# TerraFM: A Scalable Foundation Model for Unified Multisensor Earth Observation

Muhammad Sohail Danish[1]    Muhammad Akhtar Munir[1]    Syed Roshaan Ali Shah[2]
Muhammad Haris Khan[1]    Rao Muhammad Anwer[1,3]    Jorma Laaksonen[3]
Fahad Shahbaz Khan[1,4]    Salman Khan[1,5]
[1]Mohamed bin Zayed University of Artificial Intelligence    [2]University College London
[3]Aalto University    [4]Linköping University, Sweden    [5]Australian National University
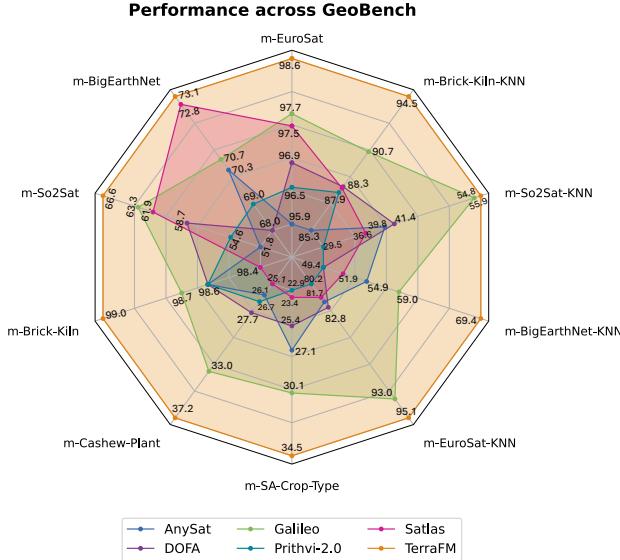
## Abstract

Modern Earth observation (EO) increasingly leverages deep learning to harness the scale and diversity of satellite imagery across sensors and regions. While recent foundation models have demonstrated promising generalization across EO tasks, many remain limited by the scale, geographical coverage, and spectral diversity of their training data, factors critical for learning globally transferable representations. In this work, we introduce **TerraFM**, a scalable self-supervised learning model that leverages globally distributed Sentinel-1 and Sentinel-2 imagery, combined with large spatial tiles and land-cover aware sampling to enrich spatial and semantic coverage. By treating sensing modalities as natural augmentations in our self-supervised approach, we unify radar and optical inputs via modality-specific patch embeddings and adaptive cross-attention fusion. Our training strategy integrates local-global contrastive learning and introduces a dual-centering mechanism that incorporates class-frequency-aware regularization to address long-tailed distributions in land cover. TerraFM achieves strong generalization on both classification and segmentation tasks, outperforming prior models on GEO-Bench and Copernicus-Bench. Our code and pretrained models are publicly available at https://github.com/mbzuai-oryx/TerraFM.

## 1 Introduction

EO provides systematic measurements of the surface of the earth, supporting a wide spectrum of critical applications such as land use monitoring [31], crop evaluation [20, 17], urban development [37], and disaster response [15, 24, 21]. These capabilities are enabled by a growing fleet of earth-observing satellites, most notably the Sentinel missions, which deliver multi-modal, multi-temporal data at a global scale [27, 7]. The rise of deep learning, particularly deep neural networks (DNNs), has fundamentally reshaped how EO data is processed and interpreted [32, 25, 28]. Modern DNNs enable automated extraction of spatial and semantic patterns from raw imagery, driving downstream tasks such as scene classification, object detection, and semantic segmentation [1, 29, 16, 28, 10]. These models offer a scalable and adaptive alternative to traditional hand-engineered pipelines by learning generalizable representations directly from the data [12]. As EO datasets continue to expand in scale, diversity, and complexity, DNNs have become the foundation for building high-capacity models capable of generalizing across geographies, modalities, and tasks [25, 1, 29].

Remote sensing data is inherently multimodal, comprising diverse sensor types such as optical, SAR, and multispectral imagery. Traditional EO pipelines often focus on single-modality inputs, typically high-resolution optical imagery, limiting the model's ability to generalize across varying sensing conditions. In contrast, multimodal and multispectral data sources, such as Sentinel-1 SAR and Sentinel-2 Level-1C/Level-2A optical bands, capture complementary structural and spectral information, enabling richer scene understanding [13, 10]. Foundation models that embrace this

**Performance across GeoBench**

Figure 1: Performance comparison across GEO-Bench classification tasks using supervised and kNN-based evaluation protocols. Specifically, supervised fine-tuning assesses the adaptability of the learned base representations and the frozen-encoder with kNN head assesses the generalization of representations. Five recent EO foundation models (AnySat [1], DOFA [36], Galileo [28], Prithvi-2.0 [25], Satlas [2]) performances are shown compared to our proposed TerraFM. TerraFM consistently outperforms existing models across modalities and evaluation settings, showing strong generalization.

diversity have demonstrated superior transferability across tasks and geographies [28, 12]. However, variation in ground sampling distance (GSD) across EO data makes tile size a critical factor; smaller tiles capture local detail but risk overfitting to texture, while larger tiles provide broader semantic context but require scale-robust architectures [22]. Recent works like AnySat and msGFM [13, 1] have shown that scale-invariant modeling and mixed-resolution pretraining lead to more robust and generalizable representations. Crucially, large-scale sampling across geographies and resolutions enables EO foundation models to learn invariant features across sensors and global conditions.

As EO foundation models scale to accommodate diverse sensor inputs and resolutions, two dominant pretraining paradigms have emerged: masked autoencoders (MAE) and contrastive learning. Although MAEs focus on reconstructing the spatial structure, their reliance on RGB-centric ViTs limits their adaptability to multispectral or SAR inputs with varying spectral dimensions [18, 25]. In contrast, contrastive approaches such as DINO [4, 19] and its adaptations to remote sensing [28, 10, 29] offer modality-agnostic training by aligning global and local views through student-teacher distillation. However, the expansive spatial coverage of EO datasets introduces new challenges: large portions of satellite imagery are semantically sparse or uninformative, and naïve sampling can lead to representation bias. This requires intelligent sampling that prioritizes semantically diverse regions, guided by land cover priors, for balanced and efficient representation learning.

To address these limitations in standard ViTs, particularly their RGB-centric design, lack of modality awareness, and unimodal self-supervision, we introduce **TerraFM**, a unified foundation model tailored for remote sensing. First, we propose a *Modality-Specific Patch Embedding* module, which replaces the shared projection in standard ViTs with modality-aware embeddings adapted to multispectral and SAR data. This enables flexible handling of sensor-specific spectral profiles while preserving spatial structure. To enhance scale-invariance and cross-view consistency, we adopt multi-crop learning within a self-supervised teacher-student framework, promoting robust representation learning through global-local alignment. Further, we interpret different aligned modalities (S1-SAR, S2-L1C, S2-L2A) as complementary views of the same scene and introduce a *Cross-Attention Fusion* module that dynamically aggregates modality-specific tokens using learnable spatial queries. This allows the model to selectively emphasize sensor contributions at each spatial location. Finally, to mitigate long-tailed land cover distribution issues prevalent in EO data, we introduce a *Dual Centering* mechanism into the distillation process. This leverages WorldCover [38] derived class statistics to compute a frequency-aware center, improving balance across dominant and rare semantic categories without requiring supervised objectives. Our key contributions are as follows.

**Contributions:** **(1)** A *modality-specific patch embedding* mechanism is introduced to generalize ViTs across heterogeneous remote sensing modalities with varying spectral dimensions. **(2)** We treat sensor modalities as natural augmentations and introduce a *cross-attention fusion* block that unifies multi-modal inputs within a shared encoder. **(3)** To address long-tailed LULC distributions, a *dual-*

Table 1: Comparison of recent remote sensing foundation models across modalities, scale, and benchmarks. **TerraFM** uniquely blends large tile size (534), WorldCover-informed metadata, and global-scale training (18.7M samples) with evaluation on both GEO-Bench and Copernicus-Bench.

| Model | Modalities | Scale | Resolution | TileSize | Metadata | Benchmarks | Pixels (~T) |
|---|---|---|---|---|---|---|---|
| SatMAE++ | S2, RGB | ~1.2 M | 10–60 m | 224, 96 | No | 6 DS | 0.12 |
| Galileo | S1, S2, NDVI, ESA WC etc | ~3–10.9 M | 10 m | 96 (flex) | Yes | GEO + 5 DS | 1.58 |
| CROMA | S1, S2 | ~1 M | 10 m | 96, 120 | No | 7 DS | 0.98 |
| SoftCon | S1, S2 | ~0.78 M | 10 m | 224 | Yes | 4 GEO + 7 DS | 0.76 |
| AnySat | Aerial, S1/S2, MODIS, etc. | 11.1 M | 0.2–250 m | 10–240 | No | 11 DS | 0.17 |
| Prithvi-2 | S2, HLS | 4.2 M | 30 m | 224 | Yes | GEO + 9 SME | 5.06 |
| DOFA | S1, S2, EnMAP, etc. | ~8 M | 1–30 m | 512, 128 | Yes | GEO + 2 DS | 6.74 |
| Panopticon | S1, S2, WV2/3, NAIP | ~2.6 M | 0.3–100 m | 96, 224 | Yes | GEO + 10 DS | 2.34 |
| MMEarth | S1, S2, DEM, etc. | ~7.2 M | 0.3–100 m | 128 | Yes | 5 GEO | 0.51 |
| msGFM | RGB, S2, SAR, DSM | ~2 M | 0.1–30 m | 192 | No | 5 DS | 0.44 |
| Copernicus-FM | S1–S5P, DEM | 18.7 M | 10 m–1 km | Mixed | Yes | Cop-Bench | 5.12 |
| **TerraFM** (Ours) | S1, S2 L1C/L2A | 18.7 M | 10–60 m | 534 | Yes | GEO + Cop-Bench | **23.32** |

*centering* strategy is incorporated to regularize representation learning using class-frequency-aware statistics. **(4)** Extensive experiments on *GEO-Bench* and *Copernicus-Bench* demonstrate leading performance across multiple downstream tasks using globally distributed data (Fig. 1).

## 2 Related Work

**Self-supervised Pretraining:** MAEs [14] have become a popular choice for self-supervised pre-training in remote sensing by reconstructing masked image regions using ViT [6]. Variants like Scale-MAE [22] and MC-MAE [11] enhance robustness across spatial scales via scale-aware encodings and convolutional tokenizers. However, MAEs struggle to scale to multisensor EO data, as their RGB-centric tokenization and reconstruction objectives limit generalization to multispectral and SAR modalities with diverse channel structures [35, 18].

Unlike MAEs, self-supervised contrastive learning focuses on learning discriminative representations by comparing semantically similar and dissimilar views. Remote sensing approaches [26, 10, 29] leverage spatial and spectral augmentations to create diverse yet consistent views. CROMA [10] combines contrastive and masked autoencoding losses, while Cross-Scale MAE [26] blends generative and contrastive objectives for multi-scale learning. Student-teacher frameworks like DINO [4, 19] scale contrastive learning via EMA-updated teachers and global-local view alignment with centering to prevent collapse. These strategies are well-suited for EO, where multimodal imagery can act as natural augmentations, enabling scalable, label-free training and broad generalization.

**Remote Sensing FMs:** Recent advances in remote sensing foundation models (FMs) have scaled self-supervised learning across architecture types, modalities, training sizes, tile resolutions, and metadata usage (Tab. 1). Multimodal integration is central to recent FMs like [12, 32, 29, 1, 28, 13]. SkySense [12] applies contrastive learning to temporal-multimodal data but requires large-scale compute. CopernicusFM [32] fuses Sentinel modalities via metadata-aware networks but faces scaling issues with heterogeneous inputs. Panopticon [29] and AnySat [1] align cross-modal views through contrastive training, while Galileo [28] uses shared embeddings for SAR and multispectral fusion. Fus-MAE [5] adopts attention-based fusion without contrastive loss, limiting generalization.

Prithvi-2 [25] is restricted to single-modal optical data with temporal-spatial modeling. DOFA [36], msGFM [13], and AnySat [1] address resolution variability using mixed tile sizes or scale-adaptive designs. Our 534px tiles capture broader spatial context than prior RSFMs. While CopernicusFM [32] and DOFA [36] incorporate metadata, we leverage land cover (LULC) priors for semantically informed learning. Both CopernicusFM and our model are trained on 18.7M samples, but ours uses over 23T pixels during pretraining, scaling the 5.1T used by Copernicus-Pretrain [32] over 4x.

## 3 TerraFM: A Scalable Multisensor Foundational Model

Unlike prior remote sensing foundation models, our approach integrates a student–teacher contrastive learning framework with dual centering (to balance long-tailed classes), modality-as-augmentation (to learn cross-modal invariances), and cross-attention fusion (to aggregate multi-sensor context), as illustrated in Fig. 2. Built on a ViT backbone and trained on 18.7M globally distributed samples
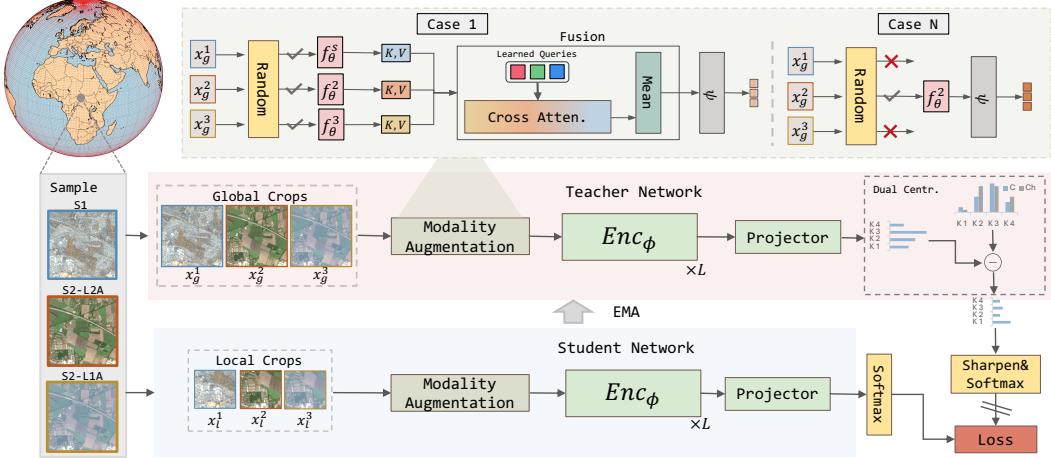
Figure 2: Overall architecture of **TerraFM**. It unifies student-teacher contrastive framework with modality augmentation with cross-attention fusion, and a new dual centering regularization. **TerraFM** is founded on ViT backbone and is trained on 18.7M globally distributed samples for pre-training and utilizes large-tile inputs for encoding broader spatial context. For illustration, RGB channels from S2-L2A and S2-L1C are selected, and S1 is visualized using a false-color RGB composite.

using $534 \times 534$ tiles, **TerraFM** captures broader spatial context and generalizes effectively across sensing modalities and geographies, achieving strong results on diverse downstream benchmarks.

## 3.1 Architecture

We use globally distributed remote sensing imagery organized over a spatial grid, partitioning the earth's surface into fixed-size tiles [9] (e.g., $5.34\,\mathrm{km} \times 5.34\,\mathrm{km}$). Each spatial unit, denoted as $s$, represents one such grid cell. For each sample, we observe a set of co-registered EO modalities:

$$\mathcal{M} = \{\text{S1},\ \text{S2-L1C},\ \text{S2-L2A}\},$$

where **S1** corresponds to Sentinel-1 SAR (Synthetic Aperture Radar), and **S2-L1C** and **S2-L2A** represent two processing levels of Sentinel-2 optical imagery: Level-1C (top-of-atmosphere reflectance) and Level-2A (bottom-of-atmosphere surface reflectance), respectively. Each modality $m \in \mathcal{M}$ provides a multi-channel image $\boldsymbol{x}^m \in \mathbb{R}^{H \times W \times C_m}$, where $H$ and $W$ denote spatial dimensions, and $C_m$ is the number of spectral channels for modality $m$. For example, Sentinel-1 contains two channels (VV and VH polarizations), therefore $C_{\text{S1}} = 2$, while Sentinel-2 modalities contain up to 13 spectral bands depending on level and resolution. These modalities are treated as complementary views of the same location, acting as natural augmentations, which support our training strategy and encourage learning modality-invariant representations.

To provide semantic grounding, each sample $s$ is assigned a high-level land use and land cover (LULC) category $y^{(s)} \in \{1, \ldots, Y\}$, derived from the ESA WorldCover product. These categories reflect coarse semantic classes at a global scale and are used to compute class-frequency-aware statistics for balanced representation learning.

**Vision Transformer Model:**

ViTs adapt the transformer architecture to visual data by treating an image as a sequence of patch tokens instead of a dense pixel grid. A typical ViT consists of two main components: a patch embedding module and a transformer encoder. Given an input image $\boldsymbol{x} \in \mathbb{R}^{H \times W \times C}$, the **patch embedding** layer $f_\theta$ divides the image into $N$ non-overlapping patches of size $P \times P$, and projects each patch into a $d$ dimensional embedding:

$$\{\boldsymbol{z}_i\}_{i=1}^N = f_\theta(\boldsymbol{x}), \qquad \boldsymbol{z}_i \in \mathbb{R}^d.$$

This projection is typically implemented using a convolutional layer with kernel size and stride equal to the patch size $P$, parameterized by weights $\mathbf{W}_\theta \in \mathbb{R}^{d \times C \times P \times P}$. To encode spatial information, the transformer encoder augments each patch token $\boldsymbol{z}_i$ with a positional vector. A learnable class token $\boldsymbol{z}_{\text{cls}}$ is added to the sequence, which yields the full input:

$$\boldsymbol{Z} = \big[\, \boldsymbol{z}_{\text{cls}};\ \{\boldsymbol{z}_i + \text{pos}_i\}_{i=1}^N \,\big].$$

4

The token sequence $\boldsymbol{Z}$ is processed by a stack of $L$ transformer layers, denoted $\text{Enc}_\phi$. For classification tasks, only the final class token $\hat{z}_{\text{cls}}$ is forwarded to a prediction head.

**Modality-Specific Patch Embedding:**

Standard patch embedding layers in ViTs are typically implemented using a shared convolutional projection across all inputs, making it unsuitable for multi-modal remote sensing data.

To better handle this heterogeneity, we adopt a modality-specific patch embedding strategy. For each modality $m \in \mathcal{M}$, we define an embedding function $f_\theta^m$ that maps the input image $\boldsymbol{x}^m \in \mathbb{R}^{H \times W \times C_m}$ to a sequence of patch tokens $\boldsymbol{Z}^m \in \mathbb{R}^{N_m \times D}$, where $C_m$ is the number of channels and $N_m$ is the number of patches. Each $f_\theta^m$ is parameterized independently to account for modality-specific dynamics. We associate each modality with a learnable embedding vector $\boldsymbol{e}^m \in \mathbb{R}^D$. This vector is added to every token from that modality via broadcasting:

$$\tilde{\boldsymbol{Z}}^m = \boldsymbol{Z}^m + \mathbf{1}_{N_m} \cdot (\boldsymbol{e}^m)^\top,$$

where $\mathbf{1}_{N_m} \in \mathbb{R}^{N_m \times 1}$ is a vector of ones. This allows the model to distinguish between modalities while preserving local spatial and spectral features. Finally, to enable shared processing in the Transformer encoder, the enriched tokens $\tilde{\boldsymbol{Z}}^m$ are linearly projected into a common latent space of dimension $d$ using a shared projection $\psi : \mathbb{R}^D \to \mathbb{R}^d$:

$$\boldsymbol{Z}^{(m)} = \psi(\tilde{\boldsymbol{Z}}^m) \in \mathbb{R}^{N_m \times d}.$$

This operation aligns all modality-specific token sequences in a unified representation space, allowing the encoder to process them jointly.

**Modality Augmentation and Cross-Attention Fusion:** Remote sensing observations of a single location are often captured using multiple sensors, each providing a unique spectral or radiometric perspective. Instead of treating these modalities as independent inputs, we interpret them as complementary views of the same scene. This allows us to use modality diversity as a form of natural augmentation, enabling the model to learn sensor-invariant representations. In our setup, each spatial sample $s$ from the Major-TOM dataset [9] is observed via a fixed set of modalities. During pretraining, we independently assign modalities to the student and teacher networks via stochastic selection (threshold = 0.5), ensuring cross-modal supervision. E.g., the teacher may observe a global crop from Sentinel-1, while the student receives local views from Sentinel-2 L2A. This modality augmentation strategy encourages the model to align features across sensors, improving robustness to sensor-specific artifacts. We consider two cases based on the number of selected modalities:

**1) Single-Modality Views:** If only one modality is selected, the input is passed through the corresponding modality-specific patch embedding layer followed by the shared transformer encoder. This follows the standard ViT pipeline but uses modality-aware embeddings to handle spectral channel differences. **2) Multi-Modality Fusion via Cross-Attention:** When multiple modalities are selected, we activate a modality fusion module based on cross-attention. For each selected modality $m \in \mathcal{M}$, we obtain a patch token sequence $\boldsymbol{Z}^{(m)} \in \mathbb{R}^{N \times D}$, where $N$ is the number of spatial positions. These are stacked into a tensor $\boldsymbol{Z}_{\text{all}} \in \mathbb{R}^{N \times M \times D}$, aligning spatial positions across modalities.

For each position $n = 1, \ldots, N$, we define shared learnable queries $\boldsymbol{q} \in \mathbb{R}^{N_q \times D}$, which attend to modality-specific keys $\boldsymbol{K}_n \in \mathbb{R}^{M \times D}$ and values $\boldsymbol{V}_n \in \mathbb{R}^{M \times D}$, yielding $N_q$ intermediate outputs:

$$\boldsymbol{z}_n' = \texttt{MultiHeadAttention}(\boldsymbol{q}, \boldsymbol{K}_n, \boldsymbol{V}_n) \in \mathbb{R}^{N_q \times D}.$$

To aggregate them, we compute a learned weighted mean using softmax-normalized attention scores:

$$\boldsymbol{w} = \texttt{Softmax}(\boldsymbol{z}_n' \cdot \boldsymbol{p}_r), \quad \boldsymbol{z}_n^{\text{fused}} = \sum_{j=1}^{N_q} w_j \boldsymbol{z}_n'[j],$$

where $\boldsymbol{p}_r \in \mathbb{R}^{D \times 1}$ is a learnable projection for scoring the query outputs. This results in a fused token $\boldsymbol{z}_n^{\text{fused}} \in \mathbb{R}^D$. The final sequence $\boldsymbol{Z}_{\text{fused}} \in \mathbb{R}^{N \times D}$ is then passed to the shared encoder. This cross-attention fusion allows the model to dynamically weigh the modality contributions at each spatial location, capturing diverse information while maintaining spatial coherence.

## 3.2 Pretraining

Our pretraining strategy builds on the DINO framework, which performs self-supervised learning. It operates using a teacher-student setup, where both networks share the same ViT backbone and a

lightweight three-layer projection head. Let $g_{\theta_s}$ and $g_{\theta_t}$ denote the student and teacher networks, respectively. While the student is trained using gradient-based optimization, the teacher is updated using EMA of the student's weights:

$$\theta_t \leftarrow \lambda_e\,\theta_t + (1-\lambda_e)\,\theta_s, \quad \lambda_e = 1 - (1-\lambda_0)\frac{1+\cos\left(\pi e/E\right)}{2}, \tag{1}$$

where $e$ is the current epoch, $E$ is the total number of training epochs, and $\lambda_0 \in [0.996, 1)$ is the initial momentum coefficient. The cosine schedule gradually increases $\lambda_e$, stabilizing the teacher updates as training progresses. This EMA mechanism allows the teacher to serve as a temporally smoothed ensemble of past student states, yielding more stable and consistent targets. Fig. 2 shows an overview of TerraFM pre-training.

**Multi-Crop Learning:** To enable scale-invariant and cross-view representation learning, we adopt a multi-crop strategy as used in DINO [4]. For each input sample, we generate two high-resolution global crops $\{\boldsymbol{x}_g^1, \boldsymbol{x}_g^2\} \subset \mathcal{X}_g$ and $J$ low-resolution local crops $\{\boldsymbol{x}_\ell^j\}_{j=1}^J \subset \mathcal{X}_\ell$. The teacher network processes only the global crops, while the student receives both global and local views. Each network produces a $K$-dim output which is temperature-scaled and normalized via the softmax function:

$$Q_s(\boldsymbol{x})^{(i)} = \frac{\exp(g_{\theta_s}(\boldsymbol{x})^{(i)}/\tau_s)}{\sum_{k=1}^{K}\exp(g_{\theta_s}(\boldsymbol{x})^{(k)}/\tau_s)}, \quad Q_t(\boldsymbol{x})^{(i)} = \frac{\exp((g_{\theta_t}(\boldsymbol{x})^{(i)} - c^{(i)})/\tau_t)}{\sum_{k=1}^{K}\exp((g_{\theta_t}(\boldsymbol{x})^{(k)} - c^{(k)})/\tau_t)},$$

where $\tau_s$ and $\tau_t$ are temperature parameters that control output sharpness, and $c$ is a centering term representing the running mean of teacher logits, used to stabilize training and avoid representation collapse. The centering term is updated using exponential moving average over the teacher outputs:

$$c \leftarrow \beta c + (1-\beta) \cdot \frac{1}{B}\sum_{i=1}^{B} g_{\theta_t}(x_i),$$

where $\beta \in [0.9, 0.999]$ controls the momentum, and $B$ is the batch size. The overall loss encourages consistency between teacher and student predictions across all distinct view pairs:

$$\sum_{\boldsymbol{x}\in\mathcal{X}_g}\sum_{\boldsymbol{x}'\in\mathcal{X};\boldsymbol{x}'\neq\boldsymbol{x}} \mathcal{L}_{\text{CE}}\left(Q_t(\boldsymbol{x}), Q_s(\boldsymbol{x}')\right),$$

where $\mathcal{X} = \mathcal{X}_g \cup \mathcal{X}_\ell$, and $\mathcal{L}_{\text{CE}}(\cdot, \cdot)$ denotes the cross-entropy loss. This loss formulation requires the student to produce consistent representations in all views.

**Dual Centering for Long-Tailed Distributions:** Remote sensing datasets often exhibit long-tailed distributions of LULC classes, with frequent categories such as Forest dominating, while classes like Urban or Bare Land remain underrepresented as shown in Fig. 3. This imbalance persists even after subsampling and poses challenges for representation learning. Standard self-supervised approaches like DINO [4] apply a single global centering term to stabilize training and avoid representation collapse, but they do not account for semantic imbalance in the data. To address this, we propose a dual-centering scheme that combines global statistics with class-frequency-aware regularization. In addition to the standard global center vector $c$, we introduce a secondary center $c_h$, computed from a subset of samples belonging to high-frequency LULC classes, such as tree cover, grassland, and open seas, based on dataset-level statistics. Given a batch of teacher logits $g_{\theta_t}(x)$, the adjusted logits for training are computed as:

$$\hat{g}(\boldsymbol{x}) = g_{\theta_t}(\boldsymbol{x}) - \alpha \cdot c - (1-\alpha) \cdot c_h,$$

where $\alpha \in [0, 1]$ balances the contribution of the global and frequency-aware centers. The vector $c_h$ is updated via exponential moving average using only frequent-class samples within each batch. This dual-centering mechanism serves two key purposes: **(i)** it preserves the stability benefits of global centering as in DINO, and **(ii)** it introduces a soft rebalancing bias that counteracts the overrepresentation of dominant classes in the feature space. In ablations (Table 5), this adjustment leads to more balanced representation learning and improved downstream performance, particularly for underrepresented LULC categories.

## 4 Pretraining Data Sampling

We utilize the Major-TOM dataset [9] as our primary EO source for pretraining. It contains 2.24 million globally distributed grid cells, each spanning approximately 10.68 km × 10.68 km ($\approx$114
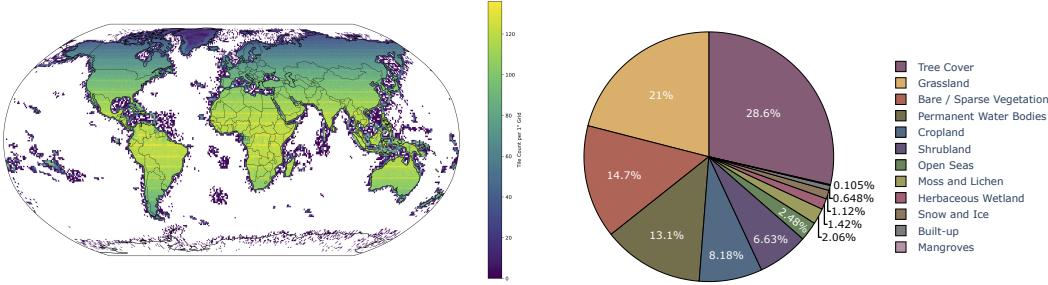
Figure 3: *Left:* Global spatial distribution of the Major-TOM training subset. Each square shows a $1° \times 1°$ cell, colored by the number of 10.68 km × 10.68 km tiles it contains. *Right:* Land-use/land-cover (LULC) breakdown across the same training tiles. A number of semantically important classes (e.g., builtup, mangroves, ice) remain underrepresented due to skewed data distribution.

$km^2$), and provides tri-modal, co-registered imagery from Sentinel-2 Level-1C, Sentinel-2 Level-2A, and Sentinel-1 RTC. Major-TOM stands out as one of the few publicly available datasets offering dense multi-modal coverage at a global scale. However, over one-third of its samples lie outside a 10 km terrestrial buffer, often within the Open Oceans class [38], limiting their relevance for land-centric tasks. Motivated by insights from [23], which emphasize the importance of semantically rich samples, and [33], which highlight the utility of structural priors, we applied a principled filtering strategy. Specifically, we removed 98% of ocean-classified tiles (retaining 2% to preserve marine representation) and sampled the terrestrial subset using global distributional priors across land cover [38], climate zones [3], and ESRI world regions [8]. This approach emphasizes meaningful land regions with ecological variety. Fig. 3 shows the global coverage of the tiles.

For pretraining, we curated a filtered subset of over 1.5 million grid cells with consistent coverage across all three modalities (S1, S2-L1C, and S2-L2A). (Need to highlight the disk space usage, it takes around 76 terabytes to store 1.5M samples in uint16). Each 10.68 km $\times$ 10.68 km grid cell was divided into four non-overlapping tiles of $534 \times 534$ pixels, resulting in more than 6 million tiles per modality. In total, this yielded 18.7 million modality-specific training tiles. During training, modalities were stochastically sampled and treated as natural augmentations to promote sensor-invariant representation learning. To mitigate spatial sampling bias and support semantically-aware learning, we enriched each grid cell with metadata from the ESRI World Regions dataset [8].

## 5 Experiments and Results

### 5.1 Pretraining Implementation Details

We pretrain our TerraFM with a $16 \times 16$ patch resolution and an input size of $224 \times 224$. The training dataset comprises around 1.53 million multi-modal samples, from which we define a virtual epoch of 300K samples to ensure frequent parameter updates and improved memory efficiency. TerraFM-B is trained for 150 epochs and TerraFM-L is trained for 200 epochs with a linear warmup over the first 30 epochs. Models are trained on 64 GPUs and the TerraFM-B training takes 92 hours with a batch size of 1024 where as the TerraFM-L use a batch size of 2048 and training time is 183 hours. The learning rate is linearly scaled with batch size, initialized as lr $= 0.0001 \times$ batch_size$/256$. Following DINO-style pretraining, we disable batch normalization in the projection head and freeze the last layer of the student for the initial 3 epochs to stabilize early training. Following DINO [4], we use two global crops with scale sampled from $[0.25, 1.0]$ and six local crops from $[0.05, 0.25]$. The output dimensionality is set to $K = 65,536$, with a teacher temperature schedule linearly increasing from 0.04 to 0.06 over the first 50 epochs. The momentum parameter for the teacher network follows a cosine schedule, starting from 0.996. A drop path rate of 0.1 is applied to regularize training. We set $N_q = 5$ and $\alpha = 0.8$ during pre-training.

### 5.2 Evaluation Implementation Details

**Linear Probing Evaluation:** To evaluate the quality of learned representations, we follow a linear probing protocol of DINOv2[19] that follows with a lightweight grid search over three key

Table 2: We evaluate image classification using k-nearest neighbors (kNN) and report Top-1 accuracy for all single-label tasks. For the multilabel BigEarthNet benchmark, we report the F1 score. Results other than Copernicus-FM and TerraFM are directly taken from [28].

| Model | Backbone | m-EuroSat Training % | | m-BigEarthNet Training % | | m-So2Sat Training % | | m-Brick-Kiln Training % | |
|---|---|---|---|---|---|---|---|---|---|
| | | 100% | 1% | 100% | 1% | 100% | 1% | 100% | 1% |
| SatMAE | ViT-Base | 84.1 | 34.8 | 50.6 | 29.0 | 36.0 | 23.1 | 86.1 | 73.5 |
| SatMAE++ | ViT-Large | 82.7 | 48.5 | 50.8 | 31.6 | 34.7 | 23.4 | 89.6 | 76.7 |
| CROMA | ViT-Base | 85.6 | 51.3 | 58.8 | 44.7 | 48.8 | 33.8 | 92.6 | 85.1 |
| SoftCon | ViT-Small | 89.8 | 27.2 | 64.7 | 43.3 | 51.1 | 31.4 | 89.2 | 77.8 |
| DOFA | ViT-Base | 82.8 | 49.6 | 49.4 | 29.9 | 41.4 | 29.4 | 88.3 | 78.3 |
| Satlas | Swin-Tiny | 81.7 | 35.8 | 51.9 | 29.6 | 36.6 | 27.1 | 88.2 | 73.0 |
| MMEarth | CNN-atto | 81.7 | 30.0 | 58.3 | 39.6 | 39.8 | 25.1 | 89.4 | 79.7 |
| DeCUR | ViT-Small | 89.0 | 46.6 | 63.8 | 49.6 | 45.8 | 30.9 | 83.7 | 74.2 |
| AnySat | ViT-Base | 82.2 | 47.1 | 54.9 | 33.7 | 39.8 | 29.0 | 85.3 | 72.0 |
| Galileo | ViT-Base | 93.0 | 56.6 | 59.0 | 36.5 | 54.8 | **43.2** | 90.7 | 78.0 |
| Prithvi-2.0 | ViT-Large | 80.2 | 48.0 | 49.4 | 28.8 | 29.5 | 26.1 | 87.9 | 80.6 |
| Copernicus-FM | ViT-Base | 76.0 | 47.4 | 53.8 | 33.3 | 38.4 | 23.3 | 93.0 | 83.2 |
| TerraFM | ViT-Base | 94.2 | 59.3 | 68.7 | 49.4 | 55.1 | 41.6 | **94.5** | **85.6** |
| | ViT-Large | **95.1** | **62.1** | **69.4** | **50.6** | **55.9** | 41.1 | 93.0 | 82.2 |

hyperparameters: (i) the learning rate, (ii) the number of transformer layers from which features are extracted, and (iii) whether to use only the [CLS] token or to concatenate it with the average-pooled patch tokens. We train the linear classifier using stochastic gradient descent (SGD) for 50 epochs. The training data is augmented using random resized cropping. Specifically, we sweep the learning rate over the set $\{10^{-5}, 2 \times 10^{-5}, 5 \times 10^{-5}, 10^{-4}, 2 \times 10^{-4}, 5 \times 10^{-4}, 10^{-3}, 2 \times 10^{-3}, 5 \times 10^{-3}, 10^{-2}, 2 \times 10^{-2}, 5 \times 10^{-2}, 1\}$ Importantly, this search is computationally efficient: features from the frozen backbone are computed once per image using a single forward pass and reused across all configurations, since each linear head only requires a simple forward pass. For each configuration, we evaluate the classifier on the validation set and report the test accuracy achieved by the best validation configuration. **UperNet Probing Evaluation:** For UperNet [34] Probing evaluation, we freeze the pretrained backbone and attach UPerNet decoder head. Specifically, we use a `Feature2Pyramid` module as the neck, followed by a UPerNet decoder and an auxiliary FCNHead. We train only the segmentation heads using the AdamW optimizer for 50 epochs without learning rate warm-up. We conduct a grid search over base learning rates $\{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. and batch size set $\{16, 32, 64\}$. **k-NN Evaluation:** To assess the quality of the learned representations without any finetuning, we apply non-parametric classification using a $k$-nearest neighbors (k-NN) classifier on the frozen features. In addition to sweeping over $k \in 3, 5, 7, 10, 15, 20, 30, 50, 100$ using validation performance, we follow the same layer selection strategy as linear probing i.e evaluating features from the last 4 transformer layers. This protocol does not require additional training or data augmentation, making it a lightweight and reliable indicator of raw feature quality in pretrained models. **Finetuning Evaluation:** For full-model finetuning, we unfreeze the backbone and jointly optimize it with the task-specific head. We perform a grid search over learning rates in the set and batch sizes. To stabilize training, we apply a reduced learning rate for the backbone, set to half of the main learning rate used for the head parameters. Once the best configuration is selected based on validation performance, we evaluate the finetuned model on the test set.

### 5.3 Evaluating Downstream Tasks

**Benchmarks:** We evaluate our model on two comprehensive remote sensing benchmarks: **GEO-Bench** and **Copernicus-Bench**, both of which include diverse downstream tasks spanning multiple domains and modalities. See more details in suppl. material.

**Discussion:** We report KNN classification accuracy on four standard GEO-Bench classification tasks to evaluate the quality of learned representations in a training-free setting. As shown in Tab. 2, TerraFM achieves the highest performance across three datasets, outperforming both modality-specific

and multimodal foundation models. Notably, our model achieves 95.1% on m-EuroSAT and 94.5% on m-Brick-Kiln, highlighting the effectiveness of the learned representations on standard scene classification tasks. On other challenging tasks such as m-So2Sat and m-BigEarthNet, our model achieves leading performance (55.9% and 69.4%, respectively), outperforming Galileo [28], despite So2Sat having fewer channels than used during pretraining, highlighting the model's robustness to missing modality information. Compared to CROMA [10] and DeCUR [30], our gains suggest that contrastive alignment combined with cross-modal fusion enhances class separability. The results across tasks of varying difficulty indicate that our model learns robust and transferable representations that generalize well across different scenarios.

Further on GEO-Bench, for classification (with fine-tuning), TerraFM achieves the improvement on m-BigEarthNet (73.1%) and m-EuroSat (98.6%), and the best-performing model on m-So2Sat (66.6%). For segmentation (with linear probing), our TerraFM-L notably outperforms existing models on m-SA-Crop-Type (34.5% mIoU) and m-Cashew-Plant (37.2% mIoU). Tab. 4 shows that TerraFM-B surpasses larger counterparts such as ViT-Large used in SatMAE++ and DOFA.

On the Copernicus-Bench [32] evaluation, our model consistently outperforms existing foundation models across tasks and modalities (Tab. 3). A comprehensive comparison of TerraFM with existing methods on Copernicus-Bench, using metrics like OA (Overall Accuracy), mAP (mean Average Precision), and mIoU (mean Intersection over Union). Notably, TerraFM consistently achieves the highest scores across most of the tasks and metrics. In particular, it achieves an OA of 99.1% on EuroSAT-S2, an mAP of 84.4% on BigEarthNet-S2, and an mIoU of 67.9% on Cloud-S2.

Table 3: Comparison of TerraFM with existing supervised and self-supervised methods on Copernicus-Bench. Metrics include OA (Overall Accuracy) for classification tasks, mAP (mean Average Precision) for multi-label classification, and mIoU (mean Intersection over Union) for segmentation.

| | Metric | Supervised | Random | SoftCon | CROMA | DOFA | Copernicus-FM | TerraFM |
|---|---|---|---|---|---|---|---|---|
| Backbone | – | ViT-B/16 | ViT-B/16 | ViT-B/14 | ViT-B/8 | ViT-B/16 | ViT-B/16 | ViT-B/16 |
| Cloud-S2 | mIoU | 59.4 | 60.4 | 66.9 | 65.0 | 65.0 | 66.7 | **67.9** |
| EuroSAT-S1 | OA | 81.5 | 75.4 | 83.6 | 83.9 | 81.7 | 87.2 | **87.8** |
| EuroSAT-S2 | OA | 97.6 | 92.5 | 96.7 | 97.0 | 97.2 | 97.9 | **99.1** |
| BigEarthNet-S1 | mAP | 70.6 | 63.8 | **78.7** | 70.8 | 70.5 | 77.9 | 76.9 |
| BigEarthNet-S2 | mAP | 80.1 | 71.6 | 83.6 | 76.4 | 75.5 | 79.0 | **84.4** |
| DFC2020-S1 | mIoU | 50.8 | 45.4 | 52.8 | 52.7 | 49.7 | 52.4 | **55.4** |
| DFC2020-S2 | mIoU | 66.2 | 62.3 | 64.1 | **66.5** | 61.8 | 64.5 | 63.8 |
| LCZ-S2 | OA | 85.3 | 77.4 | 83.6 | 84.1 | 83.0 | 84.4 | **87.0** |

Table 4: Performance comparison on GEO-Bench for both classification (Top-1 Accuracy), segmentation (mIoU), and F1 score (for m-BigEarthNet). TerraFM achieves state-of-the-art results across multiple datasets, outperforming previous FMs.

| | | Classification | | | | Segmentation | |
|---|---|---|---|---|---|---|---|
| Method | Backbone | m-EuroSat | m-BigEarthNet | m-So2Sat | m-Brick-Kiln | m-Cashew-Plant | m-SA-Crop-Type |
| SatMAE | ViT-Large | 96.6 | 68.3 | 57.2 | 98.4 | 30.8 | 24.8 |
| SatMAE++ | ViT-Large | 96.5 | 67.9 | 56.0 | 98.6 | 29.6 | 25.7 |
| CROMA | ViT-Large | 96.6 | 71.9 | 60.6 | 98.7 | 31.8 | 32.0 |
| SoftCon | ViT-Base | 97.5 | 70.3 | 61.7 | 98.7 | 29.6 | 30.8 |
| DOFA | ViT-Large | 96.9 | 68.0 | 58.7 | 98.6 | 27.7 | 25.4 |
| Satlas | Swin-Base | 97.5 | _72.8_ | 61.9 | 98.4 | 25.1 | 23.4 |
| MMEarth | CNN-atto | 95.7 | 70.0 | 57.2 | _98.9_ | 24.2 | 22.2 |
| DeCUR | ViT-Small | 97.9 | 70.9 | 61.7 | 98.7 | 26.2 | 21.5 |
| Prithvi 2.0 | ViT-Large | 96.5 | 69.0 | 54.6 | 98.6 | 26.7 | 22.9 |
| AnySat | ViT-Base | 95.9 | 70.3 | 51.8 | 98.6 | 26.1 | 27.1 |
| Galileo | ViT-Base | 97.7 | 70.7 | 63.3 | 98.7 | 33.0 | 30.1 |
| TerraFM | ViT-Base | _98.1_ | 72.6 | _64.9_ | 98.7 | _34.1_ | _33.0_ |
| | ViT-Large | **98.6** | **73.1** | **66.6** | **99.0** | **37.2** | **34.5** |

## 5.4 Ablations and Analysis

**Imapct of Components:** Tab. 5 highlights the incremental benefits of each component in our framework. We train TerraFM-B for 150 epochs on a 200k-sample subset from our full training dataset. The model was trained on a subset of the training data, and KNN classification accuracy is reported for each dataset. To measure the performance on segmentation task, we use uppernet probing on the m-Cashew-Plantation dataset from GeoBench. Adding modality as augmentation improves performance on m-EuroSat by +4.5 and m-BigEarthNet by +3.01. Incorporating fusion yields a large gain on m-Cashew-Plantation by +11.82, while dual centering provides further improvements: +3.44 on m-BigEarthNet, +7.2 on m-EuroSat, and +14.0 on m-Cashew-Plantation.

Table 5: Ablation of components: SS = Self-supervised contrastive learning, MAug = Modality Augmentation, Fus = Fusion, DC = Dual Centering. BEN = m-BigEarthNet, ES = m-EuroSat, CP = m-Cashew-Plant.

| SS | MAug | Fus | DC | BEN | ES | CP |
|----|------|-----|----|-----|-----|-----|
| ✓ | – | – | – | 54.62 | 83.20 | 50.58 |
| ✓ | ✓ | – | – | 57.63 | 87.70 | 59.17 |
| ✓ | ✓ | ✓ | – | 57.74 | 88.50 | 62.40 |
| ✓ | ✓ | ✓ | ✓ | 58.06 | 90.40 | 64.58 |

**Dual-centering Motivation and Visualization:** Here, we discuss the impact of Dual-centering on class-wise prediction behavior and representation diversity. Fig. 4 shows that models with Dual-centering exhibit higher softmax entropy across most classes, indicating more calibrated predictions, particularly benefiting rare classes like "Mangroves". Fig. 5 reveals that Dual Centering significantly increases prototype diversity, i.e., the number of distinct top-5 features activated, especially for tail classes. This suggests that the model avoids collapsing onto frequent-class prototypes and learns more diverse, semantically rich representations. These results motivate Dual-centering as an effective strategy for reducing class imbalance effects in representation learning.
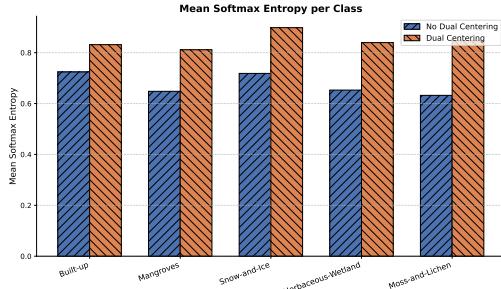


Figure 4: Mean entropy per LULC class computed on 5k uniformly sampled training samples. Logits ($K{=}65,536$) are projected to a lower-dimensional space using a fixed random Gaussian matrix before computing entropy. The baseline model (No Dual-centering) exhibits lower entropy for most classes, showing overconfident predictions biased toward frequent-class prototypes. In contrast, our Dual-centering model yields higher entropy, suggesting reduced dominance of high-frequency prototypes, especially for rare classes ("Mangroves", "Herbaceous-Wetland").
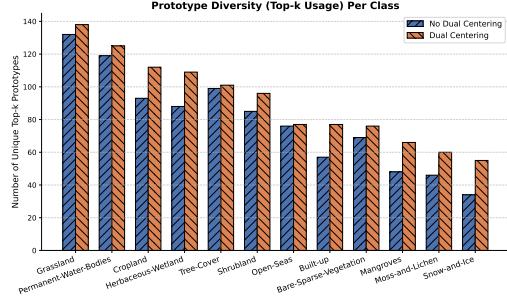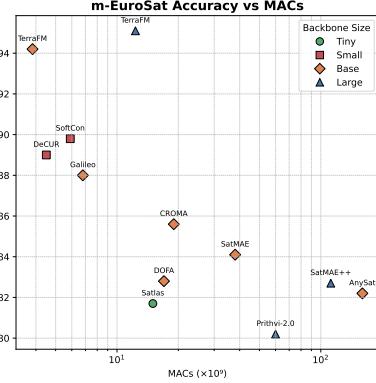
Figure 5: Prototype diversity measured as the number of unique top-5 prototypes activated across 5k uniformly sampled training samples. Dual-centering leads to greater prototype diversity in tail classes such as "Mangroves", "Herbaceous-Wetland", and "Built-up", suggesting more diverse representation learning. The baseline (No Dual-centering) tends to reuse a smaller subset of prototypes, especially for rare classes, reflecting over-reliance on dominant features from high-frequency categories.

**MACs-Performance Trade-Off:** We evaluate the compute-efficiency trade-off of various remote sensing foundation models using Multiply-Accumulate operations (MACs) as a measure of inference cost. As shown in Fig. 6, TerraFM achieves the highest accuracy on m-EuroSat while operating at significantly lower MACs compared to other large-scale models. This highlights the efficiency of our fusion design and pretraining strategy, demonstrating that strong performance can be achieved without excessive computational overhead. Models with higher MACs do not consistently translate to better accuracy, highlighting the importance of efficient and expressive architectures for scalable EO applications.

Figure 6: Model accuracy vs. efficiency on m-EuroSat dataset. Each point reports a backbone's k-NN classification accuracy (y-axis) against its inference MACs in billions (log-scaled x-axis). Marker shape and colour encode backbone size, as indicated by the inset legend. We report here two variants of the TerraFM (Base, Large), which achieve the highest accuracy while maintaining moderate computational cost relative to both lightweight and heavyweight baselines.

## 6  Conclusion

In this work, we introduced TerraFM, a unified and scalable foundation model (FM) specifically designed for multisensor EO. Given the unique nature of EO data, our approach pays special treatment to sensor heterogeneity, scale-invariance, and class-frequency imbalance which is critical for building generalizable EO FMs. Our pretraining approach leverages contrastive learning to obtain geographically and spectrally aware representations from large-scale Sentinel-1 and 2 data. Specifically, we integrate modality-specific patch embeddings, adaptive cross-attention fusion, and a dual-centering contrastive learning objective to enrich the representations on heterogeneous RS data. Our extensive evaluations on GEO-Bench and Copernicus-Bench demonstrate that TerraFM consistently outperforms SoTA self-supervised ViT models across both classification and segmentation tasks.

## References

[1] Astruc, G., N. Gonthier, C. Mallet, and L. Landrieu (2024). Anysat: An earth observation model for any resolutions, scales, and modalities. *arXiv preprint arXiv:2412.14123*.

[2] Bastani, F., P. Wolters, R. Gupta, J. Ferdinando, and A. Kembhavi (2023). Satlaspretrain: A large-scale dataset for remote sensing image understanding.

[3] Beck, H. E., N. E. Zimmermann, T. R. McVicar, N. Vergopolan, A. Berg, and E. F. Wood (2018). Present and future köppen-geiger climate classification maps at 1-km resolution. *Scientific data 5*(1), 1–12.

[4] Caron, M., H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin (2021). Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660.

[5] Chan-To-Hing, H. and B. Veeravalli (2024). Fus-mae: A cross-attention-based data fusion approach for masked autoencoders in remote sensing. In *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing Symposium*, pp. 6953–6958. IEEE.

[6] Dosovitskiy, A., L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby (2021). An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*.

[7] Drusch, M., U. Del Bello, S. Carlier, O. Colin, V. Fernandez, F. Gascon, B. Hoersch, C. Isola, P. Laberinti, P. Martimort, et al. (2012). Sentinel-2: Esa's optical high-resolution mission for gmes operational services. *Remote sensing of Environment 120*, 25–36.

[8] Esri, Global Mapping International, and U.S. Central Intelligence Agency (2025). World Regions. `https://www.arcgis.com/home/item.html?id=84dbc97915244e35808e87a881133d09`. Layer package representing boundaries for 25 commonly recognized world regions. Updated April 29, 2025. Accessed April 30, 2025.

[9] Francis, A. and M. Czerkawski (2024). Major tom: Expandable datasets for earth observation. In *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing Symposium*, pp. 2935–2940. IEEE.

[10] Fuller, A., K. Millard, and J. Green (2023). Croma: Remote sensing representations with contrastive radar-optical masked autoencoders. *Advances in Neural Information Processing Systems 36*, 5506–5538.

[11] Gao, P., T. Ma, H. Li, Z. Lin, J. Dai, and Y. Qiao (2022). MCMAE: Masked convolution meets masked autoencoders. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho (Eds.), *Advances in Neural Information Processing Systems*.

[12] Guo, X., J. Lao, B. Dang, Y. Zhang, L. Yu, L. Ru, L. Zhong, Z. Huang, K. Wu, D. Hu, et al. (2024). Skysense: A multi-modal remote sensing foundation model towards universal interpretation for earth observation imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27672–27683.

[13] Han, B., S. Zhang, X. Shi, and M. Reichstein (2024). Bridging remote sensors with multisensor geospatial foundation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27852–27862.

[14] He, K., X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick (2022). Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009.

[15] Huot, F., R. L. Hu, N. Goyal, T. Sankar, M. Ihme, and Y.-F. Chen (2022). Next day wildfire spread: A machine learning dataset to predict wildfire spreading from remote-sensing data. *IEEE Transactions on Geoscience and Remote Sensing 60*, 1–13.

[16] Kuckreja, K., M. S. Danish, M. Naseer, A. Das, S. Khan, and F. S. Khan (2024). Geochat: Grounded large vision-language model for remote sensing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27831–27840.

[17] Kussul, N., M. Lavreniuk, S. Skakun, and A. Shelestov (2017). Deep learning classification of land cover and crop types using remote sensing data. *IEEE Geoscience and Remote Sensing Letters 14*(5), 778–782.

[18] Li, X., D. Hong, and J. Chanussot (2024). S2mae: A spatial-spectral pretraining foundation model for spectral remote sensing data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24088–24097.

[19] Oquab, M., T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, et al. (2023). Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*.

[20] Prodhan, F. A., J. Zhang, F. Yao, L. Shi, T. P. Pangali Sharma, D. Zhang, D. Cao, M. Zheng, N. Ahmed, and H. P. Mohana (2021). Deep learning for monitoring agricultural drought in south asia using remote sensing data. *Remote sensing 13*(9), 1715.

[21] Rahnemoonfar, M., T. Chowdhury, A. Sarkar, D. Varshney, M. Yari, and R. R. Murphy (2021). Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access 9*, 89644–89654.

[22] Reed, C. J., R. Gupta, S. Li, S. Brockman, C. Funk, B. Clipp, K. Keutzer, S. Candido, M. Uyttendaele, and T. Darrell (2023). Scale-mae: A scale-aware masked autoencoder for multiscale geospatial representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4088–4099.

[23] Roscher, R., M. Russwurm, C. Gevaert, M. Kampffmeyer, J. A. Dos Santos, M. Vakalopoulou, R. Hänsch, S. Hansen, K. Nogueira, J. Prexl, et al. (2024). Better, not just more: Data-centric machine learning for earth observation. *IEEE Geoscience and Remote Sensing Magazine*.

[24] Sarkar, A., T. Chowdhury, R. R. Murphy, A. Gangopadhyay, and M. Rahnemoonfar (2023). Sam-vqa: Supervised attention-based visual question answering model for post-disaster damage assessment on remote sensing imagery. *IEEE Transactions on Geoscience and Remote Sensing 61*, 1–16.

[25] Szwarcman, D., S. Roy, P. Fraccaro, Þ. E. Gíslason, B. Blumenstiel, R. Ghosal, P. H. de Oliveira, J. L. d. S. Almeida, R. Sedona, Y. Kang, et al. (2024). Prithvi-eo-2.0: A versatile multi-temporal foundation model for earth observation applications. *arXiv preprint arXiv:2412.02732*.

[26] Tang, M., A. Cozma, K. Georgiou, and H. Qi (2023). Cross-scale mae: A tale of multiscale exploitation in remote sensing. *Advances in Neural Information Processing Systems 36*, 20054–20066.

[27] Torres, R., P. Snoeij, D. Geudtner, D. Bibby, M. Davidson, E. Attema, P. Potin, B. Rommen, N. Floury, M. Brown, et al. (2012). Gmes sentinel-1 mission. *Remote sensing of environment 120*, 9–24.

[28] Tseng, G., A. Fuller, M. Reil, H. Herzog, P. Beukema, F. Bastani, J. R. Green, E. Shelhamer, H. Kerner, and D. Rolnick (2025). Galileo: Learning global and local features in pretrained remote sensing models. *arXiv preprint arXiv:2502.09356*.

[29] Waldmann, L., A. Shah, Y. Wang, N. Lehmann, A. J. Stewart, Z. Xiong, X. X. Zhu, S. Bauer, and J. Chuang (2025). Panopticon: Advancing any-sensor foundation models for earth observation. *arXiv preprint arXiv:2503.10845*.

[30] Wang, Y., C. M. Albrecht, N. A. A. Braham, C. Liu, Z. Xiong, and X. X. Zhu (2024). Decoupling common and unique representations for multimodal self-supervised learning. In *European Conference on Computer Vision*, pp. 286–303. Springer.

[31] Wang, Y., Y. Sun, X. Cao, Y. Wang, W. Zhang, and X. Cheng (2023). A review of regional and global scale land use/land cover (lulc) mapping products generated from satellite remote sensing. *ISPRS Journal of Photogrammetry and Remote Sensing 206*, 311–334.

[32] Wang, Y., Z. Xiong, C. Liu, A. J. Stewart, T. Dujardin, N. I. Bountos, A. Zavras, F. Gerken, I. Papoutsis, L. Leal-Taixé, et al. (2025). Towards a unified copernicus foundation model for earth vision. *arXiv preprint arXiv:2503.11849*.

[33] Wang, Z., R. Prabha, T. Huang, J. Wu, and R. Rajagopal (2024). Skyscript: A large and semantically diverse vision-language dataset for remote sensing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Volume 38, pp. 5805–5813.

[34] Xiao, T., Y. Liu, B. Zhou, Y. Jiang, and J. Sun (2018). Unified perceptual parsing for scene understanding. In *European Conference on Computer Vision*. Springer.

[35] Xie, Z., Z. Zhang, Y. Cao, Y. Lin, Y. Wei, Q. Dai, and H. Hu (2023). On data scaling in masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10365–10374.

[36] Xiong, Z., Y. Wang, F. Zhang, A. J. Stewart, J. Hanna, D. Borth, I. Papoutsis, B. Le Saux, G. Camps-Valls, and X. X. Zhu (2024). Neural plasticity-inspired foundation model for observing the earth crossing modalities. *arXiv e-prints*, arXiv–2403.

[37] Yu, D. and C. Fang (2023). Urban remote sensing with spatial big data: A review and renewed perspective of urban studies in recent decades. *Remote Sensing 15*(5), 1307.

[38] Zanaga, D., R. Van De Kerchove, D. Daems, W. De Keersmaecker, C. Brockmann, G. Kirches, J. Wevers, O. Cartus, M. Santoro, S. Fritz, et al. (2022). Esa worldcover 10 m 2021 v200.

# Supplementary Material

This supplementary material presents additional experiments, analyses, and visualizations that complement the main paper. It includes detailed descriptions and experiments for our multimodal fusion strategies (S1), and qualitative figures (S2). We also report GPU-hour comparisons with comparable methods (S3), and visualize the land cover distribution of our dataset using global maps (S4).

## S1    Multi-Modal Fusion Strategies:

We investigate various strategies for multi-modal fusion and report results in Table A1 on two benchmark datasets: m-BigEarthNet and m-EuroSat. As a baseline, we evaluate standard DINO training using only Sentinel-2 L2A input (*DINO (S2-L2A)*), which learns unimodal representations. To enable explicit modality-aware learning, we apply a *Multi-Student-Teacher* approach where each modality has its own student and teacher networks, along with an alignment loss between student outputs to enforce cross-modal consistency. This yields consistent gains across both datasets. We also test a more expressive fusion approach, *CrossAttn (Q = 196) Global*, where 196 learned queries (standard for 224×224 image inputs) attend globally to multi-modal tokens immediately after patch embedding. However, this method does not perform well, likely due to excessive parameterization and lack of inductive bias for spatial alignment. Figure A1 visually summarizes key fusion strategies evaluated in Table A1, including (a) Multi-Student-Teacher, (b) unimodal DINO, and (c) CrossAttn (Q = 196) Global, highlighting their architectural differences and fusion mechanisms. Our proposed approach, *TerraFM-B (Q = 1)*, treats a modality as an augmentation and performs fusion using a single learned spatial query per location. This lightweight attention mechanism yields the best performance among non-ensemble methods. To further analyze architectural choices, we test a variant, *TerraFM-B (ViT PatchEmb)*, where the convolutional patch embedding is replaced by a ViT-S backbone purely for token extraction. While competitive, this setup slightly drops the performance due to increased model complexity and potential overfitting. Finally, our full model, *TerraFM-B (Q = 5)*, employs multiple learned spatial queries to achieve richer fusion between modalities. It achieves the best overall performance, validating the scalability and effectiveness of our fusion design.

Table A1: Ablation study on multi-modal fusion strategies using k-NN evaluation. TerraFM-B with multiple spatial queries (Q = 5) achieves the best performance.

|  | m-BigEarthNet | m-EuroSat |
|---|---|---|
| DINO (S2-L2A) | 54.6 | 83.2 |
| Multi-Student-Teacher | 55.8 | 87.8 |
| CrossAttn (Q = 196) Global | 52.0 | 77.1 |
| TerraFM-B (Q = 1) | 57.2 | 89.2 |
| TerraFM-B (ViT PatchEmb) | 56.9 | 87.2 |
| TerraFM-B (Q = 5) | **58.1** | **90.4** |

## S2    Qualitative:

Fig. A2 illustrates qualitative results for the cloud and cloud shadow segmentation task. TerraFM accurately outlines both cloud and shadow regions, effectively distinguishing visually similar patterns while maintaining spatial coherence across varied scenes. These results demonstrate the model's strong generalization ability under diverse and challenging atmospheric conditions.

## S3    Landslide Detection

We evaluate landslide segmentation on the Landslide4Sense (L4S) benchmark, which provides segmentation labels for landslide and non-landslide regions across diverse mountainous areas using multi-source satellite data, including Sentinel-2 bands, DEM, and slope information. Our method, TerraFM, achieves strong performance with a mean IoU of 70.8 and a landslide IoU of 43.1, outperforming the Prithvi-EO-2.0 baseline (Table A2). Both TerraFM and Prithvi-EO-2.0 are trained using focal loss with a batch size of 16, Adam optimizer with a learning rate of $1 \times 10^{-4}$. Figure A3 shows qualitative results from TerraFM, illustrating predicted landslide masks alongside the ground truth.
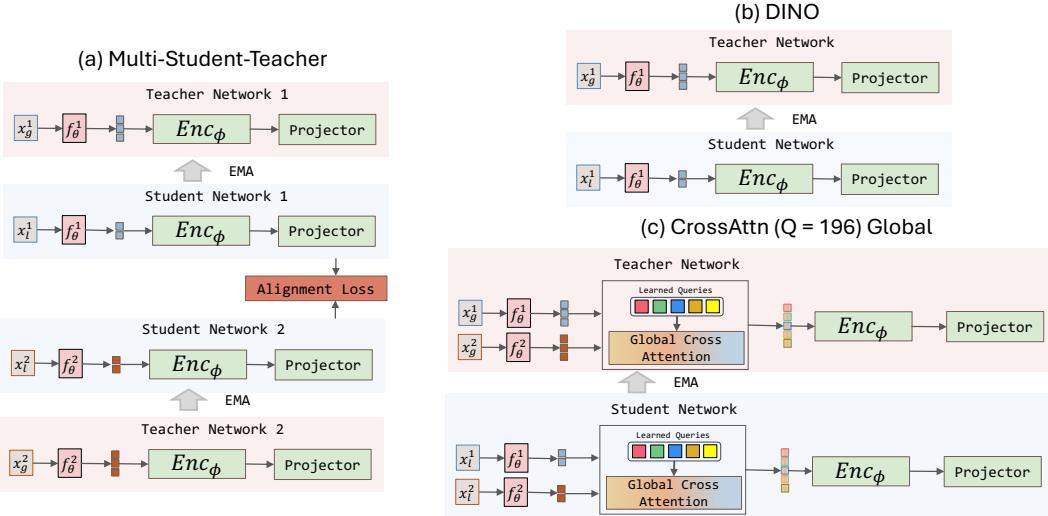
Figure A1: Architectural overview of different fusion strategies: (a) Multi-Student-Teacher with alignment loss, (b) unimodal DINO baseline, and (c) CrossAttn (Q = 196) with global learned queries.
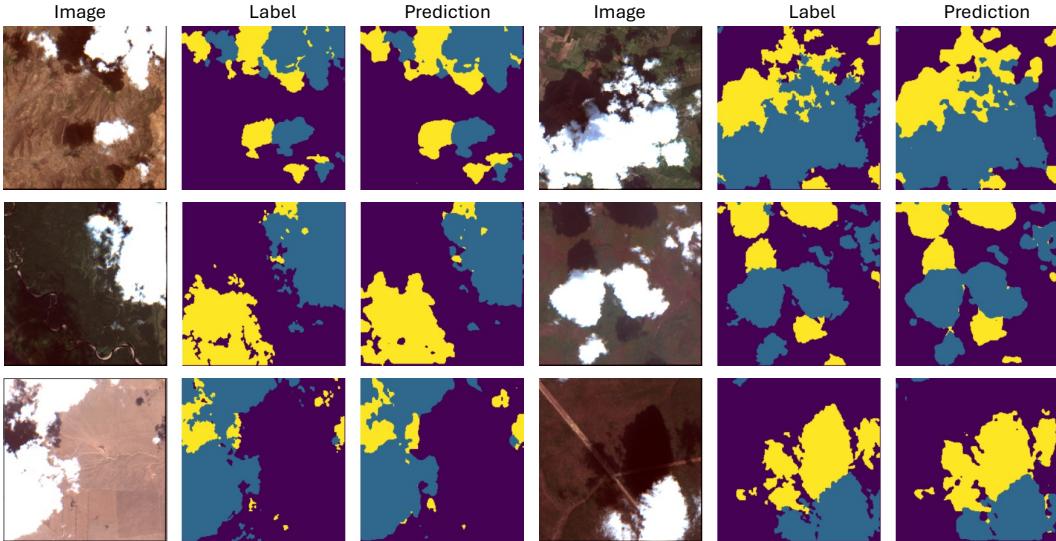


Figure A2: Qualitative results for cloud and cloud shadow segmentation. Each triplet shows the input image (left), the ground truth mask (middle), and the TerraFM prediction (right).

## S3 GPU Hour Comparison:

Compared to Prithvi-2.0, which trains ViT-L (300M) model using up to 80 GPUs for 400 epochs, consuming approximately 21,000 GPU-hours [25], our TerraFM (300M) achieves comparable scale using significantly fewer resources. Specifically, TerraFM is trained for 200 epochs on 64 GPUs, amounting to approximately 12,000 GPU-hours.

## S4 Land Cover Distribution:

Fig. A4 illustrates the global spatial coverage of our pretraining data. The selected samples span diverse ecosystems, capturing a balanced mix of urban, vegetation, sea, and arid regions. The insets demonstrate fine-grained land cover variability, ensuring semantic richness across training tiles.

Table A2: Landslide detection performance on the Landslide4Sense test set. Despite being significantly smaller (120M parameters vs. 300M for Prithvi-EO-2.0), TerraFM achieves higher overall segmentation performance, especially for landslide regions.

|                      | mIoU | IoU (Landslide) |
|----------------------|------|-----------------|
| Prithvi-EO-2.0 (300M) | 65.0 | 31.5            |
| TerraFM (120M)       | **70.8** | **43.1**    |

This diverse geographic grounding plays a crucial role in enabling the generalization capabilities of TerraFM across regions and tasks.
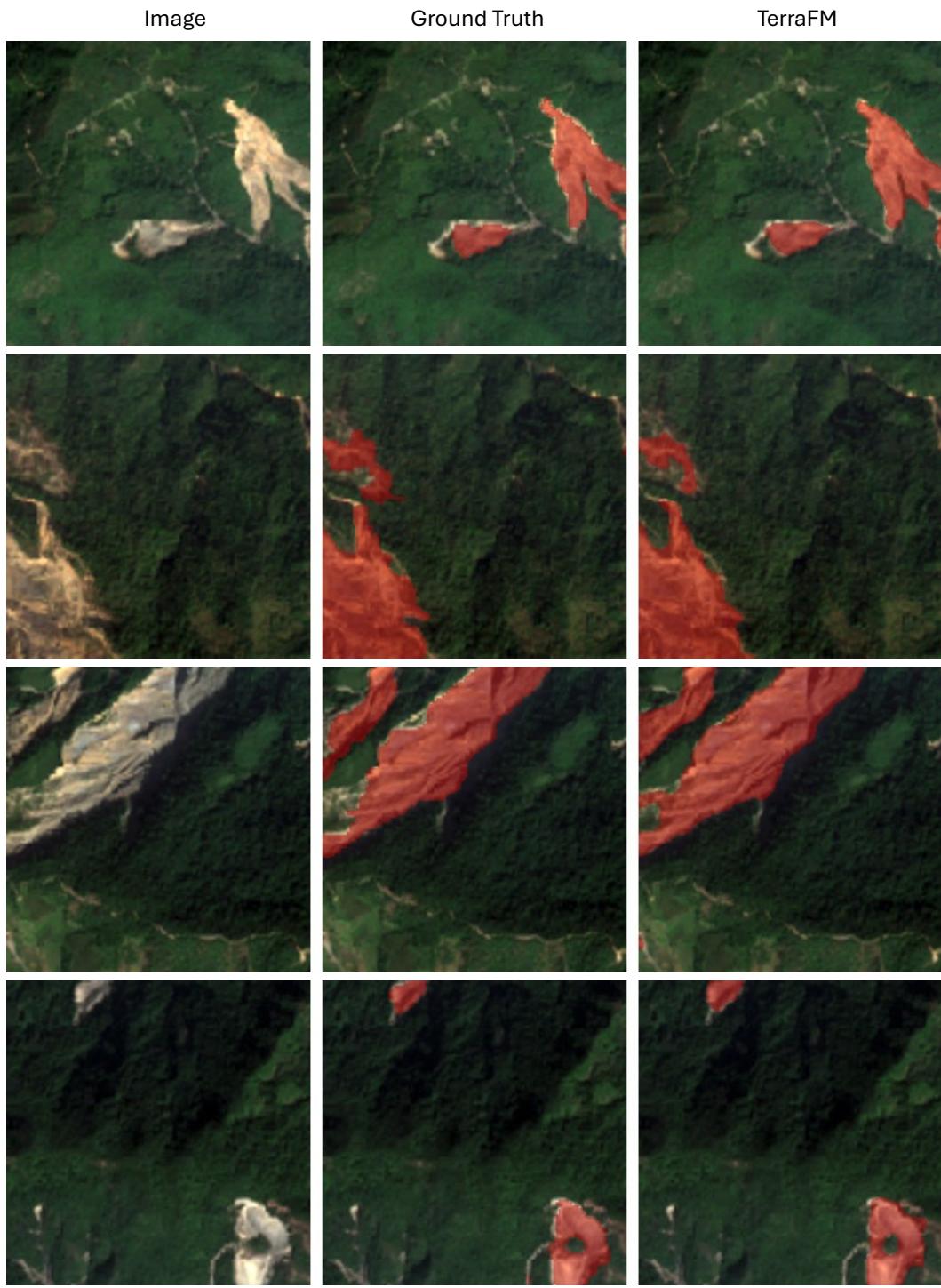
Figure A3: Qualitative results for landslide segmentation. Each triplet shows the input image (left), the ground truth mask (middle), and the TerraFM prediction (right).
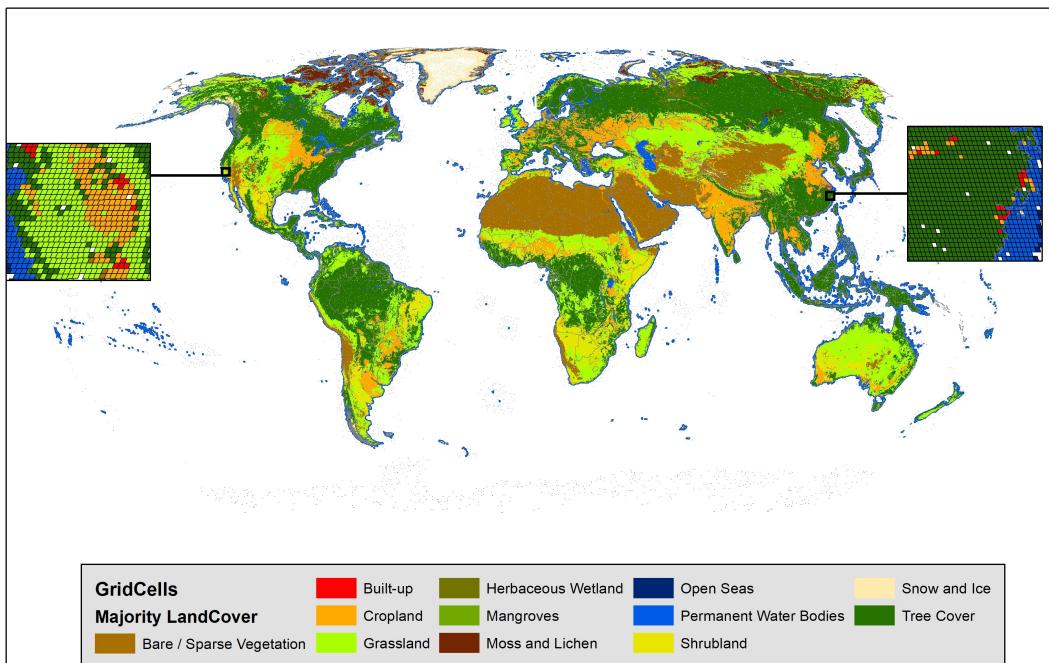
Figure A4: Global distribution of sampled training tiles by dominant land cover class, based on ESA WorldCover labels. Insets show detailed tile-level diversity, highlighting coverage across built-up, vegetation, and water classes.