

# Domain Adaptable Self-supervised Representation Learning on Remote Sensing Satellite Imagery

Muskaan Chopra<sup>γ, \*</sup>, Prakash Chandra Chhipa<sup>δ, \*</sup>, Gopal Mengi<sup>γ, \*</sup>, Varun Gupta<sup>γ</sup> and Marcus Liwicki<sup>δ</sup>

<sup>δ</sup> Machine Learning Group, EISLAB, Luleå Tekniska Universitet, Luleå, Sweden

{prakash.chandra.chhipa, marcus.liwicki}@ltu.se

<sup>γ</sup> Chandigarh College of Engineering and Technology, Punjab University, Chandigarh, India

chopramuskaan47@gmail.com, gopalmengi@gmail.com, varungupta@cet.ac.in

\* co-first authors with equal contribution

**Abstract**—This work presents a novel domain adaption paradigm for studying contrastive self-supervised representation learning and knowledge transfer using remote sensing satellite data. Major state-of-the-art remote sensing visual domain efforts primarily focus on fully supervised learning approaches that rely entirely on human annotations. On the other hand, human annotations in remote sensing satellite imagery are always subject to limited quantity due to high costs and domain expertise, making transfer learning a viable alternative. The proposed approach investigates the knowledge transfer of self-supervised representations across the distinct source and target data distributions in depth in the remote sensing data domain. In this arrangement, self-supervised contrastive learning-based pretraining is performed on the source dataset, and downstream tasks are performed on the target datasets in a round-robin fashion. Experiments are conducted on three publicly available datasets, UC Merced Landuse (UCMD), SIRI-WHU, and MLRSNet, for different downstream classification tasks versus label efficiency. In self-supervised knowledge transfer, the proposed approach achieves state-of-the-art performance with label efficiency labels and outperforms a fully supervised setting. A more in-depth qualitative examination reveals consistent evidence for explainable representation learning. The source code and trained models are published on GitHub<sup>1</sup>.

**Index Terms**—contrastive learning; self-supervised learning; representation learning, domain adaptation, remote sensing, satellite image

## I. INTRODUCTION

To formulate the policies and schemes, the region’s geographical and demographic information and its efficient representation are essential [1] [2]. Visual interpretation of aerial and space images is the most common method of producing topographic and thematic maps. Satellite images are also used to classify different types of crops using deep learning techniques [3] [4]. Today many high-resolution satellites can be relied upon to develop cartographic projects [5]. But it’s not always the case when you get a high-resolution image which makes a major concern. Satellite images are not always provided in abundance, and there may be fewer image samples, which further poses a challenge in classification and segmentation. The applications of satellite imagery classification include disaster prediction using remote sensing images, and these early predictions are used to take necessary precautions. [31]

<sup>1</sup><https://github.com/muskaan712/Domain-Adaptable-Self-Supervised-Representation-Learning-on-Remote-Sensing-Satellite-Imagery>

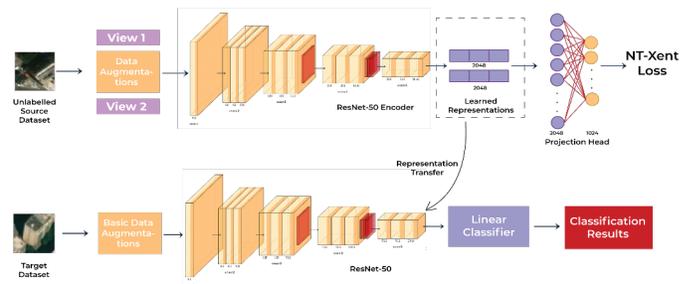


Fig. 1: Demonstrate the instance of proposed domain adaptation framework where self-supervised representation learning is performed at source dataset using contrastive learning method and downstream task performed on target datasets.

Satellite images can also be classified and segmented into areas with more wind and solar power so that adequate coverage of windmills and solar panels can be achieved to harness the power efficiently. [32]

Substantial human-labeled data is necessary to train a deep neural network successfully. Unfortunately, data collection and labeling are time-consuming and challenging in many fields. However, acquiring sufficient annotated data can be quite expensive and time-intensive. The process of cleaning, screening, labeling, evaluating, and reorganizing data by a training framework can be exceedingly time-consuming and complicated [11]. The lack of data has spawned a variety of solutions, the most prevalent of which is transfer learning. This work presents an approach that minimizes the training samples and puts less stress on the data labeling compared to the architecture modeling [6]. For most supervised learning approaches, annotated data is necessary to train a machine. This work employs self-supervised learning-based satellite image classification to deal with scarce labeled data in satellite imagery. When representations are learned from a pretext task using unlabeled input images and then used for a downstream task of interest, self-supervised learning is akin to transfer learning [12]. This work has used Domain Adaptation (DA) to prove the robustness of the model toward performance generalization on unseen data distribution. Domain Adaptation follows the concept that the model gets trained on one source

dataset and evaluated on the other target dataset, increasing the model’s reliability, re-usability, and results. The proposed work’s main contributions are outlined as-

- Establishing the adaptation of the self-supervised representation learning on remote-sensing satellite imagery by proposing a domain adaptation framework for rigorous evaluation.
- Achieving label-efficient representational knowledge transfer across multiple public datasets by obtaining state-of-the-art performance with limited labels and outperforming in fully supervised settings.
- Explaining improvement in quantitative results by qualitative analysis with significant and consistent evidence.

The rest of the article is organized as follows. Section 2 discusses the datasets. Section 3 presents a domain adaptation framework for self-supervised contrastive learning. Section 4 presents the experiments and results. Section 5 discusses the experiments and results obtained from the proposed work, Section 6 presents the related work, and Section 7 concludes the proposed work and provides the future scope of the work.

## II. SATELLITE IMAGERY DATASET DESCRIPTION

There are many applications for satellite images in meteorology, oceanography, fisheries, agriculture, biodiversity, geology, cartography, and land use planning. Instead of only having an image of a place, satellite image classification aims to transform satellite imagery into valuable information. Satellite Imagery of residential and non-residential buildups varies with objects and natural scenes captured in the image. A dataset with images labeled as a whole is required for categorizing satellite images. Three public satellite imagery datasets are used in this work, SIRI-WHU and UC Merced have equally distributed images among their classes and MLRSNet has non-uniform distribution which is demonstrated in the graph below, other details about the datasets are discussed below and summarized in Table I.

<i>Dataset</i>	<i>Total Images</i>	<i>No. of Classes</i>
<i>SIRI-WHU</i>	<i>2400</i>	<i>12</i>
<i>UC Merced Land Use Dataset</i>	<i>2100</i>	<i>21</i>
<i>MLRSNet</i>	<i>109,161</i>	<i>46</i>

TABLE I: Dataset description

### A. SIRI-WHU Dataset

The SIRI-WHU dataset<sup>2</sup> for classification has 2400 photos sorted into 12 classifications. This dataset was obtained from Google Earth and mainly included metropolitan regions in China, with the image collection developed by Wuhan University’s RS IDEA Group (SIRI-WHU). It consists of 12 classes: Agriculture, Commercial, Harbor, Idle land, Industrial, Meadow, Overpass, Park, Pond, Residential, River, and Water.

<sup>2</sup>[http://www.lmars.whu.edu.cn/prof\\_web/zhongyanfei/e-code.html](http://www.lmars.whu.edu.cn/prof_web/zhongyanfei/e-code.html)

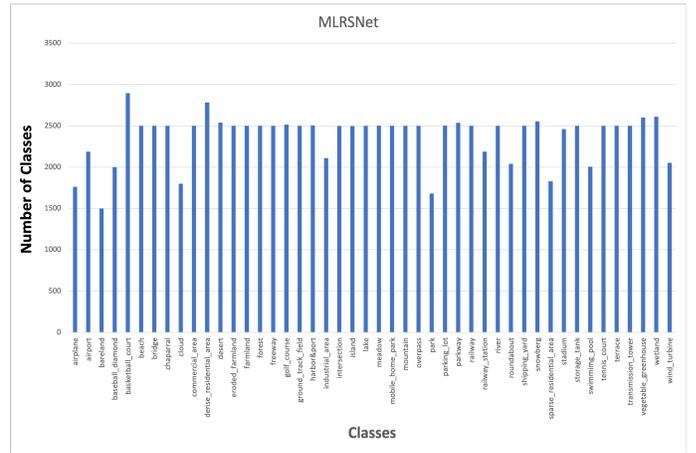


Fig. 2: Class-wise distribution of the MLRSNet dataset.

Each class comprises 200 pictures that are 200 x 200 pixels in size.

### B. UC Merced Dataset

The image data in the UC Merced dataset<sup>3</sup> were manually extracted from large-sized images in the United States Geological Survey (USGS) National Map Urban Area Imagery collection for numerous cities across the country (United States). This big ground truth picture collection consists of 21 land-use types, each with 100 pictures. The 21 classes were namely agricultural, airplane, baseball diamond, beach, buildings, chaparral, dense residential, forest, freeway, golf course, harbor, intersection, medium residential, mobile home park, overpass, parking lot, river, runway, sparse residential, storage tanks, and tennis court. This public domain imagery has a pixel resolution of 1 foot, with each image being 256x256 pixels.

### C. MLRSNet

MLRSNet<sup>4</sup> offers several satellite-based perspectives of the world. It comprises optical satellite images with great spatial resolution—between 1,500 and 3,000 example photos in every 46 categories in the 109,161 remote sensing photographs makeup MLRSNet. The photos are 256x256 pixels and have different pixel sizes (10m to 0.1m). The dataset can be used for picture segmentation, image retrieval, and classification based on multiple labels.

## III. DOMAIN ADAPTATION FRAMEWORK FOR SELF-SUPERVISED CONTRASTIVE LEARNING

The proposed framework consists of two main tasks: (i) pretext task, in which learning representations following contrastive self-supervised learning on satellite imagery datasets within the source domain is performed (ii) downstream task, in which satellite images are classified based on the representations learned in pretext task. Figure 3 depicts a schematic

<sup>3</sup><http://weege.vision.ucmerced.edu/datasets/landuse.html>

<sup>4</sup><https://data.mendeley.com/datasets/7j9bv9vwsx/2>

diagram of the proposed approach where knowledge transfer on self-supervised learnt representation is comprehensively validated

In the pretext task, various augmentations are applied to

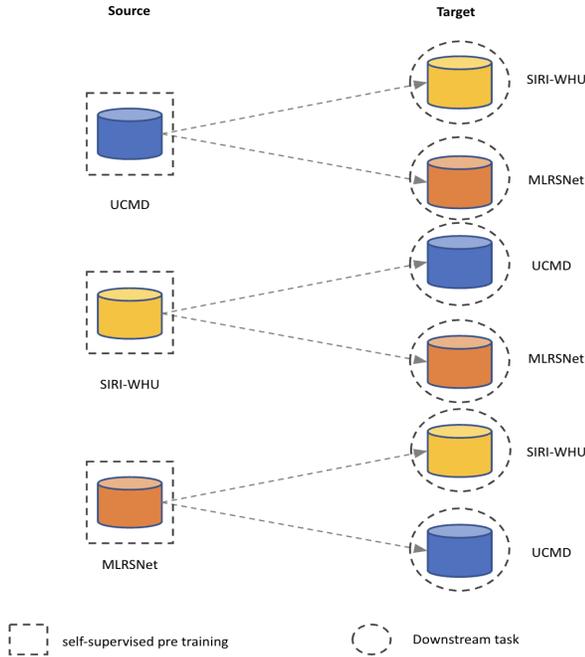


Fig. 3: Proposed approach for domain adaptation. It ensures to investigate of each dataset for self-supervised representation learning and for downstream tasks under all possible domain adaptation scenarios

the images, such as flipping, affine transformations, jitter, grayscale, etc., to create different views of the images. A positive pair is a pair of views created from the same image, whereas a negative pair is a pair of views created from different images. Then, positive and negative pairs of images are contrastively learned to form image representations. Labeled images are not required for representation learning in this task. Figure 1 depicts the contrastive learning architecture. SimCLR (simple framework for contrastive learning) [21] has been used to learn the representations. Positive and negative pairs of satellite images are created from unlabelled satellite images with augmentations such as Gaussian blur, flipping, translation, rotation, etc. These pairs of positive and negative views are fed to the encoder network. ResNet-50 encoder is a backbone for the pretext task network, followed by two fully connected layers containing 2048 and 1024 neurons each. The encoder part helps in extracting image representations for positive and negative pairs. Normalized Temperature-scaled Cross-Entropy loss (NT-Xent) is used to pull close representations and push away different representations. This loss function for a positive pair is defined below.

$$\ell(\mathbf{z}_i, \mathbf{z}_j) = -\log \frac{\exp(\text{sim}(z_i z_j) / \tau)}{\sum_{k=1}^{2n} \mathbb{1}_{k \neq i} \exp(z_i z_k / \tau)}$$

$\tau$  is the temperature parameter, where  $z_i$  and  $z_j$  in the numerator represent positive pairs, where  $z_i$  and  $z_k$  in the denominator represent all possible pairs, including positive and negative. In terms of a loss function, it comes down to the ratio of the sum of similarities between all positive pairs divided by the negative log-likelihood of these pairs being similar. A softmax function-based temperature parameter is used to normalize this loss function. It is designed to maximize an agreement between positive pairs in a mini-batch. The loss function for all the positive pairs is given below.

$$\mathcal{L} = -\frac{1}{N} \sum_{i,j \in \mathcal{MB}} \log \frac{\exp(\text{sim}(z_i, z_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{[k \neq i]} \exp(\text{sim}(z_i, z_k) / \tau)}$$

The downstream task uses the learned embeddings of images during the pretext task as input. Only a few basic augmentations like resizing and cropping are used during the downstream task. The downstream task involves binary and multi-class classification of satellite images. A binary and multi-class classification task is involved in the downstream task. For classification, fewer labeled images are now required for the downstream task. In the downstream task, the input image from the target dataset is taken as input, and primary augmentations (resizing, cropping) are applied to this image. The augmented images are fed to the encoder, initialized from the pretext task-trained model. A linear classifier having layers 512 and a number of classes has been appended to the encoder part of the network to classify satellite images.

#### IV. EXPERIMENTS AND RESULTS

Extensive experimentation is designed and performed to investigate the domain adaptation in self-supervised learning based representational knowledge transfer on three datasets, UC Merced, SIRI-WHU, and MLRSNet, covering binary and multi-class classifications downstream tasks under varying label efficiency. Table II shows the augmentations for the pretext task, and Table III shows the hyperparameters used for training the pretext task. Details of hyperparameters for fine-tuning and other configuration is available in open-source source code. The dataset follows a 70%, 20%, and 10% split for training, testing, and validation. The next subsections discuss the binary classification results and the multi-class classification results. The performance metrics are defined below.

$$\text{Precision} = \frac{\text{Total true positives}}{\text{Real actual positives} + \text{Total false positives}}$$

$$\text{Recall} = \frac{\text{Total true positives}}{\text{Total true positives} + \text{Total false-negatives}}$$

$$\text{Accuracy} = \frac{\text{true negatives} + \text{true positives}}{\text{total cases}}$$

$$\text{F1 Score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

<i>Augmentations</i>	<i>Value</i>
<i>Resize</i>	$224 \times 224$
<i>Horizontal Flip</i>	$P = 0.5$
<i>Vertical Flip</i>	$P = 0.5$
<i>Rotation</i>	$(-90, 90)$
<i>Grayscale</i>	$P = 0.2$
<i>Gaussian Blur</i>	$P = 0.51, \text{Kernel size} = [21, 21]$

TABLE II: Augmentations for pretext task

<i>Hyperparameters</i>	<i>Value</i>
<i>Batch size</i>	$256$
<i>Optimizer</i>	$SGD$
<i>Momentum</i>	$0.9, \text{nesterov} = \text{True}$
<i>Learning Rate</i>	$0.0005$
<i>Weight decay</i>	$0.0005$

TABLE III: Hyperparameters for pretext task

### A. Domain Adaptation

For domain adaptation, three different datasets have been used to evaluate the results and performance of the proposed methodology for satellite image classification. The three datasets are used without labels in the pretext task to generate representations and fine-tuning is performed on the other two remaining datasets to evaluate domain adaptation. The results of experiments are shown in Table IV.

<i>Dataset used</i>			<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-Score</i>
<i>Pretext</i>	<i>Downstream</i>	<i>% of data in downstream</i>				
<i>UC Merced</i>	<i>MLRSNet</i>	100%	96.34	96.21	96.56	96.87
		50%	95.18	94.83	94.76	94.34
		10%	92.23	91.98	91.73	91.45
<i>UC Merced</i>	<i>SIRI-WHU</i>	100%	96.87	96.32	96.34	96.87
		50%	94.99	94.12	94.76	94.12
		10%	87.50	87.43	87.24	87.97
<i>MLRSNet</i>	<i>UC Merced</i>	100%	98.50	98.54	98.21	98.32
		50%	96.01	96.80	96.79	96.85
		10%	92.32	92.32	92.56	92.81
<i>MLRSNet</i>	<i>SIRI-WHU</i>	100%	97.50	97.11	96.98	96.43
		50%	96.24	96.87	96.32	96.76
		10%	89.58	89.90	90.45	89.34
<i>SIRI-WHU</i>	<i>UC Merced</i>	100%	98.75	98.21	97.93	98.53
		50%	96.51	96.89	96.43	96.21
		10%	94.23	94.98	94.71	94.89
<i>SIRI-WHU</i>	<i>MLRSNet</i>	100%	97.87	97.43	97.54	97.32
		50%	94.40	94.87	94.91	94.51
		10%	90.02	90.83	90.26	90.73

TABLE IV: Results for Domain Adaptation on multiclassification

### B. Multi-class Classification

For the multi-class classification task, three different datasets have been used to evaluate contrastive learning for satellite image classification. The datasets are UC Merced, which has 21 classes on which our model achieved an accuracy of 99.35%, precision of 99.91%, recall value of 98.95%, and F1 score of 98.89% on the 100% dataset in the downstream task. Another dataset used is the SIRI-WHU with 12 classes,

and our model scored an accuracy of 99.68%, precision of 95.36%, recall value of 96.56%, and F1 score of 97.92% on the 100% dataset in the downstream. MLRSNET is another dataset that was used, a 46 class dataset and our model achieved an accuracy of 96.59%, precision of 96.79%, recall value of 96.545, and F1 score of 96.65% with 100% data in the downstream. The results of further experiments are shown in Table V, with fewer datasets downstream.

<i>Dataset used</i>			<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-Score</i>
<i>Pretext</i>	<i>Downstream</i>	<i>% of data in downstream</i>				
<i>UC Merced</i>	<i>UC Merced</i>	100%	99.35	99.91	98.95	98.89
		50%	95.98	91.78	93.32	93.82
		10%	89.67	87.63	85.86	85.85
<i>SIRI-WHU</i>	<i>SIRI-WHU</i>	100%	99.68	95.36	96.56	97.92
		50%	95.54	94.47	95.29	94.78
		10%	88.43	83.76	79.77	81.43
<i>MLRSNet</i>	<i>MLRSNet</i>	100%	96.59	96.79	96.54	96.65
		50%	93.87	91.37	91.88	92.68
		10%	88.66	88.54	88.91	88.44

TABLE V: Results for Multiclassification using contrastive learning

### C. Comparison with existing results

The concept of self-supervised learning and domain adaptation-based self-supervised learning applied to satellite imagery has yet to be explored. This work considered various methods applied to these datasets, including supervised learning-based methods. Table 6 compares the results obtained from the proposed work with the existing binary and multi-class classification methods on satellite imagery.

<i>Author</i>	<i>Method</i>	<i>Accuracy</i>
[23]	<i>SVM</i>	<b>98.8</b>
[24]	<i>GIST</i>	46.9
[25]	<i>ResNet 50</i>	98
[26]	<i>DCNN</i>	93.48
[16]	<i>GoogleNet</i>	97.10
[14]	<i>Semisupervised ensemble projection</i>	66.49
<b>Our Results</b>	<b><i>Self-supervised Domain Adaptation</i></b>	<b>98.75</b>
	<b><i>Self-supervised Same Dataset</i></b>	<b>99.35</b>

TABLE VI: Comparative results of multi-class classification 21 class UCMD. (Top 2 results are shown)

<i>Author</i>	<i>Method</i>	<i>Accuracy</i>
[27]	<i>AlexNet SPP SS</i>	95.07
[28]	<i>MCNN</i>	93.75
[29]	<i>Inception-LSTM</i>	<b>99.73</b>
<b>Our Results</b>	<b><i>Self-Supervised Domain Adaptation</i></b>	<b>96.87</b>
	<b><i>Self-supervised Same Dataset</i></b>	<b>99.68</b>

TABLE VII: Comparative results for multi-class classification 12 class SIRI-WHU. (Top 2 results are shown)

Author	Method	Accuracy
[30]	DenseNet201-SR-Net	87.87
[30]	ResNet101-SR-Net	87.48
[17]	Self-Supervised Learning	96
Our Results	Self-Supervised Domain Adaptation	97.87
	Self-supervised Same Dataset	96.59

TABLE VIII: Comparative results for multi-class classification 46 class MLRSNet. (Top 2 results are shown)

Based on the comparisons in Tables VI, VII, VIII the proposed work performs better than the existing works. According to the above comparative analysis, the proposed work outperforms all previous works and achieves state-of-the-art results for multi-class classification of satellite imagery.

## V. DISCUSSIONS

This section discusses the key outcomes of the proposed work and provides analysis based on the achieved results on three datasets for different scenarios.

### A. Self-supervised learnt representations are domain adaptable

Results in Table IV clearly indicate that the performance of domain adaptation with different sources and targets achieves comparable results with in-domain knowledge transfer presented in Table V. While investigating and comparing the domain adaptation results with ImageNet supervised knowledge transfer (refer Table IX & X), all the models outperform which indicates the successful domain adaptations across the datasets. Following the trend, the proposed framework consistently outperformed on the given datasets compared with the ImageNet pretrained ResNet50 in a complete range of labels from 10% to 100%, shown in Figure 4, 5, & 6.

Dataset	% of data	Accuracy	Precision	Recall	F1-Score	AUC
UCMD	100%	83.00	82.89	82.34	82.47	93.83
	50%	81.87	81.77	81.32	81.53	92.93
	10%	66.66	65.78	66.45	66.21	88.87
SIRI-WHU	100%	85.67	85.54	85.33	85.44	95.99
	50%	80.21	80.02	79.99	79.34	93.95
	10%	60.71	60.43	60.33	60.21	90.99
MLRSNet	100%	93.02	92.98	92.74	92.34	97.99
	50%	90.54	90.32	90.55	90.43	96.99
	10%	82.07	81.76	81.34	81.33	93.69

TABLE IX: Results for ResNet50 finetuning with Imagenet weights

Dataset	% of data	Accuracy	Precision	Recall	F1-Score	AUC
UCMD	100%	80.67	80.32	78.98	80.16	93.23
	50%	78.75	77.97	78.78	78.84	93.12
	10%	45.84	45.43	45.12	45.59	88.98
SIRI-WHU	100%	83.41	83.78	83.43	83.61	93.99
	50%	71.87	71.43	71.91	71.79	92.87
	10%	67.85	67.99	67.65	67.31	90.54
MLRSNet	100%	91.95	90.93	91.99	91.63	95.67
	50%	90.46	88.65	90.32	90.91	94.09
	10%	79.89	80.00	79.02	79.45	85.34

TABLE X: Results for ResNet50 linear evaluation with Imagenet weights

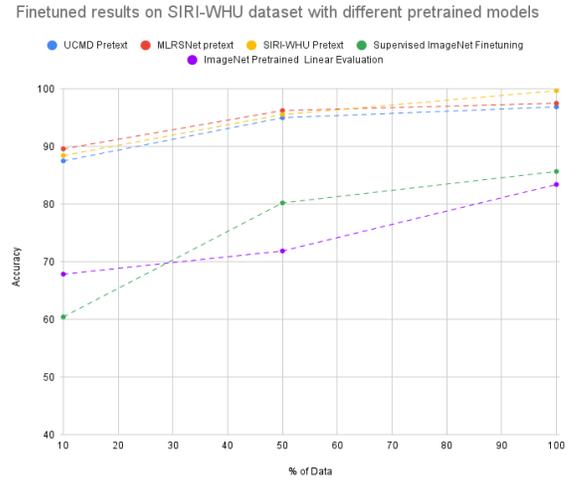


Fig. 4: Comparison of accuracy of the proposed method for SIRI-WHU with supervised learning.

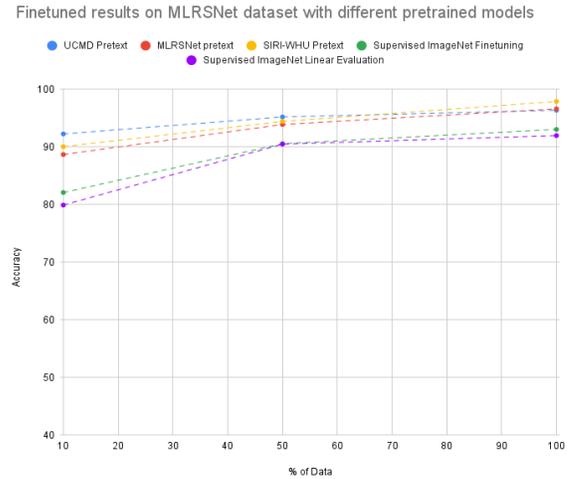


Fig. 5: Comparison of accuracy of the proposed method for MLRSNet with supervised learning.

### B. Self-supervised representation based knowledge transfer demonstrate label efficiency

Results on all three datasets indicate that the proposed framework obtains state-of-the-art results with only 10% and 50% of labels comparing previous work, which indicates that self-supervised learnt representations capture the important features of visual concepts of interest and adapt it for downstream tasks without additional efforts and parameter adjustments.

### C. Self-supervised pretrained models outperforms in fully supervised setting

Besides label efficiency, knowledge transfer in self-supervised pretrained models outperforms previous works for

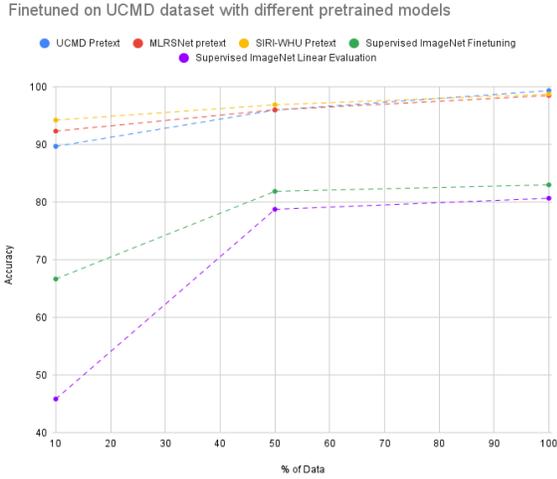


Fig. 6: Comparison of accuracy of the proposed method for UCMD with supervised learning.

all three datasets. This trend indicates that self-supervised representations are efficient end-to-end learning.

#### D. Robust and explainable representations

All the quantitative results are well described, with qualitative analysis performed on all three datasets with activation maps. Figure 7 demonstrates the explainability and attention for self-supervised pretrained models against ImageNet supervised models. This indicates that learnt representations in self-supervised manner are more efficient and thus achieve higher performance in downstream tasks.

### VI. RELATED WORK

During the past few years, self-supervised contrastive learning has emerged as a new training paradigm. Using this training paradigm, comprehensive representations can be learned without human annotation, which could solve the lack of annotated data problem. Much research has yet to be done on self-supervised learning in satellite imagery. Here, the main focus is on discussing deep learning methods applied to satellite imagery classification.

#### A. Supervised learning on satellite images

R. Naushad et al. [13] proposed a transfer learning approach to classify land use and land cover on the Eurosat dataset. For this, four CNN models were pre-trained: VGG-16 (without data augmentation), VGG-16 (with data augmentation), wide ResNet-50 (Without Data Augmentation), and wide ResNet-50 (With Data Augmentation). They achieved an accuracy of 99.17%. However, the proposed approach has been validated only on one dataset (the Eurosat dataset ) while using 100% of the data. At a large scale, CNN-based models were used by A. Albert et al. [14] to identify patterns in urban environments using satellite imagery. They used pre-trained models: VGG 16 and ResNet, for the classification task and achieved different

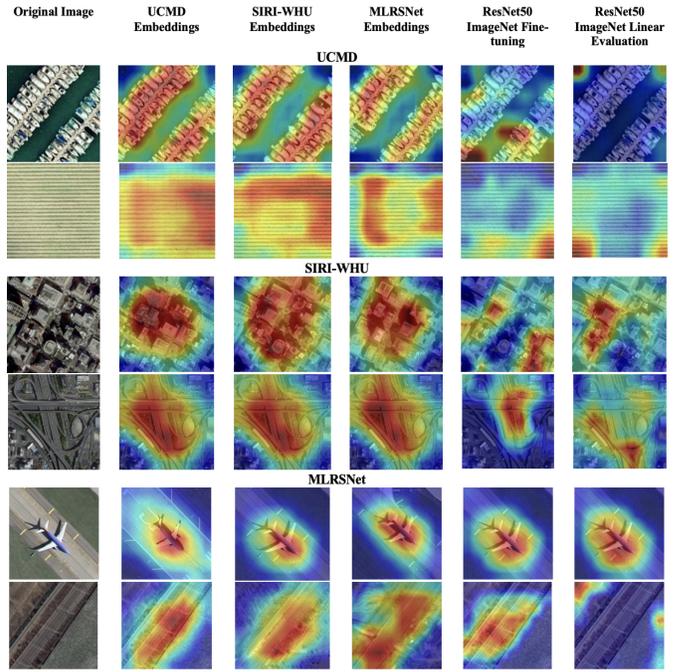


Fig. 7: Learning representation of the architecture using Class Activation Maps

accuracies for different countries, with which authors showed the highest accuracy of 83%. However, there remains scope for improvement in the results. Another method was applied by M. Castelluccio et al. [16] for land use classification in remote sensing images. They used pre-trained models: CaffeNet and GoogLeNet. These models provided an accuracy of 90.17% and 91.83%, respectively. Although authors achieve competitive results, the possibility of improvement of results and the use of better models remains.

X. Tan et al. proposed a multilabel classification to classify the MLRSNet, a benchmark dataset of 46 classes. They achieved an accuracy of 87.87% with DenseNet201-SR-Net [30]. However, there is still scope for improvement in accuracy. Furthermore, DenseNet201 is a very heavy computational method that needs more resources. S.Jog et al. [19] performed a supervised classification of satellite images using the Landsat dataset and support vector machine as the classifier, which achieved an accuracy of 92.84%. However, this approach was tested on a single dataset. Thus, the proposed method needs to be validated on other datasets as well to check the robustness of the model. M. Pritt et al. [22] used convolutional neural networks on the FMoV dataset for satellite image classification and achieved an accuracy of 83%. However, there is much scope for improvement in accuracy.

Özyurt, F. et al. [23] attempted to classify satellite images using the UC Merced dataset with a unique approach to feature extraction. For the classification, they used an SVM-based machine learning model and achieved an accuracy of 98.8%. However, their work focussed on a single dataset only. Kadhim et al. [25] used a pre-trained ResNet50 model on

the UC Merced dataset for the satellite image classification task and achieved an accuracy of 98%. However, the model has been evaluated on a single dataset only. F. P. S. Luus et. al. [26] used deep convolutional neural networks for land use classification on the UC Merced dataset and achieved an accuracy of 93.48%. Though the methodology used differs from existing methods, there is a scope for improvement in accuracy.

Han. Xiaobing et.al. [27] used a pre-trained AlexNet model on the SIRI-WHU dataset for satellite image classification and obtained an accuracy of 95.07%. However, the accuracy can be improved further using better network architectures. Y. Liu et al. [28] worked on the SIRI-WHU dataset and used multiscale convolutional neural networks for scene classification using satellite images. They achieved an accuracy of 93.75%. Although the method differs from existing approaches, this work is based on a single dataset only and needs to be tested on other datasets also. Y. Dong et al. [29] classified satellite images from the SIRI-WHU dataset using inception-based LSTM approach with an accuracy of 99.73%. However, the approach is tested on a single dataset only.

### *B. Self-supervision and domain adaptation on specialized visual domain*

Self-supervised methods on ImageNet and natural scenes have advanced in recent times. It has also been considerable advancements in other specialized visual domains to adapt self-supervised representation learning. Self-supervised methods in medical images showed progress where data and human labels are limited, Chhipa et al. [33] demonstrated self-supervised domain adaptation on histopathology images, and Azizi et al. [34] showed knowledge transfer on X-ray. Other interesting applications for self-supervised methods exploring on underwater images Tarling et al. [37] for the fish count and identifying mining materials from three-dimensional particle management sensors in Chhipa et al. [36] shown progress in a specialized domain.

### *C. Domain adaptation and self-supervised learning on satellite images*

Few non-supervised learning based methods for satellite image classification have been proposed in the literature. A semi-supervised learning based approach for satellite image classification was proposed by W. Yang et al. [15] to solve the problem of fewer images. They achieved an accuracy of 73.82% on 19 class data and 65.34% on UCMD dataset. However, there is scope for further improvements in the results obtained by the authors for satellite image classification.

A few self-supervised methods have also been used to classify remotely sensed satellite images in the literature, as proposed by V. Stojnic et al. [17]. They used self-supervised methods with a pre-trained Imagenet model on MLRSNet and achieved an accuracy of 96%. Manas et al. [?] have shown self-supervised pretraining on remote sensing data using weather information. Yi Wang et al. [18] proposed

contrastive multiview coding (CMC) based approach for satellite image classification, where one image is an anchor, and other images are neighbored around that image. They used pre-trained models for feature extraction, and the number of training samples was large. However, they did not validate the proposed approach in cross-domain settings wherein learning the representations from one dataset of satellite images and performing downstream tasks on another dataset.

From the above analysis of the existing work in the literature, it can easily be observed that most of the existing supervised learning-based satellite image classification methods require a lot of labeled data to perform satisfactorily. Only a few semi-supervised or self-supervised satellite image classification methods exist in the literature. However, these methods use the same dataset for pretext tasks, and downstream tasks, and these methods have not been evaluated in cross-domain settings. To mitigate these research gaps in the literature, this work proposes a domain adaptation-based self-supervised representation learning approach for classifying satellite images. This work proposes a domain adaptable self-supervised learning approach to reuse the representations learned on one unlabelled dataset from the source domain for classifying satellite images taken from a different target domain dataset.

## VII. CONCLUSION

This work proposed a domain-adaptable self-supervised representation learning based framework focusing on the robust evaluation of learnt representations rather than one-directional knowledge transfer, which ultimately reviews the effectiveness and applicability of such methods in the satellite imagery visual domain. One significant outcome is achieving improved performance by applying domain-adapted knowledge transfer across the datasets, outperforming the existing methods of satellite image classification, even in cross-domain settings. By applying the self-supervised representation learning, the proposed work has surpassed the existing results by 1%, with fewer training data. The proposed evaluation framework is conveniently applicable to other visual domains which are not thoroughly explored yet for the usability of self-supervised representation learning to reduce human annotation needs. In future work, i) We aim to investigate domain adaptation for other computer vision downstream tasks, i.e., segmentation and localization, ii) Investigate non-contrastive representation learning methods, and iii) Candidates for standard augmentation methods in self-supervised learning to adapt remote sensing visual domain.

## REFERENCES

- [1] Shafaey, M.A., Salem, M.A.M., Ebied, H.M., Al-Berry, M.N., Tolba, M.F. (2019). Deep Learning for Satellite Image Classification. In: Hassanien, A., Tolba, M., Shaalan, K., Azar, A. (eds) Proceedings of the International Conference on Advanced Intelligent Systems and Informatics 2018. AISI 2018. Advances in Intelligent Systems and Computing, vol 845. Springer, Cham.
- [2] Muruganandham, S. (2016). Semantic Segmentation of Satellite Images using Deep Learning (Dissertation). Retrieved from <http://urn.kb.se/resolve?urn=urn:nbn:se:ltu:diva-38558>

- [3] J. Senthilnath, S. Kulkarni, J. A. Benediktsson and X. S. Yang, "A Novel Approach for Multispectral Satellite Image Classification Based on the Bat Algorithm," in *IEEE Geoscience and Remote Sensing Letters*, vol. 13, no. 4, pp. 599-603, April 2016, doi: 10.1109/LGRS.2016.2530724.
- [4] M. Pritt and G. Chern, "Satellite Image Classification with Deep Learning," 2017 IEEE Applied Imagery Pattern Recognition Workshop (AIPR), 2017, pp. 1-7, doi: 10.1109/AIPR.2017.8457969.
- [5] Wei Zhang, Ping Tang & Lijun Zhao (2021) Fast and accurate land-cover classification on medium-resolution remote-sensing images using segmentation models, *International Journal of Remote Sensing*, 42:9, 3277-3301, DOI: 10.1080/01431161.2020.1871094
- [6] D. Ienco, Y. J. E. Gbdjo, R. Gaetano and R. Interdonato, "Weakly Supervised Learning for Land Cover Mapping of Satellite Image Time Series via Attention-Based CNN," in *IEEE Access*, vol. 8, pp. 179547-179560, 2020, doi: 10.1109/ACCESS.2020.3024133. interface," *IEEE Transl. J. Magn. Japan*, vol. 2, pp. 740-741, August 1987 [Digests 9th Annual Conf. Magnetics Japan, p. 301, 1982].
- [7] Pan, Z., Xu, J., Guo, Y., Hu, Y., & Wang, G. (2020, May 15). Deep learning segmentation and classification for urban village using a worldview satellite image based on U-Net. MDPI. Retrieved October 10, 2022, from <https://www.mdpi.com/2072-4292/12/10/1574>
- [8] Guo, Y., Liu, Y., Georgiou, T. et al. A review of semantic segmentation using deep neural networks. *Int J Multimed Info Retr* 7, 87-93 (2018). <https://doi.org/10.1007/s13735-017-0141-z>
- [9] A. Nivaggioli and H. Randrianarivo, "Weakly Supervised Semantic Segmentation of Satellite Images," 2019 Joint Urban Remote Sensing Event (JURSE), 2019, pp. 1-4, doi: 10.1109/JURSE.2019.8809060.
- [10] Iqbal, J., & Ali, M. (2020, July 29). Weakly-supervised domain adaptation for built-up region segmentation in aerial and satellite imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*. Retrieved October 10, 2022,
- [11] T. Mehari and N. Strodthoff, "Self-supervised representation learning from 12-lead ECG data," *Comput. Biol. Med.*, vol. 141, no. September 2021, p. 105114, 2022.
- [12] K. Wickstrøm, M. Kampffmeyer, K. Ø. Mikalsen, and R. Jenssen, "Mixing up contrastive learning: Self-supervised representation learning for time series," *Pattern Recognit. Lett.*, vol. 155, pp. 54-61, 2022.
- [13] Naushad, R., Kaur, T., & Ghaderpour, E. (2021, November 25). Deep Transfer Learning for Land Use and land cover classification: A comparative study. *arXiv.org*. Retrieved November 4, 2022, from <https://arxiv.org/abs/2110.02580>
- [14] Albert, Adrian & Gonzalez, Marta. (2017). Using Convolutional Networks and Satellite Imagery to Identify Patterns in Urban Environments at a Large Scale. 1357-1366. 10.1145/3097983.3098070.
- [15] W. Yang, X. Yin and G. -S. Xia, "Learning High-level Features for Satellite Image Classification With Limited Labeled Samples," in *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 8, pp. 4472-4482, Aug. 2015, doi:10.1109/TGRS.2015.2400449.
- [16] Castelluccio, Marco & Poggi, Giovanni & Sansone, Carlo & Verdoliva, Luisa. (2015). Land Use Classification in Remote Sensing Images by Convolutional Neural Networks.
- [17] V. Stojnić and V. Risojević, "Self-Supervised Learning of Remote Sensing Scene Representations Using Contrastive Multiview Coding," 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 2021, pp. 1182-1191, doi: 10.1109/CVPRW53098.2021.00129.
- [18] Wang, Yi & Albrecht, Conrad & Braham, Nassim & Mou, Lichao & Zhu, Xiao. (2022). Self-Supervised Learning in Remote Sensing: A Review. *IEEE Geoscience and Remote Sensing Magazine*. PP. 2-36.10.1109/MGRS.2022.3198244.
- [19] S. Jog and M. Dixit, "Supervised classification of satellite images," 2016 Conference on Advances in Signal Processing (CASP), 2016, pp. 93-98, doi: 10.1109/CASP.2016.7746144.
- [20] Manohar, N., Pranav, M.A., Aksha, S., Mytravarun, T.K. (2021). Classification of Satellite Images. In: Senjyu, T., Mahalle, P.N., Perumal, T., Joshi, A. (eds) *Information and Communication Technology for Intelligent Systems*. ICTIS 2020. Smart Innovation, Systems and Technologies, vol 195. Springer, Singapore.
- [21] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Article 149, 1597-1607.
- [22] M. Pritt and G. Chern, "Satellite Image Classification with Deep Learning," 2017 IEEE Applied Imagery Pattern Recognition Workshop (AIPR), 2017, pp. 1-7, doi: 10.1109/AIPR.2017.8457969.
- [23] Özyurt, F., Avcı, E., Sert, E. (2020). UC-Merced image classification with CNN feature reduction using wavelet entropy optimized with genetic algorithm. *Traitement du Signal*, Vol. 37, No. 3, pp. 347-353. <https://doi.org/10.18280/ts.370301>
- [24] G. -S. Xia et al., "AID: A Benchmark Data Set for Performance Evaluation of Aerial Scene Classification," in *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 3965-3981, July 2017, doi: 10.1109/TGRS.2017.2685945.
- [25] Kadhim, M.A., Abed, M.H. (2020). Convolutional Neural Network for Satellite Image Classification. In: Huk, M., Maleszka, M., Szczerbicki, E. (eds) *Intelligent Information and Database Systems: Recent Developments*. ACIIDS 2019. Studies in Computational Intelligence, vol 830. Springer, Cham.
- [26] F. P. S. Luus, B. P. Salmon, F. van den Bergh and B. T. J. Maharaj, "Multiview Deep Learning for Land-Use Classification," in *IEEE Geoscience and Remote Sensing Letters*, vol. 12, no. 12, pp. 2448-2452, Dec. 2015, doi: 10.1109/LGRS.2015.2483680.
- [27] Han, Xiaobing, Yanfei Zhong, Liqin Cao, and Liangpei Zhang. 2017. "Pre-Trained AlexNet Architecture with Pyramid Pooling and Supervision for High Spatial Resolution Remote Sensing Image Scene Classification" *Remote Sensing* 9, no. 8: 848. <https://doi.org/10.3390/rs9080848>
- [28] Y. Liu, Y. Zhong and Q. Qin, "Scene Classification Based on Multiscale Convolutional Neural Network," in *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 12, pp. 7109-7121, Dec. 2018, doi: 10.1109/TGRS.2018.2848473.
- [29] Y. Dong and Q. Zhang, "A Combined Deep Learning Model for the Scene Classification of High-Resolution Remote Sensing Image," in *IEEE Geoscience and Remote Sensing Letters*, vol. 16, no. 10, pp. 1540-1544, Oct. 2019, doi: 10.1109/LGRS.2019.2902675.
- [30] X. Tan, Z. Xiao, J. Zhu, Q. Wan, K. Wang and D. Li, "Transformer-Driven Semantic Relation Inference for Multilabel Classification of High-Resolution Remote Sensing Images," in *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 1884-1901, 2022, doi: 10.1109/JSTARS.2022.3145042.
- [31] Linardos, V., Drakaki, M., Tzionas, P., & Karnavas, Y. L. (2022, May 7). Machine learning in disaster management: Recent developments in methods and applications. MDPI. Retrieved November 30, 2022, from <https://www.mdpi.com/2504-4990/4/2/20>
- [32] Tawn, R., & Browell, J. (2022). A review of very short-term wind and solar power forecasting. *Renewable and Sustainable Energy Reviews*, 153, 111758. doi:10.1016/j.rser.2021.111758
- [33] Chhipa, P. C., Upadhyay, R., Pihlgren, G. G., Saini, R., Uchida, S., & Liwicki, M. (2023). Magnification prior: a self-supervised method for learning representations on breast cancer histopathological images. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (pp. 2717-2727).
- [34] Azizi, S., Mustafa, B., Ryan, F., Beaver, Z., Freyberg, J., Deaton, J., ... & Norouzi, M. (2021). Big self-supervised models advance medical image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 3478-3488).
- [35] Manas, O., Lacoste, A., Giró-i-Nieto, X., Vazquez, D., & Rodriguez, P. (2021). Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 9414-9423).
- [36] Chhipa, P. C., Upadhyay, R., Saini, R., Lindqvist, L., Nordenskjöld, R., Uchida, S., & Liwicki, M. (2022). Depth Contrast: Self-Supervised Pretraining on 3DPM Images for Mining Material Classification. *arXiv preprint arXiv:2210.10633*.
- [37] Tarling, P., Cantor, M., Clapés, A., & Escalera, S. (2022). Deep learning with self-supervision and uncertainty regularization to count fish in underwater images. *Plos one*, 17(5), e0267759.