
TESSERA: Temporal Embeddings of Surface Spectra for Earth Representation and Analysis

Zhengpeng Feng¹, Sadiq Jaffer¹, Jovana Knezevic², Silja Sormunen³,
 Robin Young¹, Madeline Lisaius¹, Markus Immitzer⁴, James Ball²,
 Clement Atzberger⁴, David A. Coomes², Anil Madhavapeddy¹, Andrew Blake⁵,
 Srinivasan Keshav^{1*}

¹ Department of Computer Science and Technology, University of Cambridge, UK

² Department of Plant Sciences, University of Cambridge, UK

³ Department of Computer Science, Aalto University, Finland

⁴ dClimate Labs, New York, USA

⁵ Clare Hall, University of Cambridge CB3 9AL, UK

*Corresponding author. Email: sk818@cam.ac.uk.

Abstract

Satellite remote sensing (RS) enables a wide array of downstream Earth observation (EO) applications, including climate modeling, carbon accounting, and strategies for conservation and sustainable land use. We present TESSERA, a novel Remote Sensing Foundation Model (RSFM) that uses Self-Supervised Learning (SSL) to generate global, robust representations at 10m scale from pixel-level satellite time series data. TESSERA combines information from only optical and SAR data streams using two parallel Transformer-based encoders: one dedicated to Sentinel-1 SAR polarizations and another to Sentinel-2 MSI data (10 selected spectral bands) to create representations that are then fused using a multilayer perceptron (MLP), resulting in a global representation map covering the years 2017 to 2024. Our precomputed representations set a new state-of-the-art performance benchmark and our open-source approach democratizes access to high-performance, high-resolution representations. We benchmark the performance of TESSERA in five diverse tasks, comparing our work with state-of-the-art task-specific models and other foundation models. Our results show that TESSERA outperforms both traditional RS baselines and the leading geospatial foundation models in these diverse downstream tasks.

Monitoring Earth’s dynamic systems through satellite Earth Observation (EO) is critical for addressing global challenges including food security, biodiversity loss, climate change, and disaster mitigation. Petabytes of EO data are available from sensors in various modalities, and declining launch costs promise exponential growth in the future. However, practitioners face a fundamental bottleneck: the scarcity of large, accurately labeled datasets required for supervised training. Direct use of EO data also presents challenges as the data quality is affected by cloud cover, atmospheric effects, sensor biases, and non-uniform temporal sampling, preventing full utilization of this valuable resource.

Users of EO data usually mitigate cloudiness by annual or seasonal compositing and then use supervised training of a neural network to map from composited images to tasks-specific classification, such as crop identification. However, compositing removes the critical temporal signal and training a task-specific neural network is computationally onerous. More recent EO approaches tackle the second problem by adopting the paradigm of *foundation models*, creating general-purpose models pre-trained on a broad set of EO data. For example, Visual Foundation Models (VFMs) have been

engineered to interpret a wide array of geospatial data types, including multispectral, hyperspectral, and SAR imagery. Their development follows several paths, including supervised pre-training on large labeled datasets for spatiotemporal analysis [13, 12, 83, 11], and various self-supervised learning (SSL) strategies that learn from the inherent structure of the data itself. These SSL techniques include generative methods, such as masked autoencoders [19, 55, 69, 32, 62, 48, 61, 28, 77, 74, 30, 72, 37, 79] and diffusion-based models [92, 42, 64], and contrastive methods which learn efficient representations by comparing data samples [7, 54, 5, 52, 43, 14, 29, 38, 85]. Subsequently, these general-purpose models are adapted and fine-tuned for specific remote sensing tasks such as high-resolution segmentation [84, 87, 88, 16, 91].

Whilst existing SSL approaches offer rich multi-dimensional ‘representations’ from unlabeled data through auxiliary tasks that reveal latent patterns [8, 51], they still do not address the problem of loss of the temporal signal during compositing. In recent work Lisaius et al [49] have demonstrated that the Barlow Twins method [90] is an elegant approach, grounded in redundancy reduction principles [10], to preserve the temporal signal in latent representations. Building on this insight, we introduce TESSERA, a foundation model for EO data that creates 128-dimensional latent representations at global scale to enable state-of-the-art performance across a diverse array of complex tasks. Our experiments demonstrate that TESSERA-learned representations outperform baseline methods in tasks as diverse as estimating canopy height in rainforests in Borneo, detecting crop types in Austria, and identifying fire scars in California. Crucially, our model is open source, enabling ease of adoption and reproducibility. Our global representation map also enables users to keep their data private rather than having to upload it to a centralized repository.

1 Our Model

TESSERA processes unlabeled time series from Sentinel-1 SAR and Sentinel-2 Multispectral Instrument (MSI) data. It is a remote sensing foundation model designed for pixel-wise feature extraction from Sentinel-1 and Sentinel-2 imagery. It comprises two main components: a dual-branch encoder and a projector network (see Fig. 1 for a detailed diagram). For each 10-meter pixel globally, it generates a compact, 128-dimensional embedding that encapsulates that pixel’s annual temporal and spectral characteristics. A primary goal of this work is the production of global, annual, 10-meter resolution representation maps covering the years 2017 to 2024. These pre-computed maps significantly lower the barrier of entry for a wide array of downstream Earth observation applications, by providing readily usable, information-rich features. The core methodology is based on a self-supervised learning paradigm, leveraging the complementary information from optical and SAR data streams.

The fundamental input unit for TESSERA is termed the “d-pixel”. For multi-spectral (from Sentinel-2) or SAR backscatter (from Sentinel-1) observations each d-pixel uniquely represents a single geographic 10-meter pixel by structuring the repeat observations from the complete annual time series of the modality into a two-dimensional array (i.e, timesteps by channels). This format explicitly preserves the temporal evolution and spectral/backscatter signatures inherent in the satellite data. The d-pixel design also intrinsically accommodates common observational gaps, such as those due to cloud cover, by masking invalid data points, which are then handled during subsequent processing stages. See Appendix A.4 and Supplementary Fig. 1.

1.1 Dual-encoder

TESSERA employs a multi-modal architecture to process and integrate information from the Sentinel-1 and Sentinel-2 observations. The model consists of specialized encoders for each modality followed by a projection head for self-supervised learning.

The model features two separate, parallel Transformer-based encoders: one for Sentinel-1 SAR backscatter data (VV and VH polarizations) and another for Sentinel-2 MSI data (10 selected spectral bands). Each encoder processes the temporal sequence of a d-pixel, utilizing multi-head self-attention mechanisms to capture complex temporal dependencies and patterns specific to each data modality. Positional encodings derived from the Day-of-Year (DOY) are incorporated to provide temporal context. An attention-pooling layer within each encoder aggregates the temporal features into a fixed-size vector representing the annual signature for that modality. See Appendix A.7 and Supplementary Figure 1 for architectural details.

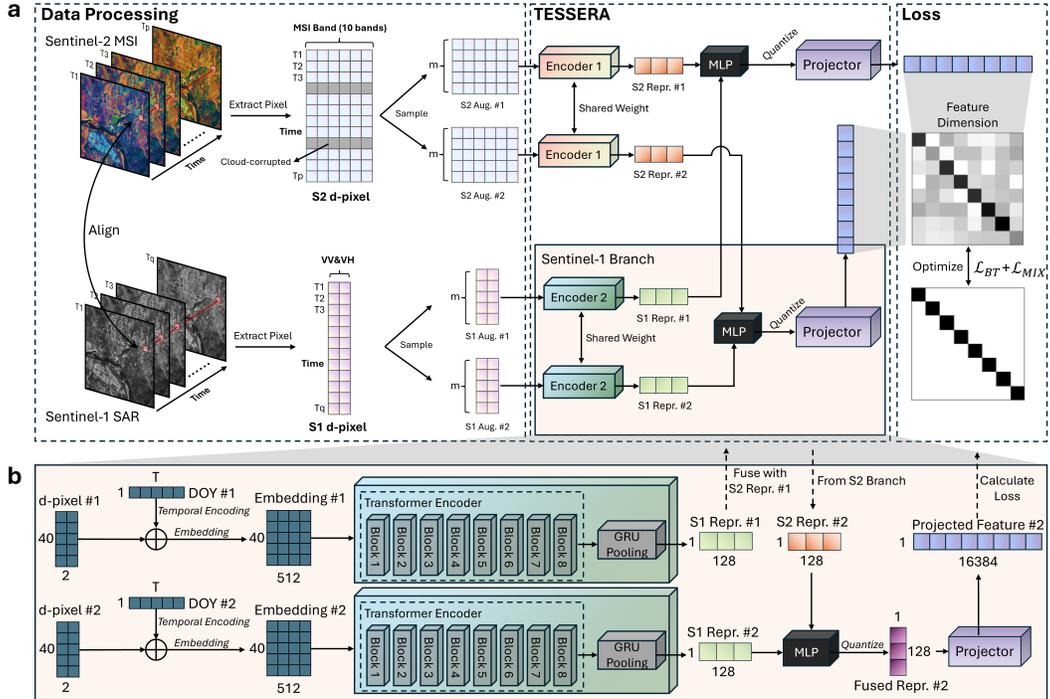


Figure 1: **Architecture and data processing pipeline of the TESSERA foundation model.** **a**, Overview of the end-to-end workflow. Spatially aligned time-series data from Sentinel-2 MSI and Sentinel-1 SAR are first extracted for each pixel location, forming modality-specific ‘d-pixels’. For self-supervised training, two distinct augmentations are created from each d-pixel and processed through a dual-branch architecture with shared weights within each branch. The resulting 128-dimensional representations from both modalities are fused by a multilayer perceptron MLP. The fused features are then expanded to 16,384 dimensions by a large projector network. A modified Barlow Twins loss ($\mathcal{L}_{BT} + \mathcal{L}_{MIX}$) is calculated on the cross-correlation matrix of these projected features, optimizing the model to learn invariant and decorrelated representations. **b**, Detailed structure of a single encoder branch. A d-pixel, represented as a matrix of temporal steps versus spectral bands (e.g., 40 time-steps \times 2 bands for Sentinel-1), is first embedded. A DOY-based temporal positional encoding is added to these embeddings. The resulting sequence is then processed by an 8-block Transformer Encoder. Finally, a GRU pooling layer aggregates the temporal features to produce a single 128-dimensional representation for that modality, which is then passed to the fusion MLP shown in panel **a**.

1.2 Projector

The modality-specific representations derived from the dual encoders are first fused using a multilayer perceptron (MLP). This fused representation, which serves as the final 128-dimensional pixel embedding for downstream tasks, is then quantized before being fed into a high-dimensional projector network. This projector, a significantly larger MLP, expands the dimensionality of the fused representation¹ to facilitate effective redundancy reduction during the self-supervised learning phase.

1.3 Model Training

The TESSERA model learns its representations through a self-supervised pretraining process on a large-scale, unlabeled global dataset. It is trained by optimizing a modified Barlow Twins loss function [90]. The core principle is to learn embeddings that are invariant to differing or incomplete views of the same pixel’s data (e.g., different sets of cloud-free observation dates) while minimizing redundancy between the learned feature dimensions. This is achieved by generating two distinct

¹This is a critical step inspired by Barlow Twins [90]

augmented views for each input d-pixel, primarily through independent sparse temporal sampling of annual Sentinel-1 and Sentinel-2 observations, rather than applying artificial distortions or augmentations. The model then learns to reconcile these partial views. The loss function includes the standard Barlow Twins terms for invariance and redundancy reduction, augmented with an additional mix-up regularization term to enhance robustness and inter-sample interaction during training (see Appendix A.7.2 for the full loss function, Equation 12).

The pretraining phase utilizes a vast global dataset comprising approximately 0.8 billion d-pixels, sampled from Sentinel-1 and Sentinel-2 imagery. A crucial element for successful generalization was a specific multi-level data sampling, shuffling, and batching strategy. This ensured that each training mini-batch contained a highly diverse set of d-pixels from various geographic locations and acquisition conditions, preventing model overfitting and promoting the learning of universally applicable features. The model was trained for one epoch over this dataset using distributed training techniques.

1.4 Model Inference and Global Representation Maps

Following pretraining, the TESSERA dual-encoder (with frozen weights and excluding the projector) is used to generate the final 128-dimensional representation for every 10-meter pixel globally, for each year from 2017 to 2024. This involves processing the annual Sentinel-1 and Sentinel-2 time series for each pixel through the trained encoders. The primary output of this inference stage is a set of annual, global, 10-meter resolution TESSERA representation maps. These maps are designed as readily usable, multi-channel geospatial data layers. This “Model-as-Data” approach significantly lowers the barrier to entry for end-users, as these rich, pre-computed features can be directly ingested by downstream models without the need for raw satellite data processing or running the TESSERA model itself.

2 Downstream Tasks

The effectiveness and generalizability of the TESSERA representations are rigorously evaluated across a diverse range of downstream remote sensing tasks. These include, but are not limited to, pixel-wise classification (e.g., crop type mapping), pixel-wise regression (e.g., canopy height estimation), and patch-based dense predictions (e.g., semantic segmentation and land change detection). For all such evaluations, the pretrained TESSERA encoder acts as a fixed feature extractor. Lightweight, task-specific model heads (e.g., shallow MLPs or UNet-style decoders) are then trained using these fixed representations. This approach typically requires substantially less labeled data and computational resources compared to training deep models from scratch, highlighting the transfer learning capabilities of TESSERA. See Appendix A.9 and Supplementary Fig. 9 for the general downstream application workflow.

2.1 Crop Type Classification Evaluation

To validate the effectiveness of TESSERA representations for agriculture, we performed a rigorous evaluation on a crop type classification task. We used the official Austrian INVEKOS dataset for the 2021-2022 growing season [1], which contains 17 aggregated crop classes.

Our evaluation methodology centered on comparing the performance of a lightweight classification head trained on TESSERA representations against two key baselines: (1) a traditional Random Forest classifier trained on engineered temporal features from raw satellite data, and (2) the PRESTO foundation model [79], where we applied the same lightweight head to its embeddings for a direct comparison. We assessed performance using multiple metrics, including F1 scores, under varying data availability scenarios, from 30% of labeled data down to one-shot learning. We also evaluated performance on a patch-based semantic segmentation task against several other foundation models. Further details on data processing, model architectures, and training protocols are provided in Section B.

2.2 Crop classification

Agricultural monitoring through EO is a cornerstone of managing global food security, providing critical data to stakeholders in the agricultural sector [86]. Crop type classification inform vital agricultural decisions by allowing accurate estimates of crop area and yield, and facilitating low-cost monitoring of diseases and pests over vast regions [31]. However, this task is notably more complex than many other land cover classification challenges, as it requires discerning subtle spectral and temporal variations that differentiate various crop classes [31]. For years, classifiers like Random Forest (RF) applied to hand-crafted composite time-series data have served as a standard baseline [31, 63].

Recent advances have shifted towards foundation models, which, unlike bespoke machine learning models that require region-specific training, offer generalizable representations that can be applied universally [4, 41]. Here, we demonstrate that TESSERA’s pre-computed representations set a new state-of-the-art. We compare TESSERA against two critical benchmarks: the traditional RF baseline and the leading pixel-based foundation model, PRESTO (Pretrained Remote Sensing Transformer) [79], a transformer-based model utilizing multi-modal inputs. Our results show TESSERA not only substantially outperforms both, but also offers a more streamlined and efficient application workflow.

To evaluate performance, we used the INVEKOS Austrian crop dataset from the 2021–2022 growing season, an extensive data set that covers 1850 km² east of Vienna with 154 different crop types, which we consolidated into 17 classes based on phenology and data availability [1]. TESSERA consistently surpasses both baselines across all training regimes, especially in low-data settings crucial for operational deployment where labeled data is scarce (Fig. 2). In pixel-wise classification tasks with training data splits from 1% to 30%, a simple multilayer perceptron (MLP) trained on TESSERA representations achieves higher average and balanced F1 scores than both RF and an identical MLP trained on PRESTO embeddings, with performance gains frequently exceeding 10% and 30%, respectively (Fig. 2a). TESSERA also maintains statistically significant advantages in one-shot and few-shot scenarios (Fig. 2b). Furthermore, TESSERA’s utility extends to patch-based semantic segmentation, where it consistently achieves higher mIoU and macro F1 scores compared to other leading foundation models across various patch sizes (Fig. 2c).

These performance gains are rooted in the superior quality of the learned representations. A UMAP analysis of the embedding space reveals that TESSERA’s representations form tighter, more semantically coherent clusters by crop class compared to those from PRESTO (Fig. 2d). This improved separability highlights TESSERA’s strength as a foundational tool for agricultural applications. Notably, using TESSERA’s compact, pre-computed 128-dimensional representations eliminates the need for the complex data preparation and feature engineering required by traditional methods like RF, democratizing access to high-performance crop classification.

2.3 Canopy height estimation

Canopy height is a key structural attribute of forests, closely linked to above-ground biomass and carbon stocks, and therefore important for climate mitigation efforts and carbon accounting frameworks like REDD+ [57, 26]. It reflects ecosystem function [56], forest age, and successional stage, while also serving as a proxy for habitat quality and vertical complexity—important determinants of biodiversity patterns and species richness [20, 15, 73]. Accurate, scalable canopy height models (CHMs) enable detection of forest degradation and disturbance arising from selective logging, disease outbreaks and storm damage.

Despite substantial progress in global canopy height mapping, current approaches exhibit limitations that restrict their reliability and applicability. GEDI, a spaceborne LiDAR system, provides spatially sparse measurements of canopy height that are used to train models for generating maps from optical and radar data; however, GEDI does not sample high latitudes [21], and produces unreliable estimates in mountainous terrain [27]. Radar-based CHMs (e.g., TanDEM-X) achieve global coverage but are affected by canopy height saturation in dense forests and depend on high-quality terrain models to isolate canopy height accurately. Optical data, while widely available, lacks inherent sensitivity to vertical structure and requires indirect proxies or fusion with information from active sensors [82, 50] - yet none of the three most widely used global CHM models currently implement sensor fusion [46, 68, 78]. Most approaches instead rely on a single sensor type and de-

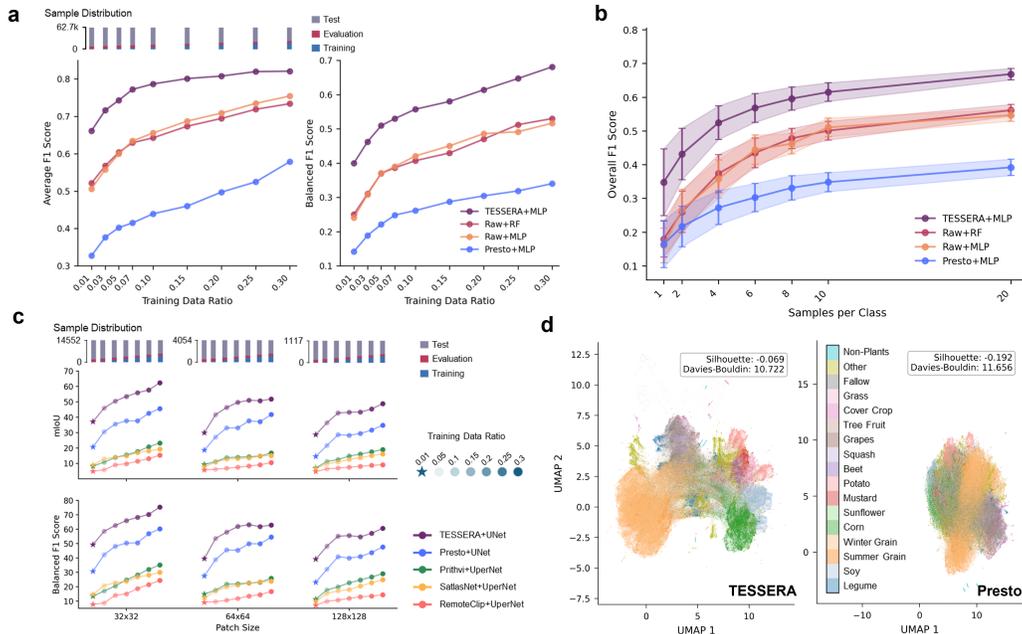


Figure 2: TESSERA representations achieve state-of-the-art performance in crop type classification. The evaluation was conducted on the 2022 Austrian INVEKOS dataset. **a**, Average and Balanced F1 scores for pixel-wise classification as a function of the training data ratio. TESSERA embeddings, coupled with a simple MLP, consistently outperform a Random Forest model trained on raw time-series data and the PRESTO foundation model. **b**, Overall F1 score in few-shot learning scenarios. TESSERA maintains a significant performance advantage even with very few training samples per class. Error bars represent the standard deviation over multiple runs. **c**, Patch-based semantic segmentation performance comparison. TESSERA representations, used with a UPerNet head, achieve higher mean Intersection over Union (mIoU) and Balanced F1 scores than other foundation models across various patch sizes and training data ratios. **d**, UMAP visualization of the embedding spaces for TESSERA and PRESTO for 17 crop classes. TESSERA’s embeddings exhibit clearer separation and more coherent clustering, as supported by superior Silhouette and Davies-Bouldin scores, indicating a more semantically meaningful representation space.

rive height predictions either from temporal composites or from large numbers of hand-engineered features (e.g., > 500 features in Potapov et al. [68]), which may generalize poorly across ecosystems. These models tend to under-predict tall canopies, especially in tropical forests where the majority of above-ground carbon is stored, and face challenges from persistent cloud cover, sensor saturation [45], and inconsistencies in input quality. As a result, significant discrepancies remain between global CHMs and airborne LiDAR benchmarks, with current products underestimating local variability and structural extremes [58]. Most CHMs are also static, failing to capture seasonal or interannual dynamics, and inconsistencies in methodology and sensor characteristics hinder cross-product integration. These limitations underscore the need for scalable, temporally resolved, and sensor-fused CHMs with quantified uncertainties to support forest monitoring, carbon accounting, and Earth system modelling [59, 65].

We evaluated TESSERA representations to predict the height of the airborne LiDAR-derived canopy at 10 m resolution within a 5×6 km area of tall old-growth tropical forest in the Danum Valley, Borneo[40]. To ensure spatial independence, we conducted four-fold spatial cross-validation, each time holding out a contiguous 50% of the region for testing and using the remaining 50% for training and validation. A 30-million-parameter UNet was trained using 64×64 pixel patches of the representations as input. To account for variability in training convergence, we performed three independent training runs per fold, each for 200 epochs, and retained the model with the best performance on the validation set. Accuracy metrics for the foundation models are reported as the average across all 12 runs. TESSERA achieved an R^2 of 0.66, a root mean squared error (RMSE) of 8.88 m, and a mean bias of -0.62 m.

We benchmarked TESSERA’s performance against state-of-the-art approaches through three comparisons (Fig. 3). First, we compared predictions against three global canopy height products: Potapov et al. [68] (‘GLAD’ map), Lang et al. [46] (‘ETH’ map), and Tolan et al. [78] (‘Meta’ map). Although local models generally outperform global products in regional evaluations [70], the pronounced saturation of these global models at much lower canopy heights—combined with RMSEs exceeding their reported global averages [68, 46, 78]—underscores the challenge of accurately mapping structurally complex, tall-canopy regions such as Danum Valley. Second, we compared TESSERA to the PRESTO [79] foundation model, using the same model architecture and training procedure but replacing the inputs of TESSERA with the representations of PRESTO. Third, we compared our results to regional wall-to-wall canopy height maps for Indonesia, Malaysia, and the Philippines generated by Lang et al. [47], which used deep convolutional neural networks trained on Sentinel-2 imagery and GEDI LiDAR reference data for regional canopy height and biomass estimation.

In all comparisons, TESSERA outperformed competing methods. Moreover, its representations do not require preprocessing, greatly enhancing reproducibility. Instead, only the region and year of interest are needed to replicate the input data pipeline. By fusing Sentinel-1 and Sentinel-2 data without requiring preprocessing or feature engineering, TESSERA’s 128-dimensional representations encode cloud-free, information-rich signals that effectively address key challenges in canopy height modelling.

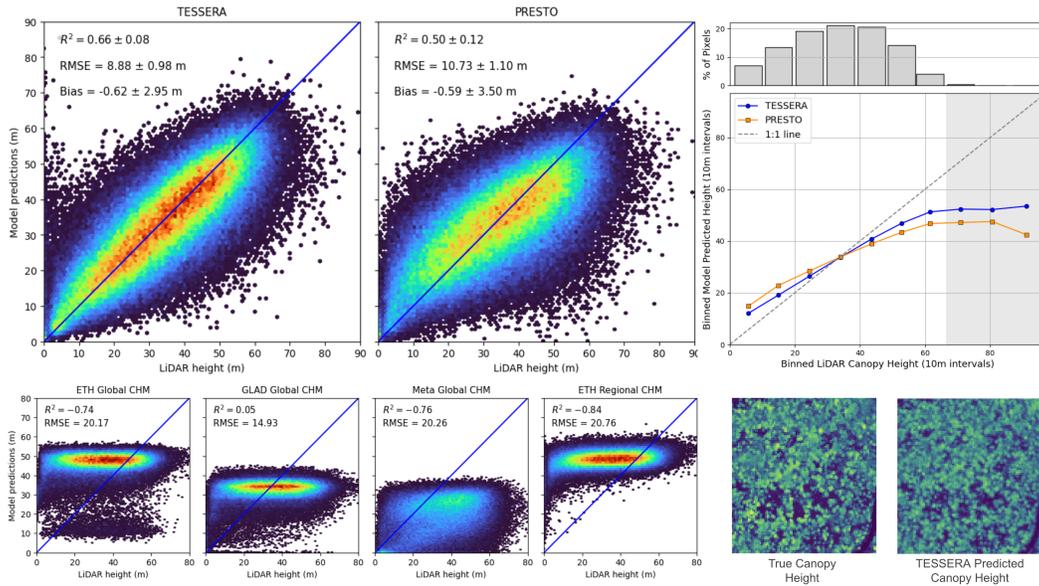


Figure 3: TESSERA representations outperform global and regional models in canopy height estimation. **Top row:** Density scatter plots compare predicted canopy heights from the TESSERA (left) and PRESTO (middle) models against airborne LiDAR-derived heights. TESSERA achieves higher accuracy than PRESTO, demonstrating superior performance in capturing structural forest attributes. The rightmost panel shows model bias across 10m canopy height bins, with LiDAR values binned on the x-axis and model predictions on the y-axis. TESSERA more closely follows the 1:1 line, particularly in mid- to high-canopy regimes, while PRESTO underestimates heights above 40m. Histogram above indicates pixel count distribution per bin. **Bottom row:** Comparison with three global canopy height products—ETH Global CHM, GLAD (Potapov) Global CHM, Meta (Tolan) Global CHM, and ETH Regional CHM—evaluated on the same test region. All four products exhibit substantially lower predictive performance ($R^2 < 0.05$, $RMSE > 14m$), including negative correlations for three out of four models. These results underscore the difficulty of accurately estimating tall tropical canopies using global models trained without local data. **Rightmost column:** Spatial comparison between true canopy height (LiDAR) and TESSERA-predicted canopy height for a representative test patch, illustrating the model’s capacity to capture fine-scale structural variation while still saturating for taller canopy values.

2.4 Burned Area Detection

Wildfire frequency and severity are projected to increase in many parts of the world due to climate change [3, 39, 66]. Accurate maps of burned areas (BAs) and burn severity are critical for monitoring wildfire trends [6], predicting wildfire occurrence [17], assessing ecosystem recovery [89], and establishing effective fire management strategies [33]. Remote sensing is widely used to detect burned areas and assess severity regionally and globally [44].

Most existing BA mapping approaches rely on time series analysis or direct comparison between pre-fire and post-fire images, often requiring careful selection of suitable cloud-free images, manually tuned thresholds, and expert input [44, 67]. These methods typically use spectral indices such as the Normalized Burn Ratio (NBR), though many studies have also explored the design of custom features and applied spectral dimensionality reduction techniques to preserve key fire-related signals [44]. The majority of existing methods use optical sensors [44], which can detect fire-specific spectral signatures but are susceptible to cloud cover [34], particularly problematic for fire detection in ecosystems where vegetation rapidly re-greens [35]. SAR offers complementary sensitivity to vegetation structure and is unaffected by cloud cover [44, 76, 81]. Recent reviews emphasize the need for sensor fusion, combining passive and active sensors to improve accuracy and robustness [44, 18]. Sentinel sensors in particular, when combined together at 10 m resolution, have shown promise in detecting BA from small fires that are often left out of coarse, global BA products [76, 81, 71]. TESSERA extends the philosophy of traditional approaches by encoding temporal dynamics, while addressing major limitations through sensor fusion and the use of abstract, data-driven features learned directly from large volumes of combined Sentinel-1 and Sentinel-2 imagery (Fig. 4).

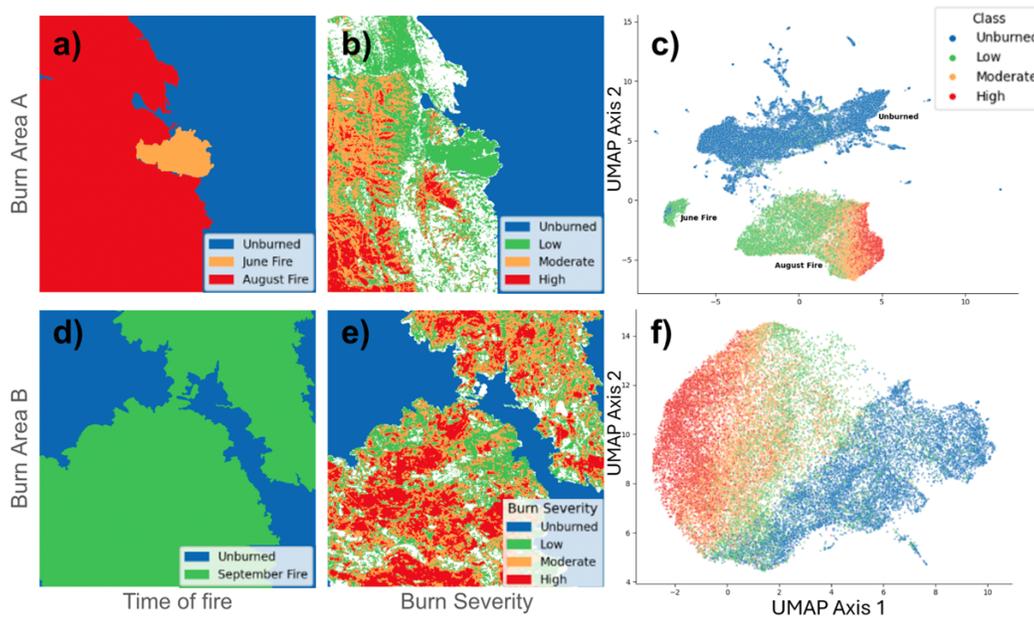


Figure 4: **TESSERA representations capture wildfire extent, timing, and severity across two burned regions.** (a–c) Burn Area A, which experienced two distinct fires in June and August 2021. (a) Time of fire map based on MTBS, showing unburned areas (blue), the June fire (orange), and the August fire (red). (b) Burn severity classification for the same region, with unburned (blue), low (green), moderate (orange), and high (red) severity. (c) UMAP projection of TESSERA embeddings from 2021, with colors corresponding to burn severity. Embeddings form well-separated clusters for unburned, June-, and August-burned areas, with a visible internal gradient reflecting severity within August fire. (d–f) Burn Area B, which experienced a single large fire in September 2021. (d) Time of fire map showing unburned (blue) and burned (green) areas. (e) Burn severity classification showing heterogeneous severity across the landscape. (f) UMAP projection for Burn Area B, revealing fuzzier boundaries between burned and unburned pixels, though the severity gradient remains evident.

We tested whether TESSERA’s temporal representations capture fire-induced changes in spectral and radar signals. We evaluated its ability to detect burned areas and burn severity using two randomly selected 15×15 km regions from the Burn Scar Benchmarking dataset, each affected by wildfires in California during 2021. TESSERA representations for 2021 were obtained and UMAP dimensionality reduction was applied to a random subsample of pixels. Fire boundaries and burn severity data were sourced from the Monitoring Trends in Burn Severity (MTBS) dataset [2]. The resulting UMAP projections demonstrated strong separability in two-dimensional space, visually highlighting TESSERA’s capacity to detect disturbances at the pixel level. Burn Area A experienced two distinct wildfires—in June and August—while Burn Area B underwent a single, large fire affecting mostly forested areas, resulting in heterogeneous burn severity. The UMAP revealed separability across three domains (Fig. 4):

- **Burned area detection:** In both regions, burned pixels formed distinct clusters from unburned ones, with clearer separation in Burn Area A.
- **Temporal fire differentiation:** In Burn Area A, embeddings separated pixels burned in June vs. August, suggesting temporal sensitivity to disturbance events.
- **Burn severity gradient:** In both regions, a gradient structure in UMAP space corresponded to burn severity, indicating that TESSERA encodes disturbance magnitude.

Because TESSERA encodes annual time series of spectral and radar signals for each pixel, it inherently captures changes in pixel trajectories without requiring prior preprocessing or the selection of cloud-free before-and-after images. This capacity would support a simplified, scalable pipeline for large-area disturbance monitoring.

2.5 Above-ground biomass estimation

Forests play a central role in the global carbon cycle, creating a pressing need for reliable and scalable AGB estimates for climate modeling, carbon accounting, and informed strategies for conservation and sustainable land use [36]. Accurate estimation of above-ground biomass (AGB) is essential for quantifying forest carbon stocks and monitoring temporal changes due to deforestation, degradation, and recovery [22].

We evaluated the performance of TESSERA in predicting AGB against two leading models on a benchmark dataset. Firstly, we considered the geospatial foundation model SpectralGPT [37], which was found in the Pangaea benchmark [53] to outperform other tested foundation models on the BioMassters dataset [60]. Secondly, as the state-of-the-art bespoke model, we considered the winner of the Biomasters competition, a UNet with a temporal attention encoder [60]. The BioMassters dataset comprises 11,462 image patches representing above-ground biomass across Finnish forests from 2017 to 2021. Ground-truth AGB values are derived from airborne LiDAR and aerial imagery, using calibrated allometric equations informed by extensive field plots by the Finnish Forest Centre and National Land Survey. Each patch covers a 2560×2560 m area, with 10×10 m spatial resolution per pixel. The fine spatial granularity of the dataset makes the estimation task more challenging, as it introduces a long, thin tail of high AGB values that would not appear in coarser-resolution data.

A comparison of RMSE across varying training label fractions (Fig. 5a) shows that a UNet trained on the TESSERA representations consistently surpasses a UPerNet trained on the SpectralGPT representations, achieving lower prediction errors on 1% of the training set than SpectralGPT on the full training set. These results underscore the robustness of TESSERA in low-label regimes. Moreover, while the best task-specific model trained on the full dataset outperforms TESSERA, the gap is moderate, even with low label availability. In general, bespoke models are inherently difficult to surpass in fully supervised settings. It is worth noting that the comparison to benchmarks is not entirely straightforward, as the satellite data used in the benchmark models differs from the TESSERA input. However, these discrepancies can be considered as an inherent part of the respective processing pipelines and hence largely unavoidable.

2.6 Stocking indices in voluntary carbon markets

The most widely applied standard in the voluntary carbon market (VCM) for reforestation and afforestation projects of degraded pastures requires project proponents to demonstrate that their chosen remotely-sensed “stocking index” is correlated with the AGB (REF). Upon acceptance, project

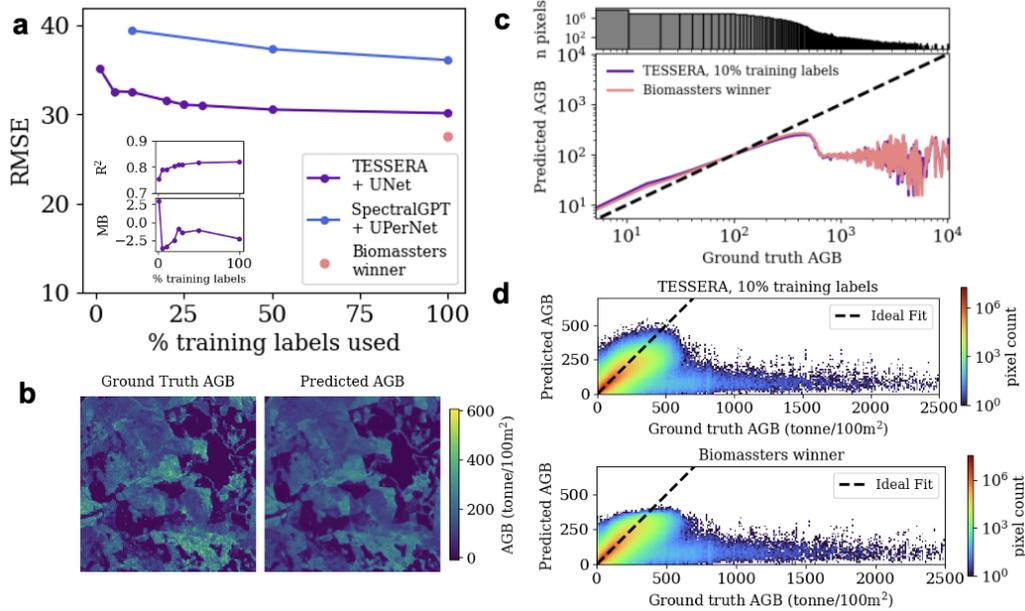


Figure 5: **TESSERA representations demonstrate robustness under limited label availability.** **a**, TESSERA achieves lower RMSE values than the foundation model benchmark SpectralGBT (RMSE values from Ref. [53]). Following the evaluation protocol of Ref. [53], RMSE is computed on a fixed test set of 2773 patches, with the remaining data randomly split into training and validation sets (80:20). The UNet model is trained on the TESSERA representations for 80 epochs, with the checkpoint achieving the lowest validation loss selected as the final model. For limited labels testing, a subset of labels is randomly chosen from the training set. The inset depicts the coefficient of determination (R^2) and mean bias (MB) as functions of label availability. **b**, Predicted AGB map for an example patch using 10% of the training labels. **c**, The mean predicted AGB values show increasing bias at higher ground-truth AGB values, which are sparse in the dataset. **d**, Pixel-wise comparison of predicted versus ground-truth AGB values for the test data with TESSERA+UNet (top) and Biomasters winner (bottom).

developers can use the index to quantify project additionality and to assess project performance in-between the 5-yearly ground-truth assessments. There is considerable freedom in choosing a stocking index, which can be either a pre-existing or proprietary canopy height or AGB product. Here, we evaluate TESSERA’s performance against five widely used, globally available canopy height and AGB products often used in the VCM. The *in situ* data used to compare the different stocking indices was collected by ICRAF in 38 sites in Para State, Brazil, that were degraded pastureland until being converted into agroforestry systems during 1980-2010 (Atzberger et al., 2025). To provide a fair comparison, stocking indices were rescaled using a linear transformation to minimize the RMSE with respect to the *in situ* data (details in the SI).

3 Conclusion

In this work, we introduce TESSERA, an innovative Remote Sensing Foundation Model that achieves worldwide coverage at 10-meter resolution through self-supervised learning applied to pixel-level satellite time series data. The representations generated by TESSERA establish new performance standards, while our open source methodology ensures widespread access to superior, high-resolution representations. We evaluated TESSERA’s capabilities across five distinct tasks, benchmarking our approach against leading task-specific models and existing foundation models. The findings demonstrate that TESSERA surpasses both conventional RF baselines and current state-of-the-art geospatial foundation models in these diverse downstream applications.

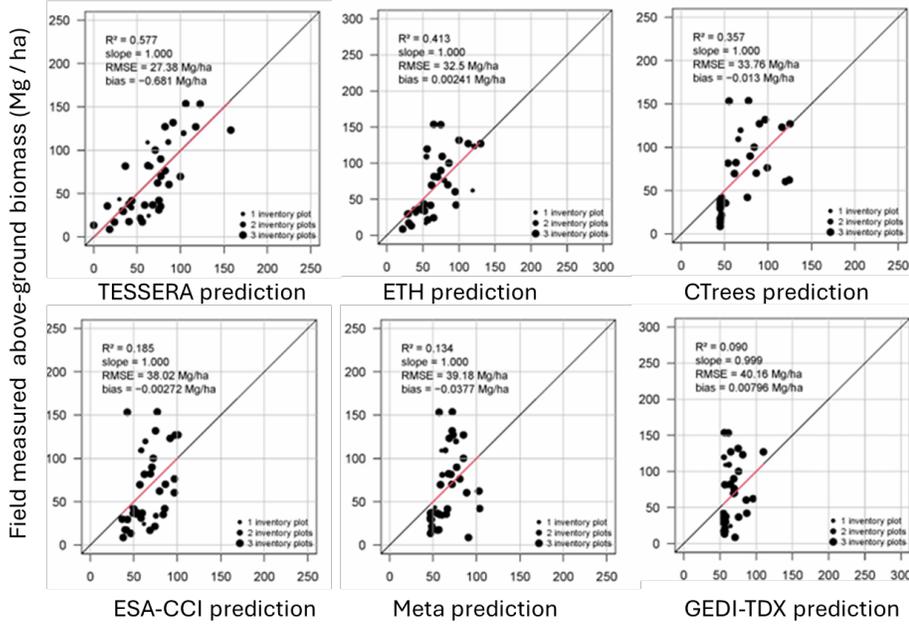


Figure 6: **TESSERA outperforms commonly used stocking indices in the Voluntary Carbon Market.** Comparison of in-situ aboveground biomass measurements in 38 agro-forestry plots in Para State, Brazil with TESSERA (top left) and five other remote sensing products. Note that the different products were originally provided in different units and converted to AGB as follows: (a) TESSERA: difference in relative heights of GEDI estimates of RH90 and RH10 converted to AGB using $X * 15.502 + 160.5$; (b) ETH canopy height converted using $3.4 X$; (c) CTrees canopy height converted using $1.806 X + 44.9$; (d) ESA: AGB converted using $0.407 X + 34$; (e) Meta canopy height converted using $3.723 X + 46.6$; (f) GEDI-Tandem-X above ground biomass prediction converted using $0.3X + 54.3$. Details of conversation functions provided in SI

We next provide supplementary materials and methods about the details of how TESSERA is implemented.

A Earth Observation data

Our input data consists of Sentinel-2 optical satellite imagery (Level-2A bottom-of-atmosphere) and Sentinel-1 Synthetic Aperture Radar (SAR, Radiometrically Terrain Corrected).

A.1 Data Representation

We consider remote sensing data with C channels (spectral bands or polarizations). Each data tile R_t at time t is represented as a 3D array with dimensions:

$$R_t \in \mathbb{R}^{W \times H \times C} \quad (1)$$

where W is the width (longitude dimension), H is the height (latitude dimension), and C is the number of spectral channels.

Each tile is accompanied by a corresponding binary mask V_t of dimensions:

$$V_t \in \{0, 1\}^{W \times H} \quad (2)$$

where $V_t(i, j) = 0$ indicates clouding or missing data for the pixel at spatial coordinates (i, j) , and $V_t(i, j) = 1$ indicates valid data.

A.2 Temporal Data Stacking

We stack spatially aligned tiles over a time period spanning T time steps ($t = 0, 1, \dots, T - 1$). The temporal data stack is defined as:

$$\mathbf{D} = [R_0, R_1, \dots, R_{T-1}] \quad (3)$$

$$\mathbf{M} = [V_0, V_1, \dots, V_{T-1}] \quad (4)$$

A.3 Time Series Extraction and d-pixel Definition

For a given spatial location (i, j) and spectral channel c , the time series $S_{i,j,c}$ represents all channel c values at coordinates (i, j) over the entire time period:

$$S_{i,j,c} = [R_0(i, j, c), R_1(i, j, c), \dots, R_{T-1}(i, j, c)] \quad (5)$$

We define a *d-pixel* $P_{i,j}$ as the collection of all spectral channels by timesteps at a given spatial location (i, j) :

$$P_{i,j} = S(i, j) \quad (6)$$

In other words, the d-pixel provides all spectral values (Sentinel-2) or backscatter values (Sentinel-1) at a given point over time. Note that d-pixels are potentially sparse and have an accompanying mask $m_{i,j}$ that indicates the timesteps in which there is valid data:

$$m_{i,j} = [v_{i,j,0}, v_{i,j,1}, \dots, v_{i,j,T-1}] \quad (7)$$

Variable	Dimensions	Description
R_t	$W \times H \times C$	Remote sensing data tile at time t
V_t	$W \times H$	Binary mask for tile at time t (1=valid, 0=invalid/clouded)
\mathbf{D}	$T \times W \times H \times C$	Complete temporal data stack
\mathbf{M}	$T \times W \times H$	Complete temporal mask stack
$S_{i,j,c}$	$T \times 1$	Time series for channel c at location (i, j)
$P_{i,j}$	$C \times T$	d-pixel: all spectral channels by timesteps at location (i, j)
$m_{i,j}$	$T \times 1$	d-pixel mask indicating valid timesteps at location (i, j)
W, H	scalar	Spatial dimensions (width, height)
C	scalar	Number of spectral channels
T	scalar	Number of time steps

Table 1: Variable definitions for remote sensing data structure

A.4 The d-pixel Representation

We represent the rich information contained within satellite image time series at the level of individual pixels using the aforementioned “d-pixel”. For a single 10-meter geographic pixel location, it consolidates its complete annual multi-spectral (from Sentinel-2) and SAR backscatter (from Sentinel-1) time series. This structure explicitly preserves the inherent spectral and temporal dimensions of the data (Supplementary Fig. 1 a illustrates the d-pixel concept).

In this array, each row corresponds to a distinct observation date throughout the year, ordered chronologically. Each column corresponds to a specific spectral band or polarization. The value at the intersection of a given row and column thus represents the measured reflectance or backscatter for that specific band/polarization on that particular date for that single pixel.

A critical feature of the d-pixel is its capacity to accommodate the frequent data gaps encountered in remote sensing due to cloud cover (for optical data) or other atmospheric interference and acquisition irregularities. Observations identified as invalid (e.g., via Sentinel-2 scene classification layers for clouds, or missing Sentinel-1 acquisitions) are masked within the d-pixel array. These missing entries are explicitly handled by downstream model components, for instance, by being ignored during attention computation or managed during the temporal sampling process described in Appendix A.7.3.

A.5 Dual-Encoder Architecture

Given the distinct nature of Sentinel-1 SAR and Sentinel-2 MSI data, TESSERA employs two separate, parallel Transformer-based encoder branches.

- **Sentinel-2 MSI Encoder:** This branch processes time series of 10 spectral bands from Sentinel-2. We used blue (B2), green (B3), red (B4), red edges 1–3 (B5, B6, B7), near-infrared (B8, B8A), and shortwave infrared (B11, B12).
- **Sentinel-1 SAR Encoder:** This branch processes time series of 2 polarizations from Sentinel-1 (VV and VH).

Each encoder begins by linearly embedding the input features (spectral bands or polarizations) for each time step. To preserve sequence order and incorporate temporal context, learnable positional encodings based on the Day-of-Year (DOY) of each observation are added to these embeddings. The core of each encoder consists of a stack of 8 standard Transformer blocks [80], featuring multi-head self-attention and feed-forward layers to learn temporal patterns within the data streams.

To derive a single vector summarizing the entire time series for each modality, an attention-pooling layer weighs the importance of different time steps before aggregation. The resulting modality-specific representations (one from the S1 encoder, one from the S2 encoder) are then fused using a multi-layer perceptron.

A.6 Projector Network

The fused representation from the dual-encoder stage is subsequently fed into a large projector network. This projector is a five-layer MLP, with each hidden layer having 16,384 dimensions. This significant expansion in dimensionality is crucial, as suggested by the original Barlow Twins work [90], to enable effective redundancy reduction during the self-supervised loss computation. The final output of the projector for each input d-pixel is an embedding, which is then used in the loss calculation. For downstream tasks, we typically use the 128-dimensional output from the fusion MLP (before the projector) as the final pixel representation. The TESSERA encoder (up to the fusion MLP) has approximately 40 million parameters, while the projector accounts for the majority of the model’s ~ 1.4 billion parameters.

A.7 Self-Supervised Training

A.7.1 Augmented View Generation

The TESSERA model is trained using a modified Barlow Twins objective function [90]. For this objective, two distorted views, denoted as Y_A and Y_B , are generated for each input d-pixel. In TESSERA, these views are created by independently running the temporal sampling and preprocessing pipeline twice for the Sentinel-1 and Sentinel-2 data associated with a given d-pixel. This process involves:

1. For each view, independently sampling a fixed number of valid observation dates from the annual Sentinel-2 time series (10 spectral bands).
2. For each view, independently sampling a fixed number of valid observation dates from the annual Sentinel-1 time series (2 polarizations).

These views represent different, valid, but inherently incomplete glimpses of the pixel’s true temporal-spectral evolution, akin to observing the same location through intermittent cloud cover or from different satellite passes at different times. The model learns by reconciling these partial views. The inherent differences between the Sentinel-1 SAR and Sentinel-2 MSI modalities further provide diverse perspectives on the same underlying physical processes. Thus, our augmentations are fundamentally about sampling from the available, inherently incomplete information streams, rather than artificially distorting a complete input.

A.7.2 Loss Function

The network processes these two views (Y_A, Y_B) through the dual-encoder and the projector to produce batch-normalized embeddings Z_A and Z_B . The standard Barlow Twins loss function, \mathcal{L}_{BT} , is defined as [90]:

$$\mathcal{L}_{BT} = \sum_i (1 - C_{ii})^2 + \lambda_{BT} \sum_i \sum_{j \neq i} C_{ij}^2 \quad (8)$$

Here, C is the cross-correlation matrix computed between the batch-normalized embeddings Z_A and Z_B . The indices i and j iterate over the dimensions of the embedding vectors. The first term (invariance term) encourages similar representations for different views of the same input ($C_{ii} \rightarrow 1$). The second term (redundancy reduction term) promotes informational efficiency by minimizing correlation between different embedding dimensions ($C_{ij} \rightarrow 0$ for $i \neq j$), weighted by λ_{BT} .

To further enhance model robustness and mitigate overfitting, TESSERA incorporates an additional mix-up regularization term, \mathcal{L}_{MIX} , inspired by Bandara et al. [9]. This involves shuffling one set of views (e.g., Y_B) along the batch dimension to create $Y_S = \text{Shuffle}(Y_B)$, then generating mixed views $Y_M = \alpha_{mix} Y_A + (1 - \alpha_{mix}) Y_S$, where $\alpha_{mix} \sim \text{Beta}(\beta_p, \beta_p)$. The embeddings Z_M and Z_S are obtained. The mix-up loss penalizes deviations from the assumption that a linear interpolation in input space corresponds to a linear interpolation in embedding space:

$$C_{target}^{MA} = \alpha_{mix} (Z_A)^T Z_A + (1 - \alpha_{mix}) (Z_S)^T Z_A \quad (9)$$

$$C_{target}^{MS} = \alpha_{mix} (Z_A)^T Z_S + (1 - \alpha_{mix}) (Z_S)^T Z_S \quad (10)$$

$$\mathcal{L}_{MIX} = \frac{1}{2} (\|C^{MA} - C_{target}^{MA}\|_F^2 + \|C^{MS} - C_{target}^{MS}\|_F^2) \quad (11)$$

where $C^{MA} = (Z_M)^T Z_A$ and $C^{MS} = (Z_M)^T Z_S$ are the actual cross-correlation matrices from the model’s outputs. The total loss function optimized during the training of TESSERA is a weighted sum:

$$\mathcal{L}_{total} = \mathcal{L}_{BT} + \lambda_{mix} \mathcal{L}_{MIX} \quad (12)$$

where λ_{mix} controls the strength of the mix-up regularization. We found $\lambda_{BT} = 5 \times 10^{-3}$ and $\lambda_{mix} = 1.0$ to be effective.

A.7.3 Pretraining Details

The TESSERA model (Alpha version) was pretrained using approximately 0.8 billion d-pixel samples derived from 3,012 globally distributed Military Grid Reference System (MGRS) tiles. For pretraining, these d-pixels were generated from Sentinel-1 and Sentinel-2 data that were spatially downsampled by a factor of 400 from their native 10-meter resolution. Each d-pixel for Sentinel-2 contained time series for 10 spectral bands, and for Sentinel-1, 2 polarizations (VV and VH). The sequence length for the Transformer encoders was fixed at 40 timesteps for both modalities.

For each pixel location, after filtering invalid observations (e.g., due to cloud cover for Sentinel-2), we performed sparse temporal sampling. This involves randomly selecting a fixed number of 40 valid observation dates from the year’s data. If fewer than 40 valid dates exist, sampling is performed with replacement. This strategy standardizes the input sequence length and serves as a key data augmentation mechanism, building invariance to data gaps and teaching the model that the underlying signal persists regardless of the specific dates observed. The temporal context of each sampled observation was encoded by transforming its normalized Day-of-Year (DOY) into sine and cosine features, which were then concatenated with the corresponding spectral or backscatter measurements. Finally, these values were standardized using global statistics to stabilize training.

The model was trained for 1 epoch over the entire dataset, which corresponded to approximately 3,000 GPU hours on 8 AMD MI300X GPUs (192GB memory each). We used PyTorch with Fully

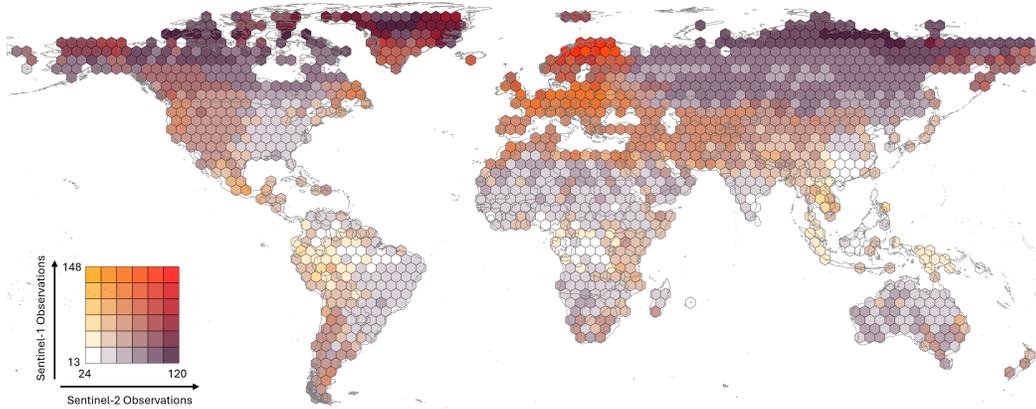


Figure 7: **Geographical distribution and data density of the training dataset.** TESSERA was trained on over 3,000 MGRS tiles distributed globally from 2017 to 2024. The color of each hexagon in the map corresponds to the number of valid observation days for Sentinel-1 (S1) and Sentinel-2 (S2), as defined by the bivariate color legend. This visualization highlights the density of combined S1 and S2 observations available for training across different regions, ensuring the model learns from a diverse range of geographical and environmental conditions.

Sharded Data Parallel (FSDP) and Automatic Mixed Precision (AMP) enabled. The AdamW optimizer was employed with a base learning rate of 0.002 and weight decay of 1×10^{-6} . The learning rate schedule included a linear warmup over the initial 10% of steps, a plateau for the next 20%, followed by a cosine decay. The global batch size was 8192.

A crucial aspect of our training methodology is a data shuffling strategy, essential for learning globally representative features from a vast and geographically diverse dataset. Given that d-pixels within an individual MGRS tile exhibit high spatial autocorrelation, a naive sequential or locally-shuffled data loading process would expose the model to strong geographic biases in each batch. To overcome this, we developed a custom data processing pipeline to implement a truly global shuffle across all ≈ 0.8 billion training samples, which constituted over 2TB of initial d-pixel data. This pipeline is illustrated conceptually in Fig. 8a.

The impact of this approach on training stability is empirically demonstrated in Fig. 8b. Compared to a conventional, localized shuffling strategy which results in a highly volatile loss curve (top plot), our global shuffling strategy yields a markedly smoother and more stable convergence (bottom plot). This enhanced stability is fundamental for robust convergence and for preventing the model from overfitting to regional characteristics.

Operationally, the process begins with the aggregation of d-pixels from all MGRS tiles into a single, comprehensive pool. A global shuffling operation is performed on this pool, a critical step to break the spatial contiguity of data from individual tiles and ensure each training batch contains a diverse mix of geographical and environmental contexts. Following this global shuffle, the data augmentation required by the Barlow Twins framework is applied. As detailed in Appendix 5.1, this involves generating two distinct augmented views (e.g., Y^A and Y^B) for each d-pixel.

To manage the significant I/O demands of shuffling and augmenting such a large volume of data, we developed the pipeline as a custom Rust binary. This high-performance binary handles the reading of raw d-pixel data, executes the global shuffle, and prepares the data for augmentation. The resulting pairs of augmented d-pixels are then serialized into a compact, pickle-like file format. These files are organized into manageable chunks and loaded by PyTorch *DataLoader* workers, which stream the data and assemble the final training batches. This end-to-end pipeline ensures that each batch presented to the model is a well-shuffled, globally diverse representation of Earth’s surface characteristics, which is fundamental for training a robust pixel-wise foundation model like TESSERA.

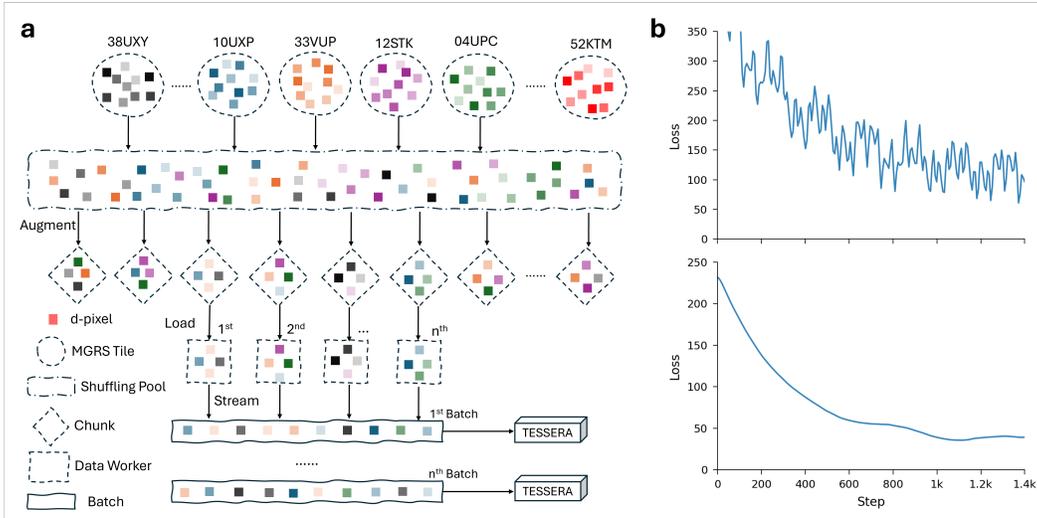


Figure 8: **Overview of the data shuffling pipeline and its impact on training stability.** **a**, Schematic of the data shuffling and loading process. D-pixels (colored squares) from thousands of MGRS tiles are first aggregated into a global pool. A custom Rust binary performs a global shuffle on this multi-terabyte dataset before applying augmentations. The processed data is then organized into chunks and streamed by data workers to form well-shuffled, globally diverse training batches. **b**, Comparison of training loss curves. The top plot shows the volatile loss progression typical of a localized shuffling strategy, which is susceptible to geographic bias. The bottom plot shows the significantly smoother and more stable loss curve achieved with our global shuffling pipeline, demonstrating more effective and robust model convergence.

A.8 Global Representation Map Generation

A primary output of the TESSERA project is the generation of annual global representation maps at 10-meter resolution for the years 2017-2024.

A.8.1 Model Inference

To generate these maps, the pretrained and frozen TESSERA dual-encoder (excluding the projector) is used. For each 10-meter pixel on the globe and for each year:

1. The full Sentinel-1 and Sentinel-2 time series data at 10-meter resolution are acquired and preprocessed to form d-pixels. Unlike pretraining, no spatial downsampling is performed at this stage.
2. A fixed number of 40 timesteps are sampled from the valid observations within the year for both Sentinel-1 and Sentinel-2 data, along with their DOY positional encodings.
3. These sampled time series are fed into their respective frozen TESSERA encoders.
4. The outputs from the S1 and S2 encoders are fused by the MLP, producing a 128-dimensional representation vector for that pixel for that year.

This process is repeated for all land pixels globally to create an annual representation map of shape (H, W, 128), where H and W are the dimensions of the global 10-meter grid.

A.8.2 Data Product and Accessibility

The resulting product consists of eight such global representation maps, one for each year from 2017 to 2024. These maps are intended to be released publicly, alongside the open-sourced TESSERA model parameters and generation code. This “Model-as-Data” approach significantly lowers the entry barrier for users, who can treat these representation maps as conventional multi-channel images. They can be directly ingested by downstream models without requiring users to process raw

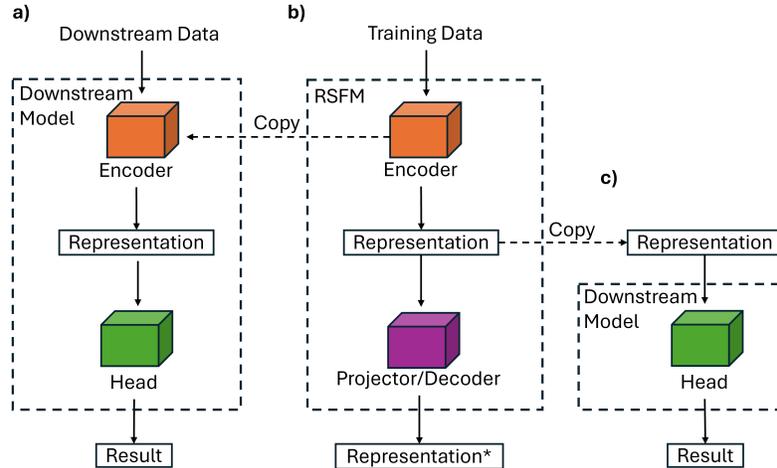


Figure 9: **Illustration of RSFM training and downstream task applications.** (b) shows the training process of a typical RSFM, where both the encoder and the projector (or decoder) are updated. (a) depicts a common application method: the RSFM encoder is extracted and combined with a task-specific head to form a new model for downstream tasks. Users can then fine-tune the entire model or just the head. (c) presents the emerging paradigm championed by TESSERA, where pre-generated representations from the RSFM encoder are used directly as input. This results in extremely lightweight downstream models that consist only of a task-specific head.

satellite data or run the TESSERA model themselves, thus democratizing access to advanced RSFM capabilities.

A.9 Downstream Task Application Methodology

A core motivation for self-supervised learning with foundation models is the creation of task-agnostic feature representations that can be effectively transferred to various downstream tasks, particularly in scenarios with limited labeled data. Having pretrained TESSERA on large unlabeled datasets, we evaluated the utility of its learned representations across different remote sensing applications. The evaluation methodology involves using the pretrained TESSERA encoders as fixed feature extractors.

The pipeline for applying TESSERA representations to downstream tasks is as follows (illustrated in Supplementary Fig. 9):

- 1. Load Pretrained Encoder:** The weights from the saved pretraining checkpoint are loaded into the TESSERA dual-encoder architecture. The parameters of these encoders are frozen and are not updated during downstream model training. This ensures that downstream performance directly reflects the quality of the fixed representations.
- 2. Prepare Labeled Downstream Data:** The specific labeled dataset for the target task (e.g., pixel-level crop type labels, canopy height measurements, or land change polygons) is prepared.
- 3. Extract Representations:** For each input sample (e.g., pixel, object, or patch) in the labeled dataset, its corresponding Sentinel-1 and Sentinel-2 time series data for the relevant year undergo the same preprocessing and d-pixel creation pipeline used during TESSERA inference (i.e., using 10m resolution data and temporal sampling to 40 timesteps). The preprocessed sequences are passed through the respective frozen encoders, and their outputs are fused by the MLP to generate the final 128-dimensional TESSERA representation for that sample. To enhance stability, this extraction can be repeated multiple times with different random temporal samplings, and the resulting representations averaged, although for many tasks, a single extraction is sufficient.

4. **Design Task-Specific Head:** A lightweight, task-specific neural network module (the "head") is designed. This head takes the extracted TESSERA representations as input.
 - For pixel-wise classification (e.g., crop classification), the head is typically a shallow MLP (e.g., 1-3 layers) ending in a softmax output layer.
 - For pixel-wise regression (e.g., canopy height regression), the head is usually an MLP ending in a single linear output neuron.
 - For tasks requiring spatial context from the representations (e.g., canopy height mapping over an area, semantic segmentation), the input to the head can be a patch of TESSERA representations (e.g., $64 \times 64 \times 128$). The head might then be a convolutional architecture, such as a UNet, that processes these spatial feature maps to produce dense predictions. For land change detection, a simple approach involves computing the dot product between representations from two different years.
5. **Train Downstream Head:** Only the parameters of this newly defined task head are trained using the extracted TESSERA representations as input features and the corresponding labels. Standard supervised learning techniques, optimizers (e.g., Adam), and task-appropriate loss functions (e.g., Cross-Entropy for classification, Mean Squared Error for regression) are used. This training typically requires significantly less labeled data and computational power compared to training a deep model from scratch.
6. **Evaluation:** Once the head is trained, inference is performed on a test set by extracting TESSERA representations for the test samples and passing them through the trained head. Performance is evaluated using standard metrics relevant to the task.

This standardized workflow allows for robust assessment of TESSERA representations across diverse applications, demonstrating their value as foundational features for geospatial analysis.

B Downstream Task: Austrian Crop Classification

The following section details the experimental setup for the crop classification task presented in the main text.

B.1 Dataset and Preprocessing

We used the publicly available INVEKOS dataset for Austria, focusing on the 2022 growing season [1]. The dataset originally contained 154 crop types, which we grouped into 17 broader classes (e.g., merging different varieties of wheat) based on phenological similarity and sample availability to ensure robust training and evaluation. For each pixel in the dataset, we extracted its corresponding 128-dimensional TESSERA representation from our generated 2022 global representation map.

B.2 Pixel-wise Classification Baselines

To provide a comprehensive performance comparison, we implemented three distinct models for pixel-wise classification:

- **TESSERA + MLP:** The primary model, where the frozen 128-dimensional TESSERA representations were used as input to a simple MLP. The MLP consisted of two hidden layers with 256 and 128 neurons, respectively, using ReLU activation functions, followed by a softmax output layer for the 17 classes.
- **PRESTO + MLP:** For a direct and fair comparison with the closest foundation model, we used the official pre-trained PRESTO model [79] to generate its pixel embeddings. These embeddings were then fed into an identical MLP head as the one used for TESSERA.
- **Random Forest:** As a traditional baseline, we trained a Random Forest classifier directly on raw time-series data. For each pixel, we utilized all available Sentinel-1 (2 polarizations) and Sentinel-2 (10 spectral bands) observations throughout the year. The temporal and spectral/polarization dimensions were flattened and concatenated to form a single 1256-dimensional feature vector. The RF model consisted of 200 trees, with other parameters set to standard values.

For the experiments shown in the main text’s figure, we trained these models on randomly selected subsets of the data (from 1% to 30% for panel a; specified samples-per-class for panel b), using a fixed validation set for hyperparameter tuning and a held-out test set for final evaluation.

B.3 Patch-wise Semantic Segmentation

To assess spatial-contextual performance, we conducted a semantic segmentation experiment. The approach varied based on the foundation model’s architecture:

- For pixel-based foundation models like **TESSERA** and **PRESTO**, which do not explicitly model spatial context in their representations, we first constructed representation patches (e.g., of size $64 \times 64 \times 128$) from the pixel-wise embeddings. These patches were then fed into a standard **UNet** architecture, which learns to model the spatial relationships between the representations to produce a segmentation map.
- For other foundation models that are inherently patch-based (e.g., Prithvi [75], Satlas [11]), their encoders already process image patches. Therefore, for these models, we only needed to attach and train a **UPerNet decoder head** to their frozen encoders to generate the final segmentation outputs.

Performance was measured using mean Intersection over Union (mIoU) and macro F1 scores.

B.4 Embedding Space Analysis

The 2D visualizations shown in the main text were generated by applying the Uniform Manifold Approximation and Projection (UMAP) algorithm to the representations of every pixel within the Austrian study area. Specifically, the entire representation map (e.g., an array of shape $H \times W \times 128$, where H and W are the height and width of the region) was used as direct input to UMAP. This was performed for both TESSERA and PRESTO representations.

To quantitatively measure the quality of the clustering in the original 128-dimensional space, we calculated the Silhouette score and the Davies-Bouldin Index.

- The **Silhouette score**, $s(i)$, for a single data point i measures how similar it is to its own cluster compared to other clusters. It is defined as:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \tag{13}$$

where $a(i)$ is the mean distance between i and all other points in the same cluster, and $b(i)$ is the mean distance from i to all points in the nearest neighboring cluster. The score ranges from -1 to 1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.

- The **Davies-Bouldin Index** (DBI) evaluates clustering quality by computing the ratio of within-cluster scatter to between-cluster separation. For a set of k clusters, it is defined as:

$$DBI = \frac{1}{k} \sum_{i=1}^k \max_{j \neq i} \left(\frac{\sigma_i + \sigma_j}{d(c_i, c_j)} \right) \tag{14}$$

where σ_i is the average distance of all points in cluster i to their centroid c_i , and $d(c_i, c_j)$ is the distance between the centroids of clusters i and j . Lower DBI values indicate better clustering, with a score of 0 representing the ideal case where clusters are compact and well-separated.

References and Notes

- [1] INVEKOS Schläge Österreich 2022 - Open Government Data Austria.
- [2] Monitoring Trends in Burn Severity (ver. 12.0, April 2025) - ScienceBase-Catalog. <https://www.sciencebase.gov/catalog/item/5e541969e4b0ff554f753113>.
- [3] J. T. Abatzoglou and A. P. Williams. Impact of anthropogenic climate change on wildfire across western US forests. *Proceedings of the National Academy of Sciences*, 113(42):11770–11775, Oct. 2016.
- [4] A. A. Adegun, S. Viriri, and J.-R. Tapamo. Review of deep learning methods for remote sensing satellite images classification: experimental survey and comparative analysis. *Journal of Big Data*, 10(1):93, June 2023.
- [5] P. Akiva, M. Purri, and M. Leotta. Self-supervised material and texture representation learning for remote sensing tasks. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8193–8205, 2022.
- [6] N. Andela, D. C. Morton, L. Giglio, Y. Chen, G. R. van der Werf, P. S. Kasibhatla, R. S. DeFries, G. J. Collatz, S. Hantson, S. Kloster, D. Bachelet, M. Forrest, G. Lasslop, F. Li, S. Mangeon, J. R. Melton, C. Yue, and J. T. Randerson. A human-driven decline in global burned area. *Science*, 356(6345):1356–1362, June 2017.
- [7] K. Ayush, B. Uz Kent, C. Meng, K. Tanmay, M. Burke, D. Lobell, and S. Ermon. Geography-aware self-supervised learning. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10161–10170, 2021.
- [8] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes, A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. G. Wilson, J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum. A cookbook of self-supervised learning, 2023.
- [9] W. G. C. Bandara, C. M. D. Melo, and V. M. Patel. Guarding barlow twins against overfitting with mixed samples, 2023.
- [10] H. B. Barlow. Redundancy reduction revisited. *Network: Computation in Neural Systems*, 12:241 – 253, 2001.
- [11] F. Bastani, P. Wolters, R. Gupta, J. Ferdinando, and A. Kembhavi. SatlasPretrain: A Large-Scale Dataset for Remote Sensing Image Understanding . In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16726–16736, Los Alamitos, CA, USA, Oct. 2023. IEEE Computer Society.
- [12] K. Bi, L. Xie, H. Zhang, X. Chen, X. Gu, and Q. Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619:533–538, 2023.
- [13] C. Bodnar, W. P. Bruinsma, A. Lucic, M. Stanley, A. Allen, J. Brandstetter, P. Garvan, M. Riechert, J. A. Weyn, H. Dong, J. K. Gupta, K. Thambiratnam, A. T. Archibald, C.-C. Wu, E. Heider, M. Welling, R. E. Turner, and P. Perdikaris. A foundation model for the earth system. *Nature*, 641:1180–1187, 2025.
- [14] J. Bourcier, G. Dashyan, K. Alahari, and J. Chanussot. Learning representations of satellite images from metadata supervision. In A. Leonardis, E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, editors, *Computer Vision – ECCV 2024*, volume 15085 of *Lecture Notes in Computer Science*, page 4. Springer, Cham, 2025.
- [15] R. Cazzolla Gatti, A. Di Paola, A. Bombelli, S. Noce, and R. Valentini. Exploring the relationship between canopy height and terrestrial plant diversity. *Plant Ecology*, 218(7):899–908, July 2017.
- [16] K. Chen, C. Liu, H. Chen, H. Zhang, W. Li, Z. Zou, and Z. Shi. Rsprompter: Learning to prompt for remote sensing instance segmentation based on visual foundation model. *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [17] Y. Chen, D. C. Morton, N. Andela, L. Giglio, and J. T. Randerson. How much global burned area can be forecast on seasonal time scales using sea surface temperatures? *Environ. Res. Lett.*, 11(4):045001, Mar. 2016.

- [18] E. Chuvieco, F. Mouillot, G. R. van der Werf, J. San Miguel, M. Tanase, N. Koutsias, M. García, M. Yebra, M. Padilla, I. Gitas, A. Heil, T. J. Hawbaker, and L. Giglio. Historical background and current developments for mapping burned area from satellite Earth observation. *Remote Sensing of Environment*, 225:45–64, May 2019.
- [19] Y. Cong, S. Khanna, C. Meng, P. Liu, E. Rozi, Y. He, M. Burke, D. B. Lobell, and S. Ermon. Satmae: pre-training transformers for temporal and multi-spectral satellite imagery. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [20] A. B. Davies and G. P. Asner. Advances in animal ecology from 3D-LiDAR ecosystem mapping. *Trends in Ecology & Evolution*, 29(12):681–691, Dec. 2014.
- [21] R. Dubayah, J. B. Blair, S. Goetz, L. Fatoyinbo, M. Hansen, S. Healey, M. Hofton, G. Hurtt, J. Kellner, S. Luthcke, J. Armston, H. Tang, L. Duncanson, S. Hancock, P. Jantz, S. Marselis, P. L. Patterson, W. Qi, and C. Silva. The Global Ecosystem Dynamics Investigation: High-resolution laser ranging of the Earth’s forests and topography. *Science of Remote Sensing*, 1:100002, June 2020.
- [22] L. Duncanson, J. Armston, M. Disney, V. Avitabile, N. Barbier, K. Calders, S. Carter, J. Chave, M. Herold, T. Crowther, M. Falkowski, J. Kellner, N. Labriere, R. Lucas, N. MacBean, R. Mcroberts, V. Meyer, E. Næsset, J. Nickeson, and M. Williams. The importance of consistent validation of global forest aboveground biomass products. *Surveys in Geophysics*, 40, 07 2019.
- [23] European Space Agency. Sentinel-1 data. Copernicus Data Space Ecosystem, 2014–present.
- [24] European Space Agency. Sentinel-2 data. Copernicus Data Space Ecosystem, 2015–present.
- [25] European Space Agency (ESA). SNAP - ESA Sentinel Application Platform, 2024.
- [26] T. R. Feldpausch, J. Lloyd, S. L. Lewis, R. J. W. Brienen, M. Gloor, A. Monteagudo Mendoza, G. Lopez-Gonzalez, L. Banin, K. Abu Salim, K. Affum-Baffoe, M. Alexiades, S. Almeida, I. Amaral, A. Andrade, L. E. O. C. Aragão, A. Araujo Murakami, E. J. M. M. Arets, L. Arroyo, G. A. Aymard C., T. R. Baker, O. S. Bánki, N. J. Berry, N. Cardozo, J. Chave, J. A. Comiskey, E. Alvarez, A. de Oliveira, A. Di Fiore, G. Djagbletey, T. F. Domingues, T. L. Erwin, P. M. Fearnside, M. B. França, M. A. Freitas, N. Higuchi, E. H. C. Y. Iida, E. Jiménez, A. R. Kassim, T. J. Killeen, W. F. Laurance, J. C. Lovett, Y. Malhi, B. S. Marimon, B. H. Marimon-Junior, E. Lenza, A. R. Marshall, C. Mendoza, D. J. Metcalfe, E. T. A. Mitchard, D. A. Neill, B. W. Nelson, R. Nilus, E. M. Nogueira, A. Parada, K. S.-H. Peh, A. Pena Cruz, M. C. Peñuela, N. C. A. Pitman, A. Prieto, C. A. Quesada, F. Ramírez, H. Ramírez-Angulo, J. M. Reitsma, A. Rudas, G. Saiz, R. P. Salomão, M. Schwarz, N. Silva, J. E. Silva-Espejo, M. Silveira, B. Sonké, J. Stropp, H. E. Taedoumg, S. Tan, H. ter Steege, J. Terborgh, M. Torello-Raventos, G. M. F. van der Heijden, R. Vásquez, E. Vilanova, V. A. Vos, L. White, S. Willcock, H. Woell, and O. L. Phillips. Tree height integrated into pantropical forest biomass estimates. *Biogeosciences*, 9(8):3381–3403, Aug. 2012.
- [27] L. Fu, Q. Shu, Z. Yang, C. Xia, X. Zhang, Y. Zhang, Z. Li, and S. Li. Accuracy assessment of topography and forest canopy height in complex terrain conditions of Southern China using ICESat-2 and GEDI data. *Frontiers in Plant Science*, 16, Mar. 2025.
- [28] A. Fuller, K. Millard, and J. R. Green. Croma: remote sensing representations with contrastive radar-optical masked autoencoders. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- [29] X. Guo, J. Lao, B. Dang, Y. Zhang, L. Yu, L. Ru, L. Zhong, Z. Huang, K. Wu, D. Hu, H. He, J. Wang, J. Chen, M. Yang, Y. Zhang, and Y. Li. Skysense: A multi-modal remote sensing foundation model towards universal interpretation for earth observation imagery. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27662–27673, 2024.
- [30] A. Gupta, A. Sheshadri, S. Roy, V. Gaur, M. Maskey, and R. Ramachandran. Machine learning global simulation of nonlocal gravity wave propagation. *arXiv preprint arXiv:2406.14775*, 2024.

- [31] C. Gómez, J. C. White, and M. A. Wulder. Optical remotely sensed time series data for land cover classification: A review. *ISPRS Journal of Photogrammetry and Remote Sensing*, 116:55–72, June 2016.
- [32] B. Han, S. Zhang, X. Shi, and M. Reichstein. Bridging remote sensors with multisensor geospatial foundation models. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27852–27862, 2024.
- [33] Y. Hashida, D. J. Lewis, and K. Cummins. Prescribed fires as a climate change adaptation tool. *Journal of Environmental Economics and Management*, 130:103081, Mar. 2025.
- [34] T. J. Hawbaker, V. C. Radeloff, A. D. Syphard, Z. Zhu, and S. I. Stewart. Detection rates of the MODIS active fire product in the United States. *Remote Sensing of Environment*, 112(5):2656–2664, May 2008.
- [35] T. J. Hawbaker, M. K. Vanderhoof, Y.-J. Beal, J. D. Takacs, G. L. Schmidt, J. T. Falgout, B. Williams, N. M. Fairaux, M. K. Caldwell, J. J. Picotte, S. M. Howard, S. Stitt, and J. L. Dwyer. Mapping burned areas using dense time-series of Landsat data. *Remote Sensing of Environment*, 198:504–522, Sept. 2017.
- [36] M. Herold, S. Carter, V. Avitabile, A. Espejo, I. Jonckheere, R. Lucas, R. Mcroberts, E. Næsset, J. Nightingale, R. Petersen, J. Reiche, E. Romijn, A. Rosenqvist, D. Rozendaal, F. M. Seifert, M.-J. Sanz-Sanchez, and V. De Sy. The role and need for space-based forest biomass-related measurements in environmental management and policy. *Surveys in Geophysics*, 40, 07 2019.
- [37] D. Hong, B. Zhang, X. Li, Y. Li, C. Li, J. Yao, N. Yokoya, H. Li, P. Ghamisi, X. Jia, A. Plaza, P. Gamba, J. A. Benediktsson, and J. Chanussot. Spectralgpt: Spectral remote sensing foundation model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5227–5244, Aug. 2024.
- [38] A. V. Huynh, L. E. Gillespie, J. Lopez-Saucedo, C. Tang, R. Sikand, and M. Expósito-Alonso. Contrastive ground-level image and remote sensing pre-training improves representation learning for natural world imagery. In A. Leonardis, E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, editors, *Computer Vision – ECCV 2024*, volume 15138 of *Lecture Notes in Computer Science*, page 10. Springer, Cham, 2025.
- [39] V. Iglesias, J. K. Balch, and W. R. Travis. U.S. fires became larger, more frequent, and more widespread in the 2000s. *Science Advances*, 8(11):eabc0020, Mar. 2022.
- [40] T. D. Jackson, F. J. Fischer, G. Vincent, E. B. Gorgens, M. Keller, J. Chave, T. Jucker, and D. A. Coomes. Tall Bornean forests experience higher canopy disturbance rates than those in the eastern Amazon or Guiana shield. *Global Change Biology*, 30, 2024.
- [41] A. Joshi, B. Pradhan, S. Gite, and S. Chakraborty. Remote-Sensing Data and Deep-Learning Techniques in Crop Mapping and Yield Prediction: A Systematic Review. *Remote Sensing*, 15(8):2014, Jan. 2023. Number: 8 Publisher: Multidisciplinary Digital Publishing Institute.
- [42] S. Khanna, P. Liu, L. Zhou, C. Meng, R. Rombach, M. Burke, D. B. Lobell, and S. Ermon. Diffusionsat: A generative foundation model for satellite imagery. In *The Twelfth International Conference on Learning Representations*, 2024.
- [43] K. Klemmer, E. Rolf, C. Robinson, L. Mackey, and M. Rußwurm. Satclip: Global, general-purpose location embeddings with satellite imagery. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(4):4347–4355, Apr. 2025.
- [44] E. Kurbanov, O. Vorobev, S. Lezhnin, J. Sha, J. Wang, X. Li, J. Cole, D. Dergunov, and Y. Wang. Remote Sensing of Forest Burnt Area, Burn Severity, and Post-Fire Recovery: A Review. *Remote Sensing*, 14(19):4714, Jan. 2022.
- [45] K. Lahssini, N. Baghdadi, G. le Maire, I. Fayad, and L. Villard. Canopy height mapping in French Guiana using multi-source satellite data and environmental information in a U-Net architecture. *Frontiers in Remote Sensing*, 5, Nov. 2024.
- [46] N. Lang, W. Jetz, K. Schindler, and J. D. Wegner. A high-resolution canopy height model of the Earth. *Nature Ecology & Evolution*, 7(11):1778–1789, Nov. 2023.
- [47] N. Lang, K. Schindler, and J. D. Wegner. High carbon stock mapping at large scale with optical satellite imagery and spaceborne LIDAR, July 2021.

- [48] Z. Li, B. Hou, S. Ma, Z. Wu, X. Guo, B. Ren, and L. Jiao. Masked angle-aware autoencoder for remote sensing images. In *European Conference on Computer Vision*, pages 260–278. Springer, 2025.
- [49] M. C. Lisaius, A. Blake, S. Keshav, and C. Atzberger. Using barlow twins to create representations from cloud-corrupted remote sensing time series. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 17:13162–13168, 2024.
- [50] A. Liu, Y. Chen, and X. Cheng. Improving Tropical Forest Canopy Height Mapping by Fusion of Sentinel-1/2 and Bias-Corrected ICESat-2–GEDI Data. *Remote Sensing*, 17(12):1968, Jan. 2025.
- [51] X. Liu, F. Zhang, Z. Hou, L. Mian, Z. Wang, J. Zhang, and J. Tang. Self-supervised learning: Generative or contrastive. *IEEE Transactions on Knowledge and Data Engineering*, page 1–1, 2021.
- [52] G. Mai, N. Lao, Y. He, J. Song, and S. Ermon. Csp: self-supervised contrastive spatial pre-training for geospatial-visual representations. In *Proceedings of the 40th International Conference on Machine Learning, ICML’23*. JMLR.org, 2023.
- [53] V. Marsocci, Y. Jia, G. L. Bellier, D. Kerekes, L. Zeng, S. Hafner, S. Gerard, E. Brune, R. Yadav, A. Shibli, H. Fang, Y. Ban, M. Vergauwen, N. Audebert, and A. Nascetti. Pangaea: A global and inclusive benchmark for geospatial foundation models, 2025.
- [54] O. Mañas, A. Lacoste, X. Giró-i Nieto, D. Vazquez, and P. Rodríguez. Seasonal contrast: Un-supervised pre-training from uncurated remote sensing data. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9394–9403, 2021.
- [55] M. Mendieta, B. Han, X. Shi, Y. Zhu, and C. Chen. Towards geospatial foundation models via continual pretraining. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16760–16770, 2023.
- [56] M. Migliavacca, T. Musavi, M. D. Mahecha, J. A. Nelson, J. Knauer, D. D. Baldocchi, O. Perez-Priego, R. Christiansen, J. Peters, K. Anderson, M. Bahn, T. A. Black, P. D. Blanken, D. Bonal, N. Buchmann, S. Caldararu, A. Carrara, N. Carvalhais, A. Cescatti, J. Chen, J. Cleverly, E. Cremonese, A. R. Desai, T. S. El-Madany, M. M. Farella, M. Fernández-Martínez, G. Filippa, M. Forkel, M. Galvagno, U. Gomasasca, C. M. Gough, M. Göckede, A. Ibrom, H. Ikawa, I. A. Janssens, M. Jung, J. Kattge, T. F. Keenan, A. Knohl, H. Kobayashi, G. Kraemer, B. E. Law, M. J. Liddell, X. Ma, I. Mammarella, D. Martini, C. Macfarlane, G. Matteucci, L. Montagnani, D. E. Pabon-Moreno, C. Panigada, D. Papale, E. Pendall, J. Penuelas, R. P. Phillips, P. B. Reich, M. Rossini, E. Rotenberg, R. L. Scott, C. Stahl, U. Weber, G. Wohlfahrt, S. Wolf, I. J. Wright, D. Yakir, S. Zaehle, and M. Reichstein. The three major axes of terrestrial ecosystem function. *Nature*, 598(7881):468–472, Oct. 2021.
- [57] A. L. Mitchell, A. Rosenqvist, and B. Mora. Current remote sensing approaches to monitoring forest degradation in support of countries measurement, reporting and verification (MRV) systems for REDD+. *Carbon Balance and Management*, 12(1):9, Apr. 2017.
- [58] V. Moudrý, L. Gábor, S. Marselis, P. Pracná, V. Barták, J. Prošek, B. Navrátilová, J. Novotný, M. Potůčková, K. Gdulová, P. Crespo-Peremarch, J. Komárek, M. Malavasi, D. Rocchini, L. A. Ruiz, J. Torralba, M. Torresani, R. Cazzolla Gatti, and J. Wild. Comparison of three global canopy height maps and their applicability to biodiversity modeling: Accuracy issues revealed. *Ecosphere*, 15(10):e70026, 2024.
- [59] V. Moudrý, R. Remelgado, M. Forkel, M. Torresani, G. V. Laurin, E. Sarovcova, V. E. G. Millan, F. J. Fischer, T. Jucker, M. Gally, P. Kacic, C. R. Hakkenberg, Ž. Kokalj, K. Stereńczak, Y. Erfanifard, D. Rocchini, J. Prošek, S. Roilo, K. Gdulova, A. F. Cord, M. Perrone, J. A. Molina-Valero, P. Surovy, Z. Melichová, M. Malavasi, R. Urban, M. Štroner, D. Seidel, S. Szabó, L. Bertalan, A. Etner, R. C. Gatti, and V. Barták. Harmonised airborne laser scanning products can address the limitations of large-scale spaceborne vegetation mapping. Dec. 2024.
- [60] A. Nascetti, R. Yadav, K. Brodt, Q. Qu, H. Fan, Y. Shendryk, I. Shah, and C. Chung. Biomasssters: A benchmark dataset for forest biomass estimation using multi-modal satellite time-series. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [61] V. Nedungadi, A. Kariryaa, S. Oehmcke, S. J. Belongie, C. Igel, and N. Lang. Mmearth: Exploring multi-modal pretext tasks for geospatial representation learning. In A. Leonardis,

- E. Ricci, S. Roth, O. Russakovsky, T. Sattler, and G. Varol, editors, *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXIV*, volume 15122 of *Lecture Notes in Computer Science*, pages 164–182. Springer, 2024.
- [62] M. Noman, M. Naseer, H. Cholakkal, R. M. Anwar, S. Khan, and F. S. Khan. Rethinking Transformers Pre-training for Multi-Spectral Satellite Imagery . In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27811–27819, Los Alamitos, CA, USA, June 2024. IEEE Computer Society.
- [63] A. O. Ok, A. , Ozlem, , and O. Gungor. Evaluation of random forest method for agricultural crop classification. *European Journal of Remote Sensing*, 45(1):421–432, Jan. 2012. Publisher: Taylor & Francis _eprint: <https://doi.org/10.5721/EuJRS20124535>.
- [64] L. Pang, D. Tang, S. Xu, D. Meng, and X. Cao. Hsigene: A foundation model for hyperspectral image generation, 2024.
- [65] J. Pauls, M. Zimmer, B. Turan, S. Saatchi, P. Ciais, S. Pokutta, and F. Gieseke. Capturing Temporal Dynamics in Large-Scale Canopy Tree Height Estimation. <https://arxiv.org/abs/2501.19328v1>, Jan. 2025.
- [66] O. Pechony, D. T. Shindell, and F. S. Chapin. Driving forces of global wildfires over the past millennium and the forthcoming century. *Proceedings of the National Academy of Sciences of the United States of America*, 107(45):19167–19170, 2010.
- [67] J. J. Picotte, K. Bhattarai, D. Howard, J. Lecker, J. Epting, B. Quayle, N. Benson, and K. Nelson. Changes to the Monitoring Trends in Burn Severity program mapping production procedures and data products. *Fire Ecology*, 16(1):16, June 2020.
- [68] P. Potapov, X. Li, A. Hernandez-Serna, A. Tyukavina, M. C. Hansen, A. Kommareddy, A. Pickens, S. Turubanova, H. Tang, C. E. Silva, J. Armston, R. Dubayah, J. B. Blair, and M. Hofton. Mapping global forest canopy height through integration of GEDI and Landsat data. *Remote Sensing of Environment*, 253:112165, Feb. 2021.
- [69] C. J. Reed, R. Gupta, S. Li, S. Brockman, C. Funk, B. Clipp, K. Keutzer, S. Candido, M. Uyttendaele, and T. Darrell. Scale-MAE: A Scale-Aware Masked Autoencoder for Multiscale Geospatial Representation Learning . In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4065–4076, Los Alamitos, CA, USA, Oct. 2023. IEEE Computer Society.
- [70] E. Rolf, L. Gordon, M. Tambe, and A. Davies. Contrasting local and global modeling with machine learning and satellite data: A case study estimating tree canopy height in African savannas, Nov. 2024.
- [71] E. Roteta, A. Bastarrika, M. Padilla, T. Storm, and E. Chuvieco. Development of a Sentinel-2 burned area algorithm: Generation of a small fire database for sub-Saharan Africa. *Remote Sensing of Environment*, 222:1–17, Mar. 2019.
- [72] J. Schmude, S. Roy, W. Trojak, J. Jakubik, D. S. Civitarese, S. Singh, J. Kuehnert, K. Ankur, A. Gupta, C. E. Phillips, R. Kienzler, D. Szwarcman, V. Gaur, R. Shinde, R. Lal, A. D. Silva, J. L. G. Diaz, A. Jones, S. Pfreundschuh, A. Lin, A. Sheshadri, U. Nair, V. Anantharaj, H. Hamann, C. Watson, M. Maskey, T. J. Lee, J. B. Moreno, and R. Ramachandran. Prithvi wxc: Foundation model for weather and climate, 2024.
- [73] A. K. Skidmore, N. C. Coops, E. Neinavaz, A. Ali, M. E. Schaepman, M. Paganini, W. D. Kissling, P. Vihervaara, R. Darvishzadeh, H. Feilhauer, M. Fernandez, N. Fernández, N. Gorelick, I. Geijzendorffer, U. Heiden, M. Heurich, D. Hobern, S. Holzwarth, F. E. Muller-Karger, R. Van De Kerchove, A. Lausch, P. J. Leitão, M. C. Lock, C. A. Múcher, B. O’Connor, D. Rocchini, C. Roeoesli, W. Turner, J. K. Vis, T. Wang, M. Wegmann, and V. Wingate. Priority list of biodiversity metrics to observe from space. *Nature Ecology & Evolution*, 5(7):896–906, July 2021.
- [74] X. Sun, P. Wang, W. Lu, Z. Zhu, X. Lu, Q. He, J. Li, X. Rong, Z. Yang, H. Chang, Q. He, G. Yang, R. Wang, J. Lu, and K. Fu. Ringmo: A remote sensing foundation model with masked image modeling. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–22, 2023.

- [75] D. Szwarcman, S. Roy, P. Fraccaro, P. E. Gislason, B. Blumenstiel, R. Ghosal, P. H. de Oliveira, J. L. de Sousa Almeida, R. Sedona, Y. Kang, S. Chakraborty, S. Wang, C. Gomes, A. Kumar, M. Truong, D. Godwin, H. Lee, C.-Y. Hsu, A. A. Asanjan, B. Mujeci, D. Shidham, T. Keenan, P. Arevalo, W. Li, H. Alemohammad, P. Olofsson, C. Hain, R. Kennedy, B. Zadrozny, D. Bell, G. Cavallaro, C. Watson, M. Maskey, R. Ramachandran, and J. B. Moreno. Prithvi-eo-2.0: A versatile multi-temporal foundation model for earth observation applications, 2025.
- [76] M. A. Tanase, M. A. Belenguer-Plomer, E. Roteta, A. Bastarrika, J. Wheeler, Á. Fernández-Carrillo, K. Tansey, W. Wiedemann, P. Navratil, S. Lohberger, F. Siegert, and E. Chuvieco. Burned Area Detection and Mapping: Intercomparison of Sentinel-1 and Sentinel-2 Based Algorithms over Tropical Africa. *Remote Sensing*, 12(2):334, Jan. 2020.
- [77] M. Tang, A. Cozma, K. Georgiou, and H. Qi. Cross-scale mae: a tale of multi-scale exploitation in remote sensing. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- [78] J. Tolan, H.-I. Yang, B. Nosarzewski, G. Couairon, H. V. Vo, J. Brandt, J. Spore, S. Majumdar, D. Haziza, J. Vamaraju, T. Moutakanni, P. Bojanowski, T. Johns, B. White, T. Tiedecke, and C. Couprie. Very high resolution canopy height maps from RGB imagery using self-supervised vision transformer and convolutional decoder trained on aerial lidar. *Remote Sensing of Environment*, 300:113888, Jan. 2024.
- [79] G. Tseng, R. Cartuyvels, I. Zvonkov, M. Purohit, D. Rolnick, and H. Kerner. Lightweight, Pre-trained Transformers for Remote Sensing Timeseries, Feb. 2024. arXiv:2304.14065 [cs].
- [80] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc.
- [81] A. Verhegghen, H. Eva, G. Ceccherini, F. Achard, V. Gond, S. Gourlet-Fleury, and P. O. Cerutti. The Potential of Sentinel Satellites for Burnt Area Mapping and Monitoring in the Congo Basin Forests. *Remote Sensing*, 8(12):986, Dec. 2016.
- [82] C. Wang, C. Song, T. A. Schroeder, C. E. Woodcock, T. M. Pavelsky, Q. Han, and F. Yao. Interpretable Multi-Sensor Fusion of Optical and SAR Data for GEDI-Based Canopy Height Mapping in Southeastern North Carolina. *Remote Sensing*, 17(9):1536, Jan. 2025.
- [83] D. Wang, J. Zhang, B. Du, G.-S. Xia, and D. Tao. An empirical study of remote sensing pretraining. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–20, 2023.
- [84] D. Wang, J. Zhang, B. Du, M. Xu, L. Liu, D. Tao, and L. Zhang. Samrs: Scaling-up remote sensing segmentation dataset with segment anything model. In *Advances in Neural Information Processing Systems*, volume 36, pages 8815–8827, 2023.
- [85] Y. Wang, N. A. A. Braham, Z. Xiong, C. Liu, C. M. Albrecht, and X. X. Zhu. Ssl4eo-s12: A large-scale multi-modal, multi-temporal dataset for self-supervised learning in earth observation. *ArXiv*, abs/2211.07044, 2022.
- [86] M. Weiss, F. Jacob, and G. Duveiller. Remote sensing for agricultural applications: A meta-review. *Remote Sensing of Environment*, 236:111402, Jan. 2020.
- [87] A. Xiao, W. Xuan, H. Qi, Y. Xing, R. Ren, X. Zhang, L. Shao, and S. Lu. Cat-sam: Conditional tuning for few-shot adaptation of segment anything model. *arXiv preprint arXiv:2402.03631*, 2024.
- [88] Z. Yan, J. Li, X. Li, R. Zhou, W. Zhang, Y. Feng, W. Diao, K. Fu, and X. Sun. Ringmo-sam: A foundation model for segment anything in multimodal remote-sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–16, 2023.
- [89] J. Zang, F. Qiu, and Y. Zhang. A global dataset of forest regrowth following wildfires. *Sci Data*, 11(1):1052, Sept. 2024.
- [90] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow twins: Self-supervised learning via redundancy reduction. *ArXiv*, abs/2103.03230, 2021.
- [91] X. Zhang, Y. Liu, Y. Lin, Q. Liao, and Y. Li. Uv-sam: adapting segment anything model for urban village identification. In *Proceedings of the Thirty-Eighth AAAI Conference on Artificial Intelligence and Thirty-Sixth Conference on Innovative Applications of Artificial*

Intelligence and Fourteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'24/IAAI'24/EAAI'24. AAAI Press, 2024.

- [92] Z. Zheng, S. Ermon, D. Kim, L. Zhang, and Y. Zhong. Changen2: Multi-temporal remote sensing generative change foundation model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47(2):725–741, 2025.

Acknowledgments

We gratefully acknowledge help from AMD Inc., Tarides, Jane Street, the Dawn supercomputing team at Cambridge, the Aalto University Science-IT project and the Aalto Scientific Computing team, and the UKRI STFC AIRR Programme for providing us access to compute resources.

Funding: Z. F. and S. K. were funded by a charitable donation made by Dr Robert Sansom to the University of Cambridge. S. J. was funded by a charitable donation from John Bernstein. J. B. received support from the UKRI CRCRM scheme (grant MR/Y024354). S. S. received funding from the Jenny and Antti Wihuri Foundation and the Helsinki Institute for Information Technology.

Author contributions: System design: CA, ZF, SJ, SK
System implementation: ZF, RY,
Data curation: ZF, RY, JK, SJ, AM, DAC, SK
Downstream tasks: ML, CA, SS, DAC, CA, MI, ZF, JK
Paper conceptualization: SK, AM, DAC, AB, SJ
Paper draft: ZF, ML, RY, SS, SK, JB, JK
Paper review and refinement: AB, AM, DAC, SJ, SK

Competing interests: There are no competing interests to declare.

Data and materials availability: Our framework integrates data from the Copernicus Sentinel-2 (S2) [24] optical and Sentinel-1 (S1) [23] SAR missions, leveraging the strengths of each sensor. We utilized S2 Level-2A surface reflectance products, providing detailed spectral information across 10 bands². A challenge with optical data is frequent cloud contamination; we filtered the time series using the associated Sen2Cor scene classification map (SCL) to mask cloudy observations. For SAR data, we processed Sentinel-1 Ground Range Detected imagery to derive calibrated backscatter coefficients in VV and VH polarizations, applying standard corrections using the ESA SNAP toolbox [25].

²We used red, blue, green, NIR, NIR8, red-edge1, red-edge2, red-edge3, SWIR16, and SWIR22.