# Vision Foundation Models in Remote Sensing: A Survey

Siqi Lu[1], Junlin Guo[1], James R Zimmer-Dauphinee[2],
Jordan M Nieusma[4], Xiao Wang[3], Parker VanValkenburgh[5],
Steven A Wernke[2], Yuankai Huo[1,6]
[1]Department of Electrical and Computer Engineering, Vanderbilt University,
[2]Department of Anthropology, Vanderbilt University,
[3]Oak Ridge National Laboratory,
[4]Data Science Institute, Vanderbilt University,
[5]Department of Anthropology, Brown University
[6]Department of Computer Science, Vanderbilt University,

*Abstract*—**Artificial Intelligence (AI) technologies have profoundly transformed the field of remote sensing, revolutionizing data collection, processing, and analysis. Traditionally reliant on manual interpretation and task-specific models, remote sensing research has been significantly enhanced by the advent of foundation models—large-scale, pre-trained AI models capable of performing a wide array of tasks with unprecedented accuracy and efficiency. This paper provides a comprehensive survey of foundation models in the remote sensing domain. We categorize these models based on their architectures, pre-training datasets, and methodologies. Through detailed performance comparisons, we highlight emerging trends and the significant advancements achieved by those foundation models. Additionally, we discuss technical challenges, practical implications, and future research directions, addressing the need for high-quality data, computational resources, and improved model generalization. Our research also finds that pre-training methods, particularly self-supervised learning techniques like contrastive learning and masked autoencoders, remarkably enhance the performance and robustness of foundation models. This survey aims to serve as a resource for researchers and practitioners by providing a panorama of advances and promising pathways for continued development and application of foundation models in remote sensing.**

*Index Terms*—**Remote sensing, Machine learning, Artificial intelligence, Image processing, Computer vision, Transformers.**

## I. INTRODUCTION

ARTIFICIAL Intelligence (AI) technologies have profoundly transformed the field of remote sensing, revolutionizing how data is collected, processed, and analyzed. Traditionally, remote sensing projects relied heavily on manual interpretation and task-specific models that required extensive labeled datasets and significant computational resources. However, the advent of AI and deep learning (DL) has ushered in a new era in which large-scale, pre-trained models, known as foundation models, are capable of performing a wide array of tasks with unprecedented accuracy and efficiency. These advancements have not only enhanced the potential applications of remote sensing but have also opened new avenues for its usage across various domains.

In recent years, numerous vision foundation models have emerged, demonstrating remarkable performance in handling diverse remote sensing tasks. These models have shown the potential to significantly improve performance on multiple downstream tasks such as scene classification, semantic segmentation, object detection, and more. By leveraging vast amounts of pre-training data and sophisticated architectures, these foundation models have set new benchmarks in the field, making them indispensable tools for researchers and engineers alike.

This paper aims to provide a comprehensive survey of vision foundation models in the remote sensing domain ad rem, and is limited to foundation models released between June 2021 and June 2024. This timeframe marks a surge in the development of modern foundation models, including vision transformers and advanced self-supervised learning techniques. Although early models like Tile2Vec [47] and others laid the groundwork for representation learning in remote sensing, they were typically limited in scale and generalization capabilities. Furthermore, numerous review papers have already provided comprehensive overviews of these pre-2021 models. Our review, therefore, focuses on recent developments to highlight the unique contributions and innovations that have emerged in the past few years.

In figure 1, **58** vision foundation models are listed in chronological order. To facilitate navigation and en-

**Foundation Models in Remote Sensing**

2021
Dec — Akiva et al.
Oct — Mañas et al., Li et al.
Jun — Stojnic et al.

2022
Mar. — Ayush et al.
May — Wang D. et al.
Jun. — Wang Y. et al., Scheibenreif et al.
Jul. — Sun et al.
Aug. — Li et al.
Sep. — Jain et al.
Nov. — Wang et al., Zhang et al.

2023
Jun — Mall et al., Prexl et al., Stewart et al.
Apr — Muhtar et al., Tang et al.
Jan — Cong et al.
Aug. — Mendieta et al., Bastani et al.
Sep. — Yao et al., Reed et al., Wang et al., Wang et al.
Oct. — Feng et al., Wang et al
Nov. — Jakubik et al., Fuller et al.
Dec. — Irvin et al.

2024
Jan. — Tang et al., Dumeur et al., Smith et al., Huang et al., Tian et al., Dong et al.
Feb. — Hong et al., Tseng et al.
Mar. — Noman et al., Li et al., Bountos et al., Guo et al., Dong et al.
Apr. — Han et al., Wanyan et al., Wang et al.
May — Xiong et al., Wang et al., Cha et al., Nedungadi et al., Zhang et al., Wang et al., Jiang et al.
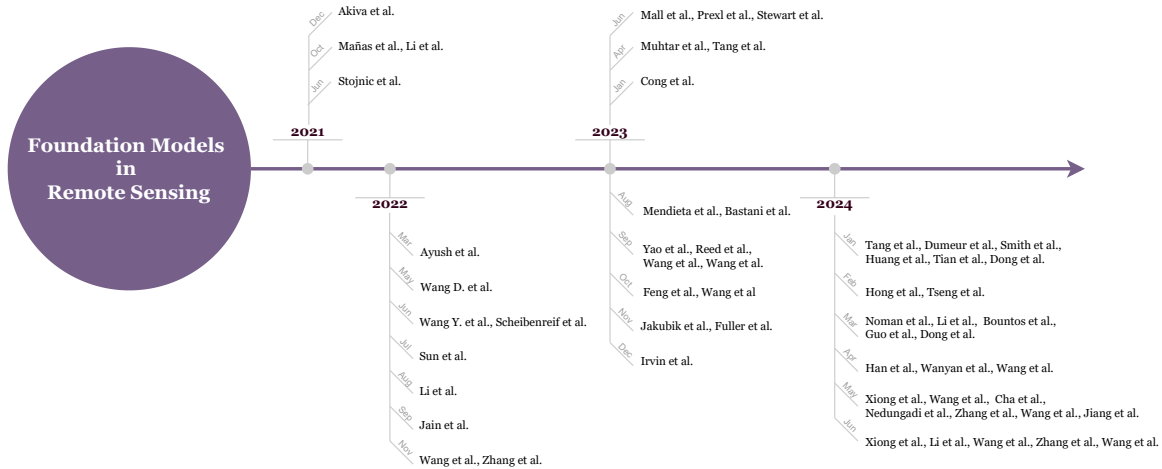Jun. — Xiong et al., Li et al., Wang et al., Zhang et al., Wang et al.

Fig. 1: Overview of some well-known foundation models for remote sensing from 2021 June to 2024 June.

hance utility for researchers, we categorized existing models based on their perception levels (e.g., image-level, region-level, pixel-level). This organization helps clarify which models have been tested for general image-based challenges or specialized applications such as environmental monitoring, land cover mapping, archaeological exploration, disaster management, and more. It is essential to distinguish between applications that models have been explicitly tested on and those for which they could potentially be effective for. In this review, the fact that a model has not been tested on a particular application does not mean it won't perform well. Foundation models, especially convolutional neural network (CNN) backbones like residual networks (ResNet) [36] and vision transformers (ViT) [25], may still be suitable for various downstream tasks, even if prior work has not yet demonstrated this.

Our contributions include:

1) An exhaustive review of current state of vision foundation models proposed in the field of remote sensing, starting from the background and methodologies of these models to specific applications across different domains and tasks in a hierarchical and structured manner.
2) Categorization and analysis of the models based on their application in both image analysis (table I) and practical applications (table VI-B). We discuss the architecture, pre-training datasets, pre-training methods, and performance of each model.
3) Discussion of challenges and unresolved aspects related to foundation models in remote sensing. We pinpoint new trends, raise important questions, and proposed future directions for further exploration.

## II. BACKGROUND

### A. Remote Sensing

Remote sensing (RS) refers to the process of acquiring information about objects or areas from a distance, typically using satellite or airborne sensors. These technologies and techniques serve vital roles in diverse fields, enabling the collection of data over geographic areas without physical contact. Applications of remote sensing include earth observation, digital archaeology, urban planning and development, and disaster management. The field of remote sensing has developed rapidly since the mid 20th century. Initially, remote sensing predominately consisted of analog photographic techniques via aerial and satellite platforms, which provided limited spectral and spatial resolution. The launch of early Earth observation satellites, such as Landsat program commenced in 1967 [112], marked a significant advancement, enabling consistent and wide-ranging data collection for environmental monitoring.

Modern remote sensing employs a variety of sensors suited for specific types of data collection, including optical, thermal, and radar.. Optical sensors capture a wide variety of spectral bands, including visible and near-infrared light, allowing for detailed imaging of land cover and vegetation health. Thermal sensors detect heat emitted or reflected from the Earth's surface, useful for monitoring volcanic activity, forest fires, and climate change monitoring. Radar sensors can penetrate clouds and vegetation, providing critical information in all-weather conditions and for applications such as soil moisture estimation and urban infrastructure mapping [17], [71].

In recent years, remote sensing has found applications in many fields. With regard to environmental monitoring,

**Data**   |   **Downstream Tasks**

Panchromatic   True Color   SAR

Hyperspectral   Multispectral

...

Segmentation   Object Detection
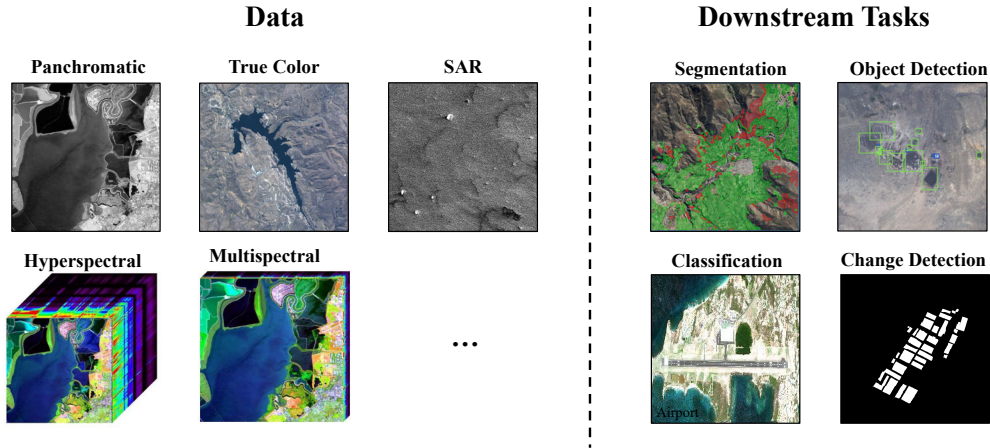
Classification   Change Detection

Airport

Fig. 2: Examples of data types used in those foundation models and downstream tasks that can be done by foundation models. Data: (1) Panchromatic [4], (2) True Color, (3) SAR [94], (4) Hyperspectral [4], (5) Multispectral [4]. Downstream tasks: (1) Segmentation, (2) Object Detection, (3) Classification [15], (4) Change Detection [76]. [1]

it is used to track deforestation, to monitor air and water quality, and to assess the impacts of climate change [30], [39]. In agriculture, remote sensing helps in crop health monitoring, yield estimation, and efficient resource management [71]. Urban planning and development benefit from remote sensing through the monitoring of urban sprawl, infrastructure development, and land-use planning [17], [48]. Furthermore, in disaster management, remote sensing is crucial for assessing the damage caused by natural disasters, aiding in the planning and execution of relief operations [1], [30].

The integration of remote sensing data with Geographic Information Systems (GIS) has further enhanced its utility. GIS provides a framework for capturing, storing, analyzing, and visualizing spatial and geographic data. When combined with remote sensing data, GIS can be used to create detailed and dynamic maps and models for various applications. This synergy is particularly valuable in resource management, urban planning, and disaster response, where accurate and timely information is critical [17], [30], [71].

### B. Foundation Models for Remote Sensing

Foundation models (FMs) refer to large-scale, pretrained models that provide a robust starting point for various downstream tasks across different domains [50]. These models leverage extensive datasets and advanced architectures, enabling them to capture complex patterns and features that can be fine-tuned for specific

applications with minimal additional training. In remote sensing, FMs are particularly valuable due to the diverse and complex nature of the data, including multispectral and multi-temporal imagery. Techniques such as self-supervised learning (SSL) [51] and transformers [93] have significantly enhanced the performance and efficiency of tasks such as image classification, object detection, and change detection, addressing the unique challenges posed by remote sensing data [19].

A major strength of these models lies in their ability to utilize SSL to learn effective representations from largely unlabeled data, which is often abundant in remote sensing scenarios [126]. By integrating advanced architectures like transformers [93], FMs in remote sensing can handle the unique characteristics of geospatial data, such as varying spatial resolutions and temporal dynamics, without requiring separate task-specific models.

The evolution of FMs has been driven by advancements in deep learning and the availability of large datasets. Initially, convolutional neural networks (CNNs) like ResNet [37] paved the way for improved image recognition and classification tasks [65]. The introduction of transformers, which use self-attention mechanisms to model long-range dependencies, has further advanced the capabilities of FMs in handling large-scale image data [16]. Vision transformers (ViTs) [25] extend the transformer architecture to process image data by treating image patches as sequences of tokens, enabling models to learn both local and global relationships. This capability makes transformers particularly effective for semantic segmentation and change detection tasks, where

---

[1]True Color, Segmentation, and Object detection images © MAXAR 2024, provided through the NextView License Agreement.

| Year-Month | Architecture | Model Name | Image-Level | Pixel-Level | Region-Level | Spatial-Temporal | Contrastive Learning | Predictive Coding |
|---|---|---|---|---|---|---|---|---|
| 2021 Jun | ResNet-50 | CMC-RSSR [84] | ✓ | | | | ✓ | |
| 2021 Oct | ResNet-50 | SeCo [66] | ✓ | | | ✓ | | |
| 2021 Oct | ResNet-50 | GeoKR [56] | ✓ | ✓ | ✓ | | | |
| 2021 Dec | ResNet-34 | MATTER [2] | ✓ | ✓ | | ✓ | | ✓ |
| 2022 Mar | ResNet-50 | GASSL [6] | ✓ | ✓ | ✓ | | ✓ | |
| 2022 May | ViTAEv2-S | RSP [96] | ✓ | ✓ | ✓ | ✓ | | |
| 2022 Jun | ViT-S/8 | DINO-MM [105] | ✓ | | | | ✓ | |
| 2022 Jun | Swin Transformer | Scheibenreif, et al. [79] | ✓ | ✓ | | | ✓ | |
| 2022 Jul | ViT/Swin Transformer | RingMo [87] | ✓ | ✓ | ✓ | ✓ | | ✓ |
| 2022 Aug | ResNet-50 | GeCO [57] | ✓ | ✓ | ✓ | | | ✓ |
| 2022 Sep | BYOL | RS-BYOL [45] | ✓ | ✓ | | | ✓ | |
| 2022 Nov | ViT-B | CSPT [124] | ✓ | | ✓ | | | ✓ |
| 2022 Nov | ViT | RVSA [100] | ✓ | ✓ | ✓ | | | ✓ |
| 2023 Jan | MAE-based Framework | SatMAE [16] | ✓ | ✓ | | | | ✓ |
| 2023 Apr | TOV | TOV [89] | ✓ | ✓ | ✓ | | | ✓ |
| 2023 Apr | Teacher-student Self-distillation | CMID [70] | ✓ | ✓ | ✓ | ✓ | | |
| 2023 Jun | CACo | CACo [67] | ✓ | ✓ | | ✓ | ✓ | |
| 2023 Jun | ResNet-18 | IaI-SimCLR [77] | ✓ | | | | ✓ | |
| 2023 Jun | ResNet | SSL4EO-L [83] | | ✓ | | | ✓ | |
| 2023 Aug | Teacher-Student | GFM [69] | ✓ | ✓ | | ✓ | | ✓ |
| 2023 Aug | Swim Transformer | SatLasPretrain [7] | ✓ | ✓ | | | | |
| 2023 Sep | Multi-Branch | RingMo-Sense [119] | | ✓ | | | | ✓ |
| 2023 Sep | ViT | Scale-MAE [78] | ✓ | ✓ | | | | ✓ |
| 2023 Sep | CNN-Transformer | RingMo-lite [109] | ✓ | ✓ | ✓ | ✓ | | ✓ |
| 2023 Sep | Multimodel SSL | DeCUR [102] | ✓ | ✓ | | | | ✓ |
| 2023 Oct | MSFE+MMFH | Feng et al. [27] | ✓ | ✓ | ✓ | ✓ | | ✓ |
| 2023 Oct | ViT | FG-MAE [108] | ✓ | ✓ | | | | ✓ |
| 2023 Nov | ViT | Prithvi [46] | | ✓ | | | | ✓ |
| 2023 Nov | Multimodal Encoder | CROMA [28] | ✓ | ✓ | | | ✓ | ✓ |
| 2023 Dec | ViT | USat [44] | ✓ | | | | | ✓ |
| 2024 Jan | ViT-B | Cross-Scale MAE [88] | ✓ | ✓ | | | | ✓ |
| 2024 Jan | Unet+Transformer | U-BARN [26] | ✓ | ✓ | | | | |
| 2024 Jan | Autoregressive Transformer | EarthPT [82] | ✓ | | | | | ✓ |
| 2024 Jan | Teacher-Student Network | GeRSP [42] | ✓ | ✓ | ✓ | | ✓ | ✓ |
| 2024 Jan | Dual-Branch | SwiMDiff [91] | ✓ | | | | | |
| 2024 Jan | Generative ConvNet | SMLFR [22] | | ✓ | ✓ | | | ✓ |
| 2024 Feb | 3D GPT | SpectralGPT [40] | ✓ | ✓ | | ✓ | | ✓ |
| 2024 Feb | MAE-based Framework | Presto [92] | | ✓ | | | ✓ | ✓ |
| 2024 Mar | SatMAE | SatMAE++ [73] | ✓ | | | | | ✓ |
| 2024 Mar | Joint-Embedding Predictive Architecture | SAR-JEPA [58] | ✓ | | | | | ✓ |
| 2024 Mar | ViT | FoMo-Bench [8] | ✓ | ✓ | ✓ | | | ✓ |
| 2024 Mar | Factorized Multi-Modal Spatiotemporal Encoder | SkySense [32] | ✓ | ✓ | ✓ | ✓ | | ✓ |
| 2024 Mar | Multi-Modules | UPetu [24] | ✓ | ✓ | | ✓ | | ✓ |
| 2024 Apr | Swim Transformer | msGFM [33] | ✓ | ✓ | | | | ✓ |
| 2024 Apr | DINO | DINO-MC [111] | ✓ | | | ✓ | ✓ | |
| 2024 May | OFA-Net | OFA-Net [118] | ✓ | ✓ | | | | ✓ |
| 2024 May | Shared Encoder, Task-Specific Decoders | MTP [99] | ✓ | ✓ | ✓ | ✓ | | |
| 2024 May | ViT | BFM [11] | | ✓ | ✓ | | | ✓ |
| 2024 May | MP-MAE | MMEarth [72] | ✓ | ✓ | | | | ✓ |
| 2024 May | ViT | CtxMIM [123] | ✓ | ✓ | ✓ | | | ✓ |
| 2024 May | HiViT | SARATR-X [54] | ✓ | | ✓ | | | ✓ |
| 2024 May | Transformer | SoftCon [106] | ✓ | ✓ | | ✓ | ✓ | |
| 2024 May | ViT | LeMeViT [49] | | ✓ | ✓ | ✓ | | |
| 2024 Jun | Masked Autoencoder | S2MAE [59] | ✓ | | | ✓ | | ✓ |
| 2024 Jun | CNN - Transformer | RS-DFM [110] | | ✓ | ✓ | | | |
| 2024 Jun | MAE-based | A2-MAE [122] | ✓ | ✓ | | ✓ | | |
| 2024 Jun | ViT | HyperSIGMA [95] | ✓ | ✓ | ✓ | ✓ | | ✓ |
| 2024 Jun | Dynamic OFA | DOFA [117] | ✓ | ✓ | | | | ✓ |

TABLE I: Summary of the pretraining methods utilized and image analysis tasks evaluated across different models. Image-level, pixel-level, region-level, and spatial-temporal classify the tasks in image analysis, while contrastive learning and predictive coding indicates the different self-supervised pretraining strategies that each study used.

capturing long-range dependencies is crucial, especially in high-resolution satellite imagery.

Notable foundation models in remote sensing include SatMAE [16], which pre-trains transformers for temporal and multi-spectral satellite imagery; Scale-MAE [78], a scale-aware masked autoencoder for multiscale geospatial representation learning; and DINO-MC [111], which extends global-local view alignment for SSL with remote sensing imagery. These models have shown remarkable performance in various remote sensing tasks such as scene classification, object detection, and change detection.

Despite their success, FMs face several challenges, including the need for high-quality and diverse training data, significant computational resources, and effective domain adaptation to specific remote sensing tasks [73]. Addressing these challenges will be crucial for the continued advancement of FMs in remote sensing.

## III. RELATED REVIEW PAPERS

Artificial intelligence in remote sensing has been a growing area of research, with numerous review papers providing insights into AI advancements and their applications. In this section, we summarize the most influential reviews on foundation models in remote sensing.

Zhang et al. (2016), in their foundational review "*Deep Learning for Remote Sensing Data: A Technical Tutorial on the State of the Art*" [121], introduced deep learning techniques to RS, focusing on convolutional neural networks (CNNs) for tasks such as image classification and object detection. This work highlighted both the promise and challenges of early AI integration in RS, setting the stage for subsequent advancements.

In 2017, Zhu et al.'s "*Deep Learning in Remote Sensing: A Comprehensive Review and List of Resources*" [129] delved into diverse AI applications, including hyperspectral analysis and synthetic aperture radar (SAR) interpretation. It also provided an extensive resource list, capturing the rapid adoption of deep learning in addressing complex RS challenges, paving the way for more advanced AI models in the following years.

More recent reviews have focused on advanced AI models and methods. Wang et al.'s 2022 review, "*Self-Supervised Learning in Remote Sensing*" [103], highlighted the ability of self-supervised learning (SSL) methods to utilize large volumes of unlabeled data, significantly reducing dependence on labeled datasets while maintaining high performance in RS tasks. The review also identified key challenges and future directions, emphasizing SSL's potential to handle large-scale RS data complexities.

Zhang et al. (2022), in "*Artificial Intelligence for Remote Sensing Data Analysis: A Review of Challenges and Opportunities*" [120], offered a comprehensive overview of AI algorithms, synthesizing findings from over 270 studies. It emphasized ongoing challenges such as explainability, security, and integrating AI with other computational techniques, serving as a roadmap for future innovation in AI-driven RS.

Aleissaee et al.'s 2023 survey, "*Transformers in Remote Sensing*" [3], explored the impact of transformer-based models across various RS tasks, comparing them with CNNs. It identified both strengths and limitations, along with unresolved challenges, providing a detailed roadmap for future research on transformers' role in RS.

Li et al.'s 2024 review, "*Vision-Language Models in Remote Sensing*" [60], examined the increasing significance of vision-language models (VLMs), which combine visual and textual data. It highlighted VLMs' potential in applications like image captioning and visual question answering, emphasizing a shift toward richer semantic understanding in RS tasks.

Additionally, the recent work, "*On the Foundations of Earth and Climate Foundation Models*" [130], provided a comprehensive review of existing foundation models, proposing features like geolocation embedding and multisensory capability. It outlined key traits for future Earth and climate models, contributing to a broader discussion on foundational advancements in geospatial AI.

Building on these reviews, our study provides a comprehensive analysis of foundation models developed from June 2021 to June 2024, focusing on advances in self-supervised learning and transformer-based architectures. Unlike previous reviews, which focused mainly on individual techniques, we explore their combined potential in remote sensing tasks like semantic segmentation, multi-spectral analysis, and change detection. For instance, SatMAE [16] demonstrates effective use of SSL for pre-training transformers, enabling improved segmentation in complex multi-spectral imagery, while Scale-MAE employs scale-aware masked autoencoders for better handling of varied spatial resolutions in remote sensing data.

Our study also highlights new models like DINO-MC [111], which integrates global-local view alignment for SSL, making it particularly effective for identifying changes in high-resolution satellite imagery. By systematically examining these innovations, we illustrate how recent models address persistent challenges like domain adaptation and computational efficiency. For example, efficient self-attention mechanisms in Scale-MAE [78] help reduce computation costs, while enhanced geolocation embeddings in models like SatMAE improve performance in geospatial feature extraction.

In contrast to earlier reviews, which often remained theoretical, we emphasize both the theoretical advancements and practical applications of recent models. For example, DINO-MC [111] and ORBIT's [101] real-world applications in environmental monitoring and disaster response highlight its practical impact, demonstrating how new FMs can be effectively leveraged to address pressing challenges in geospatial analysis.

## IV. PRETRAINING METHODS

Pretraining serves as a critical step in developing foundation models (FM), enabling them to learn transferable and generalized representations from large-scale datasets. This process leverage self-supervised or supervised learning methods to extract domain-agnostic features that can be adapted to various downstream tasks. In this section, we explore the key pretraining methods utilized commonly in foundation models for remote sensing, explaining the mechanism of these methods and their roles in enhancing model performance and addressing challenges in this field.

### A. Self-Supervised Learning

Self-supervised learning has emerged as a cornerstone of pre-training foundation models, offering a paradigm where models learn representations by predicting parts of the input data from other parts. This approach reduces reliance on expensive and time-consuming labeled
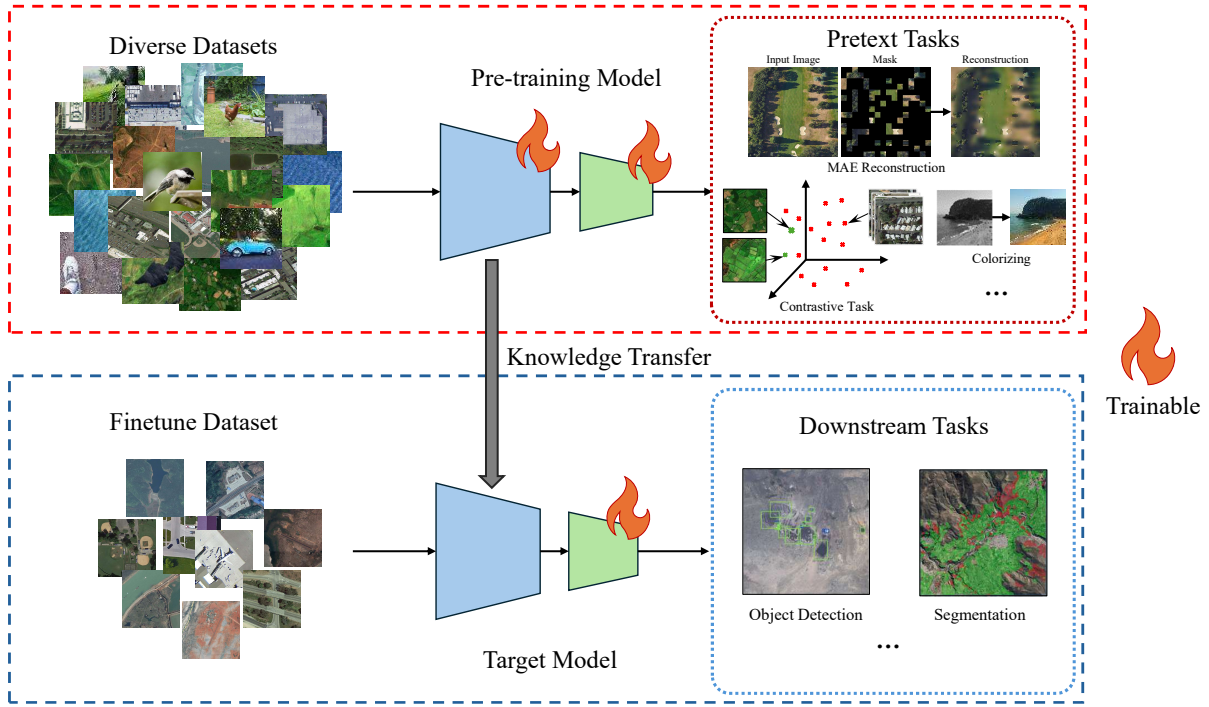
Fig. 3: General pipeline of SSL [51]. Diverse datasets images and pretext task images are acquired from ImageNet [18], BigEarthNet [85], and MillionAID [64]. Finetune dataset includes images from DIOR [55]. [2]

datasets, making it particularly advantageous in fields like remote sensing, where labeled data is often scarce or challenging to obtain.

SSL allows models to exploit vast amounts of unlabeled data, learning rich, generalizable representations that transfer well to downstream tasks such as scene classification, semantic segmentation, object detection, and change detection. By uncovering underlying data structures and patterns, SSL not only enhances model robustness but also improves adaptability across diverse domains and resolutions of remote sensing imagery [104]. Figure 3 illustrates the general pipeline of self-supervised learning.

Two SSL methods commonly used in vision foundation models for remote sensing are contrastive learning and predictive coding, each offering unique mechanisms to harness information from unlabeled data.

*1) Predictive Coding:* Predictive coding leverages a generative approach, where the model learns to predict missing or occluded parts of an image based on visible portions. This strategy helps capture spatial and contextual relationships in remote sensing imagery, which often contains diverse textures, complex scenes, and varying resolutions.

In remote sensing, predictive coding can be applied to tasks such as gap filling in satellite imagery, where the model learns to infer missing data caused by sensor limitations or occlusions like cloud cover. Popular implementations of predictive coding frameworks include autoencoder-based architectures, masked image modeling techniques like those used in MAE (Masked Autoencoders) [34], and autoregressive models. These methods are particularly effective in learning fine-grained details critical for high-resolution imagery and specialized tasks.

*2) Contrastive Learning:* Contrastive learning is another powerful SSL technique that focuses on distinguishing between similar and dissimilar samples in the data. The key idea is to bring representations of similar (positive) samples closer together while pushing apart those of dissimilar (negative) samples. This encourages the model to learn discriminative and invariant features that are crucial for remote sensing tasks.

Contrastive learning frameworks such as SimCLR [13], MoCo [35], DINO [9], and BYOL [29] have shown promise in remote sensing applications. They use augmentations like random cropping, rotations, or spectral

---

[2]Object detection and Segmentation © MAXAR 2024, provided through the NextView License Agreement.

band dropping to generate positive pairs, enabling the model to learn robust representations invariant to these transformations. For instance, in multispectral or hyperspectral imagery, contrastive learning can help models capture spectral signatures across varying conditions, improving performance in tasks like crop classification or land cover mapping [104].

Contrastive learning is especially relevant in remote sensing when labeled datasets are highly imbalanced, as it enables models to learn from underrepresented classes or regions without explicit labels.

By combining approaches like predictive coding and contrastive learning, self-supervised learning has significantly advanced the development of vision foundation models in remote sensing. These methods allow models to leverage vast, unlabeled datasets while maintaining adaptability across diverse spatial resolutions, spectral bands, and application scenarios. On the other hand, it is important to note that there are many other SSL methods that can be employed to such tasks. Other innovative methods, such as teacher-student self-distillation frameworks, have also demonstrated potential in remote sensing applications. For example, CMID [70] achieves promising performance by combining contrastive learning and masked image modeling in a teacher-student self-distillation framework. This structure enables it to capture both global and local features, making it effective for diverse remote sensing tasks.. The diversity of SSL techniques highlights the versatility and evolving nature of self-supervised learning, underscoring its critical role in unlocking the full potential of remote sensing imagery.

### B. Supervised Pretraining

Supervised pretraining is a fundamental approach in deep learning, where models are trained using labeled datasets to minimize prediction errors for specific tasks, such as image classification. This method allows models to learn direct mappings between input features and target labels, fostering the development of detailed, task-specific representations. For instance, models like ResNet [36] and VGGNet (Visual Geometry Group Network) [81] trained on large-scale datasets such as ImageNet [18] have demonstrated how supervised pretraining can capture robust feature hierarchies that are highly transferable to related tasks, including semantic segmentation or object detection.

In remote sensing, supervised pretraining has shown promise for tasks such as land cover classification and object detection using high-resolution satellite imagery [97]. However, the dependency on large-scale labeled datasets presents a major limitation. Creating labeled datasets for remote sensing tasks, particularly when involving multispectral or hyperspectral data, is resource-intensive and often requires domain expertise for annotation. For example, labeling pixel-level data for land cover classification or delineating objects in complex urban environments can be prohibitively time-consuming. Furthermore, labeled data in remote sensing is often domain-specific, limiting the generalizability of models trained on one dataset to other applications or regions [129].

These challenges highlight the need for innovative strategies to address the reliance on labeled data. Such limitations have motivated the development of alternative approaches, including self-supervised pretraining methods, which leverage the abundance of unlabeled data to learn general-purpose representations without manual annotation.

## V. IMAGE ANALYSIS METHODS

### A. Image Perception at Different Levels

Foundation models in remote sensing enable image analysis at three primary levels: image-level, region-level, and pixel-level. These levels address varying spatial, contextual, and application-specific needs, providing the foundation for a wide range of tasks such as environmental monitoring, urban planning, disaster response, and more. The following subsections outline the distinct objectives and applications at each level. A detailed summary of the models evaluated for these tasks is provided in table II. The following subsections outline the distinct objectives and applications at each level.

*1) Image-Level:* Image-level analysis focuses on classification tasks, categorizing entire images or large image segments into predefined classes, such as urban, forest, water bodies, or agricultural areas. This approach provides broad, high-level insights into geographic regions and is instrumental in large-scale applications like land use mapping, land cover classification, and resource management. By classifying entire scenes, this level of analysis enables efficient monitoring of extensive areas, supporting decision-making in environmental management and policy planning.

*2) Region-Level:* Region-level analysis identifies and localizes specific objects within an image, such as buildings, vehicles, ships, or other structures. Unlike image-level analysis, which provides holistic classifications, region-level tasks focus on object detection which is to detect individual entities and their spatial locations. This analysis is critical for targeted applications like urban planning, where the detection of infrastructure is essential, as well as disaster response and security, where identifying damaged buildings or vulnerable areas can significantly aid in timely interventions.

TABLE II: Overview of recent foundation models in remote sensing, categorized by architecture, model name, pre-training dataset, resolution, geographic coverage, image analysis lebels, visual encoder, pre-training methods, and the number of parameters. **Abbreviations for pre-training methods as specified in the original work:** SSL refers to Self-Supervised Learning, CL stands for Contrastive Learning, MIM is Masked Image Modeling, MAE is Masked Autoencoders, and FD-MIM is Feature-Distilled Masked Image Modeling.

| Model Name | Architecture | Pre-training Dataset | Resolution (m) | Geographic Coverage | Image Analysis Levels | Pretrain methods | # of Params |
|---|---|---|---|---|---|---|---|
| CMC-RSSR [84] | ResNet-50 | NWPU-DOTA [114], BigEarthNet [85], ImageNet [18] | 0.2 to 60 | Global | Image-level | Contrastive Multiview Coding | 23M |
| SeCo [66] | ResNet-50 | Sentinel-2 Imagery | 10, 20, 60 | 200k Locations Worldwide | Image-level, Spacial-temporal | CL | 23.5M |
| GeoKR [56] | ResNet-50 | Levir-KR [56] | 0.8 to 16 | Global | Image-level, Pixel-level, Region-level | Geographical Knowledge Supervision | 23.5M/138M |
| MATTER [2] | ResNet-34 | Sentinel-2 Imagery | - | Rural and Remote Regions with Little Changes | Image-level, Pixel-level | SSL | 21.3M |
| GASSL [6] | ResNet-50 | fMoW [15], GeoImageNet [18] | - | 7 Continents | Image-level, Pixel-level Region-level | CL | 23.5M |
| RSP [96] | ViTAEv2-S | MillionAID [63], [64] | 0.5 to 153 | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | Supervised Learning | 24.8M/23.5M/29M |
| DINO-MM [105] | ViT-S/8 | BigEarthNet-MM [86] | 10 | Global | Image-level | SSL | 22M |
| Scheibenreif, et al. [79] | Swin Transformer | SEN12MS [80] | 10 | Global | Image-level, Pixel-level | CL | - |
| RingMo [87] | ViT/Swin Transformer | 2 million RS images | 0.3 to 30 | 6 Continents | Image-level, Pixel-level, Region-level, Spacial-temporal | MIM | - |
| GeCo [57] | ResNet-50 | Levir-KR [56] | 0.8 to 16 | Global | Image-level, Pixel-level, Region-level | SSL | 23.5M |
| RS-BYOL [45] | BYOL | Sen12MS [80] | 10 to 20 | Global | Image-level, Pixel-level | SSL | 23.5M |
| CSPT [124] | ViT-B | ImageNet-1K [18] | - | Global | Image-level, Region-level | SSL | 86M |
| RVSA [100] | ViT | MillionAID [63], [64] | 0.5 to 153 | Global | Image-level, Pixel-level, Region-level | MAE | 100M |
| SatMAE [16] | MAE-based Framework | fMoW Sentinel-2 [15] | 10, 20, 60 | Global | Image-level, Pixel-level | MAE | 307M |
| TOV [89] | TOV | TOV-NI, TOV-RS | - | Global | Image-level, Pixel-level, Region-level | SSL | - |
| CMID [70] | Teacher-student Self-distillation | MillionAID [63], [64] | Varied | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | SSL | 25.6M/87.8M |
| CACo [67] | ResNet-18/50 | Sentinel-2 Imagery | 10 | Global | Image-level, Pixel-level, Spacial-temporal | SSL | 11.7M/23.5M |
| IaI-SimCLR [77] | ResNet-18 | SEN12MS | - | Global | Image-level | CL | 11.7M |
| SSL4EO-L [83] | ResNet/ViT | ImageNet [18], MoCo [35], SimCLR [13] | 30 | Global | Pixel-level | SSL | 11.7M/23.5M/86M |

TABLE II – continued from previous page

| Model Name | Architecture | Pre-training Dataset | Resolution (m) | Geographic Coverage | Image Analysis Levels | Pretrain methods | # of Params |
|---|---|---|---|---|---|---|---|
| GFM [69] | Teacher-Student | GeoPile [69] | - | Global | Image-level, Pixel-level | Continual Pretraining | - |
| SatlasPretrain [7] | SatlasNet | GeoPile [69] | 1, 10 | Global | Image-level, Pixel-level | Multi-task Learning | 88M |
| RingMo-Sense [119] | Multi-Branch | RS Spatiotemporal Dataset | - | Global | Pixel-level | SSL | - |
| Scale-MAE [78] | ViT-Large | FMoW [15] | - | Global | Image-level, Pixel-level | MAE | 322.9M |
| RingMo-lite [109] | CNN-Transformer | AID [115] | 0.3 to 30 | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | FD-MIM | 60% less than RingMo |
| DeCUR [102] | Multimodel SSL | SSL4EO-S12 [107], RGB-DEM/depth | Varied | Global | Image-level, Pixel-level | SSL | 23.5M |
| Feng et al. [27] | MSFE+MMFH | Multi-modal Dataset | Varied | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | SSL | - |
| FG-MAE [108] | ViT | SSL4EO-S12 [107] | 10 | Global | Image-level, Pixel-level | MAE | - |
| Prithvi [46] | ViT | Harmonized Landsat Sentinel 2 | 30 | Contiguous U.S. | Pixel-level | MAE | 100M |
| CROMA [28] | Multimodal Encoder | SSL4EO [107] | 10 | Areas Surrounding Human Settlements | Image-level, Pixel-level | CL, MAE | 86M |
| USat [44] | ViT | Satlas [7] | Varied | Global | Pixel-level | MAE | - |
| Cross-Scale MAE [88] | ViT-B | fMoW [15] | - | Global | Image-level, Pixel-level | MAE | 86M |
| U-BARN [26] | Unet+Transformer | Sentinel-2 Imagery | Varied | France | Image-level, Pixel-level | SSL | - |
| EarthPT [82] | Transformer | Sentinel-2 Imagery | 10 | UK | Image-level | Autoregressive SSL | 700M |
| GeRSP [42] | Teacher-Student Network | ImageNet [18], MillionAID [63], [64] | 0.5 to 153 | Global | Image-level, Pixel-level, Region-level | SSL, SL | - |
| SwiMDiff [91] | Dual-Branch | Sen12MS [80] | Varied | Global | Image-level, Spacial-temporal | SSL | 11.7M |
| SMLFR [22] | Generative ConvNet | GeoSense [23] | 0.05 to 150 | Multiple Continents | Pixel-level, Region-level | SSL | 88M/197M |
| SpectralGPT [40] | 3D GPT | Sentinel-2 Imagery | Varied | Global | Image-level, Pixel-level, Spacial-temporal | MAE | 100M/300M/600M |
| Presto [92] | MAE-based Framework | Presto-21.5M [92] | 10 | Global | Crop Type Segmentation | MAE | 402K |
| SatMAE++ [73] | SatMAE | fMoW [15] | Varied | Global | Image-level | Multi-Scale Pre-training | - |
| SAR-JEPA [58] | Joint-Embedding Predictive Architecture | 100K SAR Images | Varied | Global | Image-level | SSL | - |

TABLE II – continued from previous page

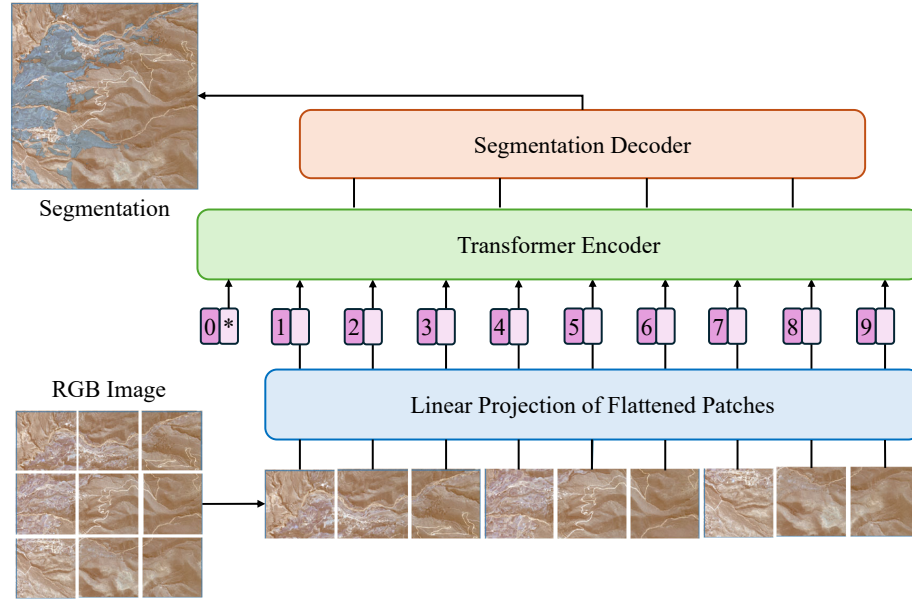| Model Name | Architecture | Pre-training Dataset | Resolution (m) | Geographic Coverage | Image Analysis Levels | Pretrain methods | # of Params |
|---|---|---|---|---|---|---|---|
| FoMo-Bench [8] | ViT | Multiple | Varied | Global | Image-level, Pixel-level, Region-level | MAE | 101M/110M |
| SkySense [32] | Factorized Multi-Modal Spatiotemporal Encoder | Multiple | Varied | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | CL | 2.06B |
| UPetu [24] | Multi-Modules | GeoSense [23] | - | Global | Image-level, Pixel-level, Spacial-temporal | SSL | 0.65M |
| msGFM [33] | Swin Transformer | GeoPile-2 [69] | 0.1 to 153 | Global | Image-level, Pixel-level | MIM | 89M |
| DINO-MC [111] | DINO | SeCo-100K [66] | 10 to 60 | Global | Image-level, Spacial-temporal | SSL | - |
| OFA-Net [118] | OFA-Net | Multi-modal Dataset | Varied | Global | Image-level, Pixel-level | MIM | - |
| MTP [99] | Shared Encoder Task-Specific Decoders | SAMRS [98] | Varied | Global | Image-level, Pixel-level, Region-level, Spacial-temporal | Multi-Task Pretraining | over 300M |
| BFM [11] | ViT | MillionAID [63], [64] | 0.5 to 153 | Global | Pixel-level, Region-level | MAE | 86M/605.26M/ 1.36B/2.42B |
| MMEarth [72] | MP-MAE | Multi-modal, Geospatial Data | - | Global | Image-level, Pixel-level | MP-MAE | 3.7M to 650M |
| CtxMIM [123] | ViT | WorldView-3 Imagery | Varied | Asia | Image-level, Pixel-level, Region-level | MIM | 88M |
| SARATR-X [54] | HiViT | SAR Datasets | 0.1 to 3 | Global | Image-level, Region-level | MIM | 66M |
| SoftCon [106] | Siamese Network with ResNet and ViT Backbones | SSL4EO-S12-ML [107] | - | Global | Image-level, Pixel-level, Spacial-temporal | Multi-label Soft Contrastive Learning | 23M, 23M, 86M |
| LeMeViT [49] | Hierarchical ViT | MillionAID [63], [64] | - | - | Image-level, Pixel-level, Region-level , Spacial-temporal | Dual Cross-Attention with Learnable Meta Token Adaptation | 8.33M to 52.61M |
| S2MAE [59] | 3D Transformer-based MAE | fMoW-Sentinel [15], BigEarthNet [85] | - | Global | Image-level, Spacial-temporal | 3D MAE | - |
| RS-DFM [110] | Multi-platform Inference Framework | AirCo-MultiTasks [110] | - | - | 3D Region-level, Pixel-level | Generalized Feature Mapping with Relative Depth Estimation | - |
| A2-MAE [122] | ViT-Large | STSSD (Spatial-Temporal-Spectral Structured Dataset) | 0.8 - 30m | Global | Image-level, Pixel-level, Spacial-temporal | Anchor-aware Masking Strategy and Geographic Encoding Module | 304M |
| HyperSIGMA [95] | ViT-based | HyperGlobal-450K [95] | 30m | Global | Image-level, Region-level, Anomaly Detection, Spacial-temporal | MAE | over 1B |
| DOFA [117] | Dynamic OFA | Multiple | 1 to 30 | Global | Image-level, Pixel-level | MIM | 111M/337M |

Fig. 4: The Vision transformer architecture.[3]

*3) Pixel-Level:* Pixel-level analysis offers the most granular form of image perception, assigning a label to every pixel within an image. This includes tasks such as semantic segmentation, where each pixel is classified into categories like vegetation, water, or buildings; it also include change detection, which identifies temporal differences between images captured at different times. Pixel-level analysis is indispensable for creating highly detailed maps used in applications like precision agriculture, deforestation tracking, and disaster management. The ability to analyze fine-grained details enables more accurate assessments and actionable insights for these critical areas.

### B. Backbone

*1) Convolutional Neural Networks (CNNs):* Convolutional Neural Networks [74] are a fundamental architecture in deep learning, designed to extract hierarchical spatial features from images through the use of convolutional layers. Each convolutional layer applies filters to the input data, detecting patterns like edges, textures, and shapes at different levels of abstraction. This makes CNNs well-suited for handling complex visual tasks in remote sensing, such as image classification, segmentation, and object detection.

Residual Neural Networks (ResNet) [36], a type of Convolutional Neural Network (CNN) , address the

degradation problem in deep neural networks by introducing residual connections, which allow gradients to bypass certain layers, facilitating the training of very deep networks. This capability is particularly beneficial in remote sensing, where deep models are often required to capture the intricate details and variations in satellite images. ResNet, as an example, is characterized by their residual blocks, which include shortcut connections that bypass one or more layers. The residual block can be described by the following equation:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, \{W_i\}) + \mathbf{x}$$

where $\mathbf{y}$ is the output, $\mathcal{F}$ represents the residual mapping to be learned, $\mathbf{x}$ is the input, and $\{W_i\}$ are the layer weights [37].

ResNet has various architectures like ResNet-50, ResNet-101, and ResNet-152, with the number indicating the total layers. These networks have shown remarkable performance in various vision tasks due to their ability to train deeper networks without degradation. In remote sensing, ResNets are widely used for image classification, object detection and change detection tasks [31] . For example, ResNet-based models can classify different land cover types [125], [128], detect objects like buildings and vehicles [31], and monitor changes [75], [128] in the landscape over time by comparing temporal sequences of satellite images.

*2) Transformers and Vision Transformers (ViTs):* Transformers, adapted for computer vision (CV) as Vision Transformers (ViT), model long-range dependencies

---

through self-attention, making them effective for complex geospatial data. Figure 4 illustrates the architecture of ViT. ViTs treat images as sequences of patches, capturing global and local patterns, which is useful for segmentation and change detection. The self-attention mechanism computes:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

where $Q$ (query), $K$ (key), and $V$ (value) are the input matrices, and $d_k$ is the dimension of the key vectors [93].

By incorporating these methodologies, foundation models for remote sensing can leverage vast amounts of data, handle complex structures, and achieve state-of-the-art performance across various applications. These methodologies enable models to effectively address the unique challenges of remote sensing, such as large image sizes, diverse data sources, and the need for high accuracy in environmental monitoring and analysis.

In the following sections, we will explore specific applications of these methodologies in different remote sensing tasks, analyze their performance, and discuss the datasets used to train and evaluate these models.

## VI. DATA AND TASKS

### A. Data

Datasets play a crucial role in remote sensing, providing the foundation for training and evaluating models. High-quality datasets enable models to learn accurate representations of the Earth's surface, improving their performance on various remote sensing tasks. In figure 2, we showcase some examples of the data used for training foundation models and their downstream tasks. In this section, we provide an overview of commonly used datasets in table VIII for remote sensing, discussing their characteristics, applications, and relevance to foundation models. These datasets, with their varying resolutions, categories, and geographic coverage, provide a rich resource for advancing remote sensing research and applications. They facilitate the development of robust models capable of addressing diverse challenges in understanding and interpreting Earth's surface through remote sensing technologies.

Datasets used in remote sensing vary significantly in size, from hundreds of thousands of samples, as in RSD46-WHU [62], [116], to over a million, as seen in MillionAID [63], [64]Generally, larger datasets contribute to model generalization by encompassing diverse geographic areas, seasonal variations, and environmental conditions. Dataset resolutions also range from high (sub-meter), suitable for tasks requiring detailed spatial analysis, to moderate (10-60 meters), as with SEN12MS [80] and SSL4EO-S12 [107], which support broader pattern recognition applications.

These datasets leverage various sensor types, including RGB, multispectral, hyperspectral, and synthetic aperture radar (SAR). For instance, SEN12MS [80] integrates both SAR and multispectral imagery, enabling models to learn from distinct data modalities. This diversity in sensor types is critical for robust model development, as each sensor type captures unique surface characteristics, supporting tasks that benefit from cross-modal information.

Foundation models, in particular, benefit from such large-scale, multimodal datasets, which support self-supervised and supervised training approaches across tasks such as scene classification, segmentation, and object detection. For further insight, the appendix includes detailed descriptions of each dataset's structure, unique characteristics, and application roles, enhancing the understanding of their impact on remote sensing advancements.

### B. Tasks

Different applications in remote sensing address particular real-world challenges by leveraging the capabilities of foundation models. These tasks include environmental monitoring, archaeology, agriculture, urban planning and development, and disaster management. To highlight the versatility of foundation models in remote sensing, we present table VI-B that categorizes models based on their applicability to various applications, as well as the different image analysis methods used. This figure serves as a quick reference for researchers to identify suitable models for their specific needs.

*1) Environmental Monitoring:* According to Himeur et al. (2022), environmental monitoring utilizes remote sensing models to observe and track environmental changes, including deforestation, desertification, and pollution. These models play a crucial role in analyzing the effects of human activities and natural phenomena on the environment [39].

*2) Agriculture:* In agriculture, remote sensing models are used to monitor crop health, estimate yields, and manage agricultural practices. According to Kamilaris et al. (2018), these models help optimize resource use and improve agricultural productivity [52].

*3) Archaeology:* In archaeology, remote sensing models have been used to identify and analyze archaeological features and sites. According to Argyrou et al. (2022), these models help detect features such as ruins, artifacts, and ancient structures from satellite imagery, leveraging technologies like Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) to process high-resolution images and capture fine details [5]. Mantovan et al. (2020) also highlight the effectiveness of AI models, particularly CNNs, in locating challenging

| Tasks | | Image Analysis by Levels | | | | Related Work |
|---|---|---|---|---|---|---|
| | | Image-Level | Pixel-Level | Region-Level | Spatial-Temporal | |
| **Environmental Monitoring** | Land Cover Change Detection | | | | ✓ | [2], [40], [66], [67], [69], [70], [87], [91], [96], [109] [24], [32], [49], [59], [95], [99], [106], [111], [122] |
| | Deforestation Monitoring | | ✓ | | | [2], [6], [16], [45], [56], [57], [79], [87], [89], [96], [100] [7], [27], [67], [69], [70], [78], [102], [109], [119] [22], [24], [26], [29], [32], [40], [42], [88], [92], [108] [33], [49], [95], [99], [106], [110], [117], [118], [122], [123] |
| | Water Body Analysis | | ✓ | ✓ | ✓ | [27], [32], [49], [70], [87], [96], [99] |
| | Forest Cover mapping | | ✓ | | ✓ | [27], [40], [49], [67], [69], [87], [96], [106], [109] [59], [95], [122] |
| | Biomass Estimation | | | | | [77] |
| | Weather/Climate Prediction | ✓ | | | | [82], [101], [119] |
| | Cloud Removal | | | | | [33] |
| | Moisture Content Measurement | | | | | [92] |
| **Agriculture** | Crop Type Mapping | ✓ | ✓ | ✓ | ✓ | [27], [32], [49], [70], [83], [87], [95], [96], [99] |
| | Weed Detection | | ✓ | | | [6], [27], [56], [57], [70], [87], [89], [96], [100], [124] [8], [22], [32], [49], [95], [99], [110], [123] |
| | Disease Monitoring | ✓ | ✓ | | ✓ | [6], [27], [56], [57], [70], [87], [89], [96], [100] [22], [32], [49], [95], [99], [110], [123] |
| | Forecasting | | | | | [82], [119] |
| | Soil Parameter Estimation | | | | | [117] |
| | Yield Estimation | ✓ | ✓ | | | [2], [6], [16], [45], [56], [57], [79], [87], [89], [96], [100] [7], [27], [67], [69], [70], [78], [102], [109] [24], [26], [32], [33], [40], [42], [88], [108], [118] [95], [99], [106], [117], [122], [123] |
| | Agricultural Pattern Segmentation | | ✓ | | | [66] |
| **Archaeology** | Artifact Classification and Recognition | ✓ | | ✓ | | [2], [6], [56], [66], [79], [84], [87], [96], [105] [7], [16], [67], [69], [70], [77], [78], [89], [100], [109], [124] [26], [27], [43], [44], [58], [82], [88], [91], [108] [8], [24], [32], [33], [72], [99], [106], [111], [118], [123] [40], [45], [57], [59], [95], [102], [117], [122] |
| | Detection of Archaeological Structures | | | ✓ | | [6], [27], [56], [57], [70], [87], [89], [96], [100], [124] [8], [22], [32], [49], [95], [99], [110], [123] |
| | Semantic Segmentation | | ✓ | | | [2], [6], [16], [45], [56], [57], [79], [87], [89], [96], [100] [7], [27], [67], [69], [70], [78], [102], [109], [119] [22], [24], [26], [28], [32], [40], [42], [88], [92], [108] [33], [49], [95], [99], [106], [110], [117], [118], [122], [123] |
| | Texture/Structural Analysis | | | | | [2] |
| | Pattern Recognition | | ✓ | ✓ | | [6], [27], [56], [57], [70], [87], [89], [96], [100] [22], [32], [49], [95], [99], [110], [123] |
| **Urban Planning & Development** | Traffic Monitoring | ✓ | | ✓ | ✓ | [27], [32], [49], [70], [87], [95], [96], [99] |
| | Land Cover/Use Classification | ✓ | | | | [3], [6], [56], [66], [79], [84], [87], [96], [105] [7], [16], [67], [69], [70], [77], [78], [89], [100], [109], [124] [26], [27], [42], [44], [58], [87], [88], [91], [108] [8], [24], [33], [53], [72], [99], [106], [111], [118], [123] [49], [45], [57], [59], [83], [95], [102], [117], [122] |
| | Road Crack Detection | | | ✓ | | [6], [27], [56], [57], [70], [87], [89], [96], [100], [124] [8], [22], [32], [49], [95], [99], [110], [123] |
| | Air Quality Monitoring | ✓ | | ✓ | | [6], [27], [33], [56], [57], [70], [87], [89], [96], [100] [22], [32], [49], [95], [99], [110], [123] |
| | Building Extraction | ✓ | ✓ | | | [27] |
| | Object/Video Tracking | | | ✓ | ✓ | [119] |
| | Infrastructure Monitoring | | | | | [44], [119] |
| **Disaster Management** | Landslide Risk Monitoring | ✓ | ✓ | ✓ | ✓ | [27], [32], [49], [70], [87], [95], [96], [99] |
| | Disaster Response | | | | | [117] |
| | Real-Time Detection and Mapping | | ✓ | ✓ | ✓ | [27], [32], [49], [70], [87], [95], [96], [99] |
| | Building Damage Assessment | ✓ | ✓ | ✓ | | [27], [32], [49], [70], [87], [95], [96], [99] |
| | Critical Infrastructure Detection | ✓ | | ✓ | | [6], [27], [56], [57], [70], [87], [89], [96], [100], [124] [8], [22], [32], [49], [95], [99], [110], [123] |
| | Flood/Fire Mapping and Prediction | ✓ | ✓ | | ✓ | [27], [40], [67], [69], [87], [96], [109] [49], [95], [106], [122] |
| | Crowd and Vehicle Detection | | | ✓ | ✓ | [2], [40], [66], [67], [69], [70], [87], [91], [96], [109] [24], [32], [49], [59], [95], [99], [106], [111], [122] |

TABLE III: This diagram illustrates various tasks in different applications for remote sensing. Key areas include Environmental Monitoring, Agriculture, Urban Planning and Development, Disaster Management, and Archaeology. Each domain comprises specific tasks in different image analysis levels like image-level, pixel-level, region-level, and spatial-temporal. The relationships between these tasks and their applications are depicted through checkmarks, emphasizing the interconnected nature of image analysis methods across different fields.

terrestrial archaeological sites and processing multispectral data [68].

*4) Urban Planning and Development:* In urban planning and development, remote sensing models are used to monitor and analyze urban expansion, infrastructure development, and land use changes. According to Jha et al. (2021), these models play a critical role in managing urban growth, planning new developments, and assessing the impact of urbanization by providing essential data for smart city planning and sustainable development [48].

*5) Disaster Management:* Remote sensing models play a crucial role in disaster management by providing timely information on affected areas. According to Abid et al. (2021), these models are used to detect and assess damage from natural disasters like earthquakes, hurricanes, and floods, enabling rapid response and recovery efforts [1].

## VII. DISCUSSION

The rapid advancement in foundation models for remote sensing underscore their transformative potential across various applications. As the field continues to evolve, it is crucial to synthesize the findings, address technical challenges, understand practical implications, and identify future research directions. In this section, we make a comprehensive analysis of these aspects, aiming to offer insights and guidance for future development and application of remote sensing foundation models.

### A. Synthesis of Findings

In our survey of foundation models for remote sensing, we identified significant advancements and trends that highlights the evolving capabilities and applications of these models. The performance metrics of various models across different downstream tasks such as scene classification, semantic segmentation, object detection, and change detection reveal the following key findings:

*1) Model Performance:* In this section, we present the performance metrics of recent foundation models in remote sensing based on results reported in the original papers. **All performance numbers mentioned here are sourced directly from the original studies to ensure accuracy and consistency in evaluating these models.** These metrics provide insights into the models' effectiveness across tasks like semantic segmentation, object detection, and change detection, highlighting their strengths and limitations under different experimental setups.

- **Image-Level** The performance of foundation models on the BigEarthNet dataset [85] for classification tasks shows variations in accuracy, as presented in table IV. Overall, msGFM [33] has the top performance of 92.90% (mAP), followed closely

by SkySense [32] with a performance of 92.09%. Other notable performers include DeCUR [102], which achieved an mAP of 89.70%, and DINO-MC [111], with an mAP of 88.75%. SeCo [66] also demonstrated strong performance with an mAP of 87.81%, while DINO-MM [105] reached an mAP of 87.10%. On the other hand, models like CACo [67] and FoMo-Bench [8] have mAP of 74.98% and F1-Score of 68.33% respectively, showing competitiveness but room for improvement.

The high mAP scores of msGFM [33] and SkySense [32] highlight their efficiency in classifcation tasks, making them suitable for applications requiring high accuracy. Other foundation models, such as DINO-MM [105] and DeCUR [102], also provide strong performance with potential for further optimization. The variety in performance metrics underscores the evolving capabilities and specialization of foundation models in handling complex classification tasks within datasets like BigEarthNet [85].

The classification advancements observed in remote sensing models stem from sophisticated pretraining techniques that capture both spatial and spectral complexity across vast datasets. SkySense, for example, shows an average improvement of 2.76% over recent models by implementing multigranularity contrastive learning on a diverse dataset of 21.5 million optical and SAR sequences [32]. This approach enables SkySense to learn nuanced spatial and temporal relationships across modalities, enhancing generalization in varied environmental conditions. Such multi-granular representation proves crucial in remote sensing, where scene classification often depends on subtle spectral differences that simpler models may overlook. Likewise, HyperSIGMA [95], pre-trained on the expansive HyperGlobal-450K hyperspectral dataset [95], leverages its sparse sampling attention mechanism to optimize spectral-spatial feature extraction in high-dimensional hyperspectral data. By selectively focusing on critical spectral bands and reducing redundancy, HyperSIGMA achieves high classification accuracy across hyperspectral scenes, a marked improvement over previous models that struggled with hyperspectral data complexity.

These models highlight the importance of designing pre-training strategies that capture multi-modal features and effectively utilize dataset diversity, as these elements directly impact the robustness and accuracy of classification in remote sensing applications.

- **Pixel-Level** For the segmentation tasks, we compared 12 foundation models which have been tested

| Dataset | Model | Performance (%) | Metrics |
|---|---|---|---|
| BigEarthNet [85] | SeCo [66] | 87.81 | mAP |
| | CMC-RSSR [84] | 82.90 | mAP |
| | DINO-MM [105] | 87.10 | mAP |
| | CACo [67] | 74.98 | mAP |
| | GFM [69] | 86.30 | mAP |
| | DINO-MC [111] | 88.75 | mAP |
| | CROMA [28] | 86.46 | mAP |
| | DeCUR [102] | 89.70 | mAP |
| | CtxMIM [123] | 86.88 | mAP |
| | FG-MAE [108] | 78.00 | mAP |
| | USat [44] | 85.82 | mAP |
| | FoMo-Bench [8] | 69.33 | F1 Score |
| | SwiMDiff [91] | 81.10 | mAP |
| | SpectralGPT [40] | 88.22 | mAP |
| | SatMAE++ [73] | 85.11 | mAP |
| | msGFM [33] | **92.90** | mAP |
| | SkySense [32] | 92.09 | mAP |
| | MMEarth [72] | 78.6 | mAP |
| | **Shallow CNN\* [85]** | **70.98** | **F1-Score** |

TABLE IV: This table provides an overview of the performance metrics for various models applied to the BigEarthNet dataset [85] for **image-level** tasks. The performance is measured using mAP (Mean Average Precision) and F1 Score. *Performance for the shallow CNN model are sourced from the original BigEarthNet [85] paper.

| Dataset | Model | Performance (%) | Metrics |
|---|---|---|---|
| ISPRS Potsdam | GeoKR [56] | 70.48 | mIoU |
| | RSP [96] | 65.30 | mIoU |
| | RingMo [87] | 91.74 | OA |
| | RVSA [100] | 91.22 | OA |
| | TOV [89] | 60.34 | mIoU |
| | CMID [70] | **87.04** | mIoU |
| | RingMo-lite [109] | 90.96 | OA |
| | Cross-Scale MAE [88] | 76.17 | mIoU |
| | SMLFR [22] | **91.82** | OA |
| | SkySense [32] | **93.99** | mF1 |
| | UPetu [24] | 83.17 | mIoU |
| | BFM [11] | 92.58 | OA |
| | **R-SegNet\* [127]** | **91.37** | **OA** |

TABLE V: This table provides an overview of the performance metrics for various models applied to the ISPRS Potsdam [43] dataset for **pixel-level** tasks. The performance is measured using Mean Intersection over Union (mIoU) and Overall Accuracy (OA). *Non-FM

on the ISPRS Potsdam dataset. As shown in table V, SkySense [32] has the better performance out of all 12 models, with a mF1 Score of 93.99%. CMID [70] stands out with the highest mIoU of 87.04%, demonstrating its superior capability in accurately segmenting different regions within the dataset. For Overall Accuracy performance, BFM [11] has the highest OA score of 91.82%. Cross-Scale MAE [88], UPetu [24] and RSP [96] have mIoU scores of 76.17%, 83.17% and 65.30% respectively, showing competitive segmentation capabilities. GeoKR [56] reaches an mIoU of 70.48%, indicating robust segmentation performance but with room for improvement compared to CMID [70]. TOV scores the lowest mIoU at 60.34%, suggesting it may struggle with finer segmentation tasks compared to the other models.

The performance metrics for the models applied to the ISPRS Potsdam dataset reveal significant variations in their effectiveness in segmentation tasks. SkySense [32] and CMID [70] emerge as top performers in mF1 score and mIoU, respectively, while SMLFR [22], RingMo [87], and RingMo-lite [109] demonstrate strong overall accuracy. These insights can guide the selection and optimization of models for specific remote sensing applications, ensuring the best possible performance for the task at hand.

For the change detection tasks, we compared the performance of foundation models on the OSCD and LEVIR-CD datasets. The models were evaluated based on their F1 Scores, which provide a balanced measure of precision and recall. As shown in the table, the performance varies significantly across different models and datasets.

SkySense [32] achieves the highest F1 Score of 60.06% on the OSCD dataset, demonstrating its superior ability to accurately detect changes. GFM [69] follows with an F1 Score of 59.82%, indicating strong performance in change detection tasks. SpectralGPT [40] also performs well with an F1 Score of 54.29%. Other notable models include DINO-MC [111] with an F1 Score of 52.71% and CACo [67] with an F1 Score of 52.11%. SeCo [66] records the lowest F1 Score at 46.94%, suggesting it may require further optimization to enhance its change detection capabilities.

In contrast, the LEVIR-CD dataset reveals higher performance metrics across the models. MTP [99] achieves the highest F1 Score of 92.67% and Sky-Sense [32] follows closely with an F1 Score of 92.58%, demonstrating their robust performance. SWiMDiff reaches a lower F1 Score of 80.90% compared to its peers but still indicates effective performance in the LEVIR-CD [12] dataset.

- **Region-Level** In table VI, the performance of foundation models on the DOTA, DIOR, and DIOR-R datasets for object detection are evaluated based on their Mean Average Precision (mAP) and Average Precision at 50% (AP50).
  On the DOTA dataset, RVSA [100] achieves the highest mAP of 81.24% in accurately detecting objects, followed by SMLFR [22] and RSP [96] with mAP of 79.33% and 77.72%. CMID [70],

GeRSP [42], and BFM [11] also demonstrate moderate performance with mAPs of 72.12%, 67.40% and 58.69%. For DIOR and DIOR-R dataset, MTP [99] and SkySense [32] are the top performers with an AP50 of 78% and an mAP of 78.73% respectively, showcasing their superior object detection capabilities. These insights can guide the selection and optimization of models to ensure the best possible performance for specific remote sensing applications.

- **Influence of Pre-training Methods.** Various pre-training methods have a substantial impact on the performance of foundation models in remote sensing. Models pre-trained using SSL techniques, such as contrastive learning (CL) and masked autoencoders (MAE), consistently exhibit superior performance compared to those pre-trained with traditional supervised learning. For instance, Sky-Sense which uses a multi-granularity contrastive learning approach, outperforms other models by approximately 3.6% in scene classification and object detection tasks [32]. Similarly, Seco, based on seasonal contrast learning, yields superior performance for land-cover classification, improving metrics by up to 7% over ImageNet-pre-trained models [66] . In handling multi-temporal and multi-spectral data, models like SatMAE [16] and Scale-MAE [78], using masked autoencoding, achieve improvements in change detection, with SatMAE showing up to a 14% performance gain in land cover classification [16] and Scale-MAE offering a 1.7% mIoU improvement for segmentation across varied resolutions [78]. These findings highlight the critical role of innovative pre-training methods in maximizing the effectiveness of foundation models and suggest that continued exploration and refinement of these techniques are essential for advancing remote sensing capabilities.

Foundation models like SatMAE, RingMo, A2-MAE, and ORBIT each demonstrate strong performance, but practical trade-offs are essential to consider, especially for application-specific constraints [16], [87], [101], [122]. SatMAE, based on a transformer architecture, effectively leverages temporal and multi-spectral embeddings to capture complex spatiotemporal patterns in satellite imagery. This strength, however, comes at the cost of significant computational requirements, which may not be feasible for real-time monitoring applications in resource-constrained environments.

In contrast, RingMo provides a more lightweight vision transformer architecture, offering efficient model inference and a balance between performance and computational demands. This makes RingMo particularly suitable for rapid-inference tasks like disaster response monitoring, where real-time processing is critical. [87] A2-MAE introduces an anchor-aware masking strategy, optimizing spatial-temporal-spectral representations and allowing effective integration of multi-source data. This design enhances its adaptability to varied data resolutions and modalities, yet the model's complex encoding techniques add to its computational load, suggesting a fit for applications that require high accuracy over efficiency [122].

Finally, ORBIT, designed with 113 billion parameters, is exceptionally scalable, achieving high-throughput performance for Earth system predictability tasks. While it excels in large-scale predictive tasks, the model's considerable resource requirements limit its deployment to specialized high-performance computing environments. [101] These trade-offs highlight the importance of selecting a model that aligns with specific operational goals, whether for maximizing accuracy or minimizing computational overhead.

Furthermore, recent studies comparing SSL approaches highlight the distinct advantages of generative methods like Masked Autoencoders (MAE) over contrastive methods for time-series data, especially when labeled data is limited [61]. Unlike contrastive approaches that emphasize distinguishing between similar and dissimilar pairs, generative methods such as MAE reconstruct data from masked segments, allowing them to capture complex underlying structures and relationships within the data. This reconstruction-based learning proves particularly advantageous for time-series and multi-spectral applications in remote sensing, where temporal and spectral dependencies are essential. Consequently, MAE-based models can achieve stronger representations under sparse labeling conditions, positioning them as powerful tools for remote sensing tasks that require nuanced temporal analysis.

*2) Practical Implications:* Foundation models offer transformative capabilities in remote sensing by building upon established applications like multi-spectral and time-series data analysis. While these applications have traditionally relied on machine learning and deep learning, foundation models reduce the need for labeled data and enable rapid adaptation to new tasks, providing robust solutions in areas previously limited by data constraints and task-specific architectures. Consequently, the advancements in foundation models have significant practical implications across various areas:

- **Environmental Monitoring.** Models like GASSL [6] and SatMAE [16] offer detailed assessments of environmental changes, aiding in conservation

| Dataset | Model | Performance (%) | Metrics |
|---|---|---|---|
| DOTA | RSP [96] | 77.72 | mAP |
| | RVSA [100] | **81.24** | mAP |
| | TOV [89] | 26.10 | mAP50 |
| | CMID [70] | 72.12 | mAP |
| | GeRSP [42] | 67.40 | mAP |
| | SMLFR [22] | 79.33 | mAP |
| | BFM [11] | 58.69 | mAP |
| | **YOLOv2-D\* [21]** | **60.51** | **AP** |
| DIOR | RingMo [87] | 75.80 | mAP |
| | CSPT [124] | 69.80 | mAP |
| | RingMo-lite [109] | 73.40 | mAP |
| | GeRSP [42] | 72.20 | mAP |
| | MTP [99] | **78.00** | AP50 |
| | **Faster R-CNN\* [55]** | **74.05** | **mAP** |
| DIOR-R | RVSA [100] | 71.05 | mAP |
| | SMLFR [22] | 72.33 | mAP |
| | SkySense [32] | **78.73** | mAP |
| | MTP [99] | 74.54 | mAP |
| | BFM [11] | 73.62 | mAP |
| | **AOPG\* [14]** | **64.41** | **mAP** |

TABLE VI: This table provides an overview of the performance metrics for various models applied to the DOTA [20], [21], [113] dataset, DIOR [55] and DIOR-R [14] dataset for **region-level** task. The performance is mainly measured using Mean Average Precision (mAP). *Model performance are aquired from original dataset papers. AOPG is Anchor-free Oriented Proposal Generator.

| Dataset | Model | F1 Score |
|---|---|---|
| OSCD [10] | SeCo [66] | 46.94 |
| | MATTER [2] | 49.48 |
| | CACo [67] | 52.11 |
| | GFM [69] | 59.82 |
| | SWiMDiff [91] | 49.60 |
| | SpectralGPT [40] | 54.29 |
| | SkySense [32] | **60.06** |
| | DINO-MC [111] | 52.71 |
| | HyperSIGMA [95] | 59.28 |
| | MTP [99] | 53.36 |
| | **CNNs\* [10]** | **89.66 (OA)** |
| LEVIR-CD [12] | RSP [96] | 90.93 |
| | RingMo [87] | 91.86 |
| | RIngMo-lite [109] | 91.56 |
| | SwiMDiff [91] | 80.90 |
| | SkySense [32] | 92.58 |
| | UPetu [24] | 88.50 |
| | **STANet\* [12]** | **85.4** |

TABLE VII: This table provides an overview of the F1 Score for various models applied to the Onera Satellite Change Detection (OSCD) dataset [10] and the LEVIR-CD dataset [12] for **spacial-temporal** downstream tasks. *Performance for the models are sourced from the original dataset papers. STANet is Spatial-Temporal Attention Network.

efforts and policy-making. These models excel in monitoring deforestation, desertification, and pollution levels, providing actionable insights for environmental management. By integrating multispectral and temporal data, these models can track changes over time, allowing for early detection of environmental degradation and the formulation of timely interventions. This capability is particularly important for the sustainable management of natural resources, as well as reducing the impacts of climate change.

- **Agriculture and Forestry.** Foundation models such as EarthPT [82] and GeCo [57] delivers valuable insights into crop health, yield predictions, and land use management, optimizing agricultural practices and resource allocation. For instance, RSP [96], leveraging multi-spectral data, enhances precision agriculture by accurately monitoring crop conditions and predicting yields. These models can detect early signs of crop stress, diseases, and pest infestations, enabling farmers to take proactive measures. Additionally, they aid in forestry management by providing detailed maps of forest cover, biomass estimation, and monitoring deforestation activities, thereby supporting conservation efforts and sustainable forestry practices.

- **Archaeology.** The use of foundation models in archaeology revolutionizes the way archaeological features and sites are discovered, mapped, and analyzed. Models, such as GeoKR [56], RingMo [87], etc., can process high-resolution satellite imagery and multi-spectral data to enhance the detection and mapping of archaeological features that might be difficult to discern with the naked eye. Others like MATTER [2] can accomplish texture and material analysis to help identify various surface. They enable large-scale surveys, allowing archaeologists to identify potential sites of interest over vast areas efficiently. Although thorough exploration still requires on-site visits and excavations or other terrestrial investigations, these significantly improve the initial identification and mapping process. Additionally, these models can track changes over time, helping archaeologists monitor environmental and human impacts, providing crucial information for preservation and restoration. This enhances the efficiency and accuracy of surveys and opens new possibilities for discovering unknown sites.

- **Urban Planning and Development.** Remote sensing models like CMID [70] and SkySense [32] are pivotal for monitoring urban expansion, infrastructure development, and land use changes. These models facilitate sustainable urban growth and development planning by providing high-resolution

data analysis and trend forecasting. They enable city planners to assess the impact of urbanization on natural habitats, optimize land use, and plan infrastructure projects more effectively.

- **Disaster Management.** Models such as OFA-Net [118], DOFA [117] and Prithvi [46] are instrumental in flood mapping, as well as fire detection. These models provide critical real-time data that helps in identifying affected areas quickly, enabling timely and effective response measures. This capability supports emergency responders in prioritizing resource allocation and implementing evacuation plans, thereby reducing the impact of natural disasters. Additionally, these models assist in post-disaster recovery by assessing damage and monitoring the recovery process over time. By integrating various data sources, they enhance the ability to make informed decisions, coordinate response efforts, and plan for future disaster mitigation strategies.

The improvements in accuracy across the models discussed have profound implications for real-world remote sensing applications. In deforestation monitoring, for instance, models like GFM achieve high pixel-level accuracy in semantic segmentation, showing up to a 4.5% improvement over baseline models, which enhances the precision of mapping forest cover changes, supporting conservation efforts [101]. Similarly, HyperSIGMA achieves an impressive 6.2% accuracy boost in hyperspectral vegetation monitoring, providing invaluable data for assessing forest health and biodiversity [95].

In urban planning, models like UPetu excel in infrastructure mapping by integrating multi-modal data, such as optical and radar imagery, achieving over 5% higher accuracy compared to single-modality models, which allows urban planners to make more informed land-use decisions [24]. Additionally, RingMo enhances object detection accuracy by 3.7% over traditional supervised models, effectively identifying dense urban features critical for disaster management and urban infrastructure assessment [87].

Finally, ORBIT demonstrates exceptional scalability, processing large climate datasets with a scaling efficiency of up to 85%, which supports applications in long-term environmental monitoring, such as climate change prediction and seasonal forecasting. This scalability not only advances traditional remote sensing workflows but also enables complex multi-temporal analyses and predictive modeling previously challenging with conventional methods [101].

While remote sensing has long benefited from multispectral and temporal data, the adaptability, scalability, and efficiency of foundation models unlock a new level of precision and accessibility in these applications. This advancement opens up opportunities to tackle complex and evolving challenges across domains—from environmental conservation to urban planning—that traditional models have struggled to address at scale.

### B. Future Direction

Future research should prioritize several key areas:

- **Efficient Model Development:** Exploring techniques such as model distillation, pruning, and quantization to reduce computational requirements without compromising performance is crucial. Additionally, developing scalable architectures that efficiently handle ultra-high-resolution images is essential. For instance, applying pruning techniques to models like SatMAE [16] could maintain performance while reducing computational load. Model adaptation techniques such as LoRA (Low-Rank Adaptation) [41] have emerged as effective methods for fine-tuning large-scale models with minimal computational overhead. By decomposing weight updates into low-rank matrices, LoRA [41] enables efficient adaptation without the need to modify the entire set of model parameters, making it suitable for resource-constrained environments or when frequent re-training is required. Incorporating methods like LoRA [41] can further enhance the applicability of foundation models across diverse tasks and domains.
- **Multi-Modal Data Integration:** Enhancing methods for integrating and processing multi-modal data (e.g., combining optical and radar imagery) will provide more comprehensive insights. Research on advanced SSL techniques capable of leveraging multi-modal data is necessary. The OFA-Net [118] framework, which integrates multi-model data, serves as a promising direction for future models to emulate and improve upon.
- **Interdisciplinary Collaboration:** Promoting collaboration between remote sensing experts, AI researchers, and domain specialists can address complex challenges and drive innovation. For example, partnerships between AI researchers and environmental scientists can refine models like GASSL [6] for better environmental monitoring and conservation efforts.

Looking ahead, the consistent success of self-supervised learning methods in foundation models marks an exciting frontier for future research. These models' ability to learn from unlabeled data and adapt to diverse remote sensing tasks with minimal fine-tuning suggests that advancements in unsupervised learning techniques could greatly reduce reliance on large labeled datasets, which remain a significant bottleneck in many remote sensing

applications. However, as these models grow in size and complexity, balancing computational demands with the need for efficiency will become increasingly crucial. Future work may focus on developing more resource-efficient versions of foundation models that maintain high performance, particularly for deployment in real-time monitoring systems or environments with limited computational resources.

*C. Limitation*

This survey has several limitations:

- **Scope and Coverage.** The review focuses on foundation models released between June 2021 and June 2024. While the scope of this review is extensive and covers many significant developments, it is not exhaustive. Some recent advancements and innovations in the field may not be included due to their release timing or the lack of sufficient evaluation metrics at the time of writing. Consequently, certain cutting-edge models that have emerged in the latter part of this period or that have not yet been thoroughly evaluated might be omitted. This limitation underscores the need for readers to seek out the most current research and updates beyond the scope of this survey. Additionally, while foundation models have been empirically tested on a specific set of downstream applications, their robust architectures and general-purpose training paradigms, such as convolutional networks (e.g., ResNet) and vision transformers (e.g., ViT), indicate their potential to perform well across a much broader range of tasks. The limited testing observed in current literature should not be seen as a constraint on their applicability, but rather as an indication of the focus of existing research efforts. Given their design, these models are expected to generalize effectively to a wide variety of remote sensing tasks, even beyond those explicitly tested. Future work should aim to explore and validate their performance across more diverse applications to unlock their full potential.
- **Evolving Field.** The field of AI and remote sensing is rapidly evolving, with continuous advancements and breakthroughs occurring at a fast pace. This dynamic nature necessitates ongoing reviews and updates to ensure the relevance and comprehensiveness of the survey. New techniques, methodologies, and models are constantly being developed, which can significantly impact the state of the art. Therefore, it is essential to recognize that this survey represents a snapshot in time and that continuous monitoring of the literature is required to capture the latest advancements and emerging trends. This approach will help maintain an up-to-date under-standing of the field and incorporate new findings as they become available.

## VIII. CONCLUSION

In this comprehensive survey, we have reviewed the recent advancements in foundation models for remote sensing. We categorized these models based on their pretraining methods, image analysis techniques, and applications across different areas, highlighting their unique methodologies and capabilities.

Our analysis covered various advanced techniques, including self-supervised learning, vision transformers, and residual neural networks. These models have significantly improved performance on different image perception levels like region-level, pixel level, and image-level, as well as in applications like environmental monitoring, digital archaeology, agriculture, urban planning, and disaster management.

While significant progress has been made, several challenges persist, such as the need for more diverse and high-quality datasets, high computational requirements, and difficulties for different applications. Addressing these challenges will require further research and collaboration across disciplines.

In summary, this survey provides a detailed overview of the current state of foundation models in remote sensing, offering valuable insights and identifying future research directions. We recommend continued efforts in developing efficient model architectures, enhancing multi-modal data integration, and expanding dataset diversity to fully realize the potential of these models in remote sensing.

## APPENDIX
### COMMONLY USED PRE-TRAIN DATASET FOR REMOTE SENSING

RSD46-WHU [62], [116] dataset, introduced in 2017, is sourced from Google Earth and Tianditu. It contains

117,000 images with a patch size of 256 pixels and spatial resolutions ranging from 0.5 to 2 meters per pixel. Covering 46 categories globally, this dataset is primarily used for scene classification. Similarly, the Functional Map of the World (fMoW) [15], released in April 2018, comprises over 1 million images from Digital Globe. Spanning 63 categories across 207 countries, it includes multispectral images used for both scene classification and object detection.

In May 2019, the DOTA [114] dataset was proposed, known for its large-scale aerial image object detection capabilities. It includes 11,268 images of various resolutions from Google Earth, GF-2 Satellite, and aerial sources, covering 18 categories globally. Another significant dataset, SEN12MS [80], released in June 2019, contains 541,986 images from Sentinel-1, Sentinel-2, and MODIS Land Cover. With a patch size of 256x256 pixels, it supports land cover classification and change detection tasks.

BigEarthNet [85], also from June 2019, consists of 590,326 images with varying sizes from 20x20 to 120x120 pixels, sourced from Sentinel-2. It covers 43 categories across Europe and is used for scene classification and object detection. The SeCo [66] dataset, another June 2019 release, contains approximately 1 million images with a resolution of 2.65x2.65km from Sentinel-2. It is designed for seasonal change detection and land cover classification over seasons.

The MillionAID dataset [63], [64], introduced in March 2021, includes over 1 million images of various sizes from Google Earth. Covering 51 categories globally, it is used for scene classification. Levir-KR, released in July 2021, contains 1,431,950 images from Gaofen-1, Gaofen-2, and Gaofen-6 satellites, supporting change detection and scene classification applications.

SoundingEarth [38], introduced in August 2021, comprises 50,545 images of 1024-pixel size from Google Earth, combining RGB and audio data for remote sensing. The TOV-RS-Balanced dataset [90] from April 2022 includes 500,000 images with a 600-pixel size from Google Earth, covering 31 categories globally, and is used for scene classification, object detection, and semantic segmentation.

SeasoNet [53], released in July 2022, features 1,759,830 images from Sentinel-2 with patch size from 20 to 120 pixels, supporting seasonal scene classification, segmentation, and retrieval over Germany. Lastly, the SSL4EO-S12 dataset [107] from November 2022 contains over 3 million images from Sentinel-1 and Sentinel-2, with patch size of 264x264 pixels. Since this dataset dose not contain any labels, it is commonly used for self-supervised learning.

In recent years, additional datasets have further enriched the resources available for remote sensing research. The SAMRS dataset [98], released in October 2023, offers a high-resolution collection of images sourced from datasets like HRSC2016 and FAIR1M-2.0, tailored for advanced segmentation tasks. With over 105,000 images and resolutions up to 1024x1024 pixels, SAMRS supports semantic and instance segmentation as well as object detection, contributing to the development of scalable segmentation models for remote sensing.

Focusing on change-aware learning, CACo [67], launched in June 2023, provides a variable patch-size dataset sourced from Sentinel-2. This dataset is optimized for change detection and contrastive learning, specifically addressing urban and rural landscapes. By prioritizing contrastive and self-supervised tasks, CACo aids in developing models that can adapt to changes in satellite imagery across various environments.

The SatlasPretrain [7] dataset, introduced in October 2023, is a large-scale collection with over 856,000 images combining Sentinel-2 and NAIP high-resolution sources. With multispectral and high-resolution imagery, SatlasPretrain supports applications such as land cover classification, segmentation, and change detection, further advancing research in high-resolution satellite image analysis.

The SSL4EO-L [83] dataset, released in October 2023, represents a vast resource with over 5 million images from Landsat, designed for self-supervised learning in cloud detection and land cover classification. By focusing on multi-year Landsat imagery, SSL4EO-L enables robust training for applications that benefit from long-term temporal coverage and cloud-resilient classification.

Finally, MMEarth [72], introduced in July 2024, combines data from Sentinel-1, Sentinel-2, and Aster DEM, providing over 1.2 million images for multimodal applications. This dataset supports land cover classification and semantic segmentation, enabling researchers to leverage multiple sensor types and climate data for improved geospatial representation learning.

These datasets, with their varying resolutions, categories, and geographic coverage, provide a rich resource for advancing remote sensing research and applications. They facilitate the development of robust models capable of addressing diverse challenges in understanding and interpreting Earth's surface through remote sensing technologies.

## REFERENCES

[1] Sheikh Kamran Abid, Noralfishah Sulaiman, Shiau Wei Chan, Umber Nazir, Muhammad Abid, Heesup Han, Antonio Ariza-Montes, and Alejandro Vega-Muñoz. Toward an integrated disaster management approach: How artificial intelligence can boost disaster management. *Sustainability*, 13(22), 2021. 3, 14

[2] Peri Akiva, Matthew Purri, and Matthew Leotta. Self-supervised material and texture representation learning for remote sensing tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8203–8215, June 2022. 4, 8, 13, 17

| Month, Year | Dataset | Title | Patch Size | Size | Resolution (m) | Sensor | Categories | Geographic Coverage | Image Type | Application |
|---|---|---|---|---|---|---|---|---|---|---|
| 2017 | RSD46-WHU [62], [116] | - | 256 x 256 | 117,000 | 0.5 - 2 | Google Earth, Tianditu | 46 | Global | RGB | Scene Classification |
| Apr, 2018 | fMoW [15] | Functional Map of the World | - | 1,047,691 | - | Digital Globe | 63 | 207 of 247 countries | Multispectral | Scene Classification, Object Detection |
| May, 2019 | DOTA [114] | DOTA: A Large-scale Dataset for Object Detection in Aerial Images | 800 x 800 to 20,000 x 20,000 | 11,268 | Various | Google Earth, GF-2 Satellite, and aerial images | 18 | Global | RGB | Object Detection |
| Jun, 2019 | SEN12MS [80] | SEN12MS – A Curated Dataset of Georeferenced Multi-Spectral Sentinel-1/2 Imagery for Deep Learning and Data Fusion | 256 x 256 | 541,986 | 10 | Sentinel-1, Sentinel-2, MODIS Land Cover | 33 | Globally distributed | SAR/Multispectral | Land Cover Classification, Change Detection |
| Jun, 2019 | BigEarthNet [85] | BigEarthNet: A Large-Scale Benchmark Archive For Remote Sensing Image Understanding | 20 x 20 to 120 x 120 | 590,326 | Various | Sentinel-2 | 43 | Europe | Multispectral | Scene Classification, Object Detection |
| Jun, 2019 | SeCo [66] | Seasonal Contrast: Unsupervised Pre-Training from Uncurated Remote Sensing Data | 264 x 264 | ~1M | 10 - 60 | Sentinel-2 | - | Global | Multispectral | Seasonal Change Detection, Land Cover Classification over Seasons |
| Mar, 2021 | MillionAID [63], [64] | Million-AID | 110 - 31,672 | 1,000,848 | Various | Google Earth | 51 | Global | RGB | Scene Classification |
| Jul, 2021 | Levir-KR [56] | Geographical Knowledge-driven Representation Learning for Remote Sensing Images | - | 1,431,950 | Various | Gaofen1, Gaofen-2, Gaofen-6 | 8 | Global | Multispectral | Change Detection, Scene Classification |
| Apr, 2022 | TOV-RS-Balanced [90] | TOV: The Original Vision Model for Optical Remote Sensing Image Understanding via Self-supervised Learning | 600 x 600 | 500,000 | 1 - 20 | Google Earth | 31 | Global | RGB | Scene Classification, Object Detection, Semantic Segmentation |
| Jul, 2022 | SeasoNet [53] | SeasoNet: A Seasonal Scene Classification, Segmentation and Retrieval dataset for satellite Imagery over Germany | up to 120 x 120 | 1,759,830 | 10 - 60 | Sentinel-2 | 33 | Germany | Multispectral | Scene Classification, Scene Segmentation |
| Nov, 2022 | SSL4EO-S12 [107] | SSL4EO-S12: A Large-Scale Multi-Modal, Multi-Temporal Dataset for Self-Supervised Learning in Earth Observation | 264 x 264 | 3,012,948 | 10 - 60 | Sentinel-1, Sentinel-2 | - | Global | SAR/Multispectral | Self-Supervised Learning |
| Oct, 2023 | SAMRS [98] | SAMRS: Scaling-up Remote Sensing Segmentation Dataset with Segment Anything Model | 600 x 600 to 1024 x 1024 | 105,090 | Various | HRSC2016, DOTA-V2.0, DIOR, FAIR1M-2.0 | 37 | Global | High-resolution | Semantic Segmentation, Instance Segmentation, Object Detection |
| Jun, 2023 | CACo [67] | Change-Aware Sampling and Contrastive Learning for Satellite Images | Variable | - | 10 | Sentinel-2 | - | Urban and Rural Areas | Multispectral | Semantic Segmentation, Change Detection, Self-Supervised Learning |
| Oct, 2023 | SatlasPretrain [7] | SatlasPretrain: A Large-Scale Dataset for Remote Sensing Image Understanding | 512 x 512 | 856,000 | 1 - 10 (Sentinel-2), 0.5 - 2 (NAIP) | Sentinel-1, Sentinel-2, Landsat, and NAIP | 137 | Global | Multispectral, High-resolution | Land Cover Classification, Segmentation, Change Detection |
| Oct, 2023 | SSL4EO-L [83] | SSL4EO-L: Datasets and Foundation Models for Landsat Imagery | 264 x 264 | 5,000,000 | 30 | Landsat 4-5 TM, Landsat 7 ETM+, Landsat 8-9 OLI/TIRS | - | Global | Multispectral | Cloud Detection, Land Cover Classification, Semantic Segmentation |
| Jul, 2024 | MMEarth [72] | MMEarth: Exploring Multi-Modal Pretext Tasks For Geospatial Representation Learning | 128 x 128 | 1,200,000 | 10 | Sentinel-2, Sentinel-1, Aster DEM | 46 | Global | Multispectral, SAR, Climate | Land Cover Classification, Semantic Segmentation |

TABLE VIII: This table summarizes a set of commonly used pre-trained datasets for remote sensing, including details on dataset, sensor type, geographic coverage, and related applications.

[3] Abdulaziz Amer Aleissaee, Amandeep Kumar, Rao Muhammad Anwer, Salman Khan, Hisham Cholakkal, Gui-Song Xia, and Fahad Shahbaz khan. Transformers in remote sensing: A survey, 2022. 5

[4] Arbeck. English: Mono, Multi and Hyperspectral Cube and corresponding Spectral Signatures, Mar. 2013. 3

[5] Argyro Argyrou and Athos Agapiou. A review of artificial intelligence and remote sensing for archaeological research. *Remote Sensing*, 14(23), 2022. 12

[6] Kumar Ayush, Burak Uzkent, Chenlin Meng, Kumar Tanmay, Marshall Burke, David Lobell, and Stefano Ermon. Geography-aware self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10181–10190, October 2021. 4, 8, 13, 16, 18

[7] Favyen Bastani, Piper Wolters, Ritwik Gupta, Joe Ferdinando, and Aniruddha Kembhavi. Satlaspretrain: A large-scale dataset for remote sensing image understanding, 2023. 4, 9, 13, 20, 21

[8] Nikolaos Ioannis Bountos, Arthur Ouaknine, and David Rolnick. Fomo-bench: a multi-modal, multi-scale and multi-task forest monitoring benchmark for remote sensing foundation models, 2024. 4, 10, 13, 14, 15

[9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021. 6

[10] Rodrigo Caye Daudt, Bertrand Le Saux, Alexandre Boulch, and Yann Gousseau. Oscd - onera satellite change detection, 2019. 17

[11] Keumgang Cha, Junghoon Seo, and Taekyung Lee. A billion-scale foundation model for remote sensing images, 2024. 4, 10, 15, 16, 17

[12] 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. 15, 17

[13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020. 6, 8

[14] 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. 17

[15] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *CVPR*, 2018. 3, 8, 9, 10, 20, 21

[16] Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall Burke, David Lobell, and Stefano Ermon. Satmae: Pre-training transformers for temporal and multi-spectral satellite imagery. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 197–211. Curran Associates, Inc., 2022. 3, 4, 5, 8, 13, 16, 18

[17] Russell G. Congalton. Remote sensing: An overview. *GIScience & Remote Sensing*, 47(4):443–459, 2010. 2, 3

[18] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 6, 7, 8, 9

[19] Philipe Dias, Abhishek Potnis, Sreelekha Guggilam, Lexie Yang, Aristeidis Tsaris, Henry Medeiros, and Dalton Lunga. An agenda for multimodal foundation models for earth observation. In *IGARSS 2023 - 2023 IEEE International Geoscience and Remote Sensing Symposium*, pages 1237–1240, 2023. 3

[20] Jian Ding, Nan Xue, Yang Long, Gui-Song Xia, and Qikai Lu. Learning roi transformer for detecting oriented objects in aerial images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 17

[21] Jian Ding, Nan Xue, Gui-Song Xia, Xiang Bai, Wen Yang, Michael Yang, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. Object detection in aerial images: A large-scale benchmark and challenges. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. 17

[22] Zhe Dong, Yanfeng Gu, and Tianzhu Liu. Generative convnet foundation model with sparse modeling and low-frequency reconstruction for remote sensing image interpretation. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2024. 4, 9, 13, 15, 17

[23] Zhe Dong, Yanfeng Gu, and Tianzhu Liu. Generative convnet foundation model with sparse modeling and low-frequency reconstruction for remote sensing image interpretation. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–16, 2024. 9, 10

[24] Zhe Dong, Yanfeng Gu, and Tianzhu Liu. Upetu: A unified parameter-efficient fine-tuning framework for remote sensing foundation model. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–13, 2024. 4, 10, 13, 15, 17, 18

[25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. 2, 3

[26] Iris Dumeur, Silvia Valero, and Jordi Inglada. Self-supervised spatio-temporal representation learning of satellite image time series. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 17:4350–4367, 2024. 4, 9, 13

[27] Yingchao Feng, Peijin Wang, Wenhui Diao, Qibin He, Huiyang Hu, Hanbo Bi, Xian Sun, and Kun Fu. A self-supervised cross-modal remote sensing foundation model with multi-domain representation and cross-domain fusion. In *IGARSS 2023 - 2023 IEEE International Geoscience and Remote Sensing Symposium*, pages 2239–2242, 2023. 4, 9, 13

[28] Anthony Fuller, Koreen Millard, and James R. Green. Croma: Remote sensing representations with contrastive radar-optical masked autoencoders, 2023. 4, 9, 13, 15

[29] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning, 2020. 6

[30] Shengxi Gui, Shuang Song, Rongjun Qin, and Yang Tang. Remote sensing object detection in the deep learning era—a review. *Remote Sensing*, 16(2), 2024. 3

[31] Shengxi Gui, Shuang Song, Rongjun Qin, and Yang Tang. Remote sensing object detection in the deep learning era—a review. *Remote Sensing*, 16(2), 2024. 11

[32] 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, 2024. 4, 10, 13, 14, 15, 16, 17

[33] Boran Han, Shuai Zhang, Xingjian Shi, and Markus Reichstein. Bridging remote sensors with multisensor geospatial foundation models, 2024. 4, 10, 13, 14, 15

[34] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners, 2021. 6

[35] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning, 2020. 6, 8

[36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 2, 7, 11

[37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 3, 11

[38] Konrad Heidler, Lichao Mou, Di Hu, Pu Jin, Guangyao Li, Chuang Gan, Ji-Rong Wen, and Xiao Xiang Zhu. Self-supervised audiovisual representation learning for remote sensing data, 2021. 20

[39] Yassine Himeur, Bhagawat Rimal, Abhishek Tiwary, and Abbes Amira. Using artificial intelligence and data fusion for environmental monitoring: A review and future perspectives. *Information Fusion*, 86-87:44–75, 2022. 3, 12

[40] 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*, page 1–18, 2024. 4, 9, 13, 15, 17

[41] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 18

[42] Ziyue Huang, Mingming Zhang, Yuan Gong, Qingjie Liu, and Yunhong Wang. Generic knowledge boosted pre-training for remote sensing images, 2024. 4, 9, 13, 16, 17

[43] International Society for Photogrammetry and Remote Sensing (ISPRS). 2d semantic labeling contest – potsdam, 2024. Accessed: 2024-07-08. 15

[44] Jeremy Irvin, Lucas Tao, Joanne Zhou, Yuntao Ma, Langston Nashold, Benjamin Liu, and Andrew Y. Ng. Usat: A unified self-supervised encoder for multi-sensor satellite imagery, 2023. 4, 9, 13, 15

[45] Pallavi Jain, Bianca Schoen-Phelan, and Robert Ross. Self-supervised learning for invariant representations from multi-spectral and sar images. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 15:7797–7808, 2022. 4, 8, 13

[46] Johannes Jakubik, Sujit Roy, C. E. Phillips, Paolo Fraccaro, Denys Godwin, Bianca Zadrozny, Daniela Szwarcman, Carlos Gomes, Gabby Nyirjesy, Blair Edwards, Daiki Kimura, Naomi Simumba, Linsong Chu, S. Karthik Mukkavilli, Devyani Lambhate, Kamal Das, Ranjini Bangalore, Dario Oliveira, Michal Muszynski, Kumar Ankur, Muthukumaran Ramasubramanian, Iksha Gurung, Sam Khallaghi, Hanxi, Li, Michael Cecil, Maryam Ahmadi, Fatemeh Kordi, Hamed Alemohammad, Manil Maskey, Raghu Ganti, Kommy Weldemariam, and Rahul Ramachandran. Foundation models for generalist geospatial artificial intelligence, 2023. 4, 9, 18

[47] Neal Jean, Sherrie Wang, Anshul Samar, George Azzari, David Lobell, and Stefano Ermon. Tile2vec: Unsupervised representation learning for spatially distributed data, 2018. 1

[48] Avinash Kumar Jha, Awishkar Ghimire, Surendrabikram Thapa, Aryan Mani Jha, and Ritu Raj. A review of ai for urban planning: Towards building sustainable smart cities. In *2021 6th International Conference on Inventive Computation Technologies (ICICT)*, pages 937–944, 2021. 3, 14

[49] Wentao Jiang, Jing Zhang, Di Wang, Qiming Zhang, Zengmao Wang, and Bo Du. Lemevit: Efficient vision transformer with learnable meta tokens for remote sensing image interpretation, 2024. 4, 10, 13

[50] Licheng Jiao, Zhongjian Huang, Xiaoqiang Lu, Xu Liu, Yuting Yang, Jiaxuan Zhao, Jinyue Zhang, Biao Hou, Shuyuan Yang, Fang Liu, Wenping Ma, Lingling Li, Xiangrong Zhang, Puhua Chen, Zhixi Feng, Xu Tang, Yuwei Guo, Dou Quan, Shuang Wang, Weibin Li, Jing Bai, Yangyang Li, Ronghua Shang, and Jie Feng. Brain-inspired remote sensing foundation models and open problems: A comprehensive survey. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16:10084–10120, 2023. 3

[51] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey, 2019. 3, 6

[52] Andreas Kamilaris and Francesc X. Prenafeta-Boldú. Deep learning in agriculture: A survey. *Computers and Electronics in Agriculture*, 147:70–90, 2018. 12

[53] Dominik Koßmann, Viktor Brack, and Thorsten Wilhelm. Seasonet: A seasonal scene classification, segmentation and retrieval dataset for satellite imagery over germany, 2022. 20, 21

[54] Weijie L, Wei Yang, Yuenan Hou, Li Liu, Yongxiang Liu, and Xiang Li. Saratr-x: A foundation model for synthetic aperture radar images target recognition, 2024. 4, 10

[55] Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han. Object detection in optical remote sensing images: A survey and a new benchmark. *ISPRS Journal of Photogrammetry and Remote Sensing*, 159:296–307, Jan. 2020. 6, 17

[56] Wenyuan Li, Keyan Chen, Hao Chen, and Zhenwei Shi. Geographical knowledge-driven representation learning for remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–16, 2022. 4, 8, 13, 15, 17, 21

[57] Wenyuan Li, Keyan Chen, and Zhenwei Shi. Geographical supervision correction for remote sensing representation learning. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–20, 2022. 4, 8, 13, 17

[58] Weijie Li, Yang Wei, Tianpeng Liu, Yuenan Hou, Yuxuan Li, Zhen Liu, Yongxiang Liu, and Li Liu. Predicting gradient is better: Exploring self-supervised learning for sar atr with a joint-embedding predictive architecture, 2024. 4, 9, 13

[59] Xuyang Li, Danfeng Hong, and Jocelyn Chanussot. 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 (CVPR)*, pages 24088–24097, June 2024. 4, 10, 13

[60] Xiang Li, Congcong Wen, Yuan Hu, Zhenghang Yuan, and Xiao Xiang Zhu. Vision-language models in remote sensing: Current progress and future trends. *IEEE Geoscience and Remote Sensing Magazine*, 12(2):32–66, 2024. 5

[61] Ziyu Liu, Azadeh Alavi, Minyi Li, and Xiang Zhang. Self-supervised learning for time series: Contrastive or generative?, 2024. 16

[62] Yang Long, Yiping Gong, Zhifeng Xiao, and Qing Liu. Accurate object localization in remote sensing images based on convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing*, PP:1–13, 01 2017. 12, 19, 21

[63] Yang Long, Gui-Song Xia, Shengyang Li, Wen Yang, Michael Ying Yang, Xiao Xiang Zhu, Liangpei Zhang, and Deren Li. On creating benchmark dataset for aerial image interpretation: Reviews, guidances and million-aid. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:4205–4230, 2021. 8, 9, 10, 12, 20, 21

[64] Yang Long, Gui-Song Xia, Liangpei Zhang, Gong Cheng, and Deren Li. Aerial scene parsing: From tile-level scene classification to pixel-wise semantic labeling, 2022. 6, 8, 9, 10, 12, 20, 21

[65] Yuchi Ma, Shuo Chen, Stefano Ermon, and David B. Lobell. Transfer learning in environmental remote sensing. *Remote Sensing of Environment*, 301:113924, 2024. 3

[66] Oscar Mañas, Alexandre Lacoste, Xavier Giró-i Nieto, David Vazquez, and Pau Rodríguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9414–9423, October 2021. 4, 8, 10, 13, 14, 15, 16, 17, 20, 21

[67] Utkarsh Mall, Bharath Hariharan, and Kavita Bala. Change-aware sampling and contrastive learning for satellite images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5261–5270, June 2023. 4, 8, 13, 14, 15, 17, 20, 21

[68] Lorenzo Mantovan and Loris Nanni. The computerization of archaeology: Survey on artificial intelligence techniques. *SN Computer Science*, 1(5), Aug. 2020. 14

[69] Matias Mendieta, Boran Han, Xingjian Shi, Yi Zhu, and Chen Chen. Towards geospatial foundation models via continual pretraining, 2023. 4, 9, 10, 13, 15, 17

[70] Dilxat Muhtar, Xueliang Zhang, Pengfeng Xiao, Zhenshi Li, and Feng Gu. Cmid: A unified self-supervised learning framework for remote sensing image understanding. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–17, 2023. 4, 7, 8, 13, 15, 17

[71] Ranganath Navalgund, Jayaraman V, and Parth Roy. Remote sensing applications: An overview. *Current science*, Vol. 93, 12 2007. 2, 3

[72] Vishal Nedungadi, Ankit Kariryaa, Stefan Oehmcke, Serge Belongie, Christian Igel, and Nico Lang. Mmearth: Exploring multi-modal pretext tasks for geospatial representation learning, 2024. 4, 10, 13, 15, 20, 21

[73] Mubashir Noman, Muzammal Naseer, Hisham Cholakkal, Rao Muhammad Anwar, Salman Khan, and Fahad Shahbaz Khan. Rethinking transformers pre-training for multi-spectral satellite imagery, 2024. 4, 9, 15

[74] Keiron O'Shea and Ryan Nash. An introduction to convolutional neural networks, 2015. 11

[75] Eleonora Jonasova Parelius. A review of deep-learning methods for change detection in multispectral remote sensing images.

*Remote Sensing*, 15(8), 2023. 11

[76] Daifeng Peng, Lorenzo Bruzzone, Yongjun Zhang, Haiyan Guan, Haiyong Ding, and Xu Huang. Semicdnet: A semisupervised convolutional neural network for change detection in high resolution remote-sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 59(7):5891–5906, 2021. 3

[77] Jonathan Prexl and Michael Schmitt. Multi-modal multi-objective contrastive learning for sentinel-1/2 imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 2136–2144, June 2023. 4, 8, 13

[78] Colorado J. Reed, Ritwik Gupta, Shufan Li, Sarah Brockman, Christopher Funk, Brian Clipp, Kurt Keutzer, Salvatore Candido, Matt Uyttendaele, and Trevor Darrell. Scale-mae: A scale-aware masked autoencoder for multiscale geospatial representation learning, 2023. 4, 5, 9, 13, 16

[79] Linus Scheibenreif, Joëlle Hanna, Michael Mommert, and Damian Borth. Self-supervised vision transformers for land-cover segmentation and classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1422–1431, June 2022. 4, 8, 13

[80] Michael Schmitt, Lloyd Haydn Hughes, Chunping Qiu, and Xiao Xiang Zhu. Sen12ms – a curated dataset of georeferenced multi-spectral sentinel-1/2 imagery for deep learning and data fusion, 2019. 8, 9, 12, 20, 21

[81] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. 7

[82] Michael J. Smith, Luke Fleming, and James E. Geach. Earthpt: a time series foundation model for earth observation, 2024. 4, 9, 13, 17

[83] Adam J. Stewart, Nils Lehmann, Isaac A. Corley, Yi Wang, Yi-Chia Chang, Nassim Ait Ali Braham, Shradha Sehgal, Caleb Robinson, and Arindam Banerjee. Ssl4eo-l: Datasets and foundation models for landsat imagery, 2023. 4, 8, 13, 20, 21

[84] Vladan Stojnic and Vladimir Risojevic. Self-supervised learning of remote sensing scene representations using contrastive multiview coding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1182–1191, June 2021. 4, 8, 13, 15

[85] Gencer Sumbul, Marcela Charfuelan, Begum Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, July 2019. 6, 8, 10, 14, 15, 20, 21

[86] Gencer Sumbul, Arne de Wall, Tristan Kreuziger, Filipe Marcelino, Hugo Costa, Pedro Benevides, Mario Caetano, Begum Demir, and Volker Markl. Bigearthnet-mm: A large-scale, multimodal, multilabel benchmark archive for remote sensing image classification and retrieval [software and data sets]. *IEEE Geoscience and Remote Sensing Magazine*, 9(3):174–180, Sept. 2021. 8

[87] 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. 4, 8, 13, 15, 16, 17, 18

[88] Maofeng Tang, Andrei Cozma, Konstantinos Georgiou, and Hairong Qi. Cross-scale mae: A tale of multi-scale exploitation in remote sensing, 2024. 4, 9, 13, 15

[89] Chao Tao, Ji Qi, Guo Zhang, Qing Zhu, Weipeng Lu, and Haifeng Li. Tov: The original vision model for optical remote sensing image understanding via self-supervised learning. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 16:4916–4930, 2023. 4, 8, 13, 15, 17

[90] Chao Tao, Ji Qia, Guo Zhang, Qing Zhu, Weipeng Lu, and Haifeng Li. Tov: The original vision model for optical remote sensing image understanding via self-supervised learning, 2022. 20, 21

[91] Jiayuan Tian, Jie Lei, Jiaqing Zhang, Weiying Xie, and Yunsong Li. Swimdiff: Scene-wide matching contrastive learning with diffusion constraint for remote sensing image, 2024. 4, 9, 13, 15, 17

[92] Gabriel Tseng, Ruben Cartuyvels, Ivan Zvonkov, Mirali Purohit, David Rolnick, and Hannah Kerner. Lightweight, pre-trained transformers for remote sensing timeseries, 2024. 4, 9, 13

[93] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. 3, 12

[94] Chen Wang, Alexis Mouche, Pierre Tandeo, Justin Stopa, Nicolas Longépé, Guillaume Erhard, Ralph Foster, Douglas Vandemark, and Bertrand Chapron. A labelled ocean sar imagery dataset of ten geophysical phenomena from sentinel-1 wave mode. *Geoscience Data Journal*, 6, 07 2019. 3

[95] Di Wang, Meiqi Hu, Yao Jin, Yuchun Miao, Jiaqi Yang, Yichu Xu, Xiaolei Qin, Jiaqi Ma, Lingyu Sun, Chenxing Li, Chuan Fu, Hongruixuan Chen, Chengxi Han, Naoto Yokoya, Jing Zhang, Minqiang Xu, Lin Liu, Lefei Zhang, Chen Wu, Bo Du, Dacheng Tao, and Liangpei Zhang. Hypersigma: Hyperspectral intelligence comprehension foundation model, 2024. 4, 10, 13, 14, 17, 18

[96] Di Wang, Jing Zhang, Bo Du, Gui-Song Xia, and Dacheng Tao. An empirical study of remote sensing pretraining. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–20, 2023. 4, 8, 13, 15, 17

[97] Di Wang, Jing Zhang, Bo Du, Gui-Song Xia, and Dacheng Tao. An empirical study of remote sensing pretraining. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–20, 2023. 7

[98] Di Wang, Jing Zhang, Bo Du, Minqiang Xu, Lin Liu, Dacheng Tao, and Liangpei Zhang. Samrs: Scaling-up remote sensing segmentation dataset with segment anything model, 2023. 10, 20, 21

[99] Di Wang, Jing Zhang, Minqiang Xu, Lin Liu, Dongsheng Wang, Erzhong Gao, Chengxi Han, Haonan Guo, Bo Du, Dacheng Tao, and Liangpei Zhang. Mtp: Advancing remote sensing foundation model via multi-task pretraining, 2024. 4, 10, 13, 15, 16, 17

[100] Di Wang, Qiming Zhang, Yufei Xu, Jing Zhang, Bo Du, Dacheng Tao, and Liangpei Zhang. Advancing plain vision transformer toward remote sensing foundation model. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–15, 2023. 4, 8, 13, 15, 17

[101] Xiao Wang, Siyan Liu, Aristeidis Tsaris, Jong-Youl Choi, Ashwin Aji, Ming Fan, Wei Zhang, Junqi Yin, Moetasim Ashfaq, Dan Lu, and Prasanna Balaprakash. Orbit: Oak ridge base foundation model for earth system predictability, 2024. 5, 13, 16, 18

[102] Yi Wang, Conrad M Albrecht, Nassim Ait Ali Braham, Chenying Liu, Zhitong Xiong, and Xiao Xiang Zhu. Decur: decoupling common & unique representations for multimodal self-supervision, 2023. 4, 9, 13, 14, 15

[103] Yi Wang, Conrad M. Albrecht, Nassim Ait Ali Braham, Lichao Mou, and Xiao Xiang Zhu. Self-supervised learning in remote sensing: A review. *IEEE Geoscience and Remote Sensing Magazine*, 10(4):213–247, 2022. 5

[104] Yi Wang, Conrad M Albrecht, Nassim Ait Ali Braham, Lichao Mou, and Xiao Xiang Zhu. Self-supervised learning in remote sensing: A review, 2022. 6, 7

[105] Yi Wang, Conrad M Albrecht, and Xiao Xiang Zhu. Self-supervised vision transformers for joint sar-optical representation learning, 2022. 4, 8, 13, 14, 15

[106] Yi Wang, Conrad M Albrecht, and Xiao Xiang Zhu. Multi-label guided soft contrastive learning for efficient earth observation pretraining, 2024. 4, 10, 13

[107] Yi Wang, Nassim Ait Ali Braham, Zhitong Xiong, Chenying Liu, Conrad M Albrecht, and Xiao Xiang Zhu. Ssl4eo-s12: A large-scale multi-modal, multi-temporal dataset for self-supervised learning in earth observation, 2023. 9, 10, 12, 20, 21

[108] Yi Wang, Hugo Hernández Hernández, Conrad M Albrecht, and Xiao Xiang Zhu. Feature guided masked autoencoder for self-supervised learning in remote sensing, 2023. 4, 9, 13, 15

[109] Yuelei Wang, Ting Zhang, Liangjin Zhao, Lin Hu, Zhechao

Wang, Ziqing Niu, Peirui Cheng, Kaiqiang Chen, Xuan Zeng, Zhirui Wang, Hongqi Wang, and Xian Sun. Ringmo-lite: A remote sensing multi-task lightweight network with cnn-transformer hybrid framework, 2023. 4, 9, 13, 15, 17

[110] Zhechao Wang, Peirui Cheng, Pengju Tian, Yuchao Wang, Mingxin Chen, Shujing Duan, Zhirui Wang, Xinming Li, and Xian Sun. Rs-dfm: A remote sensing distributed foundation model for diverse downstream tasks, 2024. 4, 10, 13

[111] Xinye Wanyan, Sachith Seneviratne, Shuchang Shen, and Michael Kirley. Extending global-local view alignment for self-supervised learning with remote sensing imagery, 2024. 4, 5, 10, 13, 14, 15, 17

[112] Michael Wulder, Thomas Loveland, David Roy, Christopher Crawford, Jeffrey Masek, Curtis Woodcock, Richard Allen, Martha Anderson, Alan Belward, Warren Cohen, John Dwyer, Angela Erb, Feng Gao, Patrick Griffiths, Dennis Helder, Txomin Hermosilla, James Hipple, Patrick Hostert, M. Hughes, and Zhe Zhu. Current status of landsat program, science, and applications. *Remote Sensing of Environment*, 225:127–147, 03 2019. 2

[113] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 17

[114] 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, 2019. 8, 20, 21

[115] 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, July 2017. 9

[116] Zhifeng Xiao, Yang Long, Deren Li, Chunshan Wei, Gefu Tang, and Junyi Liu. High-resolution remote sensing image retrieval based on cnns from a dimensional perspective. *Remote Sensing*, 9(7), 2017. 12, 19, 21

[117] Zhitong Xiong, Yi Wang, Fahong Zhang, Adam J. Stewart, Joëlle Hanna, Damian Borth, Ioannis Papoutsis, Bertrand Le Saux, Gustau Camps-Valls, and Xiao Xiang Zhu. Neural plasticity-inspired multimodal foundation model for earth observation, 2024. 4, 10, 13, 18

[118] Zhitong Xiong, Yi Wang, Fahong Zhang, and Xiao Xiang Zhu. One for all: Toward unified foundation models for earth vision, 2024. 4, 10, 13, 18

[119] Fanglong Yao, Wanxuan Lu, Heming Yang, Liangyu Xu, Chenglong Liu, Leiyi Hu, Hongfeng Yu, Nayu Liu, Chubo Deng, Deke Tang, Changshuo Chen, Jiaqi Yu, Xian Sun, and Kun Fu. Ringmo-sense: Remote sensing foundation model for spatiotemporal prediction via spatiotemporal evolution disentangling. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–21, 2023. 4, 9, 13

[120] Lefei Zhang and Liangpei Zhang. Artificial intelligence for remote sensing data analysis: A review of challenges and opportunities. *IEEE Geoscience and Remote Sensing Magazine*, 10(2):270–294, 2022. 5

[121] Liangpei Zhang, Lefei Zhang, and Bo Du. Deep learning for remote sensing data: A technical tutorial on the state of the art. *IEEE Geoscience and Remote Sensing Magazine*, 4(2):22–40, 2016. 5

[122] Lixian Zhang, Yi Zhao, Runmin Dong, Jinxiao Zhang, Shuai Yuan, Shilei Cao, Mengxuan Chen, Juepeng Zheng, Weijia Li, Wei Liu, Wayne Zhang, Litong Feng, and Haohuan Fu. A$^2$-mae: A spatial-temporal-spectral unified remote sensing pre-training method based on anchor-aware masked autoencoder, 2024. 4, 10, 13, 16

[123] Mingming Zhang, Qingjie Liu, and Yunhong Wang. Ctxmim: Context-enhanced masked image modeling for remote sensing image understanding, 2024. 4, 10, 13, 15

[124] Tong Zhang, Peng Gao, Hao Dong, Yin Zhuang, Guanqun Wang, Wei Zhang, and He Chen. Consecutive pre-training: A knowledge transfer learning strategy with relevant unlabeled data for remote sensing domain. *Remote Sensing*, 14(22), 2022. 4, 8, 13, 17

[125] Yi Zhao, Xinchang Zhang, Weiming Feng, and Jianhui Xu. Deep learning classification by resnet-18 based on the real spectral dataset from multispectral remote sensing images. *Remote Sensing*, 14(19), 2022. 11

[126] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, Hao Peng, Jianxin Li, Jia Wu, Ziwei Liu, Pengtao Xie, Caiming Xiong, Jian Pei, Philip S. Yu, and Lichao Sun. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt, 2023. 3

[127] Huming Zhu, Chendi Liu, Qiuming Li, et al. Deep convolutional encoder–decoder networks based on ensemble learning for semantic segmentation of high-resolution aerial imagery. *CCF Transactions on High Performance Computing*, 6:408–424, August 2024. 15

[128] Hao Zhu, Mengru Ma, Wenping Ma, Licheng Jiao, Shikuan Hong, Jianchao Shen, and Biao Hou. A spatial-channel progressive fusion resnet for remote sensing classification. *Information Fusion*, 70:72–87, 2021. 11

[129] Xiao Xiang Zhu, Devis Tuia, Lichao Mou, Gui-Song Xia, Liangpei Zhang, Feng Xu, and Friedrich Fraundorfer. Deep learning in remote sensing: A comprehensive review and list of resources. *IEEE Geoscience and Remote Sensing Magazine*, 5(4):8–36, 2017. 5, 7

[130] Xiao Xiang Zhu, Zhitong Xiong, Yi Wang, Adam J. Stewart, Konrad Heidler, Yuanyuan Wang, Zhenghang Yuan, Thomas Dujardin, Qingsong Xu, and Yilei Shi. On the foundations of earth and climate foundation models, 2024. 5

## Biography Section

**Siqi Lu** (Student Member, IEEE) is a second-year Master's student in Electrical and Computer Engineering at Vanderbilt University, supervised by Dr. Yuankai Huo and Dr. Mitchell M Wilkes. Her research interests include deep learning, medical image analysis, and software engineering. She received the B.S. degree in Electrical Engineering from the University of Illinois, Urbana-Champaign in 2023.

**Junlin Guo** is currently working toward the Ph.D. degree in Electrical and Computer Engineering at Vanderbilt University. His research interests include medical image analysis, deep learning, computer vision, and brain study. Before joining Vanderbilt University, he received the B.S. degree in telecommunication engineering from Northeastern University in 2017, and the M.S. degree from the Department of Electrical and Computer Engineering at Vanderbilt University in 2020, focusing on fMRI brain activation study.

**James R Zimmer-Dauphinee** received the B.A. degree in Anthropology and the B.S. degree in Mathematics from Georgia Southern University in 2011, the M.A. degree in Anthropology from the University of Arkansas in 2014, and the Ph.D. degree in Anthropology from Vanderbilt University in 2023. He is currently a postdoctoral fellow in the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University, funded by the GeoPACHA 2.0 Grant from the National Endowment for the Humanities. His research interests include developing deep learning models for large-scale autonomous archaeological satellite imagery surveys, geophysical methods, and spatial modeling to understand the impact of colonization on indigenous peoples.

**Jordan M Nieusma** is a research assistant for the Vanderbilt University Spatial Analysis Research Laboratory. She received her M.S degree in Data Science at Vanderbilt University in May 2024 and holds a Bachelor of Arts degree in English with a French Minor from Haverford College.

**Xiao Wang** received the B.S. degrees in Mathematics and Computer Science from Saint John's University, MN, in 2012, the M.S. degree in Electrical and Computer Engineering from Purdue University, West Lafayette, IN, in 2016, and the Ph.D. degree in Electrical and Computer Engineering from Purdue University in 2017. He pursued postdoctoral research at Harvard Medical School and Boston Children's Hospital until 2021. He is currently a research staff scientist at Oak Ridge National Laboratory. His research interests include applying machine learning, medical physics, image processing, and high-performance computing to various imaging problems, including CT reconstruction, electron tomography imaging, and MRI. He was the 2022 AAPM Truth CT reconstruction challenge winner and a 2017 ACM Gordon Bell Prize finalist.

**Parker VanValkenburgh** is an Associate Professor of Anthropology and Interim Director of Latin American and Caribbean Studies at Brown University. His research focuses on the impacts of colonialism and imperialism on Indigenous people and environments in the Peruvian Andes. He utilizes diverse materials and digital methodologies, including GIS, to understand the transformation of relationships in imperial histories. He co-directs the Paisajes Arqueológicos de Chachapoyas (PACha) project and GeoPACHA (Geospatial Platform for Andean Culture, History, and Archaeology). VanValkenburgh received his Ph.D. from Harvard University.

**Steven A Wernke** is Associate Professor and Chair of Anthropology at Vanderbilt University, director of the Spatial Analysis Research Laboratory, and director of the Vanderbilt Institute for Spatial Research. Prof. Wernke is an archaeologist and historical anthropologist of the Andean region of South America. His research takes place at the intersection of several disciplines: archaeology and history, prehispanic and colonial studies, anthropology and cultural geography. Prof. Wernke's interests center on the lived experiences of indigenous communities across the Spanish invasion of the Andes–especially how new kinds of communities, landscapes, and religious practice emerged out of successive attempts by the Inkas and the Spanish to subordinate and remake Andean societies in their own self-image. Methodologically, his work brings together analyses of archaeological and documentary datasets in geospatial frameworks.

**Yuankai Huo** is an Assistant Professor in Computer Science at Vanderbilt University, TN, USA. He received his B.S. degree in Electrical Engineering from Nanjing University of Posts and Telecommunications (NJUPT) in 2008, and Master degree in Electrical Engineering from Southeast University in 2011. After graduation, He worked in Columbia University and New York State Psychiatric Institute as a staff engineer and research officer from 2011 to 2014. He received his Master degree in Computer Science from Columbia University in 2014, and Ph.D. degree in Electrical Engineering from Vanderbilt University in 2018. Then, he had worked as a Research Assistant Professor at Vanderbilt University, and later, a Senior Research Scientist at PAII Labs. Since 2020, he has been a faculty member at the Department of Electrical Engineering and Computer Science, and Data Science Institute, Vanderbilt University.