained for a maximum of 100 epochs using a batch size of 64 . The training process employs cross-entropy loss for optimization. The best fine-tuned models and its accuracy were selected based on best validation loss.

4.3. Data splitting strategy

Train-test split ratio is a critical factor in model performance on downstream tasks. If dataset size is large, $80/20$ or higher ratios are commonly used (Rácz et al., 2021; Joseph, 2022). Based on our imbalanced dataset and previous research (Blickensdörfer et al., 2024), a $70/30$ (approximately $2:1$ ratio) train-test split was applied for fine-tuning and RF classification to

increase model performance and reduce standard deviation on tree species classification over large areas for NFI. During fine-tuning and subsequent RF classification for features extracted from the fine-tuned model, the same training and test data were used. This practice ensures that extracted features are not trained on test data, which would cause data leakage. By adopting a class-per-split strategy, each class is equally split into train and test data at this ratio.

In remote sensing machine learning, ensuring that the train-test split strategy does not introduce spatial autocorrelation is crucial for obtaining unbiased model performance estimates Karasiak et al. (2022). Our NFI data accommodates a sampling strategy which selects each plot with a density of 1 point per 100ha (Schelhaas et al., 2014). This means that each random point is designated per square kilometer. Additionally, the Francini data removed adjacent pixels with a minimum distance of 15m (Hermosilla et al., 2022; Francini et al., 2024). These sampling strategies ensure that the train-test split does not introduce strong spatial autocorrelation in the data; therefore, no additional train-test split strategy was applied in this study.

4.4. Performance evaluation

Model performance was evaluated using a comprehensive set of established accuracy metrics: per-class precision and recall, F1 scores, confusion matrices, and overall accuracy. While overall accuracy provides a general assessment of classification performance, it exhibits reduced sensitivity to class imbalance issues that are prevalent in ecological datasets. Therefore, the F1 score, which represents the harmonic mean of precision and recall, was employed as a supplementary metric to provide a more balanced evaluation of model performance when addressing the inherent class imbalances in tree species distribution data. Confusion matrices were used to visualize the model’s performance at the class level. These metrics are widely used in tree species classification studies (Goutte and Gaussier; Hermosilla et al., 2022; Francini et al., 2024; Blickensdörfer et al., 2024) and provide comprehensive assessment of classification performance across all species classes.

4.5. Statistical validation

To evaluate variability of model performance, each model was evaluated five times with different random seeds, which affect to train-test split and model initialization. The mean and standard deviation of the accuracy metrics were calculated for each model. This approach provides a more robust evaluation of model performance.

5. Results

5.1. Influence of multi-source input features

Following the approach of Francini et al. (2024), we first evaluated the contribution of Sentinel-1 radar data to optical Sentinel-2 classification. Using RF with harmonic features, the addition of Sentinel-1 provided modest improvements: 0.4% for the 13-class COMB dataset and 1.4% for the 7-class SIMB

dataset. These results justify the inclusion of Sentinel-1 features in our subsequent comparisons.

We first evaluated whether combining optical and radar satellite data improves tree species classification using hand-crafted features in a RF classifier. Feature sets were derived from seasonal medoid composites and harmonic fits, using either Sentinel-2 (S2) alone or in combination with Sentinel-1 (S1). The goal was to assess whether radar backscatter, which captures structural and moisture-related canopy properties, offers complementary information to optical reflectance data. As shown in Table 2 the inclusion of S1 provided modest but consistent improvements across both datasets. For the 13-class dataset (COMB), accuracy increased from 66.34% to 66.71% when adding S1 to S2-derived features. In the 7-class dataset (SIMB), accuracy improved from 72.87% to 74.22%. These results suggest that the radar signal contributes marginally to distinguishing between species, especially under a more simplified class grouping. While the absolute gains remain small, incorporating S1 ensures a more robust feature representation, particularly in settings where optical signals may be limited by cloud cover or seasonal variability. The modest accuracy gains observed here are also useful to contextualize the more substantial differences in performance reported in the next sections, which stem from the feature representations rather than the data sources themselves.

5.2. Deep vs Harmonic features

Table 3 demonstrates the superior performance of deep features over hand-crafted features when both are evaluated using RF classification. Deep features consistently outperform traditional harmonic and seasonal features across all three datasets, with F1 score improvements of 6-8 percentage points for the imbalanced datasets (COMB and SIMB) and over 10 percentage points for the balanced dataset (SIBA).

For direct comparison using identical S1+S2 sensor inputs (detailed in Appendix Table A.1), deep features show F1 score improvements of 12.3 percentage points (SIMB), 9.7 percentage points (COMB), and 10.7 percentage points (SIBA), confirming consistent performance advantages across different input configurations.

We compared classification performance between traditional hand-crafted features and deep features extracted from a pre-trained transformer model. In both cases, a RF classifier was used to isolate the effect of the feature representation. As shown in Table 3, deep features extracted from the pre-trained Presto model consistently outperformed the harmonic and medoid predictors across all dataset.

On the 7-class dataset (SIMB), models trained with deep features achieved an overall accuracy of 74.62% (± 1.93), compared to 70.58% (± 1.78) for the harmonic+medoid approach – an improvement of approximately 4 percentage points. The difference was even more pronounced when examining F1 scores, where deep features yielded 60.54% (± 3.08) versus 51.57% (± 3.36) for traditional features, representing an improvement of nearly 9 percentage points. In the more complex 13-classes case (COMB), deep features kept the advantage over traditional hand-crafted features, achieving an overall accuracy of 66.77%

Table 2: Classification accuracy comparison across different feature combinations and satellite data sources in NFI data (in %). Bold values indicate best performance for each label type.

Label	Features	Sentinel-1		Sentinel-2		S1-S2 Combined	
		Band	Acc.	Band	Acc.	Band	Acc.
COMB (13 classes)	Seasonal (S)	S1	51.03	S2	64.03	S1-S2	64.52
	Harmonic (H)	S1	55.04	S2	63.79	S1-S2	64.64
	All (S+H)	S1	57.72	S2	66.34	S1-S2	66.71
SIMB (7 classes)	Seasonal (S)	S1	55.60	S2	70.18	S1-S2	71.93
	Harmonic (H)	S1	60.19	S2	71.39	S1-S2	71.52
	All (S+H)	S1	62.75	S2	72.87	S1-S2	74.22

Table 3: Classification performance comparison (in %) across feature types (deep vs. harmonic+medoid) and classifiers (MLP vs. RF). Each model was run five times and the mean and standard deviation (\pm) are reported. Bold values indicate the best performance among MLP and RF per dataset type.

Dataset type			Classification Results		
			MLP Classifier (Deep features)	RF Classifier (Deep features)	RF Classifier (Harm.+med. Francini et al. (2024))
Overall Accuracy	SIMB (7 classes)	1479	77.13 \pm 1.17	74.62 \pm 1.93	70.58 \pm 1.78
	COMB (13 classes)	1462	67.99 \pm 1.33	66.77 \pm 1.51	62.71 \pm 1.51
	SIBA (7 classes)	13,790	98.21 \pm 0.20	95.28 \pm 0.60	84.31 \pm 0.99
F1 Score	SIMB (7 classes)	1479	68.24 \pm 2.85	60.54 \pm 3.08	51.57 \pm 3.36
	COMB (13 classes)	1462	52.45 \pm 2.25	46.84 \pm 2.65	38.54 \pm 2.14
	SIBA (7 classes)	13,790	98.20 \pm 0.20	95.27 \pm 0.60	84.26 \pm 0.99

(± 1.51) compared to 62.71% (± 1.51). Similarly, the F1 score showed a substantial improvement from 38.54% (± 2.14) with traditional features to 46.84% (± 2.65) with deep features – a gain of more than 8 percentage points. This indicates that the deep learning approach handles the increased complexity of multi-class classification more effectively. The highest performance gains were observed in the Francini dataset (SIBA), where deep features achieved a remarkable overall accuracy of 95.28% (± 0.60), representing an improvement of nearly 11 percentage points over the traditional features’ accuracy of 84.31% (± 0.99). The F1 scores demonstrated comparable improvements, with deep features achieving 95.27% (± 0.60) compared to 84.26% (± 0.99) for traditional features.

These results suggest that deep features extracted from the pre-trained transformer model are more effective at capturing species-specific spectral and temporal signatures than traditional hand-crafted features. The improvement is most pronounced in the SIBA dataset – where class separation is clearer due to the balanced classes and a higher sample size per class – but performance gains are evident across all dataset configurations examined in this study.

To evaluate the influence of classifier architecture, we compared the Multi Layer Perceptron (MLP) classifier for tree species classification Mouret et al. (2024) with a RF, both trained on the same deep features extracted from the pre-trained Presto model. This isolates the effect of the classifier while keeping the input representation fixed.

Across all datasets, the MLP classifier consistently outperformed the RF classifier in both overall accuracy and F1 score (Table 3). In the the 7-class dataset (SIMB), the MLP achieved an overall accuracy of 77.13% (± 1.17) compared to 74.62% (± 1.93) for the RF with a difference of approximately 2.5 percentage points. The performance gap was even more pronounced in terms of F1 score, with the MLP achieving 68.24% (± 2.85) versus 60.54% (± 3.08) for RF, representing a substantial improvement of nearly 8 percentage points. For the more challenging 13-class dataset (COMB), the MLP classifier maintained its advantage with an overall accuracy of 67.99% (± 1.33) compared to 66.77% (± 1.51) for the RF classifier. Similarly, the F1 score showed the MLP’s superior performance with 52.45% (± 2.25) versus 46.84% (± 2.65) for RF. The largest performance gain was observed in the balanced dataset (SIBA), where MLP reached 98.21% accuracy and 98.20% F1, compared to 95.28% and 95.27% with RF.

Notably, the MLP classifier not only achieved higher mean performance values but also demonstrated more consistent results across repeated runs, as evidenced by the lower standard deviations in most cases. These results indicate that while both classifiers can effectively leverage deep features, the MLP classifier with pre-trained Presto offers superior performance for forest type classification tasks across datasets of varying complexity and size. The improved class separation suggests that MLPs better exploit species-specific spectral and temporal signatures embedded in the deep features—patterns that may be

less accessible to tree-based models like RF, especially in more balanced or data-rich scenarios.

The S1+S2 only comparison (detailed in Appendix Table A.2) confirms these findings with F1 score improvements of 6.7 percentage points for SIMB (70.57% vs 63.83%), 6.7 percentage points for COMB (54.91% vs 48.24%), and 3.1 percentage points for SIBA (98.07% vs 94.94%), demonstrating consistent MLP advantages across different input configurations.

5.3. Class level performances

To better understand class-level behavior, we visualized the confusion matrices of the best-performing model – an MLP fine-tuned on deep features – for all three datasets (COMB, SIMB, SIBA). These matrices are shown in Figure 3, with each cell reporting the number of predicted samples and the corresponding precision (i.e., percentage of correct predictions among all predicted samples for that class). In the 7-class NFI dataset (SIMB), the model performs strongly across most classes. *Pinus* shows high precision with 1,668 correctly predicted samples and a precision of 84.7%. *Dark Conifer* and *Populus* also achieve high precision with 1,448 samples at 73.5% and 1,833 samples at 93.0%, respectively. Broadleaved classes exhibit more overlap: *Quercus* (oak) has 1,478 samples with 75.0% precision, while *Other Broadleaves* and *Fagus* reach 71.6% and 91.6%, respectively. Confusion is mostly observed between *Quercus*, *Populus*, and *Other Broadleaves*, reflecting the spectral similarity within this group.

In the more complex 13-class dataset (COMB), misclassifications are more pronounced, but key species maintain strong separation. *Pinus sylvestris* achieves 77.4% precision, *Pseudotsuga menziesii* reaches 83.3%, and *Picea spp.* also stands at 83.3%. *Populus spp.* improves to 73.9%, and *Betula spp.* achieves 100% precision, though on a smaller sample size. In contrast, *Quercus robur/petraea* reaches only 53.9%, and *Other Broadleaves* suffers from overlapping predictions with *Quercus* and *Populus*. Despite the increased class count, precision remains relatively high for most coniferous classes and for broadleaves with distinct phenology.

In the balanced dataset (SIBA), the model achieves near-perfect class separation. All classes have 1,970 samples, and precision ranges between 84.0% (*Larix*) and 93.0% (*Populus*). Other classes include *Fagus* (91.6%), *Pines* (84.7%), *Oak* (75.0%), and *Other Broadleaves* (71.6%). The confusion is minimal and symmetrically distributed, confirming that the model benefits substantially from balanced data and larger sample size.

6. Discussion

6.1. Multi-modal sensor integration

Our findings on the modest but consistent benefits of combining Sentinel-1 radar with Sentinel-2 optical data align well with recent European tree species mapping studies. Blickensdörfer et al. (2024) demonstrated that Sentinel-1 SAR data effectively fills potential data gaps in cloud-affected regions, particularly important for national-scale applications where clear-sky observation density varies substantially. Similarly, Liu et al.

(2023) confirmed that Sentinel-1 provides supplementary structural information to Sentinel-2 optical data, contributing to increased reliability of large-scale tree species mapping. Our results showing 0.4-1.4% accuracy improvements when adding Sentinel-1 to Sentinel-2 features corroborate these findings and provide confidence in the robustness of multi-modal approaches for operational NFI applications.

6.2. Deep vs hand-crafted features: an emerging paradigm

The transition from hand-crafted to learned feature representations marks a significant shift in remote sensing methodology. While deep learning has proven transformative in crop type mapping and land cover classification (Jakubik et al., 2023; Szwarcman et al., 2025; Wang et al., 2023), its application to tree species classification at NFI scale represents a relatively recent development. Our results demonstrate consistent improvements in F1 scores of 6-12 percentage points across datasets of varying class balance and complexity, suggesting that transformer-derived features capture phenological and structural variability more effectively than engineered metrics. This performance gain supports the ongoing shift in remote sensing from task-specific feature engineering toward transferable, learned representations that can adapt to diverse ecological contexts without domain-specific tuning.

6.3. Neural network architectures for learned features

The systematic superiority of MLP classifiers over RF when both were trained on deep features reflects broader patterns observed in high-dimensional, learned feature spaces. Neural architectures tend to better exploit the latent structure present in deep representations, where complex non-linear relationships between features may not be effectively captured by tree-based models. This finding aligns with recent advances in representation learning, where the pairing of learned features with compatible downstream architectures has proven critical for optimal performance (Mouret et al., 2024).

6.4. Computational efficiency and operational scalability

From an operational perspective, this approach offers significant advantages for national-scale implementations. The computational efficiency is evidenced by Presto’s lightweight architecture, which enables fine-tuning on a single GPU or CPU (Tseng et al., 2024), making it accessible to organizations with limited computational resources. The pre-training process itself demonstrates remarkable efficiency, processing 21.5 million pixel time-series in only 2 hours 12 minutes per epoch (43 hours 15 minutes for 20 epochs), indicating that similar approaches could be readily applied to other countries’ NFI systems. For The Netherlands specifically, with its 36.4 million 10×10m pixels, large-scale pixel-label mapping appears computationally feasible using this methodology.

6.5. Data quality and class balance considerations

Our comparative analysis across datasets of varying quality and balance reveals critical insights for operational NFI implementations. The stark performance differences between the NFI

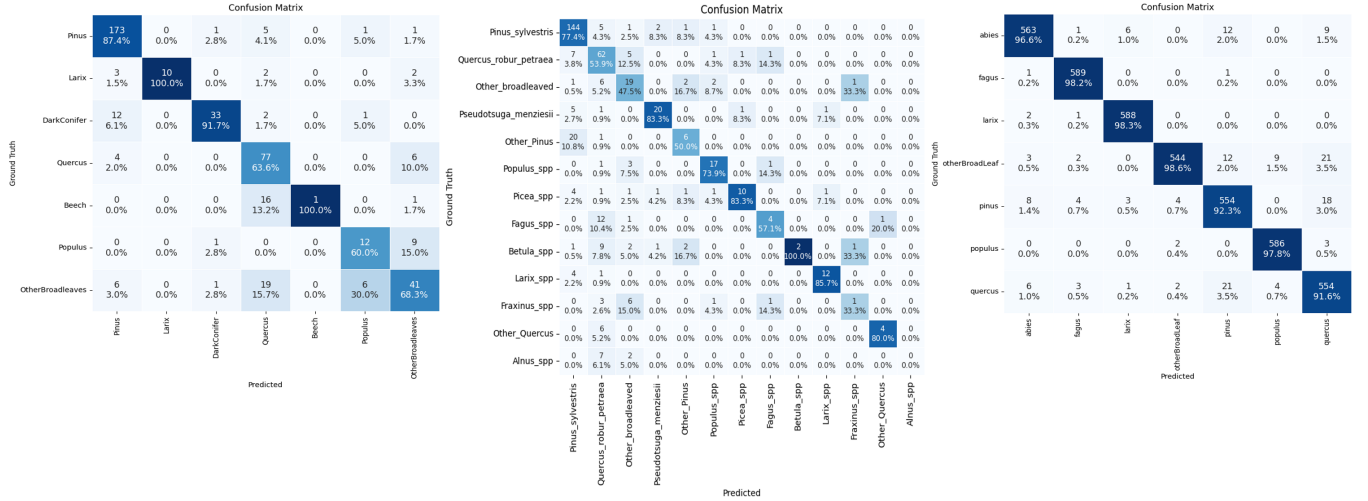


Figure 3: Confusion matrices for the best performing models with deep features by MLP classifier. *Simplified & Imbalanced (SIMB)*(left), *Complex & Imbalanced (COMB)*(center), *Simplified & Balanced (SIBA)* (right)

datasets (COMB and SIMB) and the Francini dataset (SIBA) underscore the fundamental importance of training data characteristics for deep learning approaches.

The Francini dataset’s superior performance (98.2% F1 score) compared to the imbalanced NFI datasets (52.5-68.2% F1 scores) primarily reflects the impact of class balance rather than dataset size alone. While the Francini dataset benefits from both larger sample size (13,790 vs 1,462-1,479) and better class balance, the magnitude of improvement suggests that addressing class imbalance represents a critical optimization opportunity for operational systems.

This finding has important implications for NFI data collection strategies. Resources may be more effectively allocated to targeted collection for underrepresented species classes. The approach demonstrated by Francini et al. (2024), which supplements ground-truth NFI data with additional samples through visual interpretation of high-resolution satellite imagery, offers a practical pathway for addressing these imbalances. Tools such as Collect Earth Online (Saah et al., 2019) provide accessible platforms for additional data collection.

6.6. Limitations and future directions

Several limitations merit consideration for operational deployment. First, deep features were extracted from a model pre-trained on a limited temporal and spatial range (2020-2021 global data), and further domain adaptation might improve performance for underrepresented classes to country’s forests. Second, while MLPs demonstrated superior performance in our controlled experiments, their generalization in operational settings with smaller or more imbalanced training sets requires further investigation.

Future research directions include: (i) adaptation to new regions or tree species not present in the training data, testing the transferability of learned representations across different forest ecosystems; (ii) integration of additional modalities such as hyperspectral or LiDAR data where available, potentially capturing structural metrics that complement spectral-temporal pat-

terns; and (iii) development of computational efficient uncertainty quantification methods to assess classification confidence for operational decision-making.

7. Summary and conclusions

This study examined the efficacy of deep learning approaches for tree species classification in National Forest Inventories, yielding several significant findings. First, deep features extracted from pre-trained Presto models consistently outperformed traditional methods employing harmonic and medoid predictors across varying dataset sizes and classification complexities. This performance advantage was particularly evident with larger and more balanced training datasets, as demonstrated by the results from the Francini dataset, indicating superior scalability of deep learning approaches with increased high-quality data availability. Second, the MLP classifier demonstrated superior performance compared to RF in our experimental setting, corroborating previous findings (Mouret et al., 2024). Collectively, these results demonstrate the potential for implementing computationally efficient, large-scale tree species mapping within NFI systems.

This research significantly advances our understanding of deep learning applications in forestry by having demonstrated the effectiveness of fine-tuning pre-trained deep learning models for tree species classification, establishing a computationally efficient approach for improving NFI accuracy.

This insight lays the groundwork for future developments in automated forest inventory systems and broader applications in environmental monitoring.

CRedit authorship contribution statement

Takayuki Ishikawa: Writing – review & editing, Writing – original draft, Visualization, Methodology, Investigation, Formal analysis, Conceptualization. **Carmelo Bonannella:** Writing – review & editing, Supervision, Conceptualization. **Bas J.**

W. Lerink Writing – review & editing, Data curation. **Marc Rußwurm** Writing – review & editing, Supervision, Conceptualization.

Declaration of competing interest

C.B. reports employment with OpenGeoHub Foundation. Other authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgement

This research is based on data from the Dutch National Forest Inventory (NFI), commissioned by the Ministry of Agriculture, Nature and Food Quality of the Netherlands. We thank the Ministry and the agencies responsible for collecting and providing these nationally representative forest data, which underpin our analysis. C.B. (Carmelo Bonannella) was supported by the Open-Earth-Monitor Cyberinfrastructure project, funded by the European Union's Horizon Europe research and innovation programme [grant agreement No.101059548].

Data availability

The source code is available on:
<https://github.com/wildflowers315/tree-species-classification-with-presto>.

References

- Ahlswede, S., Schulz, C., Gava, C., Helber, P., Bischke, B., Förster, M., Arias, F., Hees, J., Demir, B., Kleinschmit, B., 2023. TreeSatAI Benchmark Archive: a multi-sensor, multi-label dataset for tree species classification in remote sensing. *Earth System Science Data* 15, 681–695. , doi:10.5194/essd-15-681-2023.
- Arets, E., van Baren, S., Hendriks, C., Kramer, H., Lesschen, J., Schelhaas, M., 2023. Greenhouse gas reporting of the LULUCF sector in the Netherlands : methodological background, update 2023. Number 238 in WOT-technical report, WOT Natuur & Milieu, Netherlands. doi:10.18174/588942.
- Basu, S., Ganguly, S., Mukhopadhyay, S., DiBiano, R., Karki, M., Nemani, R., 2015. DeepSat: a learning framework for satellite imagery, in: *Proceedings of the 23rd SIGSPATIAL International Conference on Advances in Geographic Information Systems*, ACM, Seattle Washington. pp. 1–10. , doi:10.1145/2820783.2820816.
- Blickensdörfer, L., Oehmichen, K., Pflugmacher, D., Kleinschmit, B., Hostert, P., 2024. National tree species mapping using Sentinel-1/2 time series and German National Forest Inventory data. *Remote Sensing of Environment* 304, 114069. , doi:10.1016/j.rse.2024.114069.
- Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J.Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D.E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P.W., Krass, M., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X.L., Li, X., Ma, T., Malik, A., Manning, C.D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Niebles, J.C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J.S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y., Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin, A., Taori, R., Thomas, A.W., Tramèr, F., Wang, R.E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S.M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K., Liang, P., 2022. On the Opportunities and Risks of Foundation Models. , doi:10.48550/arXiv.2108.07258.
- Bonannella, C., 2024. Spatiotemporal modeling of vegetation dynamics in a changing environment: combining earth observation and machine learning. internal phd, wu. Wageningen University. Netherlands. doi:10.18174/655208.
- Breiman, L., 2001. Random Forests. *Machine Learning* 45, 5–32. , doi:10.1023/A:1010933404324.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houselby, N., 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. , doi:10.48550/arXiv.2010.11929.
- FAO, 2020. Global Forest Resources Assessment 2020: Main report. Rome. , doi:10.4060/ca9825en.
- Farr, T.G., Rosen, P.A., Caro, E., Crippen, R., Duren, R., Hensley, S., Kobrick, M., Paller, M., Rodriguez, E., Roth, L., Seal, D., Shaffer, S., Shimada, J., Umland, J., Werner, M., Oskin, M., Burbank, D., Alsdorf, D., 2007. The Shuttle Radar Topography Mission. *Reviews of Geophysics* 45. , doi:10.1029/2005RG000183.
- Fassnacht, F.E., Latifi, H., Stereńczak, K., Modzelewska, A., Lefsky, M., Waser, L.T., Straub, C., Ghosh, A., 2016. Review of studies on tree species classification from remotely sensed data. *Remote Sensing of Environment* 186, 64–87. , doi:10.1016/j.rse.2016.08.013.
- Francini, S., Schelhaas, M.J., Vangi, E., Lerink, B., Nabuurs, G.J., McRoberts, R.E., Chirici, G., 2024. Forest species mapping and area proportion estimation combining Sentinel-2 harmonic predictors and national forest inventory data. *International Journal of Applied Earth Observation and Geoinformation* 131, 103935. , doi:10.1016/j.jag.2024.103935.
- Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., Moore, R., 2017. Google Earth Engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment* 202, 18–27. , doi:10.1016/j.rse.2017.06.031.
- Goutte, C., Gaussier, E., . A probabilistic interpretation of precision, recall and f-score, with implication for evaluation, in: Losada, D.E., Fernández-Luna, J.M. (Eds.), *Advances in Information Retrieval*, Springer. pp. 345–359. doi:10.1007/978-3-540-31865-1_25.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R., 2022. Masked Autoencoders Are Scalable Vision Learners, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16000–16009. doi:10.1109/CVPR52688.2022.01553.
- Heaton, J., 2016. An empirical analysis of feature engineering for predictive modeling, in: *SoutheastCon 2016*, pp. 1–6. , doi:10.1109/SECON.2016.7506650.
- Hermosilla, T., Bastyr, A., Coops, N.C., White, J.C., Wulder, M.A., 2022. Mapping the presence and distribution of tree species in Canada's forested ecosystems. *Remote Sensing of Environment* 282, 113276. , doi:10.1016/j.rse.2022.113276.
- Jactel, H., Bauhus, J., Boberg, J., Bonal, D., Castagnérol, B., Gardiner, B., Gonzalez-Olabarria, J.R., Koricheva, J., Meurisse, N., Brockerhoff, E.G., 2017. Tree Diversity Drives Forest Stand Resistance to Natural Disturbances. *Current Forestry Reports* 3, 223–243. , doi:10.1007/s40725-017-0064-1.
- Jakubik, J., Roy, S., Phillips, C.E., Fraccaro, P., Godwin, D., Zadrozny, B., Szwarcman, D., Gomes, C., Nyirjesy, G., Edwards, B., Kimura, D., Simmumba, N., Chu, L., Mukkavilli, S.K., Lambhate, D., Das, K., Bangalore, R., Oliveira, D., Muszynski, M., Ankur, K., Ramasubramanian, M., Gurung, I., Khallaghi, S., Hanxi, Li, Cecil, M., Ahmadi, M., Kordi, F., Alemohammad, H., Maskey, M., Ganti, R., Weldemariam, K., Ramachandran, R., 2023. Foundation Models for Generalist Geospatial Artificial Intelligence. , arXiv:2310.18660.
- Joseph, V.R., 2022. Optimal ratio for data splitting. *Statistical Analysis and Data Mining: An ASA Data Science Journal* 15, 531–538. , doi:10.1002/sam.11583.
- Kang, Q., Chen, X., Li, S., Zhou, M., 2017. A Noise-Filtered Under-Sampling Scheme for Imbalanced Classification. *IEEE Transactions on Cybernetics* 47, 4263–4274. , doi:10.1109/TCYB.2016.2606104.
- Karasiak, N., Dejoux, J.F., Monteil, C., Sheeren, D., 2022. Spatial depen-

- dence between training and test sets: another pitfall of classification accuracy assessment in remote sensing. *Machine Learning* 111, 2715–2740. , doi:10.1007/s10994-021-05972-1.
- Li, J., Hong, D., Gao, L., Yao, J., Zheng, K., Zhang, B., Chanussot, J., 2022. Deep learning in multimodal remote sensing data fusion: A comprehensive review. *International Journal of Applied Earth Observation and Geoinformation* 112, 102926. , doi:10.1016/j.jag.2022.102926.
- Liu, X., Frey, J., Munteanu, C., Still, N., Koch, B., 2023. Mapping tree species diversity in temperate montane forests using Sentinel-1 and Sentinel-2 imagery and topography data. *Remote Sensing of Environment* 292, 113576. doi:10.1016/j.rse.2023.113576.
- Lu, S., Guo, J., Zimmer-Dauphinee, J.R., Nieusma, J.M., Wang, X., VanValkenburgh, P., Wernke, S.A., Huo, Y., 2024. AI Foundation Models in Remote Sensing: A Survey. , doi:10.48550/arXiv.2408.03464.
- Mouret, F., Morin, D., Planells, M., Vincent-Barbaroux, C., 2024. Tree species classification at the pixel-level using deep learning and multispectral time series in an imbalanced context. , doi:10.48550/arXiv.2408.08887.
- Pasquarella, V.J., Brown, C.F., Czerwinski, W., Rucklidge, W.J., 2023. Comprehensive quality assessment of optical satellite imagery using weakly supervised video learning, in: 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 2125–2135. doi:10.1109/CVPRW59228.2023.00206.
- Rácz, A., Bajusz, D., Héberger, K., 2021. Effect of Dataset Size and Train/Test Split Ratios in QSAR/QSPR Multiclass Classification. *Molecules (Basel, Switzerland)* 26, 1111. doi:10.3390/molecules26041111.
- Saah, D., Johnson, G., Ashmall, B., Tondapu, G., Tenneson, K., Patterson, M., Poortinga, A., Markert, K., Quyen, N.H., San Aung, K., Schlichting, L., Matin, M., Uddin, K., Aryal, R.R., Dilger, J., Lee Ellenburg, W., Flores-Anderson, A.I., Wiell, D., Lindquist, E., Goldstein, J., Clinton, N., Chishtie, F., 2019. Collect Earth: An online tool for systematic reference data collection in land cover and use applications. *Environmental Modelling & Software* 118, 166–171. , doi:10.1016/j.envsoft.2019.05.004.
- Schelhaas, M., Clerkx, A., Daamen, W., Oldenburger, J., Velema, G., Schnitger, P., Schoonderwoerd, H., Kramer, H., 2014. Zesde Nederlandse bosinventarisatie : methoden en basisresultaten. Number 2545 in Alterra-rapport, Alterra, Netherlands. .
- Schelhaas, M., Teeuwen, S., Oldenburger, J., Beerkens, G., Velema, G., Kramers, J., Lerink, B., Paulo, M., Schoonderwoerd, H., Daamen, W., Dolstra, F., Lusink, M., van Tongeren, K., Scholten, T., Pruijsten, I., Voncken, F., Clerkx, A., 2022. Zevende Nederlandse Bosinventarisatie: Methoden en resultaten. Number 142 in WOt-rapport, WOT Natuur & Milieu, Netherlands. doi:10.18174/571720.
- Szwarcman, D., Roy, S., Fraccaro, P., Gíslason, T.E., Blumenstiel, B., Ghosal, R., de Oliveira, P.H., de Sousa Almeida, J.L., Sedona, R., Kang, Y., Chakraborty, S., Wang, S., Gomes, C., Kumar, A., Truong, M., Godwin, D., Lee, H., Hsu, C.Y., Asanjan, A.A., Mujeci, B., Shidham, D., Keenan, T., Arevalo, P., Li, W., Alemohammad, H., Olofsson, P., Hain, C., Kennedy, R., Zadrozny, B., Bell, D., Cavallaro, G., Watson, C., Maskey, M., Ramachandran, R., Moreno, J.B., 2025. Prithvi-eo-2.0: A versatile multi-temporal foundation model for earth observation applications. , arXiv:2412.02732.
- Tomppo, E., Gschwantner, T., Lawrence, M., McRoberts, R.E. (Eds.), 2010. National Forest Inventories: Pathways for Common Reporting. Springer Netherlands, Dordrecht. , doi:10.1007/978-90-481-3233-1.
- Tseng, G., Cartuyvels, R., Zvonkov, I., Purohit, M., Rolnick, D., Kerner, H., 2024. Lightweight, Pre-trained Transformers for Remote Sensing Time-series. , doi:10.48550/arXiv.2304.14065.
- UNFCCC, 2015. Paris agreement. . uNTS XXVII 7.d; FCCC/CP/2015/L.9/Rev.1; T.I.A.S. No. 16-1104.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2023. Attention Is All You Need. , doi:10.48550/arXiv.1706.03762.
- Wang, H., Chang, W., Yao, Y., Yao, Z., Zhao, Y., Li, S., Liu, Z., Zhang, X., 2023. Cropformer: A new generalized deep learning classification approach for multi-scenario crop classification. *Frontiers in Plant Science* 14. doi:10.3389/fpls.2023.1130659.
- Wang, Y., Albrecht, C.M., Braham, N.A.A., Mou, L., Zhu, X.X., 2022. Self-Supervised Learning in Remote Sensing: A review. *IEEE Geoscience and Remote Sensing Magazine* 10, 213–247. , doi:10.1109/MGRS.2022.3198244.
- Yu, A., Liu, B., Cao, X., Qiu, C., Guo, W., Quan, Y., 2022. Pixel-Level Self-Supervised Learning for Semi-Supervised Building Extraction From Remote Sensing Images. *IEEE Geoscience and Remote Sensing Letters* 19, 1–5. , doi:10.1109/LGRS.2022.3207465.

Appendix A. Appendix

Appendix A.1. Additional results

Appendix A.1.1. Results of models using only S1 and S2 data for fine-tuning

Our primary models incorporated all available satellite data as input for both the second stage of pre-training and fine-tuning to extract deep features. In this section, we present results from models using only S1 and S2 data as input for fine-tuning and feature extraction, excluding SRTM and ERA5 data.

Deep features model performance.

Table A.1 demonstrates that deep features extracted by pre-trained Presto from only S1 and S2 data outperformed harmonic and medoid features in the RF classifier. These results were comparable to models that utilized all satellite data as input.

Table A.1: Classification performance comparison between deep features and harmonic+medoid features (in %). Each model was run five times and the mean and standard deviation (\pm) are reported. Bold values indicate the best performance between deep features extracted from Presto and harmonic+medoid features per dataset type.

			RF	
	Dataset type	Data size	Deep feat. ((Tseng et al., 2024))	Harm.+med. ((Francini et al., 2024))
Overall Accuracy	NFI data (7 classes)	1480	75.47 \pm 1.15	70.58 \pm 1.78
	NFI data (13 classes)	1462	67.77 \pm 2.38	62.71 \pm 1.51
	Francini data (7 classes)	13,790	94.96 \pm 0.58	84.31 \pm 0.99
F1 Score	NFI data (7 classes)	1480	63.83 \pm 1.21	51.57 \pm 3.36
	NFI data (13 classes)	1462	48.24 \pm 4.31	38.54 \pm 2.14
	Francini data (7 classes)	13,790	94.94 \pm 0.58	84.26 \pm 0.99

Table A.2 shows that the MLP classifier consistently achieved higher accuracy than the RF classifier, even when using only S1 and S2 data as input for the pre-trained Presto model.

Table A.2: Classification performance comparison between MLP and RF classifiers with deep features (in %). Each model was run five times and the mean and standard deviation (\pm) are reported. Bold values indicate the best performance between MLP and RF classifier per dataset type.

			Results	
	Dataset type	Data size	MLP classifier	RF classifier
Overall Accuracy	NFI data (7 classes)	1480	78.79 \pm 1.01	75.47 \pm 1.15
	NFI data (13 classes)	1462	69.39 \pm 1.65	67.77 \pm 2.38
	Francini data (7 classes)	13,790	98.08 \pm 0.09	94.96 \pm 0.58
F1 Score	NFI data (7 classes)	1480	70.57 \pm 2.20	63.83 \pm 1.21
	NFI data (13 classes)	1462	54.91 \pm 4.29	48.24 \pm 4.31
	Francini data (7 classes)	13,790	98.07 \pm 0.09	94.94 \pm 0.58

Comparison between all bands and only S1 and S2 bands.

Table A.3 compares models using all satellite data against those using only S1 and S2 data with MLP classifier. Interestingly, models with only S1 and S2 data achieved higher accuracy with NFI data, while no significant difference was observed with Francini data.

Table A.3: Classification performance comparison between using all available bands versus using only S1 and S2 bands (in %). Each model was run five times with mean and standard deviation (\pm) reported. Bold values indicate best performance between the two band configurations for each dataset type.

			MLP classifier	
	Dataset type	Data size	All bands	S1S2 only
Overall Accuracy	NFI data (7 classes)	1480	77.13 \pm 1.17	78.79 \pm 1.01
	NFI data (13 classes)	1462	67.99 \pm 1.33	69.39 \pm 1.65
	Francini data (7 classes)	13,790	98.21 \pm 0.20	98.08 \pm 0.09
F1 Score	NFI data (7 classes)	1480	68.24 \pm 2.85	70.57 \pm 2.20
	NFI data (13 classes)	1462	52.45 \pm 2.25	54.91 \pm 4.29
	Francini data (7 classes)	13,790	98.20 \pm 0.20	98.07 \pm 0.09