

# DeepAndes: A Self-Supervised Vision Foundation Model for Multi-Spectral Remote Sensing Imagery of the Andes

Junlin Guo, James R. Zimmer-Dauphinee, Jordan M. Nieusma, Siqi Lu, Quan Liu,  
Ruining Deng, Can Cui, Jialin Yue, Yizhe Lin, Tianyuan Yao,  
Juming Xiong, Junchao Zhu, Chongyu Qu, Yuechen Yang, Mitchell Wilkes,  
Xiao Wang, Parker VanValkenburgh, Steven A. Wernke, and Yuankai Huo

**Abstract**—By mapping sites at large scales using remotely sensed data, archaeologists can generate unique insights into long-term demographic trends, inter-regional social networks, and past adaptations to climate change. Remote sensing surveys complement field-based approaches, and their reach can be especially great when combined with deep learning and computer vision techniques. However, conventional supervised deep learning methods face challenges in annotating fine-grained archaeological features at scale. While recent vision foundation models have shown remarkable success in learning large-scale remote sensing data with minimal annotations, most off-the-shelf solutions are designed for RGB images rather than multi-spectral satellite imagery, such as the 8-band data used in our study. In this paper, we introduce DeepAndes, a transformer-based vision foundation model trained on three million multi-spectral satellite images, specifically tailored for Andean archaeology. DeepAndes incorporates a customized DINOv2 self-supervised learning algorithm optimized for 8-band multi-spectral imagery, marking the first foundation model designed explicitly for the Andes region. We evaluate its image understanding performance through imbalanced image classification, image instance retrieval, and pixel-level semantic segmentation tasks. Our experiments show that DeepAndes achieves superior F1 scores, mean average precision, and Dice scores in few-shot learning scenarios, significantly outperforming models trained from scratch or pre-trained on smaller datasets. This underscores the

effectiveness of large-scale self-supervised pre-training in archaeological remote sensing. Codes will be available on <https://github.com/geopacha/DeepAndes>.

**Index Terms**—Foundation Model, Self-supervised Learning, DINOv2, Multi-Spectral Imaging, Remote Sensing, Andean Archaeology

## I. INTRODUCTION

ONE of the most vexing and persistent challenges for field-based sciences such as archaeology, population biology, demography, environmental monitoring, and field geology is conducting analyses at large scales. At the level of small regions, conventional field-based survey methods have proven to be highly effective. Through such surveys, field scientists reconstruct and analyze populations, environments, and resources at regional levels and provide crucial data for modeling them at larger scales [2], [3], [6], [7], [18], [20], [41], [47]. However, generating accurate datasets that record continuous distributions (particularly of highly variable cultural phenomena, such as archaeological sites and the distribution of modern urban areas) has proven difficult to achieve at inter-regional and continental scales without large resource outlays to fund multi-year field campaigns. [17], [37], [42], [48]. These problems are especially notable in Andean South America, where the topography and inaccessibility of many areas make field surveys challenging. Since the early 2000s, field scientists have used high-resolution remote sensing satellite imagery to analyze the distribution of archaeological features, natural resources, and modern populations at larger scales. However, such “brute force” manual imagery surveys remain labor intensive, time consuming, and prone to observational fatigue and inter-observer variability in feature detection [9]. Developing effective AI-assisted approaches for autonomous information extraction will enable us to expand survey coverage at scale, providing new insights into human adaptation, settlement patterns, landscapes, and networks in the Andes.

Junlin Guo, Siqi Lu, Jialin Yue, Juming Xiong, Chongyu Qu, Mitchell Wilkes, and Yuankai Huo are with the Department of Electrical and Computer Engineering, Vanderbilt University, Nashville, TN, USA.

James R. Zimmer-Dauphinee and Steven A. Wernke are with the Department of Anthropology, Vanderbilt University, Nashville, TN, USA.

Jordan M. Nieusma and Yuankai Huo are with the Data Science Institute, Vanderbilt University, Nashville, TN, USA.

Quan Liu, Ruining Deng, Can Cui, Tianyuan Yao, Junchao Zhu, Yuechen Yang, and Yuankai Huo are with the Department of Computer Science, Vanderbilt University, Nashville, TN, USA.

Yizhe Lin is with the Department of Mathematics, Vanderbilt University, Nashville, TN, USA.

Xiao Wang is with Oak Ridge National Laboratory, Oak Ridge, TN, USA.

Parker VanValkenburgh is with the Department of Anthropology, Brown University, Providence, RI, USA.

Ruining Deng is also with the Department of Radiology, Weill Cornell Medicine, New York, NY, USA.

With the rapid advancement of artificial intelligence, deep learning has made significant contributions in a diverse series of domains [29], [30], [40], [49], [50], [53]. Recently, the emergence of “**foundation models**” (FMs) has further expanded the scope of deep learning applications [4], [12], [13], [22], [23], [51], including remote sensing [31]–[33]. As summarized in [33], foundation models pretrained on massive remote sensing datasets demonstrate strong adaptability across a wide range of downstream tasks in both earth and social sciences. As a domain-specific application, archaeological remote sensing has particularly benefited from foundation models, which have shown promise in tasks such as artifact recognition [5], detection of archaeological structures [19], and texture analysis [1], offering new opportunities for large-scale archaeological analysis. Despite these advances, most of the earth’s surface has yet to be systematically surveyed, and vast regions are underrepresented in archaeological research, including the Andes region. High levels of inter-regional variability in land cover and the diversity of archaeological features themselves have proven to be major barriers to achieving such coverage. *To the best of our knowledge, no foundation model has been specifically developed for Andean archaeology using 8-channel hyperspectral satellite imagery.* Thus, developing a new AI foundation model with broad utility across the social and earth sciences in the Andes would be highly valuable, enabling experts to contribute where they are most effective—as observers and analysts in the archaeological workflow.

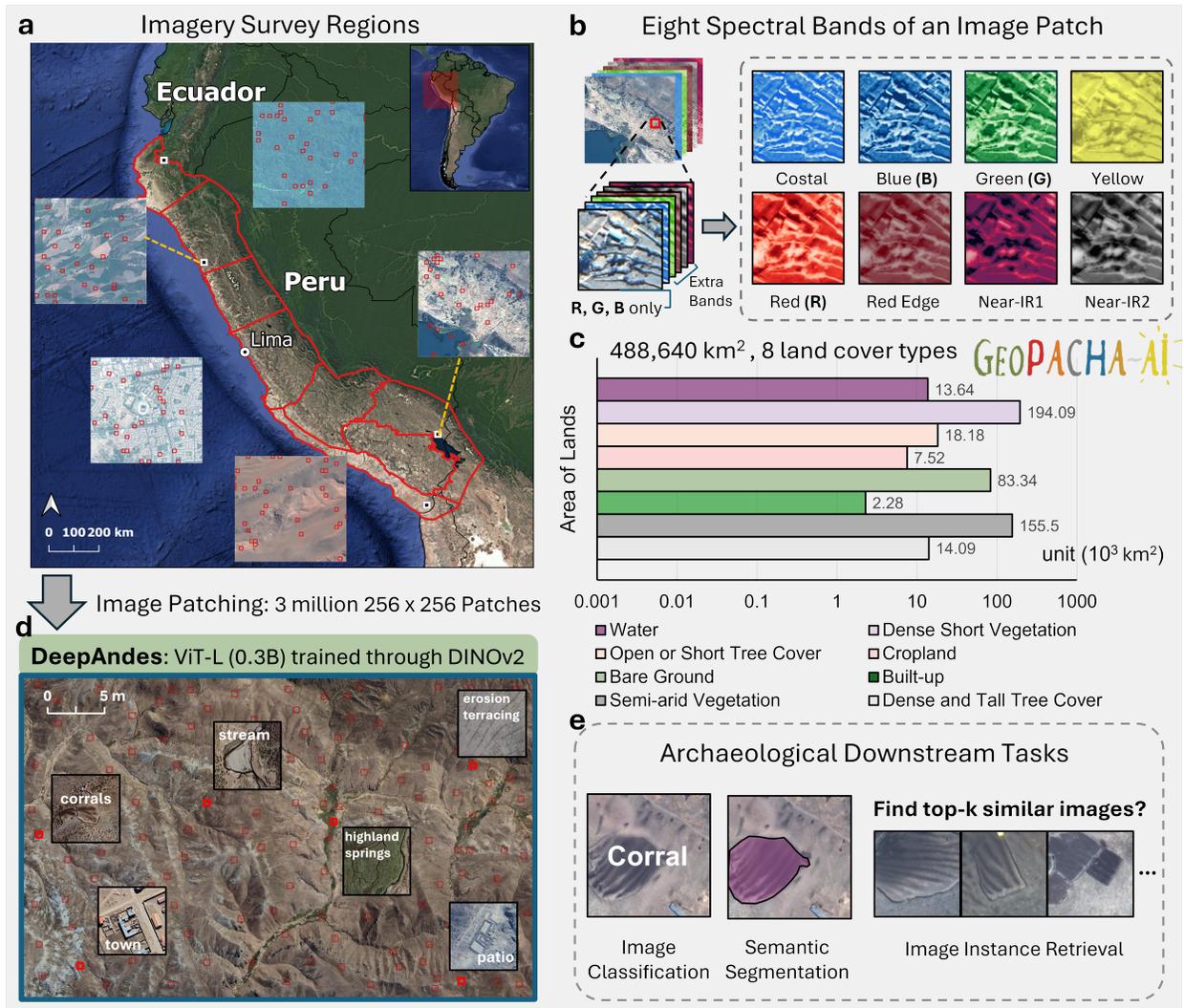
To bridge this gap, this paper proposes DeepAndes, the first vision foundation model for remote sensing of the Andean region. As shown in Fig. 1, the multi-spectral training dataset, sourced from the central Andes (Figs. 1a-c), covers 488,640 km<sup>2</sup> and encompasses 8 distinct land cover types (Fig. 1c). DeepAndes, a 307M-parameter vision transformer (ViT) model, is trained on 3 million satellite image patches. A key challenge in developing a novel vision foundation model learning framework for high-dimensional multi-spectral remote sensing images is that most off-the-shelf vision foundation models are primarily designed for natural images with only RGB bands (red, green, and blue). To address this problem, we customize the DINOv2 [36] self-supervised learning (SSL) pipelines to support our 8-band multi-spectral imagery. During pre-training, we also adjusted the scale of global-local image views to align with the size of architectural and archaeological features in the dataset. As shown in Fig. 1e, to evaluate DeepAndes, we investigated its image pattern recognition and few-shot learning capabilities through three key vision-based remote sensing tasks: imbalanced image classification, image instance retrieval, and image segmentation. These downstream tasks are crucial for large-scale remote

sensing datasets, where data are often highly heterogeneous and features are sparsely distributed. Similarly to [39], DeepAndes-based image instance retrieval can help identify patterns in the distribution of cultural and natural features within vast datasets, improving the efficiency of large-scale archaeological analysis.

Through extensive experiments employing archaeological features as a test case, our results demonstrate that self-supervised pre-training enhances DeepAndes’s performance in both image-level and pixel-level downstream tasks. Among the imbalanced archaeological loci (discrete archaeological features of interest) classification task, DeepAndes achieves an F1 score of 0.886. Even when fine-tuned on only 10% of the original classification data, it maintains a strong F1 score of 0.83, substantially outperforming models trained from scratch using supervised methods. For the archaeological image instance retrieval task, features extracted from the frozen DeepAndes backbone achieve a mean average precision (mAP) of 0.869 within the top- $k$  retrieved samples (mAP@ $k$ ,  $k=5$ ), demonstrating the effectiveness of the feature representations. For pixel-level dense feature recognition, we evaluated three challenging archaeological loci segmentation tasks, where the diversity of locus objects adds complexity. Our few-shot learning results show that DeepAndes excels in both transfer learning and fine-tuning settings when using a simple linear segmentation head. Interestingly, the scaling law is observed from the DeepAndes model as depicted in Fig. 2. In short, DeepAndes, pre-trained on 3 million diverse remote sensing satellite image patches, outperforms the same ViT models pre-trained on smaller-scale datasets (e.g., 30K and 300K images) across all downstream tasks, underscoring the effectiveness of large-scale self-supervised learning in this work.

The contributions of this paper are three-fold:

- We propose DeepAndes, the first AI-driven vision foundation model for the Andean region, leveraging large-scale pre-training on multi-spectral high-resolution satellite images for downstream analysis in field sciences.
- We revise the DINOv2 framework for multi-spectral pre-training using 8-band WorldView-2 and Worldview-3 satellite imagery. This modified framework leverages the flexibility of DINOv2 with data preprocessing, data augmentation, and network architecture, enabling seamless adaptation to other multi-channel remote sensing data in the future, regardless of the number of input channels.
- We demonstrate the few-shot adaptability of DeepAndes to downstream applications, showcasing its effectiveness across three prevalent archaeological remote sensing tasks: imbalanced image classification, image instance retrieval, and pixel-



**Fig. 1: Overview of DeepAndes.** This figure shows the training dataset (**a-d**) and three domain-specific downstream tasks (**e**) using DeepAndes — a vision foundation model designed for multi-spectral satellite imagery in the Andes region. Particularly, (**a**) shows a large-scale map of the imagery used to train DeepAndes, highlighting various land cover types, with their area distribution shown in (**c**). (**b**) presents the unit sample patch (red box in **a**, **b**, **d**) with eight spectral bands. (**d**) illustrates image patching for DINOv2 training, with geospatial sampling densely covering different archaeological sites.

level image segmentation.

## II. RELATED WORK

In this section, we describe the background of this study, covering foundation models in remote sensing and summarizing two major SSL strategies for foundation model pre-training.

### A. Foundation Models For Remote Sensing

The emergence of foundation models has revolutionized the field of remote sensing through their capacity

to serve as versatile pre-trained frameworks for various downstream applications [33]. These models excel in processing complex remote sensing data, including multi-spectral and multi-temporal imagery, by leveraging extensive pre-training on large-scale datasets [27]. The integration of SSL approaches [28] and transformer architectures [43] has substantially improved performance across tasks such as image classification and change detection [15]. A distinctive advantage of FMs in remote sensing lies in their capability to extract meaningful representations from unlabeled data through SSL techniques

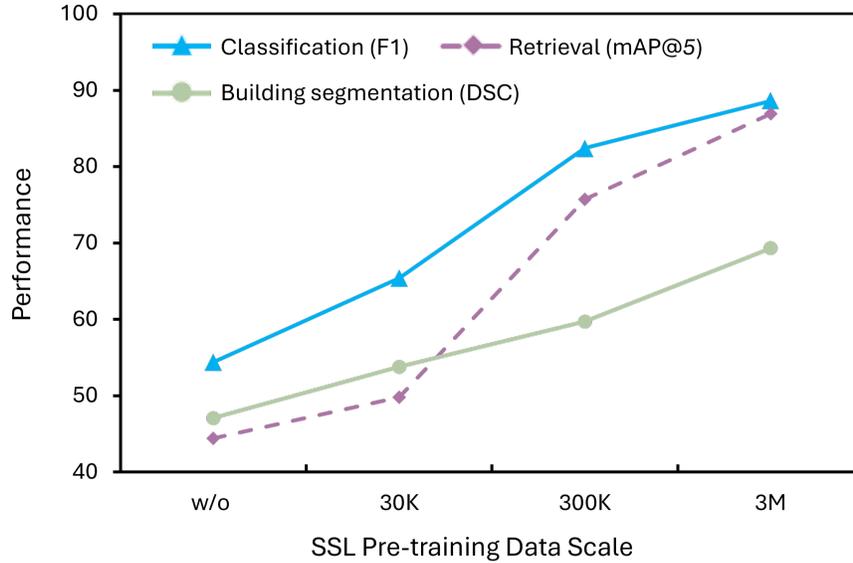


Fig. 2: **Scaling law is observed in DeepAndes.** This figure illustrates the model’s performance across three key downstream tasks: site classification, image retrieval, and building segmentation. The results are presented for models trained with no pretraining (w/o), 30K, 300K, and 3 million images. The findings highlight the scalability of DeepAndes, indicating that its performance can be further improved with larger training datasets.

[51]. The incorporation of transformer architectures [43] enables these models to effectively process geospatial data’s unique characteristics, including variable spatial resolutions and temporal patterns, building upon earlier findings [33]. The development trajectory of FMs has been shaped by both deep learning advances and data availability. While early progress centered on CNN architectures like ResNet [26] for image analysis tasks [34], the advent of transformer models has significantly enhanced the processing of large-scale imagery [11]. ViT models have particularly advanced the field by processing images as sequences of patches, enabling comprehensive analysis of both local details and global patterns. Recent advancements in FMs are transforming remote sensing by improving representation learning for diverse geospatial tasks. Notable examples include SatMAE [11], designed for temporal and multi-spectral satellite imagery analysis; Scale-MAE [38], which focuses on multiscale geospatial representation learning; SkySense [24], a billion-scale pre-trained universal model for earth observation imagery; and DINO-MC [45], which enhances SSL capabilities for remote sensing applications. However, these advances face ongoing challenges, including data quality requirements, computational demands, and domain adaptation needs [35].

### B. Self-Supervised Representation Learning Strategies

In the field of remote sensing and computer vision, self-supervised learning represents a fundamental component in foundation model pre-training [33]. During the pre-training stage, SSL allows models to learn meaningful representations without relying on labeled datasets. This capability is particularly beneficial in remote sensing, where labeled data is often scarce. Models pre-trained using SSL are adept at identifying patterns and features in large volumes of unlabeled remote sensing data, making them highly effective for various downstream applications. As summarized in [33], some commonly used SSL techniques in remote sensing FMs can be categorized into two main types: contrastive learning and predictive coding.

**Contrastive Learning.** Contrastive learning method focuses on learning data representations by comparing different augmented versions of the same data point. It pulls similar (positive) pairs closer and push dissimilar (negative) pairs apart in the representation space. The methodology depends on advanced data augmentation techniques, such as random cropping, rotation, and color jittering, to create diverse views of the same image. Notable implementations that have shown success in remote sensing applications include DINO [8], which employs self-distillation, SimCLR [10], which uses a simple contrastive framework, and MOCO (Momentum Contrast) [25], which leverages a dynamic dictionary

approach [33].

**Predictive Coding.** Predictive coding approaches, as discussed in [33], focus on training models to reconstruct missing input data components from available observations. This methodology has proven particularly effective in capturing both spatial and temporal relationships within remote sensing data, especially when dealing with multi-spectral imagery and cloud-covered regions [33]. The strategy commonly employs architectures such as autoencoders (AE), which learn compressed representations through encoding-decoding processes, and masked autoencoders (MAE), which specifically focus on reconstructing masked portions of input data [33]. These approaches have demonstrated significant success in developing robust internal representations of complex remote sensing data structures.

This work employs one of the SOTA SSL pre-training strategies, DINOv2 [36], on million-scale satellite image datasets from the Central Andes. The DINOv2 algorithm leverages contrastive learning concepts with knowledge distillation. Additionally, the “masking” concept is utilized with iBOT [52] loss to further encourage patch-level representation learning.

### III. METHODOLOGIES

This section details the methodologies employed in this work. It begins with the construction of the million-scale pre-training dataset, followed by an overview of the DINOv2 framework and our pre-training strategies. Finally, it presents the downstream analysis using the pre-trained DeepAndes backbone.

#### A. Million-Scale Training Dataset Construction

As described previously in Figure 1, the dataset used for SSL pre-training was sourced from the entire Peruvian Andes, surveyed from 6 teams, spanning 488,640  $km^2$ , and including 8 distinct land cover types. Our high-resolution satellite imagery is from the WorldView-2 (WV-2) and WorldView-3 (WV-3) satellites, both of which are composed of eight multi-spectral bands, plus a panchromatic band. After pansharpening (fusing the panchromatic band with other spectral bands), WV-2 imagery is of a Ground Sample Distance (GSD) of 0.46 m (at nadir), while WV-3 imagery provides 0.31 m (at nadir) GSD resolution. Our imagery processing pipeline includes steps for image orthorectification, cloud removal, and GSD commensuration. Image processing includes upsampling WV-2 imagery to match the GSD of WV-3 imagery. To construct a diverse, large-scale pre-training dataset, we randomly sampled the Central Andean area (encompassing western Peru and parts of northern Chile and northwestern Bolivia) into non-overlapping  $256 \times 256$  image patches (shown as the red

box in Fig. 1). Featureless images with entirely dark or bright pixels, such as those containing only water or clouds, were removed. This process yielded a total of 3 million  $256 \times 256$  image patches for pre-training DeepAndes.

#### B. DeepAndes

1) **DINOv2 Architecture:** DINOv2 (and its earlier variant DINO [8]) employs a contrastive learning framework with knowledge distillation, aiming to maximize the similarity of feature representations between two sets of augmented views of the same input image. Figure 3A illustrates the main architecture of DINOv2. The student network and teacher network use the same network architecture  $g$ , a vision transformer, but with different weights  $\theta_s$  and  $\theta_t$ . During pre-training, the student network  $g_{\theta_s}$  is trained to match the representation of the teacher network  $g_{\theta_t}$ . As shown, for an input image  $x$ , two sets of augmented views are generated and sent to  $g_{\theta_s}$  and  $g_{\theta_t}$ . Compared to  $g_{\theta_t}$ , the student network is information-limited because the augmented views sent to  $g_{\theta_s}$  are more noisy [8]. Unlike conventional knowledge distillation, the teacher’s weights are dynamically updated using exponential moving average (EMA) of past student weights. Depending on the image views (global or local views) matched between the student and teacher, the corresponding feature representations (class tokens and patch tokens) are further projected to obtain image-level and patch-level cross-entropy loss objectives. More details on the model design and derivation can be found in [8], [36].

2) **SSL Pre-training:** Similar to [44], Figure 3B shows the detail of the training paradigm of DINOv2-based SSL pre-training. As shown in Figure 3B, two sets of augmented views are sent to the student and teacher, respectively. The student network receives both global and local crops as inputs, while the teacher network is only provided with global crops. Specifically, The global crop encompasses most of the original image, whereas the local crop focuses on a small portion of it. As shown, during training, different augmented crops of the same image are fed into the networks to obtain the image-level and patch-level loss objectives.

- **DINO Loss: Image-level objective.** For global crops, the student attempts to generate the class token to match the teacher class token. The local crops are fed only to the student which tries to produce a representation that matches the teacher class token generated from the global crops, helping the model learn the local-global representations.
- **iBOT Loss: Patch-level objective.** The global crops are randomly masked given to the student network. iBOT head is applied to match the (un-

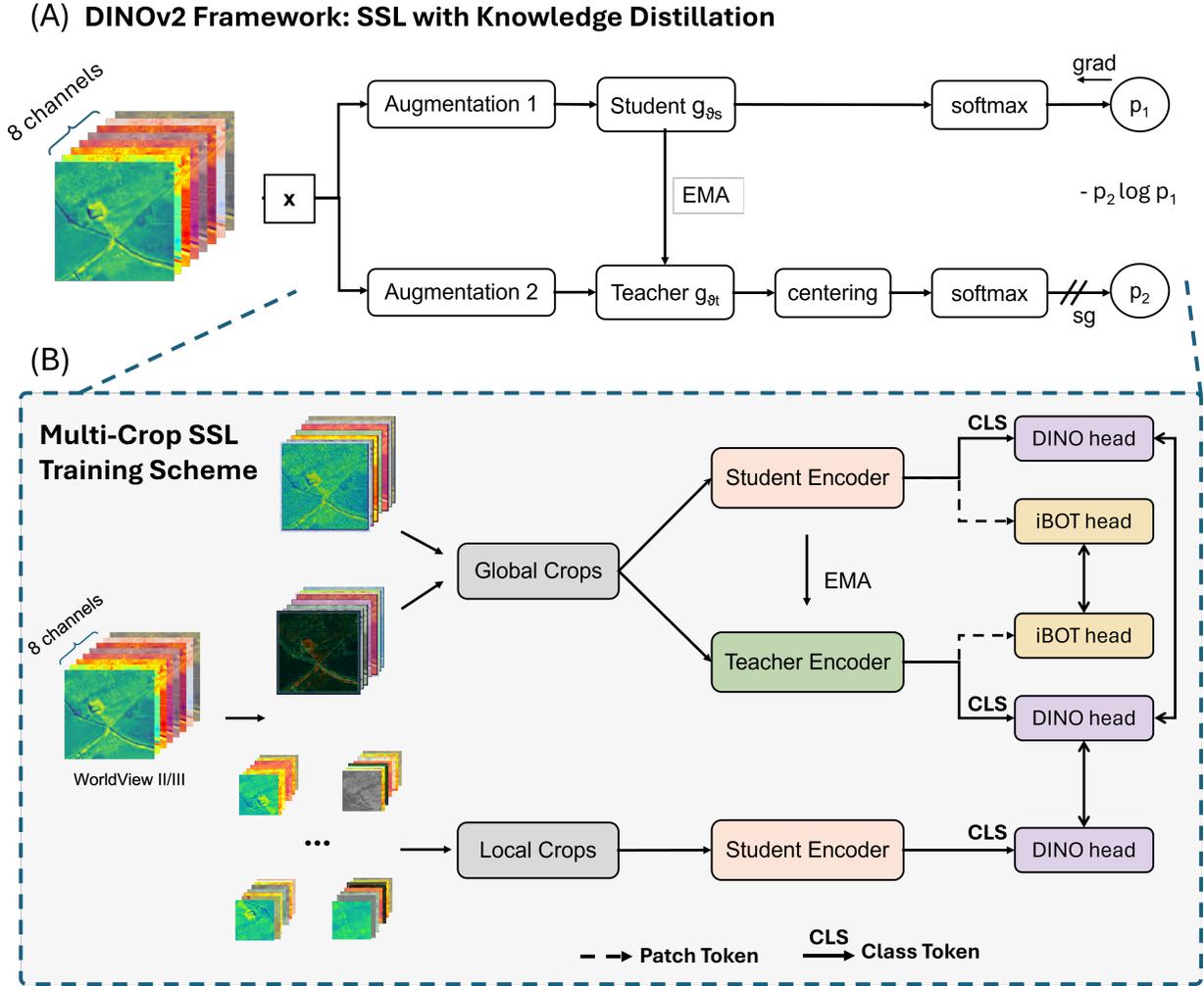


Fig. 3: **DINOv2: the self-supervised contrastive representation learning algorithm with knowledge distillation.** (A) shows an overview of the framework. (B) illustrates the details of the DINOv2 multi-crop SSL training scheme.

masked) teacher patch tokens to the student’s corresponding masked patch tokens, which encourages learning fine-grained patch-level feature representations.

In our work, we selected the ViT “large” (ViT-L/14) with 307M parameters as the backbone. The default multi-crop DINOv2 training uses a scale factor  $s = 0.32$  to define the scaling range of local crop as  $(0.05, s)$ , and  $(s, 1)$  for global views. In our study, we applied the scaling range  $(0.2, 0.5)$  for the local view and  $(0.5, 1)$  for the global view, due to the sparse image feature distribution and noise present in our remote sensing images. The other default hyperparameters for DINOv2 training algorithm were used for DeepAndes, as detailed in [36] with the following modifications: a base learning rate of 0.0002, and the entire training process covering approximately 40 times passes over the 3-million-image

dataset. DeepAndes was trained using AdamW ( $\beta_1 = 0.9, \beta_2 = 0.999$ ) with float16 precision. For ViT-L, we used 65,536 prototypes, resulting in 65,536-dimensional projection heads. Distributed data parallel (DDP) training was employed across 8 NVIDIA DGX-A100 nodes, with a batch size of 64 per GPU node. To implement complex augmentation pipelines, we adapted the default augmentation modules for 8-band image data with minor modifications.

3) **Feature Embeddings:** For an input image of size  $256 \times 256$ , it is resized to  $224 \times 224$  before being input to the model, in accordance with the default settings of DINOv2. Using ViT-L with a patch size of  $14 \times 14$ , the pre-trained DeepAndes projects the input image into vector representations (tokens) of dimension 1,024 in the feature space. With the default settings, this includes one class token and 256 image patch tokens.

These vectors generated by foundation models present meaningful feature information of the input images with a reduced dimension can be further utilized for diverse downstream tasks [24], [44], [45]. In our case, the class token can be used for downstream tasks such as image classification by concatenating a simple linear classifier. The patch tokens can be utilized for pixel-level image recognition tasks, such as semantic segmentation, by connecting them to a segmentation head.

### C. Archaeological Downstream Analysis

#### 1) Handling Imbalance in Archaeological Data:

Another major challenge in archaeological remote sensing is the inherent imbalance in archaeological data, particularly in our remote sensing satellite imagery. One of the previous studies [49] shows that the proportion of our satellite imagery that include archaeological settlement features is typically low, at less than 7% in the surveyed areas. For example, in a manually labeled dataset consisting of 5,830 labeled images sourced from a 4,000 km<sup>2</sup> survey area, the ratio of “positive” (containing settlement features) to “negative” (not containing settlement features) images can be as low as 1:100. The pre-trained model can be adapted and fine-tuned to improve performance in minority classes when the downstream task is bottlenecked by data imbalance and scarcity.

2) **Archaeological Downstream Tasks:** To support the motivation above, we evaluate DeepAndes on three image understanding downstream tasks in archaeology: image-level classification and retrieval, and pixel-level segmentation, as shown in the Figure 4.

**Imbalanced Loci Classification.** As described in the previous section, archaeological loci, such as settlement buildings and corrals, are sparsely distributed across the Andes area. Thus, we determine the first image-level downstream task as imbalanced loci classification. Specifically, images of the classification dataset are categorized into two classes: “positive” (i.e., the presence of loci) and “negative” (i.e., no presence of loci), denoted as background images. The class token of the pre-trained ViT-L backbone was connected to a simpler linear classifier with two fully connected layers. Given an input image, the fine-tuned classifier predicts either “positive” or “negative.” To represent the data imbalance in this task, the ratio of “positive” to “negative” images was set to 1:10. Additionally, to evaluate the few-shot learning capability of the pre-trained transformer backbones, we assess the classification task at different downstream dataset scales (10%, 30%, 50%, and 100% of the original classification dataset).

**Image Instance Retrieval.** Another image-level downstream task we perform is the image instance retrieval task. Unlike classification, image retrieval does

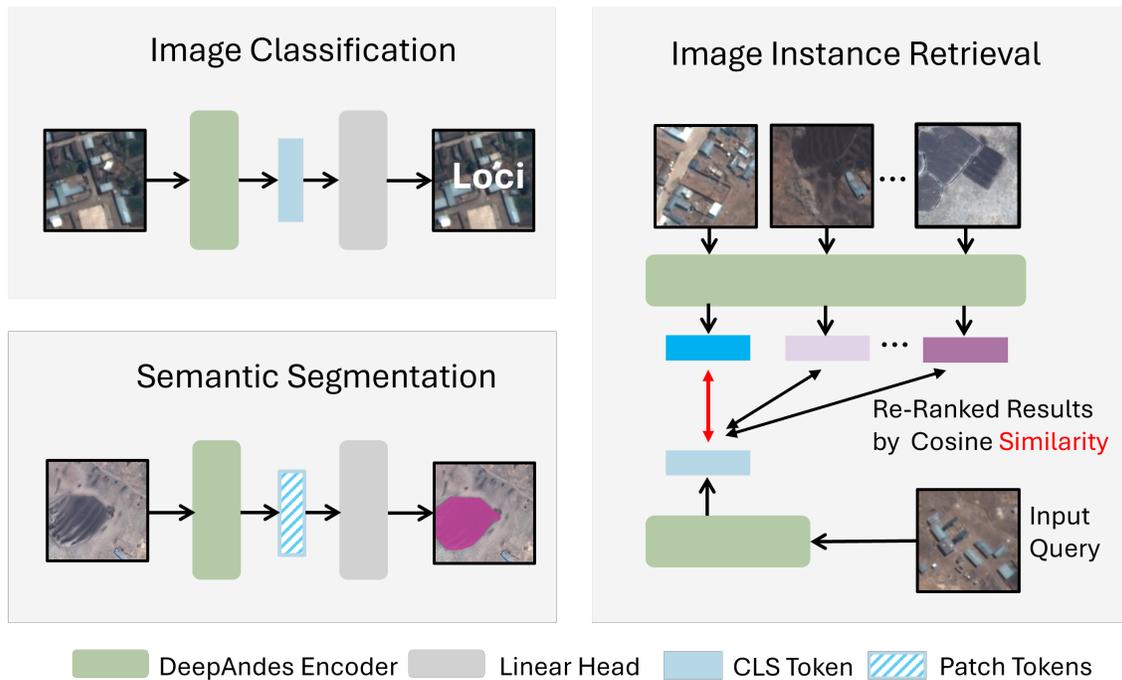


Fig. 4: **Image Understanding Downstream Analysis.** Both image-level tasks (e.g., image classification and retrieval) and pixel-level tasks (e.g., segmentation) are included.

not require fine-tuning on DeepAndes. Instead, we use an archaeological loci image as the input query and retrieve the top- $k$  most similar images from the database, ranking them based on their cosine similarity to the query loci image (via the class token). The performance of the model is then evaluated by the mean average precision of the matched image class (“positive” or “negative”) in the retrieved results.

**Few-shot Loci Segmentation** In this work, we also consider the task of recognizing dense features, such as semantic segmentation at the pixel level. Specifically, we perform few-shot semantic segmentation tasks for three types of loci: active buildings, active corrals, and archaeological corrals. Patch-level features are extracted from the images via patch tokens and then concatenated with a simple linear segmentation head. Unlike the loci classification task, pixel-level semantic segmentation enables more precise localization of archaeological features within the image. As described in [36], we freeze the pre-trained transformer backbone and only fine-tune the linear segmentation head to generate output logits from the patch tokens.

#### IV. DATA AND EXPERIMENTS

##### A. Data Pre-processing

Firstly, the raw satellite imagery is pre-processed into digitized archaeological images for the computer system in preparation for training deep learning models. The remote sensing satellite images used in this work are collected by the WV-2 and WV-3 satellites provided by the Digital Globe Foundation, following color correction and orthographic correction using a coarse digital elevation model (DEM). The data are then pan-sharpened using the Bayesian fusion algorithm from Orfeo-Toolbox [21] to increase the spatial resolution of the multi-spectral imagery to 0.31 m for WV-2 imagery to match the native resolution of WV-3 imagery. In this work, all 8 spectral bands (four standard colors—red, green, blue, and near-infrared 1—and four new bands—coastal, yellow, red edge, and near-infrared 2) are used. Lastly, the imagery is re-sampled from 32 bits to 8 bits to reduce storage size and computational load.

##### B. SSL with Different Pre-training Scales

As introduced in Section III-A, after data pre-processing, we construct the self-supervised pre-training dataset sourced from the entire Peruvian Andes. In total, the imagery covers survey regions of approximately 488,640  $km^2$  and includes 8 distinct land cover types. We densely sample 3 million image patches of size  $256 \times 256$  from all surveyed regions to ensure a diverse, large-scale pre-training dataset. In this work, we perform

SSL on pre-training datasets of varying scales to robustly evaluate the impact of pre-training on downstream tasks. Specifically, the ViT-L backbone is pre-trained using 30K, 300K, and 3 million image patches, and the corresponding pre-trained models are then evaluated respectively. The same procedure is followed to create these pre-training datasets, with the only difference being the number of patches sampled from each land cover type.

##### C. Archaeological Downstream Tasks

**1) Imbalanced Loci Classification:** For the archaeological loci classification task, the dataset is imbalanced, with a 1:10 ratio of “positive” (containing loci) to “negative” (not containing loci) images. Due to this imbalance and limited labeled data, we employ K-fold cross-validation (K=5) for robust evaluation. Specifically, for each train-test split, four folds (positive: 729, negative: 7,290) are used for training, and one fold (positive: 183, negative: 1,830) is used for testing. A simple linear classifier with two fully connected layers is concatenated to the DeepAndes backbone. Prior to the five-fold cross-validation experiments, we ran several random seed trials to identify optimized hyperparameters for model fine-tuning. Additionally, to evaluate the few-shot learning capability of the pre-trained transformer backbones, we assess the classification task at different downstream dataset scales (10%, 30%, 50%, and 100% of the original classification dataset).

**Evaluation Metric.** The F1 score is used as the evaluation metric for this imbalanced classification task, as it combines both sensitivity and precision,

$$F1 \text{ score} = 2 * \frac{Precision * Sensitivity}{Precision + Sensitivity}$$

where the Sensitivity (or Recall) determines the accuracy of the minority class classification and Precision indicates the probability of its correct detection. Particularly, the Precision is calculated as the ratio of true positives (TP) to the total number of positive predictions, which is the sum of true positives and false positives (FP):

$$Precision = \frac{TP}{TP + FP}$$

The Recall is the ratio of TP to the total number of actual positive samples, which is the sum of true positives and false negatives (FN):

$$Recall = \frac{TP}{TP + FN}$$

**2) Image Instance Retrieval:** For the image instance retrieval task, we use all available labeled data (both **training** and **testing**) from the loci classification task and probe the frozen backbones through the class token

to form a database. This can be implemented using FAISS library [16]. In our database, for each positive sample, we use it as the input query and retrieve the top  $k$  most similar instances from the database, ranking them based on their cosine similarity to the query image. The precision of retrieving image instances from the same class as the query image (positive samples) is then evaluated for different pre-trained backbones.

**Evaluation Metric.** To evaluate the performance of image instance retrieval task, we use  $mAP@k$  (mean Average Precision within top- $k$  retrieved samples). The derivation is following. When performing image retrieval, we rank the images based on their similarity to the query image, where rank 1 corresponds to the most similar image, rank 2 to the second most similar image, and so on. The precision at a given rank is defined as the proportion of relevant images (those belonging to the same class as the query image) retrieved up to that rank:

$$Precision@i = \frac{\text{Number of images retrieved at rank } i}{i}$$

Then, the Average Precision (AP) for a given query is the average of precision values at different ranks (up to  $k$ ) where relevant images are retrieved:

$$AP = \frac{1}{N} \sum_{i=1}^N Precision@i$$

where  $N$  is the number of relevant images for that query. Finally,  $mAP@k$  is calculated by averaging the AP over all queries in the dataset:

$$mAP@k = \frac{1}{Q} \sum_{q=1}^Q AP_q$$

where  $Q$  is the total number of queries (number of positive samples in the database). In this work, we use multiple  $k$  values for evaluation, ranging from 5, 20, and 50, to 100.

3) **Few-shot Loci Segmentation:** To evaluate the foundation model’s few-shot learning performance on pixel-level feature recognition, we focus on the semantic segmentation of three specific types of loci (i.e., active buildings, active corrals, and archaeological corrals) in this study. Three small-scale binary semantic segmentation datasets are constructed: the Active Buildings dataset contains 48 images, the Active Corrals dataset contains 55 images, and the Archaeological Corrals dataset contains 46 images. K-fold cross-validation (K=5) is employed for robust evaluation. The segmentation head is simple and directly utilizes the learned features from pre-training. Specifically, patch-level features are extracted from the images using patch tokens and concatenated with a simple one-layer linear segmentation head. Prior to the five-fold cross-validation

experiments, we conducted several random seed trials to identify optimized hyperparameters for model training. Similar to the few-shot loci classification, we also evaluate the pre-trained model’s segmentation performance at different downstream dataset scales, including 10, 20, and 30 images.

**Evaluation Metric.** The Dice Similarity Coefficient (DSC) is used to evaluate the precision of loci segmentation in this work. This metric ensures precise pixel-level overlap between the predicted and ground truth masks. It emphasizes the importance of accurate segmentation in smaller, localized areas, which aligns with our segmentation datasets, where loci objects (foreground) are typically much smaller than the background areas.

$$DSC = \frac{2 \times TP}{2 \times TP + FP + FN}$$

- $TP$  represents the number of pixels correctly classified as the loci in both the predicted and ground truth masks.
- $FP$  represents the number of pixels incorrectly classified as loci in the predicted mask but not in the ground truth mask.
- $FN$  represents the number of pixels incorrectly classified as background in the predicted mask but actually belonging to the loci in the ground truth mask.

## V. RESULTS

### A. Imbalanced Loci Classification

Both Table I and Figure 5 display the five-fold cross-validation results for imbalanced loci classification. Specifically, Table I details the mean F1 scores, Precision (P), and Recall (R). Figure 5a illustrates the Precision-Recall (PR) curves, while Figure 5b presents the confusion matrices for each model, highlighted using a representative hold-out fold.

As shown in Table I, the Scratch model struggles with limited and imbalanced data, converging only at  $N_{\text{train}} = 729$  (F1 = 0.544) and  $N_{\text{train}} = 365$  (F1 = 0.402), while failing entirely at smaller datasets. In contrast, DeepAn-des models with SSL pre-training improve consistently across all dataset sizes. At  $N_{\text{train}} = 72$ , FM3M achieves F1 = 0.83, Recall = 0.825, and Precision = 0.837, effectively balancing false positives and false negatives. Smaller pre-trained models decline in performance, with FM300K dropping to F1 = 0.728, R = 0.671, P = 0.812 and FM30K to F1 = 0.418, R = 0.342, P = 0.556, reflecting higher false negatives due to class imbalance. The same pattern can also be observed in Figure 5a and b, where the highlighted blue PR curve has the highest PR-AUC and seemingly fewer false negatives compared to the other models.

Model	$N_{\text{train}} = 729$			$N_{\text{train}} = 365$			$N_{\text{train}} = 218$			$N_{\text{train}} = 72$		
	F1	Recall	Precision	F1	Recall	Precision	F1	Recall	Precision	F1	Recall	Precision
Scratch	0.544	0.457	0.735	0.402	0.288	0.692	0.26	0.241	0.52	–	–	–
FM30K	0.654	0.595	0.747	0.596	0.522	0.701	0.516	0.434	0.672	0.418	0.342	0.556
FM300K	0.824	0.806	0.848	0.804	0.795	0.817	0.788	0.782	0.792	0.728	0.671	0.812
FM3M	<b>0.886</b>	<b>0.876</b>	<b>0.894</b>	<b>0.872</b>	<b>0.872</b>	<b>0.875</b>	<b>0.862</b>	<b>0.866</b>	<b>0.857</b>	<b>0.830</b>	<b>0.825</b>	<b>0.837</b>

TABLE I: **Performance on Imbalanced Loci Classification.** The mean F1 scores, Recall (R), and Precision (P) from the five-fold cross-validation are presented. We compare performance of four model backbones: ViT-L trained from scratch, and DeepAndes pre-trained using 30K, 300K, and 3M images.  $N_{\text{train}}$  represents the scale of “positive” images (containing loci) in the training dataset. The “positive” to “negative” ratio is **1:10** in both training and testing set. Entries marked with “–” indicate that the experiments do not converge or the values are not supported. For clarity, the highest metric values of each few-shot dataset evaluation are highlighted in **bold**.

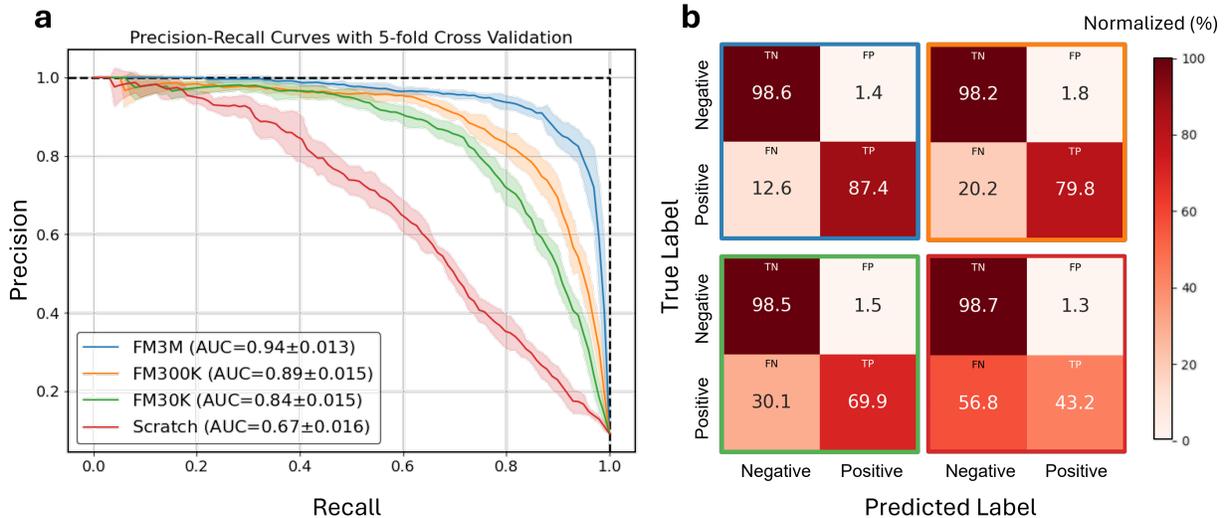


Fig. 5: **Performance on Imbalanced Loci Classification:** Precision-Recall curves from five-fold cross-validation (a). Confusion matrices (normalized) for each model shown for a representative hold-out fold (b).

Additionally, as Table I indicated, FM3M maintains strong results with few labels; at  $N_{\text{train}} = 72$ , the larger pre-trained FM3M emerges as a better few-shot learner, achieving performance comparable to FM300K at  $N_{\text{train}} = 729$ . These results underscore the advantages of large-scale SSL pre-training in managing imbalanced data and enhancing few-shot learning.

### B. Image Instance Retrieval

For loci image instance retrieval task, the performance of four frozen ViT-L backbones (Scratch, FM30K, FM300K, and FM3M) is compared in Table II. As demonstrated, for the image retrieval task, larger  $k$  values result in lower mAP@ $k$  scores, indicating that retrieving more samples from the database also introduces more

irrelevant ones relative to the query image class. It is also evident that pre-training improves retrieval performance. The Scratch model, without any pre-training, achieves the lowest mAP@ $k$  scores across all choices of  $k$  for this evaluation, with mAP@5 starting at 0.444. In contrast, pre-trained DeepAndes models—FM30K, FM300K, and FM3M—show progressively higher mAP values, with FM3M achieving the highest performance, with mAP@5 of 0.869. Figure 6 provides qualitative visualizations of the retrieved examples. As pre-training scales up, FM3M retrieves more relevant images (highlighted in blue boxes) that correctly match both the image class and features of the query image. In contrast, FM30K and the Scratch model retrieve incorrect image classes (highlighted in red boxes) among the top-5 retrieved ex-

amples. These results clearly demonstrate that SSL pre-training, particularly at a larger scale, enhances image-level feature representations and improves instance retrieval accuracy.

Model	mAP@5	mAP@20	mAP@50	mAP@100
Scratch	0.444	0.391	0.330	0.296
FM30K	0.498	0.435	0.384	0.350
FM300K	0.757	0.677	0.597	0.532
FM3M	<b>0.869</b>	<b>0.804</b>	<b>0.744</b>	<b>0.690</b>

TABLE II: **Performance on Image Instance Retrieval.** The mAP@ $k$  (mean Average Precision within top- $k$  retrieved samples) for different choices of  $k$  (5, 20, 50, and 100) is presented. Four frozen ViT-L backbones are compared: ViT-L (scratch) and DeepAndes pre-trained using 30K, 300K, and 3M images. For clarity, the highest metric values are highlighted in **bold**.

### C. Few-shot Loci Segmentation

For few-shot loci segmentation tasks, Table III summarizes the models’ performance on three loci types (i.e., active buildings, active corrals, and archaeological corrals). It is evident that as the pre-training scales increase—from 30K to 300K, and then to 3M—the model’s performance in both transfer learning (frozen) and fine-tuning improves compared to training from scratch. As highlighted in Table III, the best model performance across all few-shot datasets and tasks comes from the FM3M. For **active/archaeological corrals** segmentation tasks, simply training the linear segmentation head on top of the frozen FM3M backbone achieves the second highest DSC score across all experiments, demonstrating the benefits of our million-scale pre-training. On the other hand, models trained from scratch exhibit relatively low DSC scores, especially with small datasets. For example, at  $N_{train} = 10$ , the DSC score for the Scratch model (frozen) on the **Active Corrals** dataset is only 11.6. In contrast, frozen pre-trained models show an improvement of over 30% starting from FM30K (45.4), with DSC scores continuing to rise as the pre-training scale increases. Although model performance improves with fine-tuning the entire ViT-L, transfer learning on the frozen FM3M backbone still exhibits DSC scores that are either comparable to or surpass those of models with smaller pre-training scales. This improvement is particularly noticeable when downstream data is very limited (e.g.,  $N_{train} = 10$ ) across three segmentation tasks. Overall, these results highlight the effectiveness of million-scale DeepAndes in few-shot learning tasks, where the model can generalize well with

limited labeled data. This makes it especially valuable in fields like archaeology, where data annotation is scarce.

Additionally, we present the qualitative visualizations, paired with the Table III results, for the three types of loci segmentation tasks in Figure 7. An example patch and its segmentation results from different models are shown. As shown, segmenting archaeological corrals (bottom panel) is the most challenging, as there is little to no difference between the inside and outside, except for the boundary pixels. Segmenting active corrals (top panel) is the easiest, as the animals’ evidence is consistently darker. Active building segmentation (middle panel) shows medium performance overall. The variability in building colors and their small size makes it difficult for the model to converge. As shown, even for the best-performing task (top panel), transfer learning failed for the scratch model and less pre-trained backbones when the labels were limited ( $N_{train} = 10, 20, 30$ ). Fine-tuning FT-FM3M, which uses the DeepAndes (3M) backbone, effectively captures foreground pixels and achieves higher precision with a simple linear segmentation head. A similar pattern is observed in active building segmentation (middle panel), where transfer learning with FM3M outperforms the remaining frozen backbones and is comparable to FT-FM300K. For the most challenging task (bottom panel), archaeological corral segmentation, FM3M and FT-FM3M clearly outperform other models, which struggle to capture the foreground pixels and produce many false positives.

### D. DeepAndes Training and Fine-tuning Summary

Lastly, this section provides a summary of the pre-training and fine-tuning of DeepAndes (0.3B parameters vision transformer) in Figure 8. Particularly, it includes training data, model, environmental impact, and evaluation strategies utilized in this work. To observe fine-tuning efficiency, we selected the imbalanced loci image classification task as representative. The first 10 epochs of the training log are displayed because the model tends to overfit beyond this point. As demonstrated, an increased pre-training scale promotes faster training convergence by monitoring the running loss and running loss AUC (the smaller, the better). This aligns with the foundation model’s purpose to accommodate a broad range of tasks through few-shot or zero-shot learning.



Fig. 6: **Examples of Retrieved Images from Different Pre-trained Models.** The left column displays two example query images, with arrows pointing to key archaeological features. Both query images are from the positive class in the imbalanced loci classification dataset. The right columns show the top-5 retrieved images based on cosine similarity. The images with a red box indicate incorrect class retrieval, while the blue box highlights correct retrieval of both the image class and all archaeological features in the query image.

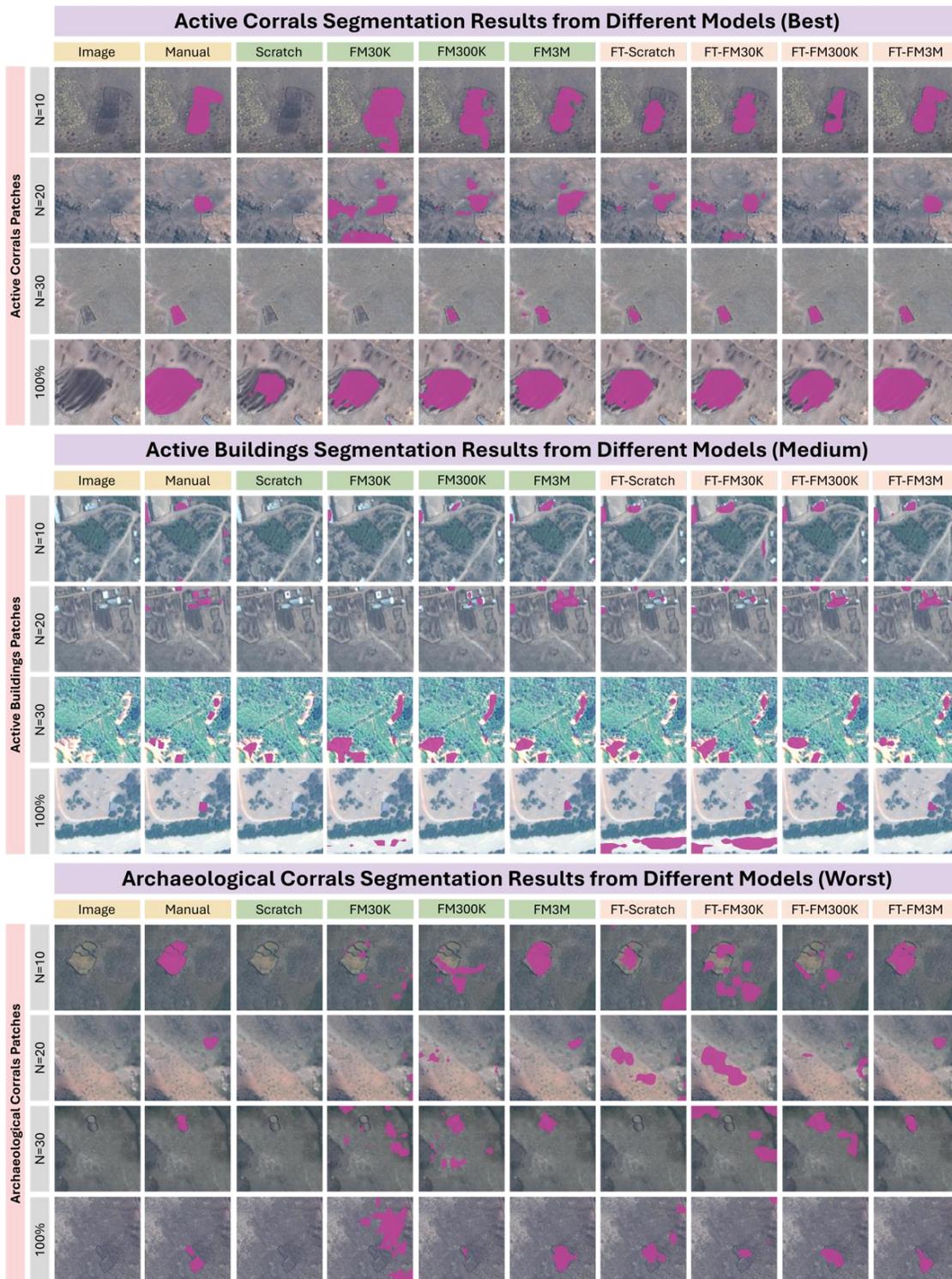


Fig. 7: **Qualitative Results of Loci Segmentation.** “FT-” denotes fine-tuning. For convenience, image patches are shown in RGB. Superior performance are from either FT-FM3M or FM3M.

Backbone	Model	Active Buildings <sup>1</sup>				Active Corrals <sup>2</sup>				Archaeological Corrals <sup>3</sup>			
		100%	$N_{\text{train}} = 10$	$N_{\text{train}} = 20$	$N_{\text{train}} = 30$	100%	$N_{\text{train}} = 10$	$N_{\text{train}} = 20$	$N_{\text{train}} = 30$	100%	$N_{\text{train}} = 10$	$N_{\text{train}} = 20$	$N_{\text{train}} = 30$
Frozen	Scratch	29.1	8.2	21.3	20.6	38.3	11.6	14.6	21.0	8.4	–	–	–
	FM30K	27.0	20.0	23.0	24.7	45.7	45.4	44.8	45.5	11.9	9.7	11.9	12.1
	FM300K	55.2	41.4	47.3	51.4	60.9	55.5	56.5	59.3	27.8	19.0	17.8	24.8
	FM3M	58.3	<u>52.3</u>	56.2	58.2	<u>70.4</u>	<u>65.5</u>	<u>67.1</u>	<u>68.2</u>	<u>67.3</u>	<b>56.3</b>	<u>61.3</u>	<u>63.0</u>
Finetuned	Scratch	47.1	42.1	47.8	45.9	64.2	57.9	62.0	62.5	11.6	15.1	15.4	15.4
	FM30K	53.8	44.0	52.8	53.2	66.4	61.3	63.7	64.8	17.7	16.2	15.6	18.4
	FM300K	<u>59.7</u>	51.4	<u>57.2</u>	<u>58.8</u>	69.2	59.3	64.6	67.8	33.9	15.3	17.9	29.4
	FM3M	<b>69.3</b>	<b>57.8</b>	<b>63.7</b>	<b>66.5</b>	<b>81.1</b>	<b>70.8</b>	<b>73.5</b>	<b>75.5</b>	<b>84.8</b>	<u>51.3</u>	<b>72.5</b>	<b>81.0</b>

<sup>1</sup> **Active Buildings Dataset** contains 48 images featuring modern buildings against dense forest or vegetation backgrounds.

<sup>2</sup> **Active Corrals Dataset** contains 55 images, each showing corrals with visible evidence of animal use.

<sup>3</sup> **Archaeological Corrals Dataset** contains 46 images showing corrals where signs of use have disappeared.

TABLE III: **Performance on Few-shot Segmentation of Three Loci Types.** Mean DSC scores from five-fold cross-validation are presented, which include both transfer learning (frozen backbones) and fine-tuning (unfrozen backbones) across four models: Scratch, FM30K, FM300K, and FM3M. Specifically,  $N_{\text{train}}$  indicates the scale of the training dataset. **Bold** entries denote the highest DSC score for each few-shot setting, while **underscored** entries indicate the second highest. Entries marked with “–” indicate that the experiments do not converge or the values are not supported.

## DeepAndes Summary

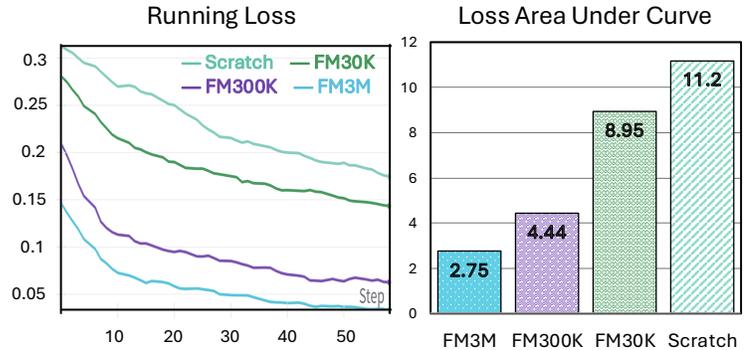
### Training

Data: 3 million 256 x 256 images, 0.5/0.3 m  
 Model: ViT-L (0.3B params), Patch size 14  
 embedding 1024, 16 heads, MLP FFN

### Environment Impact

Hardware: 8 x NVIDIA A100-80GB  
 GPU Hours: 156.92  
 Carbon Emitted: 0.2t CO<sub>2</sub>e

### Fine-tuning Efficiency



**Downstream** Classification (Segmentation): CLS (Patch) Token + Linear Classifier (Segmentation) Head

**Evaluation** Image Retrieval: CLS Token + Similarity or Distance Comparison

Fig. 8: **Summarization of DeepAndes Model and Fine-Tuning Efficiency Comparison.** The right panel illustrates the running loss experimental logs (10 epochs) from an imbalanced loci classification experiment, highlighting DeepAndes (FM3M)’s rapid convergence

## VI. DISCUSSION

The evaluation of downstream tasks indicates that DeepAndes enhances visual feature representations, which improve as pre-training scales increase. With

the 3-million scale pre-training, DeepAndes achieves over 83% accuracy in few-shot imbalance classification, surpassing other models, including those trained from scratch in a supervised manner, across various fine-

tuning scales. Similar patterns are evident in image retrieval and segmentation tasks, where the feature embeddings from the pre-trained backbone show potential for both image-level and pixel-level remote sensing tasks when labeled data is limited.

There are also some limitations to this work. First, the data for self-supervised pre-training should undergo more thorough data curation. In [36], the training data undergo a comprehensive data de-duplication process and employ image retrieval to refine the uncurated data source with curated sources such as ImageNet-22 [14]. In our work, we try to sample from diverse archaeological regions while scaling up the datasets, yet we lack curated datasets and detailed prior knowledge of the survey regions, unlike established benchmarks. The archaeological regions are not as distinct as these natural imaging datasets, suggesting a need for further de-duplication beyond excluding featureless samples. Additionally, as we refine the DeepAndes model, it could be utilized to analyze and help reduce duplicate samples in the pre-training dataset.

For the future and in progress work, the DeepAndes will be integrated into the the GeoPACHA geospatial platform [46] as a web application, collaborating regional experts and their teams. Meanwhile, depending on the remote sensing tasks, the human-in-the-loop is also included as experts will also verify and compare the autonomous model prediction with their data annotations. The curated data will be used for next round of foundation model fine-tuning to acquire more a specialized model, DeepAndesArch, for remote sensing tasks. From a model perspective, the proposed foundation model offers broad utility and can be further fine-tuned for advanced object detection tasks. By incorporating complex transformer heads, such as Co-DETR [54], it can provide valuable support to archaeological survey teams.

## VII. CONCLUSION

In this work, we present DeepAndes, the first AI vision foundation model for multi-spectral remote sensing data for social and earth science applications, utilizing the SOTA DINOv2 framework. DeepAndes is pre-trained on 3 million WorldView-2 and WorldView-3 satellite images with eight spectral bands. Through extensive downstream experiments on three prevalent computer vision tasks in archaeology—imbalanced image classification, image instance retrieval, and semantic segmentation—DeepAndes demonstrates its effectiveness in both image-level and pixel-level feature representation, as well as in few-shot learning capabilities. The pre-trained DeepAndes will be integrated into the GeoPACHA web app to expand the scale of our archaeological surveys in the Andes with human-in-the-loop

verification. The broad utility of the proposed foundation model can be further fine-tuned for more advanced object detection tasks using complex transformer heads, such as Co-DETR. As we collect more data with the AI-assisted GeoPACHA tool, experts will be able to contribute more effectively as observers and analysts in the archaeological workflow.

## ACKNOWLEDGEMENTS

This work was supported by the following grants: National Endowment for the Humanities Level III Digital Enhancement Grant (Award ID Number HAA-293452-23), National Science Foundation IIS-III: Medium: Collaborative Research (Award number 2106717); National Science Foundation Collaborative Research: Research Infrastructure: HNDS-I (Award Numbers 2419793 and 2419794); Vanderbilt University Scaling Success Grant, Vanderbilt University Discovery Grant. The project also benefited from the infrastructural support of the Vanderbilt University Data Science Institute (GPU computational infrastructure) and the Vanderbilt University HLRB Lab (computational infrastructure), and the Spatial Analysis Research Laboratory (geospatial computational infrastructure).

## REFERENCES

- [1] 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*, pages 8203–8215, 2022. 2
- [2] Susan Alcock and John Cherry. *Side-by-Side Survey : Comparative Regional Studies in the Mediterranean World*. Oxbow Books, 2016. 1
- [3] K Jerry Allwine, Joseph H Shinn, Gerald E Streit, Kirk L Clawson, and Mike Brown. Overview of urban 2000: A multiscale field study of dispersion through an urban environment. *Bulletin of the American Meteorological Society*, 83(4):521–536, 2002. 1
- [4] Muhammad Awais, Muzammal Naseer, Salman Khan, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, and Fahad Shahbaz Khan. Foundation models defining a new era in vision: a survey and outlook. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 2
- [5] 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*, pages 10181–10190, 2021. 2
- [6] Andrew Bevan and James Conolly. *Mediterranean islands, fragile communities and persistent landscapes: Antikythera in long-term perspective*. Cambridge University Press, 2013. 1
- [7] Brian R Billman and Gary M Feinman. Settlement pattern studies in the americas: fifty years since virú. (*No Title*), 1999. 1
- [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers, 2021. 4, 5
- [9] Jesse Casana. Regional-scale archaeological remote sensing in the age of big data: Automated site discovery vs. brute force methods. *Advances in Archaeological Practice*, 2(3):222–233, 2014. 1
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020. 4

- [11] 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. 4
- [12] Can Cui, Ruining Deng, Junlin Guo, Quan Liu, Tianyuan Yao, Haichun Yang, and Yuankai Huo. Enhancing physician flexibility: Prompt-guided multi-class pathological segmentation for diverse outcomes. In *2024 IEEE EMBS International Conference on Biomedical and Health Informatics (BHI)*, pages 1–8. IEEE, 2024. 2
- [13] Can Cui, Ruining Deng, Junlin Guo, Quan Liu, Tianyuan Yao, Haichun Yang, and Yuankai Huo. Pfps: Prompt-guided flexible pathological segmentation for diverse potential outcomes using large vision and language models. *arXiv preprint arXiv:2407.09979*, 2024. 2
- [14] 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. 15
- [15] Philippe 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
- [16] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. *arXiv preprint arXiv:2401.08281*, 2024. 9
- [17] Timothy K Earle et al. Archaeological field research in the upper mantaro, peru, 1982-1983: Investigations of inka expansion and exchange. (*No Title*), 1987. 1
- [18] Gary M Feinman, Stephen A Kowalewski, Laura Finsten, Richard E Blanton, and Linda Nicholas. Long-term demographic change: A perspective from the valley of oaxaca, mexico. *Journal of Field Archaeology*, 12(3):333–362, 1985. 1
- [19] 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. IEEE, 2023. 2
- [20] RM Fuller, GB Groom, S Mugisha, P Ipulet, D Pomeroy, A Katende, R Bailey, and R Ogutu-Ohwayo. The integration of field survey and remote sensing for biodiversity assessment: a case study in the tropical forests and wetlands of sango bay, uganda. *Biological conservation*, 86(3):379–391, 1998. 1
- [21] Manuel Grizonnet, Julien Michel, Victor Poughon, Jordi Inglada, Mickaël Savinaud, and Rémi Cresson. Orfeo toolbox: Open source processing of remote sensing images. *Open Geospatial Data, Software and Standards*, 2(1):15, 2017. 8
- [22] Junlin Guo, Siqi Lu, Can Cui, Ruining Deng, Tianyuan Yao, Zhewen Tao, Yizhe Lin, Marilyn Lionts, Quan Liu, Juming Xiong, et al. How good are we? evaluating cell ai foundation models in kidney pathology with human-in-the-loop enrichment. *arXiv preprint arXiv:2411.00078*, 2024. 2
- [23] Junlin Guo, Siqi Lu, Can Cui, Ruining Deng, Tianyuan Yao, Zhewen Tao, Yizhe Lin, Marilyn Lionts, Quan Liu, Juming Xiong, et al. Assessment of cell nuclei ai foundation models in kidney pathology. In *Medical Imaging 2025: Image Perception, Observer Performance, and Technology Assessment*, volume 13409, pages 76–82. SPIE, 2025. 2
- [24] Xin Guo, Jiangwei Lao, Bo Dang, Yingying Zhang, Lei Yu, Lixiang Ru, Liheng Zhong, Ziyuan Huang, Kang Wu, Dingxiang Hu, et al. Skysense: A multi-modal remote sensing foundation model towards universal interpretation for earth observation imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27672–27683, 2024. 4, 7
- [25] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning, 2020. 4
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 4
- [27] Licheng Jiao, Zhongjian Huang, Xiaoqiang Lu, Xu Liu, Yuting Yang, Jiakuan 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
- [28] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey, 2019. 3
- [29] Jing Ke, Junchao Zhu, Xin Yang, Haolin Zhang, Yuxiang Sun, Jiayi Wang, Yizhou Lu, Yiqing Shen, Sheng Liu, Fusong Jiang, et al. Tshfna-examiner: A nuclei segmentation and cancer assessment framework for thyroid cytology image. *Journal of Shanghai Jiaotong University (Science)*, 29(6):945–957, 2024. 2
- [30] Chenxin Li, Xinyu Liu, Wuyang Li, Cheng Wang, Hengyu Liu, Yifan Liu, Zhen Chen, and Yixuan Yuan. U-kan makes strong backbone for medical image segmentation and generation. *arXiv preprint arXiv:2406.02918*, 2024. 2
- [31] Yansheng Li, Jieyi Tan, Bo Dang, Mang Ye, Sergey A Bartalev, Stanislav Shinkarenko, Linlin Wang, Yingying Zhang, Lixiang Ru, Xin Guo, et al. Unleashing the potential of remote sensing foundation models via bridging data and computility islands. *The Innovation*, 2025. 2
- [32] Fan Liu, Delong Chen, Zhangqingyun Guan, Xiaocong Zhou, Jiale Zhu, Qiaolin Ye, Liyong Fu, and Jun Zhou. Remoteclip: A vision language foundation model for remote sensing. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. 2
- [33] Siqi Lu, Junlin Guo, James R Zimmer-Dauphinee, Jordan M Nieusma, Xiao Wang, Steven A Wernke, Yuankai Huo, et al. Vision foundation models in remote sensing: A survey. *IEEE Geoscience and Remote Sensing Magazine*, 2025. 2, 3, 4, 5
- [34] Yuchi Ma, Shuo Chen, Stefano Ermon, and David B. Lobell. Transfer learning in environmental remote sensing. *Remote Sensing of Environment*, 301:113924, 2024. 4
- [35] 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
- [36] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 5, 6, 8, 15
- [37] Jeffrey R Parsons, Charles M Hastings, and Ramiro Matos. *Prehispanic Settlement Patterns in the Upper Mantaro and Tarma Drainages, Junín, Peru: The Tarama-Chinchaycocha Region, Vol. 1, Parts 1 and 2*, volume 34. U OF M MUSEUM ANTHRO ARCHAEOLOGY, 2000. 1
- [38] 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
- [39] Masato Sakai, Akihisa Sakurai, Siyuan Lu, Jorge Olano, Conrad M Albrecht, Hendrik F Hamann, and Marcus Freitag. Ai-accelerated nazca survey nearly doubles the number of known figurative geoglyphs and sheds light on their purpose. *Proceedings of the National Academy of Sciences*, 121(40):e2407652121, 2024. 2
- [40] Neha Sharma, Reecha Sharma, and Neeru Jindal. Machine learning and deep learning applications-a vision. *Global Transitions Proceedings*, 2(1):24–28, 2021. 2
- [41] Nicholas Tripcevich and Steven A Wernke. On-site recording of excavation data using mobile gis. *Journal of Field Archaeology*, 35(4):380–397, 2010. 1
- [42] Parker VanValkenburgh and J Andrew Dufton. Big archaeology: Horizons and blindspots, 2020. 1
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polo-

- sukhin. Attention is all you need, 2023. [3](#), [4](#)
- [44] Eugene Vorontsov, Alican Bozkurt, Adam Casson, George Shaikovski, Michal Zelechowski, Siqi Liu, Kristen Severson, Eric Zimmermann, James Hall, Neil Tenenholtz, et al. Virchow: A million-slide digital pathology foundation model. *arXiv preprint arXiv:2309.07778*, 2023. [5](#), [7](#)
- [45] Xinye Wanyan, Sachith Seneviratne, Shuchang Shen, and Michael Kirley. Extending global-local view alignment for self-supervised learning with remote sensing imagery, 2024. [4](#), [7](#)
- [46] Steven A Wernke, Parker Van Valkenburgh, James Zimmer-Dauphinee, Bethany Whitlock, Giles Spence Morrow, Ryan Smith, Douglas Smit, Grecia Roque Ortega, Kevin Ricci Jara, Daniel Plekhov, et al. Large-scale, collaborative imagery survey in archaeology: the geospatial platform for andean culture, history and archaeology (geopacha). *Antiquity*, 98(397):155–171, 2024. [15](#)
- [47] Gordon R Willey. Prehistoric settlement patterns in the viru valley, peru. *Bureau of American Ethnology Bulletin*, 1953. [1](#)
- [48] David John Wilson. Prehispanic settlement patterns in the lower santa valley, peru: a regional perspective on the origins and development of complex north coast society. (*No Title*), 1988. [1](#)
- [49] Jiachen Xu, Junlin Guo, James Zimmer-Dauphinee, Quan Liu, Yuxuan Shi, Zuhayr Asad, D Mitchell Wilkes, Parker Van Valkenburgh, Steven A Wernke, and Yuankai Huo. Semi-supervised contrastive learning for remote sensing: identifying ancient urbanization in the south-central andes. *International journal of remote sensing*, 44(6):1922–1938, 2023. [2](#), [7](#)
- [50] Jialin Yue, Tianyuan Yao, Ruining Deng, Siqi Lu, Junlin Guo, Quan Liu, Juming Xiong, Mengmeng Yin, Haichun Yang, and Yuankai Huo. Glofinder: Ai-empowered qpath plugin for wsi-level glomerular detection, visualization, and curation. *Journal of Pathology Informatics*, page 100433, 2025. [2](#)
- [51] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *International Journal of Machine Learning and Cybernetics*, pages 1–65, 2024. [2](#), [4](#)
- [52] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021. [5](#)
- [53] Junchao Zhu, Ruining Deng, Tianyuan Yao, Juming Xiong, Chongyu Qu, Junlin Guo, Siqi Lu, Mengmeng Yin, Yu Wang, Shilin Zhao, et al. Asign: An anatomy-aware spatial imputation graphic network for 3d spatial transcriptomics. *arXiv preprint arXiv:2412.03026*, 2024. [2](#)
- [54] Zhuofan Zong, Guanglu Song, and Yu Liu. Detsr with collaborative hybrid assignments training. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6748–6758, 2023. [15](#)

## VIII. BIOGRAPHY SECTION

**Junlin Guo** (junlin.guo@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) and the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University. His research interests include deep learning, foundation models, and their applications in remote sensing and medical image analysis.

**James R Zimmer-Dauphinee** (james.r.zimmer-dauphinee@vanderbilt.edu) is a postdoctoral researcher with the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University. 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** (jordan.m.nieusma@vanderbilt.edu) is a data scientist and research assistant with the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University. Her research interests include data science, deep learning and efficient learning of large foundation models.

**Siqi Lu** (siqi.lu@vanderbilt.edu) is a Master’s student with the Biomedical Data Representation and Learning Lab (HRLB) and the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University. Her research interests include deep learning, medical image analysis, and software engineering.

**Quan Liu** (quan.liu@vanderbilt.edu) is a researcher with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include deep learning, artificial intelligence, and medical image analysis.

**Ruining Deng** (rud4004@med.cornell.edu) is a research fellow with Weill Cornell Medicine in New York and the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include deep learning, artificial intelligence, and medical image analysis, which aim to explore clinical knowledge and assist diagnosis in a data-driven way.

**Can Cui** (can.cui.1@vanderbilt.edu) is a researcher with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. Her research interests include deep learning, data science, and medical image analysis, which aims for multi-modal learning for disease diagnosis and prognosis.

**Jialin Yue** (jialin.yue@vanderbilt.edu) is a Master’s student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. Her research interests include medical image processing and machine learning.

**Yizhe Lin** (yizhe.lin@vanderbilt.edu) is an undergraduate research assistant with the Spatial Analysis Research Laboratory (SARL) at Vanderbilt University.

**Tianyuan Yao** (tianyuan.yao@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include medical image analysis, deep learning, computer vision and their applications in pathology and radiology imaging.

**Juming Xiong** (juming.xiong@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include medical image analysis, deep learning, and data science.

**Junchao Zhu** (junchao.zhu@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include medical image analysis, deep learning, computer vision and their applications in large-scale pathology image processing.

**Chongyu Qu** (chongyu.qu@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. His research interests include medical image analysis, deep learning, computer vision and post-training model quantization for large scale medical vision tasks.

**Yuechen Yang** (yuechen.yang@vanderbilt.edu) is a PhD student with the Biomedical Data Representation and Learning Lab (HRLB) at Vanderbilt University. Her research interests include medical image analysis, data science, and medical image processing, with a particular focus on Pathomics.

**Mitchell Wilkes** (mitch.wilkes@vanderbilt.edu) is an Associate Professor of Electrical and Computer Engineering at Vanderbilt University. Dr. Wilkes's research focuses on digital signal processing, image processing and computer vision, digital signal processing hardware, structurally adaptive systems, sonar, and signal modeling. Dr. Wilkes's intellectual neighborhoods also include Biomedical Imaging and Biophotonics, Surgery, and Engineering.

**Xiao Wang** (wangx2@ornl.gov) is 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** (parker\_vanvalkenburgh@brown.edu) is an Associate Professor of Anthropology and Archaeology at Brown University. Dr. VanValkenburgh's 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. Dr. VanValkenburgh co-directs the Paisajes Arqueológicos de Chachapoyas (PACHa) project and GeoPACHA (Geospatial Platform for Andean Culture, History, and Archaeology).

**Steven A Wernke** (s.wernke@vanderbilt.edu) is Professor and Chair of Anthropology at Vanderbilt University, director of the Spatial Analysis Research Laboratory (SARL), and director of the Vanderbilt Institute for Spatial Research. Dr. Wernke is an archaeologist and historical anthropologist of the Andean region of South America. His research combines archaeology and history, prehispanic and colonial studies, as well as anthropology and cultural geography. His interests center on the lived experiences of Indigenous communities across the Spanish invasion of the Americas, and on long-term, large-scale networks, social formations, and human-environment interactions across the Andes.

**Yuankai Huo** (yuankai.huo@vanderbilt.edu) is an Assistant Professor of Computer Science, and Electrical and Computer Engineering, as well as the Director of the Biomedical Data Representation and Learning Lab (HRLB Lab) at Vanderbilt University. Additionally, he is an Assistant Professor of Pathology, Microbiology, and Immunology at Vanderbilt University Medical Center. Dr. Huo's current research specializes in high-dimensional multi-modal data analysis, computational pathology and radiology, and medical computer vision. Dr. Huo has received prestigious awards, including the Charles E. Ives Journal Award from the Society for Imaging Science and Technology, the Early Career Achievement Award from the Society for Imaging Informatics in Medicine, and the NAIRR Pilot award from NSF. He is a Senior Member of IEEE and a lifetime member of SPIE, contributing as an organization committee member and area chair for leading medical image analysis conferences such as MICCAI, MIDL, and ISBI. His ongoing efforts are dedicated to advancing next-generation AI algorithms for ultra-high-resolution imaging and non-imaging data analysis.