# FALSE: False Negative Samples Aware Contrastive Learning for Semantic Segmentation of High-Resolution Remote Sensing Image

Zhaoyang Zhang, Xuying Wang, Xiaoming Mei, Chao Tao, Haifeng Li*

arXiv:2211.07928v1 [cs.CV] 15 Nov 2022

*Abstract*—Self-supervised contrastive learning (SSCL) is a potential learning paradigm for learning remote sensing image (RSI)-invariant features through the label-free method. The existing SSCL of RSI is built based on constructing positive and negative sample pairs. However, due to the richness of RSI ground objects and the complexity of the RSI contextual semantics, the same RSI patches have the coexistence and imbalance of positive and negative samples, which causing the SSCL pushing negative samples far away while pushing positive samples far away, and vice versa. We call this the sample confounding issue (SCI). To solve this problem, we propose a False negAtive sampLes aware contraStive lEarning model (FALSE) for the semantic segmentation of high-resolution RSIs. Since the SSCL pretraining is unsupervised, the lack of definable criteria for false negative sample (FNS) leads to theoretical undecidability, we designed two steps to implement the FNS approximation determination: coarse determination of FNS and precise calibration of FNS. We achieve coarse determination of FNS by the FNS self-determination (FNSD) strategy and achieve calibration of FNS by the FNS confidence calibration (FNCC) loss function. Experimental results on three RSI semantic segmentation datasets demonstrated that the FALSE effectively improves the accuracy of the downstream RSI semantic segmentation task compared with the current three models, which represent three different types of SSCL models. The mean Intersection-over-Union on ISPRS Potsdam dataset is improved by 0.7% on average; on CVPR DGLC dataset is improved by 12.28% on average; and on Xiangtan dataset this is improved by 1.17% on average. This indicates that the SSCL model has the ability to self-differentiate FNS and that the FALSE effectively mitigates the SCI in self-supervised contrastive learning.

*Index Terms*—Self-supervised contrastive learning (SSCL), false-negative sample (FNS), remote sensing image (RSI), semantic segmentation.

## I. Introduction

**D**EEP neural networks (DNNs) trained in a supervised learning manner have made remarkable progress in remote sensing image (RSI) scene classification [1], target detection [2], and semantic segmentation [3], [4]. The dependence of this approach on massive, high-quality labeled samples has become a bottleneck for wide-scale application [5]–[7]. The promise of self-supervised contrastive learning (SSCL) has made it possible to learn the RSI invariant features from massive unlabeled data [8]–[11].
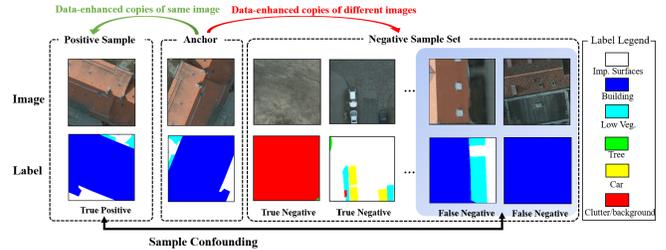
Fig. 1. Example of false negative sample (FNS) and sample confounding issue (SCI) in SSCL for semantic segmentation of high-resolution RSI. SCI will arise when model pushes negative sample image patch that contains positive samples far away.

The core idea of SSCL is to cleverly obtain copies of the same image patches as positive samples and other different images as negative samples by data augmentation methods of spatial and spectral transformations such as rotation, scaling, random color distortion, and Gaussian blur [12], and to construct self-supervised signals by pulling the positive samples closer while pushing the negative samples farther away. This is done instead of manually labeling them as supervised signals, thus forcing DNNs to obtain spatial and spectral invariant representations [11], [13]–[15].

However, due to scene complexity, ground object richness, and unbalanced distribution of samples of RSI, there is a phenomenon in which positive and negative samples coexist in the same patch and are highly unbalanced. As shown in Fig 1, the false negative samples for the selected anchor sample in the last two columns contain the same ground features as the positive sample, causing the SSCL pushing negative samples far away while pushing positive samples far away, and vice versa. We call this the sample confounding issue (SCI). The performance loss of the model due to SCI is called the sample confounding effect (SCE), while the negative sample image patch containing positive samples is called the false negative sample (FNS) because it gives the wrong feedback signal to the model.

Current methods to solve the SCE problem in SSCL are mainly considered from the perspective of samples and can be divided into two categories. One category aims to improve the quality of negative sample construction by attaching other unsupervised methods to the original SSCL and using the additional unsupervised results to guide the self-supervised model to construct higher-quality positive and negative sample pairs [16]–[18]. However, it is often difficult to improve the construction quality of positive and negative sample pairs in

RSI processing by using additional unsupervised clustering methods. Nevertheless, this may introduce defects of related unsupervised methods because the ground objects in RSI often have problems such as sample imbalance, intraclass differences, and interclass similarity, leading to the ineffectiveness of unsupervised clustering methods [19]. The second category considers abandoning the construction of negative samples [20]–[22], which means that the model's performance will depend only on the construction of positive samples. Considering that the FNS is essentially positive samples in the dataset, this type of approach completely avoids the generation of FNS. Nevertheless, it also means that the model will not use the positive samples that already exist in the dataset, which may reduce the model's ability to extract RSI invariant features.

Unlike the above methods, our observation is as follows: the SSCL of the RSI model itself can distinguish between true and false negative samples. This ability comes from the correct self-supervised signals given to the model by the true positive and true negative samples. This ability is potentially reinforced continuously as the model is trained. We refer to this as the ability of FNS self-determination (FSD). This observation motivates us to rethink the SCE problem from the perspective of the model rather than directly from the perspective of the sample.

The fundamental difficulty of using FSD to determine the FNS is that self-supervised pretraining is essentially an unsupervised process. The lack of definable criteria for the FNS leads to theoretical undecidability, so we can only approximately determine the FNS by some strategy. Approximate determination of FNS can be divided into two steps in terms of process: coarse determination of FNS and precise calibration of FNS. The former is the initial screening of FNS to ensure completeness, and the latter is the precise selection based on the former to ensure accuracy.

We propose the False negAtive sampLe aware contraStive lEarning model (FALSE), which achieves the coarse determination of FNS through the FNS self-determination (FNSD) strategy and achieves the precise calibration of FNS by designing the FNS confidence calibration (FNCC) loss function. In the FNSD strategy, the anchor sample in the closer positive sample pair in the embedding space is used as the benchmark, and the negative sample with the highest similarity to the anchor sample is determined as the possible FNS. The FNCC loss function is designed to improve the contribution of the possible FNS to the positive sample term of the original contrastive loss function [23] and reduce its contribution to the negative sample term of the loss function to mitigate SCE in the SSCL model. The contributions in this letter are as follows:

1) We propose a False negAtive sampLe aware contraStive lEarning model (FALSE) for the semantic segmentation of high-resolution RSIs. FALSE determines the approximate determination of FNS in SSCL from the perspective of the model rather than samples and mitigates the SCI in the SSCL of RSIs.

2) We designed the FNS confidence calibration (FNCC) loss function quantitatively rather than qualitatively to

characterize the strength of the ability of FNS self-determination (FSD) in the form of confidence weights.

3) The experimental results on three semantic segmentation datasets show that FALSE relative to SimCLR, PCL, and Barlow twins improves mean Intersection-over-Union (mIoU) on ISPRS Potsdam dataset by 0.7% on average on ISPRS Potsdam dataset, improves mIoU by 12.28% on average on CVPR DGLC dataset, and improves mIoU by 1.17% on average on Xiangtan dataset.

## II. METHODOLOGY

In the SSCL of RSI, the presence of true positive sample (TPS, blue dots in Fig 2) and true negative sample (TNS, red dots in Fig 2) will give the model a correct self-supervised signal about the RSI invariant features. In contrast, false negative sample (FNS, pink dots in Fig. 2) in the negative sample set will give the model an incorrect signal about the RSI invariant features, creating the SCI in SSCL.

Since SSCL pretraining is essentially an unsupervised process, the lack of definable criteria for FNS leads to theoretical undecidability, so we can only approximately determine the FNS by some strategies. The approximate determination of FNS can be divided into two steps: 1) coarse determination of FNS and 2) precise calibration of FNS. The former is the initial screening of FNS to ensure the completeness of FNS; the latter is precise on the former to ensure the accuracy of FNS selection.

### A. Coarse determination of FNS

*1) Determination benchmark:* Since the goal of the SSCL model is to bring positive samples closer and push negative samples farther, if the model projects a positive sample pair to a closer location in the embedding space, then the model is currently better at learning invariant features about that positive sample pair. Based on this, the anchor sample of the closer positive sample pair in the embedding space is selected as the benchmark for determination, maximizing the use of the feature extraction information that the model has learned, thus minimizing the model's misjudgment.

Suppose the benchmark anchor sample is denoted as $o_{key}$, its corresponding positive sample is denoted as $p$, and $sim(\cdot, \cdot)$ denotes the calculation of the feature similarity between the two samples. Then, $o_{key}$ satisfies the condition that

$$sim\left(o_{key}, \; p\right) > T \tag{1}$$

In Eq. (1), $T$ denotes the positive sample pair similarity threshold, which controls the proximity of positive sample pairs in the embedding space.

*2) Determination condition:* Based on the determination benchmark satisfying Eq. 1, we calculate the similarity between all negative samples and the benchmark anchor sample $o_{key}$ in the embedding space and determine the negative sample with the highest similarity to the benchmark anchor sample $o_{key}$ as the possible FNS. Suppose $n$ is used to denote the negative samples to be judged, and $n_{pf}$ denotes the possible FNS, the above determination condition can be simply described as:

$$\begin{aligned} &|sim(o_{key}, n_{pf}) - sim(o_{key}, p)| \\ &\rightarrow min \, |sim(o_{key}, n) - sim(o_{key}, p)| \end{aligned} \tag{2}$$
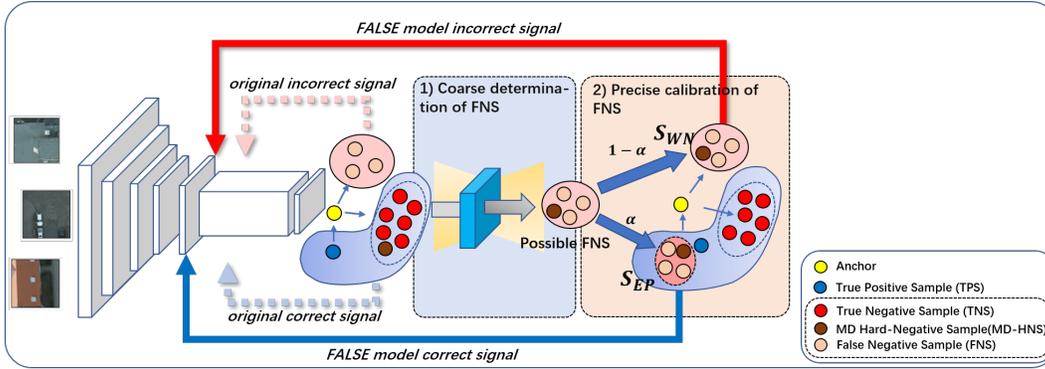
Fig. 2. Overview of FALSE model. Blue square represents FALSE model's FSD module in coarse dermination of FNS

*3) Possible FNS analysis:* Influenced by the SCI, the possible FNS obtained by the FSD is not all FNS but also contains TNS. Nevertheless, since $o_{key}$ represents the best level of the current model's ability to extract image features, it is difficult to use the model's FSD ability to eliminate this part of the TNS from the set of possible FNS. Noting that the indistinguishability of such negative samples is for the model, we refer to this part of the TNS among the possible FNS $n_{pf}$ as the model-dependent hard negative sample (MD-HNS, brown dots in Fig 2) in the SSCL. Suppose $n_f$ is used to denote the FNS and $n_h$ to denote the MD-HNS, the composition of the possible FNS obtained can be simply described as:

$$n_{pf} = n_f + n_h \qquad (3)$$

### B. Precise calibration of FNS

*1) FNS confidence calibration (FNCC) loss function:* To calibrate the obtained possible FNS and mitigate the impact of MD-HNS on the model performance, we design the FNCC loss function by introducing confidence weights $\alpha$ to calibrate the possible FNS ($n_{pf}$) as positive samples.

The original loss function of the SSCL of the RSI model mainly consist of two parts [23]: positive sample term $e^{sim(o,p)}$ and negative sample term $\sum_i^N e^{sim(o,n^i)}$.

And the FNCC multiply the similarity of possible FNS by the confidence weight $\alpha$ to get $S_{EP}$ (see Eq. (4)), and add it to the original positive sample term of SSCL loss, increases the influence of $n_{pf}$ on the positive sample term of the loss to enhance the correct signal, multiply the similarity of possible FNS by $1 - \alpha$ to get $S_{WN}$ (see Eq. (5)), and replace the similarity corresponding to the possible FNS in the original negative sample term with $S_{WN}$, reduces the influence of $n_{pf}$ on the negative sample term of the loss to weaken the incorrect signal.

When the number of negative samples corresponding to an anchor sample is $N$, the number of possible FNS determined to be obtained is $N_{pf}(N_{pf} < N)$, the FNCC loss function is defined by Eq. (6).

$$S_{EP} = \alpha \sum_j^{N_{pf}} e^{sim\left(o,n_{pf}^j\right)} \qquad (4)$$

$$S_{WN} = (1 - \alpha) \sum_j^{N_{pf}} e^{sim\left(o,n_{pf}^j\right)} \qquad (5)$$

$$L_{FNCC} = -log \frac{e^{sim(o,p)} + S_{EP}}{e^{sim(o,p)} + S_{EP} + \sum_i^{N-N_{pf}} e^{sim\left(o,n^i\right)} + S_{WN}} \qquad (6)$$

In particular, when there is no possible FNS, $N_{pf} = 0$, and Eq. (5) degenerates to original loss function and becomes the original SSCL model.

*2) Meaning of confidence weights:* The confidence weight $\alpha$ represents the degree of confidence in the model's FSD ability. When $\alpha = 0$, the positive sample signal enhancement term and the FNS signal weakening term of the FNCC loss are both 0. The model is the original SSCL model, and the FNSD strategy is not used. When $\alpha = 1$, it means that FALSE fully trusts the possible FNS obtained from the model FSD and adjusts all possible FNS to the positive sample, eliminating the contribution of these possible FNS to the negative sample term of the FNCC. When $\alpha$ takes any value less than 1 and greater than 0, it means that the model increases the contribution of possible FNS to the positive sample term of the FNCC from 0 to $\alpha$ times the original contribution to the negative sample term, and weakens its contribution to the negative sample term to $1 - \alpha$ times the original contribution.

## III. EXPERIMENTAL

### A. Datasets

The experiments were selected from the public RSI semantic segmentation dataset ISPRS Potsdam [24], competition dataset CVPR DGLC [25], and Xiangtan dataset [10] from the Gaofen-2 satellite covering Xiangtan, China. The spatial resolution and number of ground object types of the three datasets are shown in TABLE I.

TABLE I
DATASETS INTRODUCTION

| Dataset Name | Spatial Resolution | Class Num |
|---|---|---|
| ISPRS Potsdam | 0.05m | 6 |
| CVPR DGLC | 0.5m | 7 |
| Xiangtan | 2m | 8 |

### B. Experimental setup

The experiments follow the general paradigm of SSL models [9], [11], [13] and are divided into two main steps:
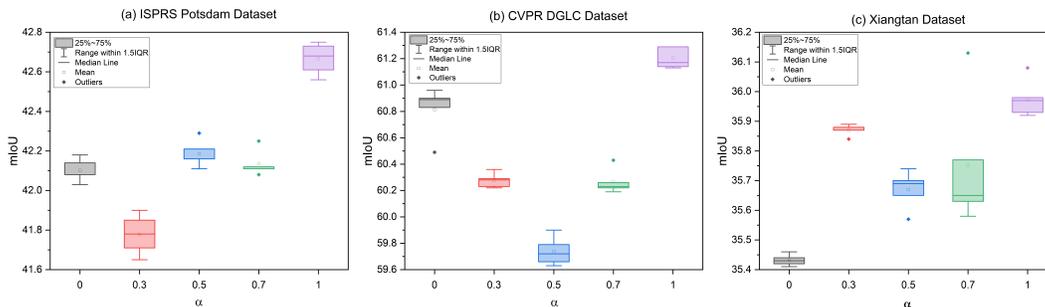
Fig. 3. mIoU (%) of 3 datasets when five different $\alpha$ are selected. Each confidence weight experiment was repeated five times.

self-supervised pretraining and supervised fine-tuning. Self-supervised pretraining uses all unlabeled training set data, pretrained for 200 epochs with the batch size set to 256. Next, the entire model encoder is frozen, and a small amount of labeled data is used to fine-tune and train the decoder, with the labeled data selected as 1% of the pretrained unlabeled data. The decoder is finally used to obtain the RSI semantic segmentation results.

To quantitatively explore the FSD of the FALSE, we conducted several experiments using the introduced confidence weights $\alpha$, with the positive sample pair similarity threshold $T$ set to 0.9. The results are presented in Experiment I.

On this basis, we selected SimCLR, which represents the original SSCL model; PCL, which represents the SSCL model that joins clusters; and the Barlow twins, which represent the SSCL model without constructing negative samples. These are used as the baseline tested on the ISPRS Potsdam, CVPR DGLC, and Xiangtan datasets, which are then compared with the FALSE model having a confidence weight of 1. The results are presented in Experiment II.

TABLE II
SEMANTIC SEGMENTATION IoU (%) OF 6 CLASSES OF GROUND OBJECTS
IN POTSDAM DATASET WHEN 5 DIFFERENT $\alpha$ ARE SELECTED

| $\alpha$ | 0 | 0.3 | 0.5 | 0.7 | 1 |
|---|---|---|---|---|---|
| Imp. Surface | 47.87 | 47.87 | 48.08 | 47.82 | **48.16** |
| Building | 47.15 | **47.72** | 46.15 | 47.42 | 46.83 |
| Low Veg. | 41.22 | 42.07 | 42.06 | 42.22 | **42.38** |
| Tree | 27.94 | 27.37 | 30.31 | 28.05 | **30.75** |
| Car | 28.22 | 24.86 | 26.42 | 27.05 | **28.29** |
| Clutter/background | 4.74 | 5.26 | 4.85 | 4.93 | **5.26** |

TABLE III
SEMANTIC SEGMENTATION IoU (%) OF 7 CLASSES OF GROUND OBJECTS
IN DGLC DATASET WHEN 5 DIFFERENT $\alpha$ ARE SELECTED

| $\alpha$ | 0 | 0.3 | 0.5 | 0.7 | 1 |
|---|---|---|---|---|---|
| Urban | 62.45 | 61.69 | 63.39 | 62.3 | **64.27** |
| Agriculture | 78.75 | 78.88 | 78.45 | 78.61 | **78.94** |
| Rangeland | **20.92** | 19.38 | 18.91 | 19.67 | 19.99 |
| Forest | 61.56 | **63.86** | 61.36 | 61.56 | 63.09 |
| Water | 57.21 | 58.15 | 57.78 | 60.41 | **61.69** |
| Barren | **46.43** | 43.56 | 43.52 | 44.43 | 45.39 |
| Unknow | **96.12** | 95.54 | 95.12 | 96.02 | 95.63 |

TABLE IV
SEMANTIC SEGMENTATION IoU (%) OF 8 CLASSES OF GROUND OBJECTS
IN XIANGTAN DATASET WHEN 5 DIFFERENT $\alpha$ ARE SELECTED

| $\alpha$ | 0 | 0.3 | 0.5 | 0.7 | 1 |
|---|---|---|---|---|---|
| Farmland | 63.82 | 64.10 | 63.82 | 64.06 | **64.45** |
| Urban | 0.00 | 0.31 | 0.53 | **2.75** | 0.96 |
| Rural areas | 16.62 | 17.16 | 17.23 | **18.25** | 14.92 |
| Water | 38.97 | 41.30 | 40.94 | 40.42 | **43.01** |
| Woodland | 78.27 | **78.49** | 78.11 | 78.12 | 78.36 |
| Grassland | 1.95 | 1.86 | 1.91 | 1.70 | **2.26** |
| Roads | 21.52 | 21.30 | 21.01 | **21.68** | 21.58 |
| Background | 97.56 | 98.08 | 98.09 | 98.2 | **98.23** |

### C. Experimental results and analysis

*1) Experiment I, Analysis of confidence weight:* We selected five values that were uniformly distributed in the interval from 0 to 1: 0, 0.3, 0.5, 0.7, and 1. The experiments were repeated five times for each confidence weight while keeping the other parameters consistent.

Fig. 3 shows the semantic segmentation mIoU results of the FALSE model on the ISPRS Potsdam, CVPR DGLC, and Xiangtan datasets with five different confidence weights. The model performed best when fully trusting the possible FNS obtained by the model's FSD ($\alpha = 1$). Compared with the original SSCL model, FALSE improves the mIoU by 0.72% on the ISPRS Potsdam dataset, 0.8% on the CVPR DGLC dataset, and 0.56% on the Xiangtan dataset. Moreover, when the model partially trusted the possible FNS obtained from the model's FSD ($0 < \alpha < 1$), the semantic segmentation performance became less stable. Considering that a confidence weight between 0 and 1 reduces the contribution of possible FNS to both positive and negative sample terms of the FNCC loss function, this phenomenon implies that the possible FNS (including FNS and MD-HNS) affects the stability of the semantic segmentation performance of the FALSE model.

TABLE II - TABLE IV show the IoU of the FALSE model for various types of ground objects on the three datasets with five different confidence weights. Compared with five different confidence weight models, the FALSE model with confidence weight is 1 achieved the best mIoU for 5 of all 6 types of ground features on the Potsdam dataset, the best mIoU for 3 and the second-ranked mIoU for 3 of all 7 types of ground features on the DGLC dataset, the best mIoU for 4 and the second-ranked mIoU for 3 of all 8 types of ground features

on the Xiangtan dataset.

TABLE V
SEMANTIC SEGMENTATION RESULT OF 4 DIFFERENT TYPES OF POSITIVE
AND NEGATIVE SAMPLE CONSTRUCTION STRATEGIES SSL MODELS

| Method | ISPRS Potsdam | | | CVPR DGLC | | | Xiangtan | | |
|---|---|---|---|---|---|---|---|---|---|
| | OA | mIoU | mAcc | OA | mIoU | mAcc | OA | mIoU | mAcc |
| SimCLR | 59.63 | 42.03 | 54.26 | 81.10 | 60.49 | 70.41 | 79.20 | 35.41 | 41.25 |
| PCL | 60.22 | 42.25 | 54.20 | 77.36 | 54.09 | 62.56 | 77.85 | 35.49 | 41.55 |
| Barlow twins | 60.03 | 41.87 | 54.06 | 66.71 | 32.45 | 39.28 | 77.57 | 33.49 | 39.33 |
| FALSE(ours) | **60.46** | **42.75** | **54.77** | **81.47** | **61.29** | **71.43** | **79.64** | **35.97** | **41.60** |

*2) Experiment II, Comparison of 4 types of SSCL models for semantic segmentation:* TABLE V shows the semantic segmentation results of the four different types of positive and negative sample construction strategy SSL models on the ISPRS Potsdam, CVPR DGLC, and Xiangtan datasets. The FALSE model achieves the best semantic segmentation results on these three datasets. Its Overall Accuracy (OA), mean Intersection-over-Union (mIoU), and mean class Accuracy (mAcc) outperform SimCLR, which represent the original SSCL model; PCL, which represents the SSCL model that joins clusters; and the Barlow twins, which represent the SSL model without constructing negative samples.

## IV. CONCLUSION AND FUTURE WORK

In this letter, we proposed the false negative sample aware contrastive learning model (FALSE) for the semantic segmentation of high-resolution RSIs. Under the restriction that self-supervised pretrained FNS are theoretically undecidable, the FALSE model achieves approximate determination of the FNS by coarse determination and precise calibration of FNS and quantitatively characterizes the ability of FNS self-determination (FSD) using confidence weights. Experiments on three RSI semantic segmentation datasets showed that FALSE effectively alleviates the SCE caused by SCI in the original SSCL of RSI. Compared with SimCLR, which represents the original SSCL model; PCL, which represents the SSCL model that joins clusters; and the Barlow twins, which represent the SSL model without constructing negative samples, FALSE improves mIoU by 0.7% on average on ISPRS Potsdam, improves mIoU by 12.28% on average on DGLC CVPR2018, and improves mIoU by 1.17% on average on Xiangtan.

The current method is only a simple implementation of the model's FSD, introducing manually set confidence weights. Through the experiment, we found that the confidence weight corresponding to the best segmentation accuracy of different ground objects is not the same, so how to adjust the confidence weight adaptively for different ground objects in the dataset and give full play to the model's ability of FSD is a further issue to be considered in the future for the FALSE model.

## REFERENCES

[1] G. Cheng, X. Xie, J. Han, L. Guo, and G.-S. Xia, "Remote sensing image scene classification meets deep learning: Challenges, methods, benchmarks, and opportunities," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 13, pp. 3735–3756, 2020.

[2] K. Li, G. Wan, G. Cheng, L. Meng, and J. Han, "Object detection in optical remote sensing images: A survey and a new benchmark," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 159, pp. 296–307, 2020.

[3] I. Kotaridis and M. Lazaridou, "Remote sensing image segmentation advances: A meta-analysis," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 173, pp. 309–322, 2021.

[4] H. Li, K. Qiu, L. Chen, X. Mei, L. Hong, and C. Tao, "Scattnet: Semantic segmentation network with spatial and channel attention mechanism for high-resolution remote sensing images," *IEEE Geoscience and Remote Sensing Letters*, vol. 18, no. 5, pp. 905–909, 2020.

[5] X. Zhao, R. Vemulapalli, P. A. Mansfield, B. Gong, B. Green, L. Shapira, and Y. Wu, "Contrastive learning for label efficient semantic segmentation," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10 623–10 633, 2021.

[6] Q. Yuan, H. Shen, T. Li, Z. Li, S. Li, Y. Jiang, H. Xu, W. Tan, Q. Yang, J. Wang *et al.*, "Deep learning in environmental remote sensing: Achievements and challenges," *Remote Sensing of Environment*, vol. 241, p. 111716, 2020.

[7] A. Vali, S. Comai, and M. Matteucci, "Deep learning for land use and land cover classification based on hyperspectral and multispectral earth observation data: A review," *Remote Sensing*, vol. 12, no. 15, p. 2495, 2020.

[8] S. Saha, M. Shahzad, L. Mou, Q. Song, and X. X. Zhu, "Unsupervised single-scene semantic segmentation for earth observation," *IEEE Transactions on Geoscience and Remote Sensing*, 2022.

[9] C. Tao, J. Qi, W. Lu, H. Wang, and H. Li, "Remote sensing image scene classification with self-supervised paradigm under limited labeled samples," *IEEE Geoscience and Remote Sensing Letters*, 2020.

[10] H. Li, Y. Li, G. Zhang, R. Liu, H. Huang, Q. Zhu, and C. Tao, "Global and local contrastive self-supervised learning for semantic segmentation of hr remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, 2022.

[11] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *International conference on machine learning*, pp. 1597–1607, 2020.

[12] O. Henaff, "Data-efficient image recognition with contrastive predictive coding," *International Conference on Machine Learning*, pp. 4182–4192, 2020.

[13] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9726–9735, 2020.

[14] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," *European conference on computer vision*, pp. 776–794, 2020.

[15] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," *Advances in neural information processing systems*, vol. 32, 2019.

[16] J. Li, P. Zhou, C. Xiong, and S. C. Hoi, "Prototypical contrastive learning of unsupervised representations," 2021.

[17] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," *Proceedings of the European conference on computer vision (ECCV)*, pp. 132–149, 2018.

[18] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9912–9924, 2020.

[19] W. Xia, C. Ma, J. Liu, S. Liu, F. Chen, Z. Yang, and J. Duan, "High-resolution remote sensing imagery classification of imbalanced data using multistage sampling method and deep neural networks," *Remote Sensing*, vol. 11, no. 21, p. 2523, 2019.

[20] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny, "Barlow twins: Self-supervised learning via redundancy reduction," *International Conference on Machine Learning*, pp. 12 310–12 320, 2021.

[21] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent-a new approach to self-supervised learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 21 271–21 284, 2020.

[22] X. Chen and K. He, "Exploring simple siamese representation learning," *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.

[23] A. Van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv e-prints*, pp. arXiv–1807, 2018.

[24] F. Rottensteiner, G. Sohn, J. Jung, M. Gerke, C. Baillard, S. Benitez, and U. Breitkopf, "The isprs benchmark on urban object classification and 3d building reconstruction," *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences I-3 (2012), Nr. 1*, vol. 1, no. 1, pp. 293–298, 2012.

[25] I. Demir, K. Koperski, D. Lindenbaum, G. Pang, J. Huang, S. Basu, F. Hughes, D. Tuia, and R. Raskar, "Deepglobe 2018: A challenge to

parse the earth through satellite images," *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018.