

Towards Privacy-preserved Pre-training of Remote Sensing Foundation Models with Federated Mutual-guidance Learning

Jieyi Tan¹, Chengwei Zhang², Bo Dang¹, Yansheng Li^{1,*}

¹Wuhan University ²University of Cambridge

tanjieyi@whu.edu.cn yansheng.li@whu.edu.cn

Abstract

Traditional Remote Sensing Foundation models (RSFMs) are pre-trained with a data-centralized paradigm, through self-supervision on large-scale curated remote sensing data. For each institution, however, pre-training RSFMs with limited data in a standalone manner may lead to suboptimal performance, while aggregating remote sensing data from multiple institutions for centralized pre-training raises privacy concerns. Seeking for collaboration is a promising solution to resolve this dilemma, where multiple institutions can collaboratively train RSFMs without sharing private data. In this paper, we propose a novel privacy-preserved pre-training framework (**FedSense**), which enables multiple institutions to collaboratively train RSFMs without sharing private data. However, it is a non-trivial task hindered by a vicious cycle, which results from model drift by remote sensing data heterogeneity and high communication overhead. To break this vicious cycle, we introduce Federated Mutual-guidance Learning. Specifically, we propose a Server-to-Clients Guidance (SCG) mechanism to guide clients updates towards global-flatness optimal solutions. Additionally, we propose a Clients-to-Server Guidance (CSG) mechanism to inject local knowledge into the server by low-bit communication. Extensive experiments on four downstream tasks demonstrate the effectiveness of our FedSense in both full-precision and communication-reduced scenarios, showcasing remarkable communication efficiency and performance gains.

1. Introduction

Recently, Remote Sensing Foundation Models (RSFMs) have gained increasing attention due to their impressive applicability and performance across various Earth Observation tasks by producing general-purpose visual features [37, 48]. Traditionally, existing RSFMs follow a data-centralized training paradigm. They are built through

*Corresponding author.

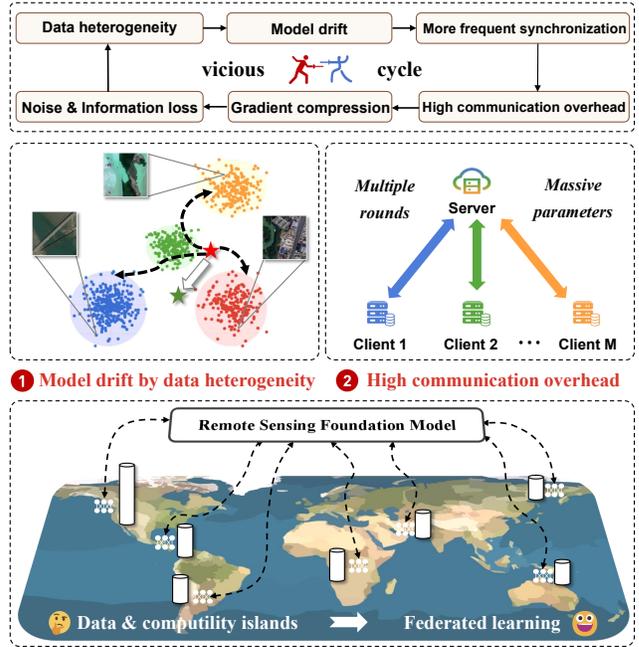


Figure 1. **Illustration of privacy-preserved pre-training of RSFMs with FL to bridge data islands.** The vicious cycle between data heterogeneity-induced model drift and communication bottlenecks reveals a critical performance-efficiency trade-off.

self-supervision on large-scale curated remote sensing data, gathered from diverse sources [10]. These data are collected by satellites, drones, or aerial platforms, and are stored in centralized archives by different institutions. A single institution could not pre-train RSFMs well due to the limited data scale and diversity [22]. At the same time, it is challenging to aggregate data from multiple institutions for data-centralized training due to geo-information security, storage bottlenecks, and industrial competitions [3, 14, 33, 38]. Such contradictory requirements urgently demand paradigm-shifting frameworks that reconcile computational synergies.

A more practical and realistic solution is to collabora-

tively inject remote sensing knowledge learned from private data owned by institutions into foundation models in a distributed manner [8, 13, 20, 34, 35]. In recent years, federated learning (FL) emerges as a promising privacy-preserving alternative, enabling collaborative model training without raw data exchange through periodic model aggregation [12, 25]. Self-supervised learning (SSL) operates on the principle of latent structure exploitation, where models learn transferable representations by solving pretext tasks derived from data’s intrinsic attributes without human annotations. However, it is a non-trivial task to train RSFMs by combining SSL with FL directly [9]. Recent studies in federated self-supervised learning (FSSL) mainly focus on natural images, showing potential for distributed visual representation learning. MocoSFL [18] focus cross-device SSL by leveraging momentum contrast on mobile devices. FedU² designs disperses local data representations uniformly using spherical Gaussian sampling and optimizes global-local model consistency. FedMKD [19] proposes a resource-adaptive FSSL to address the architecture heterogeneity and class skew issues. However, few studies have explored the **vicious cycle** challenge in FSSL, which manifests more severely in remote sensing [2, 15, 41].

To elaborate, two critical challenges are involved in this vicious cycle, as illustrated in Fig. 1. ❶ **model drift by data heterogeneity**. Remote sensing data is inherently heterogeneous due to diverse sensor types, resolutions, and geographic distribution. This heterogeneity leads to significant variability in data distributions across different institutions, causing convergence inefficiencies and degraded model performance [3]. ❷ **high communication overhead**. Foundation models are characterized by massive parameters (unlike neural networks with shallow architectures), leading to excessive communication costs and bandwidth demands [24]. The two challenges form a **vicious cycle**, where data heterogeneity necessitates more frequent model synchronization to mitigate drift, thereby amplifying communication costs. Conversely, communication compression can reduce overhead, introducing noise and information loss which further exacerbates client-side model drift. This bidirectional aggravation fundamentally undermines the efficiency and model consistency in distributed pre-training of RSFMs with FL.

In this paper, we propose a new challenging yet meaningful task: **privacy-preserved pre-training of RSFMs**. We propose FedSense, a novel FSSL framework with Federated Mutual-guidance Learning, to address the vicious cycle challenges. First, we introduce a Server-to-Client Mutual Guidance (SCG) mechanism to guide client models towards a global consensus, mitigating model drift by data heterogeneity. Second, we propose a Client-to-Server Guidance (CSG) mechanism to distill client models’ knowledge into a server-side reference model, reducing communica-

tion overhead. Our FedSense provides an integrated solution to resolve the vicious cycle challenges. As a result, our FedSense achieves state-of-the-art performance compared to existing FSSL methods with higher performance and lower communication overhead.

To sum up, this paper takes the first step towards privacy-preserved pre-training of RSFMs as far as we know. The main contributions of this work are as follows:

- We propose FedSense, establishing a new paradigm for privacy-preserved pre-training of RSFMs. To the best of our knowledge, it is the first generic FL framework that supports mainstream pre-training methods, including contrastive learning and masked image modeling.
- We resolve the vicious cycle challenges in FSSL by introducing Federated Mutual-guidance Learning, which is composed of Server-to-Client Mutual Guidance (SCG) mechanism and Client-to-Server Guidance (CSG) mechanism. The integrated solution effectively mitigates model drift by data heterogeneity and reduces communication overhead.
- We pre-train a RSFM with 10 participants with million-scale remote sensing data and evaluate the performance on four downstream tasks. Experimental results on eight datasets demonstrate that FedSense outperforms existing FSSL methods in terms of performance and communication overhead.

2. Related Work

2.1. Centralized Pre-training for RSFMs

Recent years have witnessed significant progress in data-centralized pre-training paradigms for RSFMs [21, 26]. A series of studies focus on constructing large-scale pre-training datasets and developing specialized algorithms tailored to remote sensing data characteristics. For instance, [9] introduced a rotation-variable window attention mechanism to handle large-scale and arbitrarily oriented geospatial objects, accompanied by MillionAID, a billion-parameter vision foundation model pre-trained on massive remote sensing imagery. To address dense object detection challenges in remote sensing, RingMo [36] optimized a masked autoencoder framework by incorporating multi-scale hierarchical representations and task-specific decoding strategies. For multispectral data with rich spectral information, SpectralGPT [11] treated multispectral images as 3D tensors and proposed a multi-objective reconstruction loss to jointly capture spatial-spectral correlations and spectral sequence dependencies. SkySense [10] further extended multimodal capabilities through spatiotemporal disentanglement and temporal-aware embedding mechanisms, enabling joint contrastive learning across high-resolution optical imagery, and time-series optical/SAR data. This unified framework supports diverse tasks ranging from image-

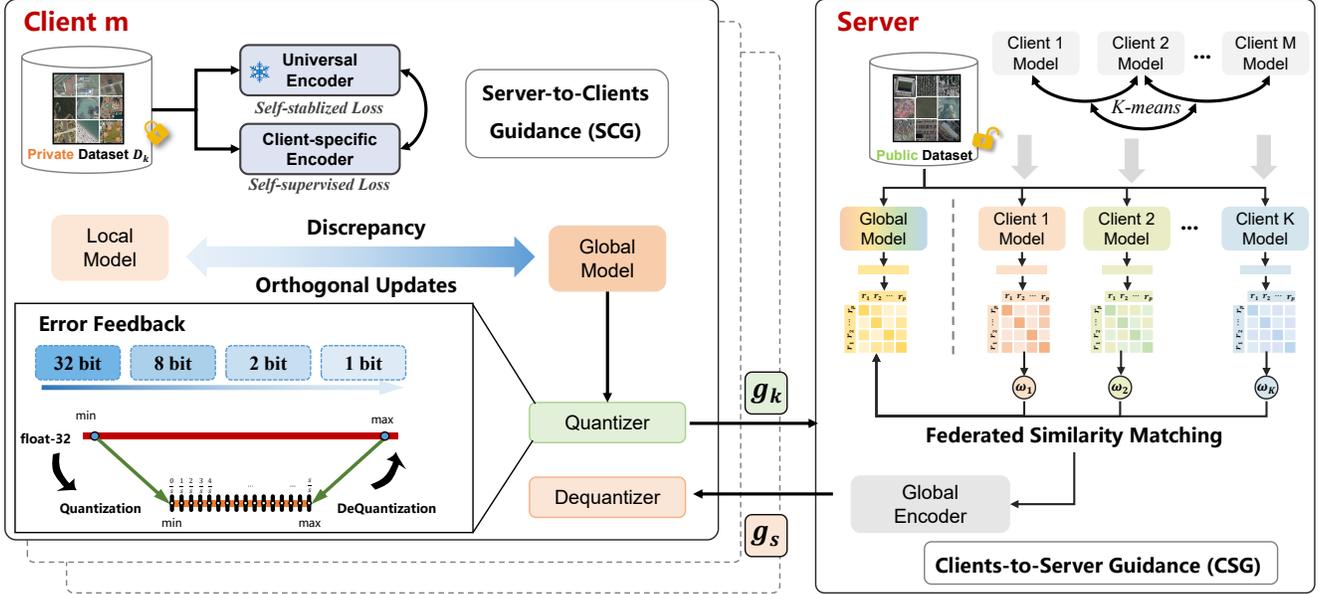


Figure 2. **Overview of FedSense.** The framework includes two components: Server-to-Clients Guidance (SCG) and Clients-to-Server Guidance (CSG). SCG guides clients’ updates towards global-flatness optimal solutions, while CSG injects local knowledge into the server by low-bit communication.

level classification to pixel-level segmentation and crop phenology monitoring. AnySat [1] is a versatile model designed to handle diverse data across resolutions, scales, and modalities.

Different from the above studies, our work takes the first step on privacy-preserved pre-training for RSFMs, and we aim to collaboratively pre-train RSFMs in a novel decentralized manner. This research is orthogonal to the existing centralized pre-training methods tailored to RSFMs, and would be complementary to scale up the performance of RSFMs in real-world applications.

2.2. Federated Self-Supervised Learning

Federated Learning (FL) [17, 29] has emerged as a promising paradigm for collaborative model training across decentralized data sources while preserving data privacy. Federated self-supervised learning (FSSL) [30, 46, 49] have demonstrated potential for privacy-preserving model training. However, existing FSSL methods predominantly address cross-device scenarios characterized by numerous resource-constrained clients (e.g., mobile devices) with homogeneous data distributions, focusing on computational efficiency and communication compression. While L-DAWA [31] mitigates data heterogeneity through hierarchical angular divergence weighting, and FedU² [23] and FedMKD [19] enhance representation consistency via unified embedding alignment and adaptive knowledge distillation, these methods remain inadequate for cross-institution pre-training of RSFMs.

Existing FSSL focuses on optimizing federated model convergence and mainly support cross-device scenarios [16, 27, 39]. Additionally, these methods are with limited applicability to support self-supervised learning frameworks (contrastive learning and masked image modeling) consistently. Thirdly, these methods ignore the vicious cycle in collaboratively pre-training RSFMs, which requires handling the drift challenges and communication efficiency simultaneously. Our work sheds lights on these limitations and proposes Federated Mutual Learning to unleash the power of FSSL for RSFM pre-training.

3. Methodology

3.1. Problem definition

An FL framework for privacy-preserving collaborative pre-training of RSFMs aims to learn general-purpose visual representations from distributed unlabeled remote sensing images. The system comprises M institutions (*i.e.* clients) and a central server. Let the m -th **client** c_m possess a private unlabeled dataset \mathcal{D}_{c_m} , and **server** s maintain a public unlabeled dataset \mathcal{D}_s . For each client c_m , the local self-supervised objective function \mathcal{F}_{c_m} with parameters θ_m can be formulated as:

$$\min_{\theta_m} \mathcal{F}_{c_m} = \mathbb{E}_{x \sim \mathcal{D}_{c_m}} \left[\psi(f_{\theta_m}(x^{(a)}), f_{\theta_m}(x^{(b)})) \right], \quad (1)$$

where $x^{(a)}$ and $x^{(b)}$ are view augmentations or masked modeling of input image x , and $\psi(\cdot, \cdot)$ denotes the self-

supervised loss function. For contrastive learning, ψ measures feature consistency between augmented views. For masked modeling, ψ computes the reconstruction error between predicted and original pixel values in masked regions. In the t -th round, upon receiving $\{\theta_m^{(t)}\}_{m=1}^M$, server s with public data \mathcal{D}_s applies aggregation $\Theta^{(t)} = \sum_{m=1}^M p_m \theta_m^{(t)}$, where $p_m = \frac{|\mathcal{D}_{c_m}|}{\sum_{m=1}^M |\mathcal{D}_{c_m}|}$. The optimization objective is:

$$\min_{\Theta} \mathcal{F}_s = \mathbb{E}_{x^s \sim \mathcal{D}_s} \left[\phi(\Theta, \{\theta_m^{(t)}\}_{m=1}^M, \dots, x^s) \right], \quad (2)$$

where x_s is unlabeled remote sensing data sampled from the public dataset, and ϕ denotes the loss function further to inject public data knowledge into the global foundation model.

3.2. Overview of our FedSense

As shown in Fig. 2, our FedSense consists two parts, which are Server-to-Clients Guidance (SCG) and Clients-to-Server Guidance (CSG). The two components are introduced sequentially. We propose a federated self-supervised learning framework for collaborative training of foundation models across multiple institutions. Note that we omit the conversion between gradients and parameters during transmission for conciseness. The algorithm is detailed in Algorithm 1.

3.3. Server-to-Clients Guidance

The core insight of SCG is to strike a balance between global knowledge preservation and local model optimization. The server guides clients to optimize their local models while restrict the discrepancy of local and global model. We design a dual loss: self-supervised loss and self-stabilized loss in Eq. (7), which are complementary to each other. The former requires models to learn orthogonality to the discrepancy of local and global model, while the latter helps to regularize the model with stabilized knowledge information in federated updates. Thus, the overall training target of the m -th client is:

$$\mathcal{L}_m^{\text{total}} = \underbrace{\mathcal{L}_m^{\text{ssl}}(\theta_m; \Theta; \mathcal{R}_C(\theta_m))}_{\text{self-supervised term}} + \underbrace{\mathcal{L}_m^{\text{sst}}(\theta_m; \theta_{\text{uni}})}_{\text{self-stabilized term}}, \quad (3)$$

where $\mathcal{L}_m^{\text{ssl}}$ denotes self-supervised loss. It corresponds to SSL objectives commonly used in the field, such as contrastive loss, and masked reconstruction loss. Note that this is a drift-aware loss incorporated with our proposed optimizing method. We conclude it as the following minimax optimization problem:

$$\min_{\theta_m} \max_{|\epsilon|_2 < \rho} \mathcal{L}_m^{\text{disc}}(\theta_m + \epsilon, f_{\theta_m}(x^{(a)}), f_{\theta_m}(x^{(b)})), \quad (4)$$

The objective of the weight perturbation on the client-side is to find the ϵ , which causes the maximum increase in the

parameter discrepancy. Given local model θ_m and global model Θ , the parameter discrepancy is calculated by:

$$\nabla_{\theta_m} \mathcal{L}_m^{\text{disc}}(\theta_m; \Theta) = \nabla_{\theta_m} (\beta(\theta_m - \Theta)), \quad (5)$$

where β is a hyperparameter to of the discrepancy term. Then, the optimal perturbation ϵ is approximated by:

$$\tilde{\epsilon} = \operatorname{argmax}_{\epsilon} \mathcal{L}_m^{\text{disc}}(\theta_m + \epsilon; \Theta) \approx \lambda \frac{\nabla_{\theta_m} \mathcal{L}_m^{\text{disc}}(\theta_m; \Theta)}{\|\nabla_{\theta_m} \mathcal{L}_m^{\text{disc}}(\theta_m; \Theta)\|_2}, \quad (6)$$

where λ is a scaling factor. The optimal perturbation ϵ will be later used to update the local model θ_m .

Moreover, the self-stabilized loss $\mathcal{L}_m^{\text{sst}}$ is designed to leverage the knowledge from a universal encoder, $f_{\theta_{\text{uni}}}$, which is pre-trained on a large-scale dataset (e.g., ImageNet-22k) and broadcast to all clients by the server at the beginning of the federated learning process. It guides the client-specific encoder to output representations that align with the universal encoder. Formally, the self-stabilized loss can be defined as:

$$\mathcal{L}_m^{\text{sst}} = \mathbb{E}_{x \sim \mathcal{D}_m} \left[-\frac{f_{\theta_m}(x)}{\|f_{\theta_m}(x)\|_2} \cdot \frac{f_{\theta_{\text{uni}}}(x)}{\|f_{\theta_{\text{uni}}}(x)\|_2} \right], \quad (7)$$

where \mathcal{D}_m denotes the local dataset of client m . $f_{\theta_{\text{uni}}}$ and f_{θ_m} are the feature extractors of the universal encoder and the client-specific encoder, respectively. Lastly, the optimal perturbation ϵ is then used to update the local model θ_m :

$$\theta_m \leftarrow \theta_m - \gamma \nabla_{\theta_m} \mathcal{L}_m^{\text{total}}(\theta_m + \epsilon; \Theta), \quad (8)$$

where γ is the learning rate. In this way, the server guides the clients to optimize their local models while preserving the global knowledge with orthogonal property and stabilized information.

3.4. Clients-to-Server Guidance

As the model size and number of clients increases, the communication overhead becomes a more serious bottleneck. The clients need to transmit the updates to the server in each communication round, and waiting for the server for aggregation and send back the updated parameters.

For the uplink communication is the most time-consuming part in this process, we propose a CSG mechanism to reduce the communication cost. Though quantization is a widely used technique to reduce the communication cost. However, it inherently introduces quantization error and information loss, which may accumulate over time and degrade the model performance. To guide the server to aggregate the quantized updates from clients, we propose a feedback error mechanism to compensate the quantization error.

Assume the feedback error at the t -th round is e_m^t , it is added to compensate updates of each client:

$$\mathcal{G}_m^{t+1} = \Delta \theta_m^t + e_m^t, \quad (9)$$

Algorithm 1: Our FedSense

Input : T : communication round; M : number of clients; E : local epochs;

Output: Weight Θ^T of the RSFM at the T -th round.

```
1 Server-side:
2 Initialize global model  $\Theta^0$ 
3 for  $t = 1$  to  $T$  do // the  $t$ -th round
4   Server broadcast  $\Theta^{t-1}$  to selected clients
5   for each client  $c_m \in \{c_m\}_{m=1}^M$  in parallel do
6     Server-to-Clients Guidance (SCG)
7   end
8    $\{\Delta\theta_m^t\}_{m=1}^M = \text{DCPR}\{\{\Delta\tilde{\theta}_m^t\}_{m=1}^M; b\}$ 
9    $\Theta^t \leftarrow \text{ServerAggregation}(\{\Delta\theta_m^t\}_{m=1}^M)$ 
10  Clients-to-Server Guidance (CSG)
11  Similarity alignment with public data Eq. (18)
12 end
13 Client-side:
14 Initialize local model  $\theta_m^{t-1} \leftarrow \Theta^{t-1}$ 
15 for  $e = 1$  to  $E$  do // the  $e$ -th epoch
16   Server-to-Clients Guidance (SCG)
17   Optimization with knowledge preservation
18    $\theta_m \leftarrow \theta_m - \gamma \nabla_{\theta_m} \mathcal{L}_m^{\text{total}}(\theta_m + \epsilon; \Theta)$  Eq. (8)
19 end
20  $\Delta\tilde{\theta}_m^t = \text{CPR}\{(\theta_m^t - \theta_m^{t-1}); b\}$ 
21 Send the updates  $\Delta\tilde{\theta}_m^t$  back to the server.
22 Downstream tasks: Institutions utilize the
    pre-trained model as a backbone, fine-tuning it on
    labeled data for specific tasks.
```

where $\Delta\theta_m^t$ is the updates of client m at round t . The quantized updates are computed as follows:

$$\tilde{\mathcal{G}}_m^{t+1}[i] = \underbrace{\|\mathcal{G}_m^{t+1}\|}_{\text{L2 norm of raw updates}} \cdot \underbrace{\text{sgn}(\mathcal{G}_m^{t+1}[i])}_{\text{sign of element (1 bit)}} \cdot \underbrace{\xi(\mathcal{G}_m^{t+1}[i]; s)}_{\text{unbiased stochastic function}(b-1 \text{ bits})}, \quad (10)$$

where $\xi(\cdot)$ is a unbiased stochastic function mapping $|\mathcal{G}_m^{t+1}[i]|/|\mathcal{G}_m^{t+1}|$ to the quantization level s .

$$e_m^t = \mathcal{G}_m^t - \text{DCPR}(\tilde{\mathcal{G}}_m^t; b), \quad (11)$$

where DCPR is the dequantization function. The feedback error is updated by:

$$e_m^t = \alpha \cdot e_m^{t-1} + (1 - \alpha) \cdot e_m^t, \quad (12)$$

where α is the momentum factor, and b is the bit-width of quantization. The feedback error is accumulated over time to preserve the gradient information across communication rounds. Inspired by the dynamic optimization characteristics of neural network training [44], we implement periodic feedback error resetting to address the fast-evolving loss

landscape. We reset the feedback error to zero at the frequency of T_{reset} rounds to ensure the feedback error remains aligned with the current optimization state.

On the server side, we propose federated similarity distillation to provide public remote sensing data guidance. Our FedSense consists of three core steps: server-side aggregation, local model clustering, and cross-model knowledge distillation. In this way, clients can leverage the public data to enhance the global model performance. The server aggregates client models $\{\theta_m\}_{m=1}^M$ using data volume weights. They are clustered into K groups via K-means to accelerate multi-model forward pass cost and mitigate the bias of some models:

$$\{\theta^{(k)}\}_{k=1}^K = \text{K-means}(\{\theta_m\}_{m=1}^M). \quad (13)$$

For each public batch $\mathcal{B} = \{x_i\}_{i=1}^p$, the group of models produce feature matrices:

$$\mathbf{Z}^{(k)} = [z_1^{(k)}, \dots, z_p^{(k)}]^\top \in \mathbb{R}^{p \times d}, \quad (14)$$

$$\mathbf{Z}^g = [z_1^g, \dots, z_p^g]^\top \in \mathbb{R}^{p \times d}, \quad (15)$$

where $z_i^{(k)} = f_{\theta^{(k)}}(x_i)$ and $z_i^g = f_{\theta^g}(x_i)$, and p and d are the batch size and feature dimension, respectively. The similarity matrices are computed as:

$$\mathbf{S}^{(k)} = \frac{\mathbf{Z}^{(k)}(\mathbf{Z}^{(k)})^\top}{\|\mathbf{Z}^{(k)}\|_F}, \quad \mathbf{S}^g = \frac{\mathbf{Z}^g(\mathbf{Z}^g)^\top}{\|\mathbf{Z}^g\|_F}. \quad (16)$$

The weighted consensus similarity combines local expertise:

$$\mathbf{S}_{\text{consensus}} = \sum_{k=1}^K \omega_k \mathbf{S}^{(k)}, \quad \omega_k = \frac{|\mathcal{D}_k|}{\sum_{i=1}^N |\mathcal{D}_i|}. \quad (17)$$

The global model is optimized by matching similarity distributions:

$$\mathcal{L}_{\text{distill}} = \frac{1}{p^2} \|\mathbf{S}^g - \mathbf{S}_{\text{consensus}}\|_F^2. \quad (18)$$

By distilling the similarity knowledge from multiple models, the global model can learn to capture the intrinsic structure of the public data, which is beneficial for enhancing the generalization ability of the global model.

4. Experiments

4.1. Federated Experimental Setups

Distributed Pre-training datasets. A distributed dataset was constructed for federated pre-training of RSFMs, comprising 10 clients with heterogeneous private remote sensing data and supplementary public datasets maintained by a server (Fig. 3). Notably, the dataset includes satellite images (e.g. WorldView-2 and JL-1) and aerial images (e.g.

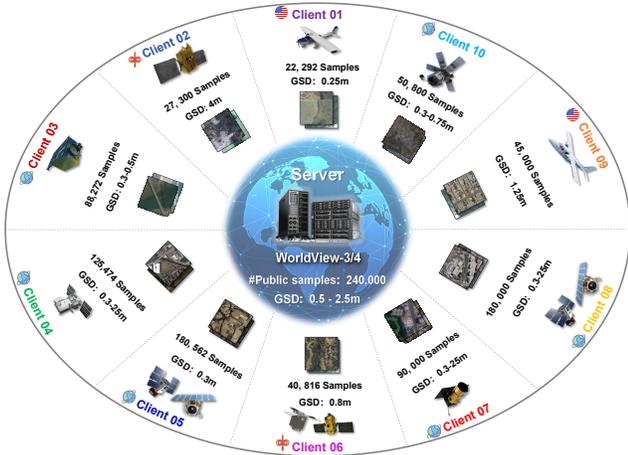


Figure 3. **Details of federated pre-training datasets.** The dataset consists of 10 clients with million-scale heterogeneous private remote sensing data and public datasets maintained by a server.

NAIP and NOAA), while some institutions possess multi-sourced collections, featuring heterogeneous sensor configurations and spatial resolutions ranging from 0.25m to 25m. The coverage spans diverse geographical regions including global, USA, and China. Such heterogeneity across data volume, geographic distribution, resolution variance, and platform diversity establishes a representative million-scale dataset that simulates real-world scenarios.

Downstream tasks. We conducted experiments on four typical downstream tasks of remote sensing image analysis to validate the effectiveness of the RSFM pretrained by our FedSense. These tasks include: scene classification (RESISC-45 [6], AID [42]); semantic segmentation (Inria [28], LoveDA [40]); object detection (DIOR-R [7], DOTA-v1.0 [43]); and change detection (LEVIR-CD+ [5], SECOND [47]). More details of the datasets are shown in Sec. C of the supplementary.

Evaluation metrics. For scene classification, we employ overall accuracy (OA) as the evaluation metric. We use mean intersection over union (mIoU) for multi-class semantic segmentation tasks and IoU for binary segmentation tasks. For object detection, we use mean Average Precision (mAP) and F1 score for change detection tasks. For semantic change detection, we use semantic change segmentation score (SCS) as the evaluation metric.

Implementation details. To systematically evaluate the universality of our framework, we establish a comprehensive evaluation framework incorporating two mainstream SSL paradigms: contrastive learning through DINO [4], and masked image modeling via SimMIM [45]. All experiments employ the tiny version of Swin Transformer (Swin-T) as the foundational backbone unless explicitly stated, with consistent training protocols (100-round pre-training,

AdamW optimizer, $1e-4$ base learning rate).

4.2. Comparison with State-of-the-Art Methods

In this section, we present a comparative analysis of performance on four downstream tasks using RSFM pre-trained by our FedSense and other state-of-the-art FSSL. The results are summarized in Tab. 1. We observe that our FedSense outperforms existing methods across almost all tasks, achieving an average improvement of 1.3% OA on scene classification, 1.1% mIoU on semantic segmentation, 0.9% mAP on object detection, and 0.8% F1 on change detection. Note that due to not using public dataset in the original paper, we conduct experiments by averagely integrate public datasets to the federated pre-training dataset for a fair comparison with existing methods with the same number of total samples. To elaborate, we provide detailed comparisons with existing methods.

Our experimental results reveal several key observations. The randomly initialized model yields inferior performance across all downstream tasks, demonstrating Transformers’ inherent limitation in learning effective representations from limited labeled data. ImageNet supervised pre-training brings significant improvements (*e.g.*, +14.36% on RESISC-45), validating the importance of pre-training for vision transformers. However, the domain gap between natural images and top-down remote sensing views limits further performance gains, motivating our FL approach.

Among FSSL methods, most approaches (FedEMA, FedMKD) achieve moderate improvements over ImageNet pretraining, confirming that collaborative pretraining with distributed data can inject domain-specific knowledge. Notably, FedU² exhibits performance degradation compared to ImageNet pretraining, suggesting its vulnerability to data heterogeneity and insufficient utilization of pre-trained knowledge. Our FedSense consistently outperforms SSL-FL, achieving absolute gains of 1.12% and 1.37% on RESISC-45 and AID for scene classification. When using DINO framework, our FedSense maintains advantages with 0.87% and 0.26% improvements, respectively.

For semantic segmentation, FedSense achieves 0.85% OA improvement on Inria and 1.07% gain on LoveDA over the best competitor. The marginal differences under DINO framework (0.12%-0.35%) suggest contrastive learning brings limited benefits for segmentation tasks that require precise pixel-level localization. In rotated object detection, FedSense with SimMIM framework surpasses SOTA by 1.50% and 0.51% on DIOR-R and DOTA-v1.0, while maintaining competitive performance (+0.11% mAP) with DINO-based approaches.

Additionally, the 1-bit quantized version of FedSense demonstrates remarkable communication efficiency with minimal performance degradation (0.23% across tasks). Quantized FedSense slightly outperforms full-precision

Framework	Method	RESISC-45	AID	Inira		LoveDA	DIOR-R	DOTA-v1.0	LEVIR-CD+			SECOND
		OA	OA	IoU	OA	mIoU	mAP	mAP	Precision	Recall	F1	SCS
Swin-Tiny	Random Init.	80.09	72.84	79.41	90.21	42.83	46.86	64.45	69.71	63.39	66.40	30.95
	ImageNet Sup.	94.45	96.36	80.65	93.16	51.52	64.37	76.90	73.26	71.44	72.34	34.79
SimMIM [45] (Swin-Tiny)	SSL-FL [46]	95.21	96.17	<u>81.33</u>	93.43	51.67	64.01	<u>76.82</u>	<u>74.31</u>	70.99	72.61	33.98
	FedEMA [49]	-	-	-	-	-	-	-	-	-	-	-
	L-DAWA [31]	-	-	-	-	-	-	-	-	-	-	-
	FedU ² [23]	-	-	-	-	-	-	-	-	-	-	-
	FedMKD [19]	-	-	-	-	-	-	-	-	-	-	-
	Our FedSense	96.33	97.54	81.66	94.28	52.74	65.51	77.33	74.82	71.67	73.21	35.43
	<u>96.01</u>	<u>96.68</u>	<u>80.97</u>	<u>93.45</u>	<u>51.91</u>	<u>64.53</u>	<u>76.76</u>	<u>73.98</u>	<u>71.59</u>	<u>72.77</u>	<u>34.87</u>	
DINO [4] (Swin-Tiny)	SSL-FL [46]	-	-	-	-	-	-	-	-	-	-	-
	FedEMA [49]	94.91	97.04	80.21	92.98	51.89	64.93	77.09	73.63	71.20	72.39	34.82
	L-DAWA [31]	94.23	96.27	80.82	93.33	51.68	65.08	77.24	<u>74.18</u>	71.82	72.98	34.97
	FedU ² [23]	93.85	95.76	80.13	92.59	51.34	64.03	75.88	73.30	70.67	71.96	33.27
	FedMKD [19]	95.34	97.17	<u>80.84</u>	<u>94.14</u>	51.99	65.44	<u>77.56</u>	74.36	73.05	<u>73.70</u>	35.11
	Our FedSense	96.21	97.43	81.96	94.23	<u>51.95</u>	<u>65.40</u>	77.64	74.98	73.13	74.04	<u>35.34</u>
	<u>95.88</u>	<u>97.23</u>	<u>80.97</u>	<u>93.88</u>	<u>51.14</u>	<u>64.62</u>	<u>77.45</u>	<u>73.88</u>	<u>73.09</u>	<u>73.48</u>	35.82	

Table 1. **Comparison results of our FedSense and previous SOTA methods.** The symbol (-) indicates unavailable results where methods are incompatible with specific SSL framework types. Results highlighted in **yellow** denote full-precision communication performance of our FedSense, while **blue** shading represents experiments with 1-bit communication-quantized transmission, demonstrating our method’s efficiency-accuracy trade-off. The best results are highlighted in **bold**, and the second-best results are underlined.

baselines on SECOND dataset (+0.06% SCS), which we conjecture stems from the error-compensated quantization acting as implicit regularization. This makes FedSense particularly suitable for bandwidth-constrained applications.

These results collectively demonstrate that our Federated Mutual-guidance Learning framework effectively coordinates distributed clients to learn transferable representations while maintaining communication efficiency. The consistent improvements across diverse tasks validate FedSense’s ability to capture domain-specific patterns from unlabeled remote sensing data through federated collaboration.

4.3. Ablation Studies

In this part, we conduct ablation studies to analyze the system scalability, model scalability, effectiveness of components, and parameter analysis of our proposed FedSense.

System Scalability Analysis. To assess our framework’s adaptability to real-world distributed scenarios, we systematically investigate how model performance scales with increasing participants and training samples. As shown in Tab. 2, expanding from 2 to 10 collaborative clients (80K to 800K samples) yields consistent performance gains across all tasks. For RESISC-45 and DIOR-R, the performance improvements scale nearly linearly with client/sample quantities, suggesting these tasks particularly benefit from diverse perspectives in FedSense. However, LoveDA exhibits marginal gains despite 10× sample growth, indicating pixel-level tasks require more sophisticated feature aggregation beyond simple data scaling. These findings confirm our framework’s effectiveness in harnessing distributed resources.

#CL	#TS	RESISC-45	LoveDA	DIOR-R	LEVIR-CD+
		OA	mIoU	mAP	F1
2	80K	94.50	51.56	64.43	72.44
4	200K	95.23	51.86	64.75	72.68
8	650K	95.85	52.33	65.16	73.00
10	800K	96.33	52.74	65.51	73.21

Table 2. **Number of participants impact analysis.** #CL means the number of participants. #TS means number of total participating samples.

Model Scalability Analysis. The impact of model capacity is systematically evaluated through progressively larger Swin Transformer variants, as detailed in Tab. 3. While all tasks benefit from increased model parameters, we observe distinct scaling patterns across task types. RESISC-45 shows diminishing returns, improving only +1.21% from Swin-T to Swin-L, suggesting vision transformers approach saturation points for scene classification. Conversely, segmentation (LoveDA mIoU +3.22%) and detection (DIOR-R mAP +3.31%) exhibit near-linear improvements with model growth, indicating complex localization tasks inherently demand higher-capacity architectures. Notably, change detection (LEVIR-CD+ F1 +2.92%) demonstrates sustained sensitivity to model size, likely requiring deeper feature hierarchies to discern subtle temporal changes. The 197M-parameter Swin-L achieves marginal gains (+0.42% mAP over Swin-B) compared to its 2.2× pa-

parameter increase, highlighting practical trade-offs between model capacity and computational costs.

Model	#Para.	RESISC-45	LoveDA	DIOR-R	LEVIR-CD+
		OA	mIoU	mAP	F1
Swin-T	28M	96.33	52.74	65.51	73.21
Swin-S	50M	96.80	54.12	67.03	74.65
Swin-B	88M	97.15	55.29	68.40	75.83
Swin-L	197M	97.54	55.96	68.82	76.13

Table 3. **Model size impact analysis.** # Para. means the number of model parameters.

Quantization Methods Comparison. We compare our proposed quantization method with uniform quantization and FedPAQ [32] on the RESISC-45 dataset. As shown in Tab. 4, our method outperforms uniform quantization and FedPAQ across all quantization bit-widths. Specifically, our method achieves 96.27% OA with 8-bit quantization, 96.13% OA with 2-bit quantization, and 96.01% OA with 1-bit quantization. Our FedSense outperforms FedPAQ by 0.14% OA with 8-bit quantization, 0.16% OA with 2-bit quantization, and 0.22% OA with 1-bit quantization. The results demonstrate the effectiveness of our proposed quantization method in enhancing the communication efficiency.

Method	32-bit	8-bit	2-bit	1-bit
Uniform Quant	96.33	96.01	95.87	95.56
FedPAQ [32]	96.33	96.13	95.97	95.79
Ours	96.33	96.27	96.13	96.01

Table 4. **Quantization bit-width comparison.**

SST	SCG	CSG	RESISC-45	LoveDA	DIOR-R	LEVIR-CD+
✓			94.52	51.87	64.83	72.43
✓	✓		94.70	51.85	64.91	72.55
✓		✓	95.25	51.97	65.12	72.78
✓	✓	✓	96.33	52.74	65.51	73.21

Table 5. **Effectiveness of proposed components.**

Effectiveness of Components. Our ablation study quantitatively validates the complementary nature of proposed components, as summarized in Tab. 5. The standalone SST mechanism achieves a 94.52% OA on RESISC-45, 51.87% mIoU on LoveDA, 64.83% mAP on DIOR-R, and 72.43% F1 on LEVIR-CD+. The SCG mechanism further boosts performance across tasks, with 0.18% OA, 0.02%

mIoU, 0.08% mAP, and 0.12% F1 improvements. The CSG mechanism contributes to a 0.83% OA gain on RESISC-45, 0.29% mAP on DIOR-R, and 0.35% F1 on LEVIR-CD+. The full model with all components achieves the best performance, demonstrating the effectiveness of our proposed components in enhancing federated learning for remote sensing tasks.

LE	Round	RESISC-45	LoveDA	DIOR-R	LEVIR-CD+
		OA	mIoU	mAP	F1
1	100	96.33	52.74	65.51	73.21
2	50	95.82	51.93	64.08	72.05
4	25	95.17	50.62	62.75	70.89
100	1 [†]	94.50	48.31	60.13	68.24

Table 6. **Local epochs (LE) impact analysis.** † means one-shot FL setting.

Parameter Analysis. The trade-off between local computation and communication frequency is systematically investigated through varying local epochs (LE), as shown in Tab. 6. Here more local epochs indicate more discrepancy between local and global models, leading to potential client drift. It is worth noting that the one-shot setting (LE=100) suffers catastrophic performance collapse, confirming remote sensing data’s inherent heterogeneity demands periodic model synchronization. However, reducing communication rounds and developing one-shot FL is becoming widely adopted for only one round of communication. We expect that our proposed FedSense can be further improved by incorporating more advanced techniques to handle the one-shot setting.

5. Conclusion & Future Work

This paper takes the first step towards developing a privacy-preserved pre-training framework (**FedSense**) for RSFMs. FedSense enables multiple institutions to collaboratively train RSFMs without sharing private data. We introduce Federated Mutual-guidance Learning, which breaks the vicious cycle caused by remote sensing data heterogeneity and high communication overhead. Specifically, we propose a SCG mechanism to guide clients updates towards global-flatness optimal solutions. Additionally, we propose a CSG mechanism to inject local knowledge into the server by low-bit communication. Extensive experiments on four downstream tasks demonstrate the effectiveness of our FedSense in both full-precision and communication-reduced scenarios, showcasing remarkable communication efficiency and performance gains. In the future, we plan to extend the current framework to support the collaborative pre-training of multi-modal RSFMs with modality heterogeneity.

References

- [1] Guillaume Astruc, Nicolas Gonthier, Clement Mallet, and Loic Landrieu. Anysat: An earth observation model for any resolutions, scales, and modalities, 2024. 3
- [2] Cosmin I Bercea, Benedikt Wiestler, Daniel Rueckert, and Shadi Albarqouni. Federated disentangled representation learning for unsupervised brain anomaly detection. *Nature Machine Intelligence*, 4(8):685–695, 2022. 2
- [3] Baris Büyüktas, Gencer Sumbul, and Begüm Demir. Federated learning across decentralized and unshared archives for remote sensing image classification: A review. *IEEE Geoscience and Remote Sensing Magazine*, pages 2–18, 2024. 1, 2
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 6, 7
- [5] Hao Chen and Zhenwei Shi. A spatial-temporal attention-based method and a new dataset for remote sensing image change detection. *Remote Sensing*, 12(10), 2020. 6, 3
- [6] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017. 6, 2
- [7] Gong Cheng, Jiabao Wang, Ke Li, Xingxing Xie, Chunbo Lang, Yanqing Yao, and Junwei Han. Anchor-free oriented proposal generator for object detection. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–11, 2022. 6, 3
- [8] Hongquan Cheng, Jie Zheng, Huayi Wu, Kunlun Qi, and Lihua He. A communication-efficient distributed deep learning remote sensing image change detection framework. *International Journal of Applied Earth Observation and Geoinformation*, 129:103840, 2024. 2
- [9] Jie Gui, Tuo Chen, Jing Zhang, Qiong Cao, Zhenan Sun, Hao Luo, and Dacheng Tao. A survey on self-supervised learning: Algorithms, applications, and future trends. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9052–9071, 2024. 2
- [10] Xin Guo, Jiangwei Lao, Bo Dang, Yingying Zhang, Lei Yu, Lixiang Ru, Liheng Zhong, Ziyuan Huang, Kang Wu, Dingxiang Hu, Huimei He, Jian Wang, Jingdong Chen, Ming Yang, Yongjun Zhang, and Yansheng Li. Skysense: A multi-modal remote sensing foundation model towards universal interpretation for earth observation imagery. In *CVPR*, pages 27672–27683, 2024. 1, 2
- [11] Danfeng Hong, Bing Zhang, Xuyang Li, Yuxuan Li, Chenyu Li, Jing Yao, Naoto Yokoya, Hao Li, Pedram Ghamisi, Xiuping Jia, Antonio Plaza, Paolo Gamba, Jon Atli Benediktsson, and Jocelyn Chanussot. Spectralgpt: Spectral remote sensing foundation model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(8):5227–5244, 2024. 2
- [12] Wenke Huang, Mang Ye, Zekun Shi, Guancheng Wan, He Li, Bo Du, and Qiang Yang. Federated learning for generalization, robustness, fairness: A survey and benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9387–9406, 2024. 2
- [13] Alex Iacob, Lorenzo Sani, Bill Marino, Preslav Aleksandrov, and Nicholas Donald Lane. Worldwide federated training of language models. *arXiv preprint arXiv:2405.14446*, 2024. 2
- [14] Sami Jaghouar, Jack Min Ong, Manveer Basra, Fares Obeid, Jannik Straube, Michael Keiblinger, Elie Bakouch, Lucas Atkins, Maziyar Panahi, et al. Intellect-1 technical report. *arXiv preprint arXiv:2412.01152*, 2024. 1
- [15] Shusen Jing, Anlan Yu, Shuai Zhang, and Songyang Zhang. FedSC: Provable federated self-supervised learning with spectral contrastive objective over non-i.i.d. data. In *ICML*, 2024. 2
- [16] Hansol Kim, Youngjun Kwak, Minyoung Jung, Jinho Shin, Youngsung Kim, and Changick Kim. Protofl: Unsupervised federated learning via prototypical distillation. In *ICCV*, pages 6470–6479, 2023. 3
- [17] Gihun Lee, Minchan Jeong, Sangmook Kim, Jaehoon Oh, and Se-Young Yun. Fedsol: Stabilized orthogonal learning with proximal restrictions in federated learning. In *CVPR*, pages 12512–12522, 2024. 3
- [18] Jingtao Li, Lingjuan Lyu, Daisuke Iso, Chaitali Chakrabarti, and Michael Spranger. MocoSFL: enabling cross-client collaborative self-supervised learning. In *ICLR*, 2023. 2
- [19] Mingyi Li, Xiao Zhang, Qi Wang, Tengfei Liu, Ruofan Wu, Weiqiang Wang, Fuzhen Zhuang, Hui Xiong, and Dongxiao Yu. Resource-aware federated self-supervised learning with global class representations. In *NeurIPS*, 2024. 2, 3, 7
- [20] Shenghui Li, Fanghua Ye, Meng Fang, Jiayu Zhao, Yun-Hin Chan, Edith C-H Ngai, and Thiemo Voigt. Synergizing foundation models and federated learning: A survey. *arXiv preprint arXiv:2406.12844*, 2024. 2
- [21] Xuyang Li, Danfeng Hong, and Jocelyn Chanussot. S2mae: A spatial-spectral pretraining foundation model for spectral remote sensing data. In *CVPR*, pages 24088–24097, 2024. 2
- [22] Yansheng Li, Jieyi Tan, Bo Dang, Mang Ye, Sergey A. Bartalev, Stanislav Shinkarenko, Linlin Wang, Yingying Zhang, Lixiang Ru, Xin Guo, Liangqi Yuan, Lei Yu, Jingdong Chen, Ming Yang, José Marcato Junior, and Yongjun Zhang. Unleashing the potential of remote sensing foundation models via bridging data and computility islands. *The Innovation*, page 100841, 2025. 1
- [23] Xinting Liao, Weiming Liu, Chaochao Chen, Pengyang Zhou, Fengyuan Yu, Huabin Zhu, Binhui Yao, Tao Wang, Xiaolin Zheng, and Yanchao Tan. Rethinking the representation in federated unsupervised learning with non-iid data. In *CVPR*, pages 22841–22850, 2024. 3, 7
- [24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021. 2
- [25] Senliang Lu, Yehang Chen, Yuan Chen, Peijun Li, Junqi Sun, Changye Zheng, Yujian Zou, Bo Liang, Mingwei Li, Qinggeng Jin, Enming Cui, Wansheng Long, and Bao Feng. General lightweight framework for vision foundation model supporting multi-task and multi-center medical image analysis. *Nature Communications*, 16(1):2097, 2025. 2
- [26] Siqi Lu, Junlin Guo, James R Zimmer-Dauphinee, Jordan M Nieuwsma, Xiao Wang, Parker VanValkenburgh, Steven A

- Wernke, and Yuankai Huo. Vision foundation models in remote sensing: A survey. *arXiv preprint arXiv:2408.03464*, 2025. 2
- [27] Ekdeep Lubana, Chi Ian Tang, Fahim Kawsar, Robert Dick, and Akhil Mathur. Orchestra: Unsupervised federated learning via globally consistent clustering. In *ICML*, pages 14461–14484, 2022. 3
- [28] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat, and Pierre Alliez. Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark. In *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2017. 6, 3
- [29] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, pages 1273–1282, 2017. 3
- [30] Yasar Abbas Ur Rehman, Yan Gao, Jiajun Shen, Pedro Porto Buarque de Gusmão, and Nicholas Lane. Federated self-supervised learning for video understanding. In *ECCV*, pages 506–522, 2022. 3
- [31] Yasar Abbas Ur Rehman, Yan Gao, Pedro Porto Buarque de Gusmao, Mina Alibeigi, Jiajun Shen, and Nicholas D. Lane. L-dawa: Layer-wise divergence aware weight aggregation in federated self-supervised visual representation learning. In *ICCV*, pages 16464–16473, 2023. 3, 7
- [32] Amirhossein Reiszadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, pages 2021–2031, 2020. 8
- [33] Lorenzo Sani, Alex Jacob, Zeyu Cao, Royson Lee, Bill Marino, Yan Gao, Dongqi Cai, Zexi Li, Wanru Zhao, Xinchu Qiu, and Nicholas D. Lane. Photon: Federated llm pre-training. *arXiv preprint arXiv:2411.02908*, 2024. 1
- [34] Lorenzo Sani, Alex Jacob, Zeyu Cao, Bill Marino, Yan Gao, Tomas Paulik, Wanru Zhao, William F Shen, Preslav Aleksandrov, Xinchu Qiu, et al. The future of large language model pre-training is federated. *arXiv preprint arXiv:2405.10853*, 2024. 2
- [35] Jun Sun, Tianyi Chen, Georgios B Giannakis, Qinmin Yang, and Zaiyue Yang. Lazily aggregated quantized gradient innovation for communication-efficient federated learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(4):2031–2044, 2020. 2
- [36] Xian Sun, Peijin Wang, Wanxuan Lu, Zicong Zhu, Xiaonan Lu, Qibin He, Junxi Li, Xuee Rong, Zhujun Yang, Hao Chang, Qinglin He, Guang Yang, Ruiping Wang, Jiwen Lu, and Kun Fu. Ringmo: A remote sensing foundation model with masked image modeling. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–22, 2023. 2
- [37] Jamie Tolan, Hung-I Yang, Benjamin Nosarzewski, Guillaume Couairon, Huy V Vo, John Brandt, Justine Spore, Sayantan Majumdar, Daniel Haziza, Janaki Vamaraju, et al. 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, 2024. 1
- [38] Devis Tuia, Konrad Schindler, Begüm Demir, Xiao Xiang Zhu, Mrinalini Kochupillai, Sašo Džeroski, Jan N. van Rijn, Holger H. Hoos, Fabio Del Frate, Mihai Datcu, Volker Markl, Bertrand Le Saux, Rochelle Schneider, and Gustau Camps-Valls. Artificial intelligence to advance earth observation: A review of models, recent trends, and pathways forward. *IEEE Geoscience and Remote Sensing Magazine*, pages 2–25, 2024. 1
- [39] Ye Lin Tun, Chu Myaet Thwal, Le Quang Huy, Minh NH Nguyen, and Choong Seon Hong. Lw-fedssl: Resource-efficient layer-wise federated self-supervised learning. *arXiv preprint arXiv:2401.11647*, 2024. 3
- [40] Junjie Wang, Zhuo Zheng, Ailong Ma, Xiaoyan Lu, and Yanfei Zhong. LoveDA: A remote sensing land-cover dataset for domain adaptive semantic segmentation. In *NeurIPS*, 2021. 6, 3
- [41] Lirui Wang, Kaiqing Zhang, Yunzhu Li, Yonglong Tian, and Russ Tedrake. Does learning from decentralized non-IID unlabeled data benefit from self supervision? In *ICLR*, 2023. 2
- [42] Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu. Aid: A benchmark data set for performance evaluation of aerial scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(7):3965–3981, 2017. 6, 2
- [43] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. Dota: A large-scale dataset for object detection in aerial images. In *CVPR*, 2018. 6, 3
- [44] Xingyu Xie, Zhijie Lin, Kim-Chuan Toh, and Pan Zhou. Loco: Low-bit communication adaptor for large-scale model training. *arXiv preprint arXiv:2407.04480*, 2024. 5
- [45] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simsim: A simple framework for masked image modeling. In *CVPR*, pages 9653–9663, 2022. 6, 7
- [46] Rui Yan, Liangqiong Qu, Qingyue Wei, Shih-Cheng Huang, Liyue Shen, Daniel L. Rubin, Lei Xing, and Yuyin Zhou. Label-efficient self-supervised federated learning for tackling data heterogeneity in medical imaging. *IEEE Transactions on Medical Imaging*, 42(7):1932–1943, 2023. 3, 7
- [47] Kunping Yang, Gui-Song Xia, Zicheng Liu, Bo Du, Wen Yang, Marcello Pelillo, and Liangpei Zhang. Semantic change detection with asymmetric siamese networks. *arXiv preprint arXiv:2010.05687*, 2020. 6, 3
- [48] Mi Zhang, Bingnan Yang, Xiangyun Hu, Jianya Gong, and Zuxun Zhang. Foundation model for generalist remote sensing intelligence: Potentials and prospects. *Science Bulletin*, 69(23):3652–3656, 2024. 1
- [49] Weiming Zhuang, Yonggang Wen, and Shuai Zhang. Divergence-aware federated self-supervised learning. In *ICLR*, 2022. 3, 7

Towards Privacy-preserved Pre-training of Remote Sensing Foundation Models with Federated Mutual-guidance Learning

Supplementary Material

A. Overview

We provide the following materials to supplement our paper and divide them into two sections.

- We provide the theoretical analysis of our proposed Fed-Sense in Sec. B.
- We provide the details of our pre-training datasets and downstream datasets in Sec. C

B. Theoretical Analysis

B.1. Assumptions

Assumption 1 (Smoothness) *The self-supervised loss \mathcal{L}_m^{ssl} is L -smooth:*

$$\|\nabla \mathcal{L}_m^{ssl}(\theta_1) - \nabla \mathcal{L}_m^{ssl}(\theta_2)\| \leq L \|\theta_1 - \theta_2\|, \quad \forall \theta_1, \theta_2 \quad (19)$$

Assumption 2 (Bounded Gradient) *Local gradients are bounded:*

$$\mathbb{E}[\|\nabla \mathcal{L}_m^{total}(\theta_m)\|^2] \leq G^2, \quad \forall m \quad (20)$$

Assumption 3 (Parameter Discrepancy) *The discrepancy between local and global models satisfies:*

$$\|\theta_m - \Theta\| \leq \delta, \quad \forall m \in [M] \quad (21)$$

where δ quantifies the maximum client drift.

B.2. Key Lemmas

Lemma 1 (Optimal Perturbation Bound) *Under Assumption 2, the optimal perturbation $\tilde{\epsilon}$ in SCG satisfies:*

$$\|\tilde{\epsilon}\| \leq \lambda \sqrt{\beta^2 \delta^2 + G^2} \quad (22)$$

Proof 1 *From the perturbation approximation:*

$$\begin{aligned} \tilde{\epsilon} &\approx \lambda \frac{\nabla \mathcal{L}_m^{disc}}{\|\nabla \mathcal{L}_m^{disc}\|} \\ \|\tilde{\epsilon}\| &\leq \lambda \sqrt{\frac{\|\nabla \mathcal{L}_m^{disc}\|^2}{\|\nabla \mathcal{L}_m^{disc}\|^2}} = \lambda \end{aligned}$$

Using the parameter discrepancy term $\nabla \mathcal{L}_m^{disc} = \beta(\theta_m - \Theta)$ and Assumption 3:

$$\|\nabla \mathcal{L}_m^{disc}\| \leq \beta \delta$$

Combining with gradient bound G via triangle inequality completes the proof.

Lemma 2 (Quantization Error Decay) *Let e_m^t be the feedback error in CSG. With momentum factor $\alpha \in (0, 1)$, the error decays geometrically:*

$$\|e_m^t\| \leq \alpha^t \|e_m^0\| + \frac{1-\alpha}{1-\alpha^{t+1}} \sum_{k=0}^t \alpha^{t-k} \|\epsilon_q^k\| \quad (23)$$

where ϵ_q^k is the quantization error at round k .

Proof 2 *Unrolling the recursive error update:*

$$\begin{aligned} e_m^t &= \alpha e_m^{t-1} + (1-\alpha) \epsilon_q^t \\ &= \alpha^t e_m^0 + (1-\alpha) \sum_{k=1}^t \alpha^{t-k} \epsilon_q^k \end{aligned}$$

Taking norms and applying triangle inequality:

$$\begin{aligned} \|e_m^t\| &\leq \alpha^t \|e_m^0\| + (1-\alpha) \sum_{k=1}^t \alpha^{t-k} \|\epsilon_q^k\| \\ &\leq \alpha^t \|e_m^0\| + \frac{1-\alpha}{1-\alpha^{t+1}} \sum_{k=0}^t \alpha^{t-k} \|\epsilon_q^k\| \end{aligned}$$

The geometric series bound completes the proof.

B.3. Main Convergence Result

Theorem 1 (Convergence Guarantee) *Under Assumptions 1-3, let learning rate $\gamma = \frac{1}{L\sqrt{T}}$. After T rounds, the averaged gradient satisfies:*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla \mathcal{L}^{total}(\Theta^t)\|^2 \leq \frac{2L(\mathcal{L}^0 - \mathcal{L}^*)}{\sqrt{T}} + \frac{C}{T} \sum_{t=1}^T (\delta^2 + \|e^t\|^2) \quad (24)$$

where C is a constant combining L, G, β, λ .

Proof 3 (Proof Sketch) *Using smoothness (Assump. 1):*

$$\mathcal{L}^{t+1} \leq \mathcal{L}^t + \langle \nabla \mathcal{L}^t, \Theta^{t+1} - \Theta^t \rangle + \frac{L}{2} \|\Theta^{t+1} - \Theta^t\|^2$$

Substituting the update rule $\Theta^{t+1} = \Theta^t - \gamma(\nabla \mathcal{L}^{total} + e^t)$:

$$\begin{aligned} \mathbb{E}[\mathcal{L}^{t+1}] &\leq \mathbb{E}[\mathcal{L}^t] - \gamma \mathbb{E} \|\nabla \mathcal{L}^t\|^2 + \gamma \mathbb{E} \langle \nabla \mathcal{L}^t, e^t \rangle \\ &\quad + \frac{L\gamma^2}{2} \mathbb{E} \|\nabla \mathcal{L}^t + e^t\|^2 \end{aligned}$$

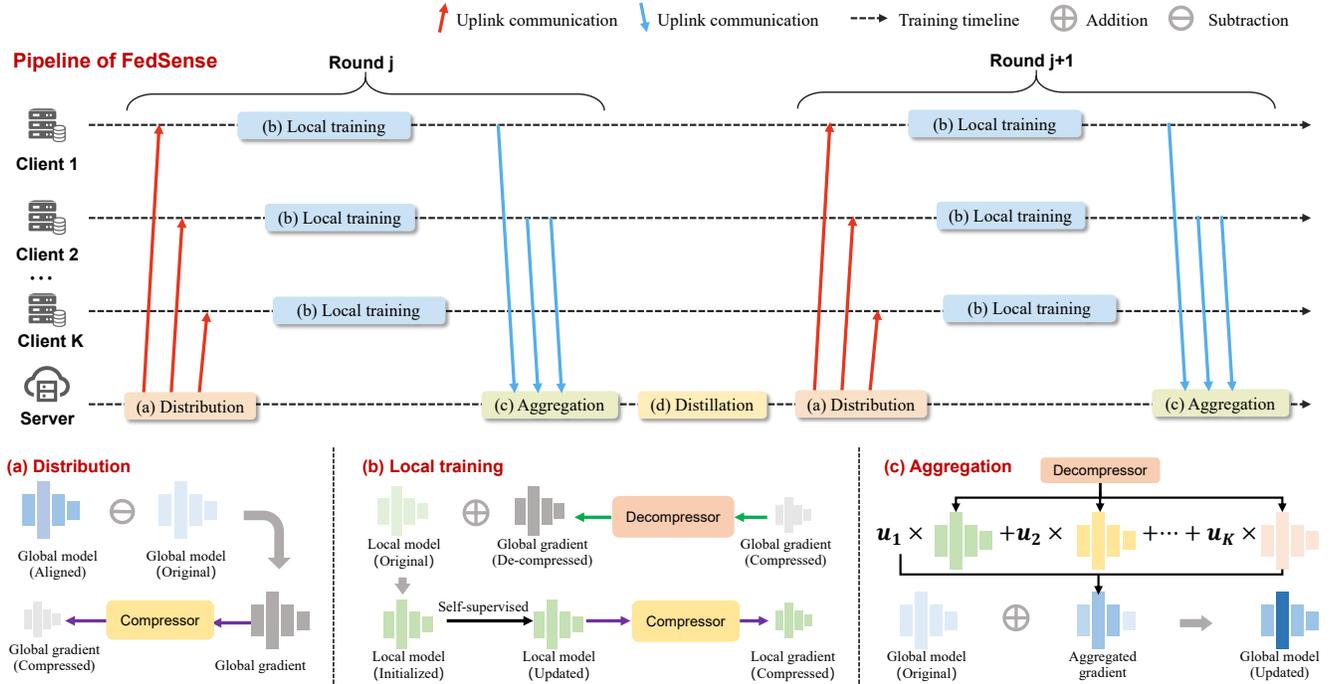


Figure 4. Pipeline of privacy-preserved pre-training of RSFMs.

ID	Source	#samples	GSD (m)	Coverage
Server	WorldView-3/4	240,000	0.5-2.5	Global
Client 01	NOAA	22,292	0.25	USA
Client 02	GF-2	27,300	4.0	China
Client 03	WorldView-2	88,272	0.3-0.5	Global
Client 04	Mixed	125,474	0.3-25	Global
Client 05	QB-2/GE-1	180,562	0.3	Global
Client 06	JL-1/GF-7	40,816	0.8	China
Client 07	Mixed	90,000	0.3-25	Global
Client 08	QB-2/GE-1	180,000	0.3-25	Global
Client 09	NAIP	45,000	1.25	USA
Client 10	Mixed	50,800	0.3-0.75	Global
Total	Multi-source	1,000,000	0.25-25	Global

Table 7. Details of the pre-training datasets.

C. Dataset details and implementation details

Scene Classification.

(1) *Aerial Image Dataset (AID)* [42]. This dataset is comprised of 10,000 images across 30 classes, all sourced from Google Earth and cropped to 600×600 pixels. It

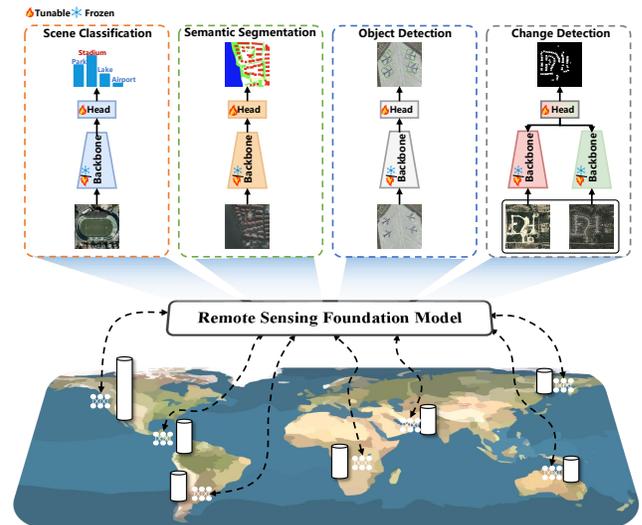


Figure 5. Illustration on downstream usage of collaboratively pre-trained RSFMs to accommodate various Earth Observation tasks.

also includes diverse resolutions from 0.5 to 8 meters per pixel and geographic variations to enhance robustness.

(2) *NWPU-RESISC45* [6]. This dataset comprises 31,500 RGB images at resolutions from 0.2m to 30m across 45 scene classes, each with 700 samples with a size of

256 times 256 pixels. It offers significant variability in translation, scale, viewpoint, illumination, and occlusion. It also has high within-class diversity and inter-class similarity.

Object Detection.

- (1) *DIOR-R* [7]. This dataset consists of 23,463 remote sensing images, with 192,472 annotated object instances spanning 20 categories. The size of each image is 800×800 pixels, and spatial resolutions range from 0.5m to 30m. With emphasis on high inter-class similarity, intra-class diversity, and object size variability, it is designed to benchmark object detection methods in diverse conditions such as different imaging times, weathers, and resolutions.
- (2) *DOTA-v1.0* [43]. This dataset consists of 2,806 aerial images, measuring from 800×800 to 4000×4000 pixels, annotated with 188,282 instances across 15 categories that include airplanes, ships, vehicles, and bridges. The objects in this dataset are presented in diverse scales, orientations and aspect ratios.

Semantic Segmentation.

- (1) *LoveDA* [40]. This dataset for domain-adaptive semantic segmentation features 5,987 images with spatial resolution of 0.3 m, each sized at 1024×1024 pixels in RGB format. Covering 536.15 km^2 , it spans urban and rural areas across Nanjing, Changzhou and Wuhan, and includes pixel-level annotations across seven land-cover categories. It addresses challenges of multi-scale objects, complex backgrounds, and inconsistent class distributions, supporting semantic segmentation and unsupervised domain adaptation.
- (2) *Inria* [28]. This dataset comprises high-resolution RGB aerial imagery, with 180 training and 180 test tiles (each 1500×1500 pixels, 0.3 m resolution). It focuses on two classes: building and non-building, covering a total of 405 km^2 of urban areas across five cities in the U.S. and Austria. The dataset emphasizes generalization challenges, supporting semantic segmentation across diverse urban landscapes.

Change Detection.

- (1) *LEVIR-CD+* [5]. This dataset is a high-resolution urban building change detection dataset comprised of 985 RGB image pairs from Google Earth, each measuring 1024×1024 pixels with a spatial resolution of 0.5 meters per pixel. Spanning 20 regions in Texas, the dataset includes building and land use change masks, covering the years 2002 to 2020 with a 5-year observation interval. It features predominantly urban areas and near-nadir imagery, making it accessible for change detection studies.
- (2) *SECOND* [47]. This dataset is a large-scale semantic change detection benchmark, comprising 4,662 image pairs, each with a size of 512×512 pixels. The images

were collected from multiple platforms across multiple cities including Hangzhou, Chengdu, and Shanghai. It focuses on six land-cover classes: non-vegetated ground surface, tree, low vegetation, water, buildings, and playgrounds. *SECOND* also offers approximately 30 change types, including changes within the same land-cover class, with pixel-level annotations ensuring high diversity and label accuracy.