



Contrastive Learning-Based Heterogeneous Network for PolSAR Land Cover Classification

Journal:	<i>Transactions on Geoscience and Remote Sensing</i>
Manuscript ID:	TGRS-2023-01902
Manuscript Type:	Regular paper
Date Submitted by the Author:	30-Apr-2023
Complete List of Authors:	Cai, Jianfeng; Xidian University, School of Artificial Intelligence Feng, Zhixi; Xidian University, School of Artificial Intelligence Yang, SY; iiip, ee
Keywords:	Synthetic Aperture Radar Data

SCHOLARONE™
Manuscripts

Contrastive Learning-Based Heterogeneous Network for PolSAR Land Cover Classification

Jianfeng Cai, Zhixi Feng, *Member, IEEE*, Shuyuan Yang, *Senior Member, IEEE*

Abstract—Polarimetric synthetic aperture radar (PolSAR) image interpretation is widely used in various fields. Recently, deep learning has made significant progress in PolSAR image classification. However, supervised learning (SL) requires a large amount of labeled PolSAR data with high quality to achieve better performance, and manually labeled data is insufficient. This causes the SL to fail into overfitting and degrades its generalization performance. Furthermore, many features in PolSAR can be used to improve methods. To solve these problems, this article proposes a Heterogeneous Network based on Contrastive Learning (HCLNet). HCLNet aims to learn high-level representation from unlabeled PolSAR data for few-shot classification according to multi-features. It introduces the heterogeneous network for the first time to utilize different PolSAR features better. Beyond the conventional CL, HCLNet develops two easy-to-use plugins to narrow the domain gap between optics and PolSAR, including beam search and superpixel-based instance discrimination. The pre-trained online network is used for the downstream task by fine-tuning. Experiments demonstrate the superiority of HCLNet on three widely used PolSAR benchmark data sets compared with state-of-the-art methods on few-shot classification. Ablation studies also verify the importance of each component. Besides, this work has implications for how to efficiently utilize the multi-features of PolSAR data to learn better high-level representation in CL and how to construct networks suitable for PolSAR data better.

Index Terms—Contrastive learning (CL), polarimetric synthetic aperture radar (PolSAR) image classification, few-shot learning, superpixel, beam search

I. INTRODUCTION

Polarimetric synthetic aperture radar (PolSAR), an active remote sensing technology, has attracted significant attention due to its ability to obtain richer information than conventional single-polarization synthetic aperture radar (SAR). By using different polarimetric combinations of transmitting and receiving backscattering waves from land covers, PolSAR can observe targets in all-weather and all-time. Therefore, PolSAR image classification [1] [2], which is the most crucial task in PolSAR image interpretation, has been widely used in various fields such as geography [3], agriculture [4], and environmental monitoring [5].

This work was supported by the National Natural Science Foundation of China (Nos. 62171357, 62276205); the Foundation of Key Laboratory of Aerospace Science and Industry Group of CASIC, China; the Foundation of Intelligent Decision and Cognitive Innovation Center of State Administration of Science, Technology and Industry for National Defense, China; the Key Project of Hubei Provincial Natural Science Foundation under Grant 2020CFA001, China. (Corresponding author: Zhixi Feng, Shuyuan Yang.)

All the authors are with the School of Artificial Intelligence, Xidian University, Xi'an 710071, China (e-mail: jfcai_1@stu.xidian.edu.cn; zxwfeng@xidian.edu.cn; syyang@xidian.edu.cn).

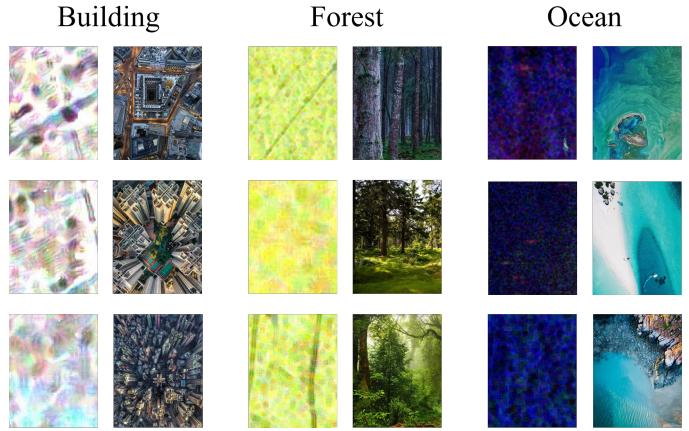


Fig. 1. Visual comparison of instance similarity between PolSAR and optical images, with PolSAR images on the left and optical images on the right.

Numerous researchers have proposed PolSAR classification methods using hand-crafted features. Two primary categories for classifying these features are inherent physical scattering and statistical features. The former is mainly based on target decomposition mechanisms: Freeman decomposition [6] decomposes the pixel into three scattering categories; H/A/ α decomposition [7] obtains entropy, anisotropy, alpha angle, and other decomposition methods, including Pauli decomposition [8], Huynen decomposition [9], Cameron decomposition [10], and Krogager decomposition [11]. The latter mainly consists of the coherency and covariance matrix, which follow the complex Wishart distribution. These methods use classifiers such as SVM [12] and MLP [13] to classify PolSAR data. However, the performance of these methods is heavily dependent on the quality of the features, and none of them can fully represent PolSAR data.

Recently, deep learning (DL) methods, especially convolutional neural networks (CNNs), have achieved magnificent success in various fields, including optical image [14]–[16], natural language processing (NLP) [17]–[19], and remote sensing image [20]–[22]. Due to the remarkable results of DL, many researchers have proposed several PolSAR deep learning methods. Zhou et al. [23] first used a CNN to replace conventional methods and achieved breakthrough results. To better adapt to the data structure of PolSAR, complex-valued CNNs (CV-CNNs) [24] were proposed. Subsequently, Wishart deep belief networks (WDBNs) [25], fully convolutional networks [26], and 3D convolution-based networks [27] have been proposed. While supervised CNN-based methods have achieved promising performance, they require a large labeled

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

training set, which is a significant expense of time and energy. When labeled data are scarce, the trained network can easily result in overfitting, leading to a lack of generalization. It means that supervised learning methods lack robustness in the case of missing labeled data, even with augmentation and regularization techniques [28]–[30].

In contrast to supervised learning, self-supervised learning (SSL) [31], where the data provides supervision, has the advantage of learning general representation from unlabeled data, which is more desirable and meaningful. Contrastive learning (CL) is the popular SSL method in optical images, which constructs simple and easy-to-use frameworks for training. From InstDisc [32], InvaSpread [33], CPC [34], and CMC [35] to MoCo [36], SimCLR [37], BYOL [38], and SimSiam [39], CL has a relatively mature architecture. However, the gap between PolSAR and optical images makes applying optical CL methods to PolSAR directly tricky. Researchers have proposed some PolSAR-tailored CL methods to address this issue: MI-SSL [40] learned the implicit multi-modal representation from unlabeled data. PCLNet [41] developed an instance discrimination proxy objective to learn representation from unlabeled data. SSPRL [42] improves CL, so no negative samples are needed; Cui et al. [43] proposed TCSPANet, the two-staged CL based on attention. However, there are still some challenges with these methods:

- Most of them fail to utilize the multi-features in PolSAR data fully. The diversity of PolSAR features makes it more advantageous in CL and can extract high-level representation through different features. However, as a standard architecture for CL of optical images, they directly utilize the Siamese network to PolSAR, making it hard to exploit these features thoroughly. In contrast, the heterogeneous network can arbitrarily combine different features for better representation learning.
- Some methods attempt to incorporate different features but fail to consider the redundancy between them. The degree of information redundancy varies for different feature combinations, and some features may even hinder model learning. Therefore, feature selection is essential.
- Some of them ignored the high similarity between pixels in PolSAR data. As shown in Fig. 1, unlike optical images, where instances are images with thousands of pixels, in PolSAR images, instances are individual pixels. As a result, traditional CL methods struggle to differentiate between different PolSAR instances. To overcome this limitation, introducing diversity between instances is necessary to enhance model learning.

Based on the preceding analysis, this article introduces a novel approach to PolSAR data classification, named Heterogeneous CL Network (HCLNet). The proposed HCLNet employs two perspectives, physical and statistical, for CL. From the physical perspective, it uses beam search to reduce feature redundancy, while from the statistical perspective, it utilizes the coherency matrix. Additionally, it constructs a heterogeneous network to learn the representation of PolSAR data using the novel superpixel-based Instance Discrimination. This approach effectively utilizes the PolSAR multi-features

and addresses the problem of pixel similarity. Furthermore, it uses two easy-to-use plugins to better adapt to PolSAR data. Specifically, the main contributions of this work can be summarized as follows:

- the *Heterogeneous Network* is proposed to learn the representation of PolSAR data using CL for the first time. The network consists of two sub-networks with different architectures, namely the online network and the target network, where the former is a 2D CNN, and the latter is a 1D CNN. This network can input different features for learning and extract high-level representations hidden between multi-features of PolSAR data without the need for labeled data.
- A novel pretext task, *Superpixel-based Instance Discrimination*, is designed to reduce the similarity between pixels and thus the model can learn representation easier and better. This task utilizes superpixel segmentation to select positive and negative samples for CL, which reduces the occurrence of highly similar pixels being negative samples of each other.
- *Beam search* is utilized to select complementary features and reduce redundancy. It created a classifier as the foundation for beam search, which implements a suitable combination of these features and removes redundant information from multi-features.
- Experimental results demonstrate the superiority of the proposed method. HCLNet is applied to three benchmark PolSAR data sets, and the results indicate that it achieves state-of-the-art classification performance, whether in few-shot or full-sample scenarios.

The rest of this article is organized as follows: A brief review of the CL for optical images and PolSAR and PolSAR multi-features are given in Section II. The details of the proposed HCLNet are described in Section III. Section IV shows the experimental results and analysis. Finally, we provide the conclusion and prospects for future research in Section V.

II. RELATED WORK

A. Contrastive Learning

As a prevalent SSL method, CL has a mature and general architecture. It usually has two networks with the same architecture; one is called the online network, which will be sent to downstream tasks for fine-tuning as the main network, and the other is called the target network. They can share parameters. Generally, CL is trained by a proxy objective, called pretext task, and usually chooses Instance Discrimination with InfoNCE loss function. Its main idea is that two representation obtained by the two networks from different perspectives of the same image should be similar and vice versa. InfoNCE loss function is a modified Cross-Entropy Loss that introduces the dot product to compute similarity. The formula is as follows:

$$L(q, k) = -\log \frac{\exp(q \times k_+ / \tau)}{\sum_{i=0}^K \exp(q \times k_i / \tau)} \quad (1)$$

where q is the output of the online network, k_+ is the output of the target network that q matches and is called the positive

sample, k_i ($i = 0 \dots K$ and $k_i \neq k_+$) is all output of the target network, which is memoried and is called the negative sample, τ is a temperature hyper-parameter that adjusts the uniformity of information distribution [32]. Intuitively, this loss tries to classify q to k_+ , and essentially, it is the log loss of a $(K+1)$ -way softmax-based classifier.

1) *The CL in Optical Images*: In optical images, according to the different ways of parameter updating and negative sample selection, CL methods can be divided into different types. Wu et al. [32] first proposed Instance Discrimination with NCE loss function and memory bank to store negative samples. It treats the two sub-images obtained from a cropped image as the positive sample and all other images in the data set as negative samples. While Ye et al. [33] select other samples in the same mini-batch as negative samples. CPC [34] is a more general architecture containing an encoder and an auto-regression model. Positive and negative samples are constructed to train the encoder autoregressively. In addition to cropping, optical images have properties such as depth that can also form positive samples with each other. So Tian et al. [35] proposed CMC using luminance (L channel), chrominance (ab channel) [44], depth, surface normal [45], and semantic labels to construct the positive sample. To make more negative samples and sample representation change smoothly, MoCo [36] introduces a dictionary named queue to increase negative samples and momentum update to update the parameter of the target network. SimCLR [37] uses a larger batch size to achieve better results. It also adds a projection head to the network's end and achieves incredible performance. To eliminate negative samples, BYOL [38] adds a prediction head to the end of the online network, then turns the similarity problem into a prediction problem, which can effectively prevent the model collapse. Chen et al. [39] summarizes the previous work and proposes a simple architecture, SimSiam, which demonstrates the importance of stop gradient.

2) *The CL in PolSAR*: To use CL based on optical images in the PolSAR domain, some researchers proposed PolSAR-tailored CL methods: MI-SSL [40] uses coherency matrix T and constructs positive samples by visual, physical, and statistical features to learn the implicit multi-modal representation with similarity and difference loss; PCLNet [41] which copy the MoCo technique, selects data set according to stimulation for interclass and intraclass diversity and uses PolSAR image rotating 180 as the positive sample. SSPRL [42] proposes two branches and dynamic convolution (DyConv) layer to improve CL, its basic architecture is similar to BYOL. TCSPANet [43] exploits unsupervised multi-scaled patch-level data sets (UsMsPD) and semi-supervised multi-scaled patch-level data sets (SsMsPD). It also proposes two CL stages in TCNet and adds attention mechanism in SPAE to get better results.

In this article, inspired by SimCLR, we construct a heterogeneous CL network based on superpixel-based instance discrimination, which selects appropriate negative samples to reduce the high similarity between positive and negative samples.

B. Multi-features within PolSAR

1) *Physical scattering features*: As mentioned in [46], most physical scattering features are based on target decomposition. Different target decompositions, which decomposes target based on Scattering matrix S , have distinct advantages for PolSAR image classification. For example, Freeman decomposition [6] decomposes the scattering matrix S into three scattering categories: surface, volume, and double bounce; entropy, anisotropy, and alpha angle are obtained by H/A/α decomposition [8]; Krogager decomposition [11] decomposed the scattering matrix S into sphere, diplane, and helix components. Other decomposition methods include Yamaguchi, Vanzyl, Neuman, Multiple-Component Scattering Model (MCSM), Huynen, Holm, Barnes, Cloude, Anned, An-Yang, Pauli decomposition, Huynen decomposition, and Cameron decomposition [8] [9] [10] [47]–[51].

2) *Statistical features*: According to the statistical characteristics, PolSAR data can become the coherency matrix T and the covariance matrix C based on Scattering matrix S , which follows the complex Wishart distance. Scattering matrix S is defined as

$$S = \begin{bmatrix} S_{HH} & S_{HV} \\ S_{VH} & S_{VV} \end{bmatrix} \quad (2)$$

where $S_{XY,X,Y} \in [H,V]$ is the scattering element of horizontal/vertical transmitting/receiving polarization. The covariance matrix C is formed by

$$h = [S_{HH} \quad \sqrt{2}S_{HV} \quad S_{VV}]^T \quad (3)$$

$$C = hh^* = \begin{bmatrix} |S_{HH}|^2 & \sqrt{2}S_{HH}S_{HV}^* & S_{HH}S_{VV}^* \\ \sqrt{2}S_{HV}S_{HH}^* & 2|S_{HV}|^2 & \sqrt{2}S_{HV}S_{VV}^* \\ S_{VV}S_{HH}^* & \sqrt{2}S_{VV}S_{HV}^* & |S_{HV}|^2 \end{bmatrix} \quad (4)$$

where the superscript “ T ” denotes the conjugate transpose. The coherent matrix T is formed by

$$k_p = [(S_{HH} + S_{VV})/\sqrt{2} \quad (S_{HH} - S_{VV})/\sqrt{2} \quad \sqrt{2}S_{HV}]^T \quad (5)$$

$$T = \langle k_p k_p^T \rangle = \begin{bmatrix} T_{11} & T_{12} & T_{13} \\ T_{21} & T_{22} & T_{23} \\ T_{31} & T_{32} & T_{33} \end{bmatrix} \quad (6)$$

As aforementioned, these features have different importances in specific scenes. Some features have functional overlap, which leads to feature redundancy. So combining these features properly to learn representation better is necessary. We propose the beam search to obtain a better complementary feature combination that is similar to [46].

III. HETEROGENEOUS NETWORK BASED ON CL

In this section, the detailed process of the proposed HCLNet is presented. The overall architecture of HCLNet is shown in Fig. 2, which includes three main components. Among them, the first component is Beam Search, which finds the appropriate combination of target decomposition features. The

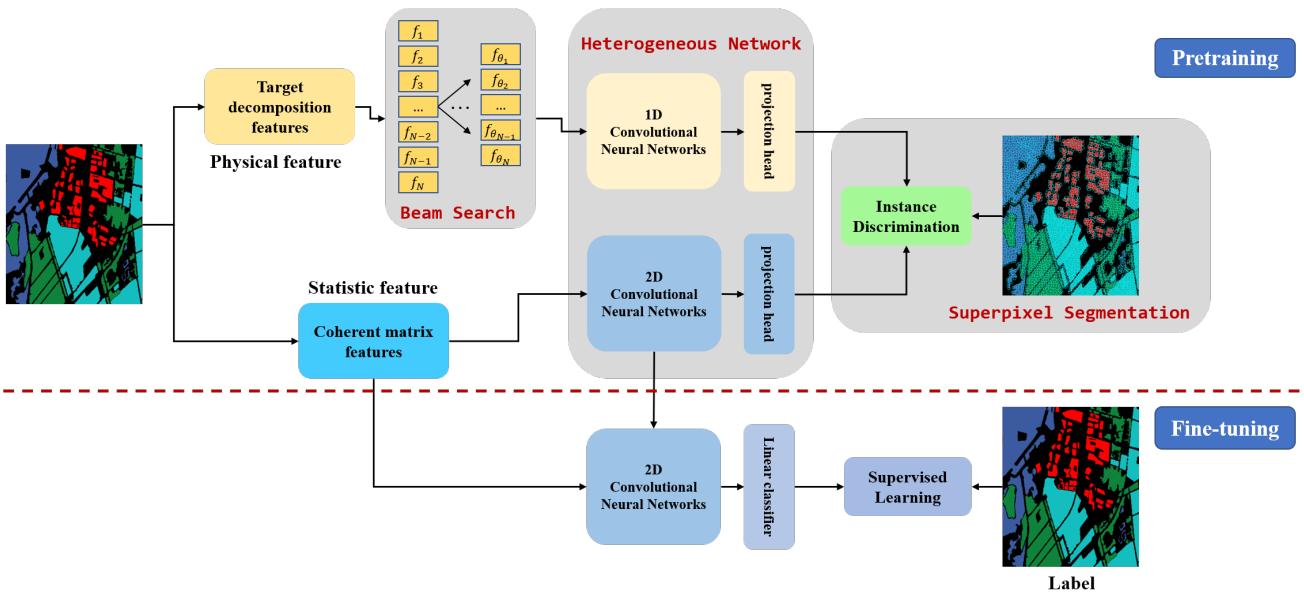


Fig. 2. The overall framework of the proposed HCLNet. It mainly contains two processes: Pretraining and Fine-tuning. In pretraining, it first uses Beam Search to combiniate features, then constructs the heterogeneous network and uses Superpixel-based Instance Discrimination to learn the high-level representation. In fine-tuning, it uses the trained online network from pretraining and fine-tunes it with a small number of labeled data to better fit the downstream distribution.

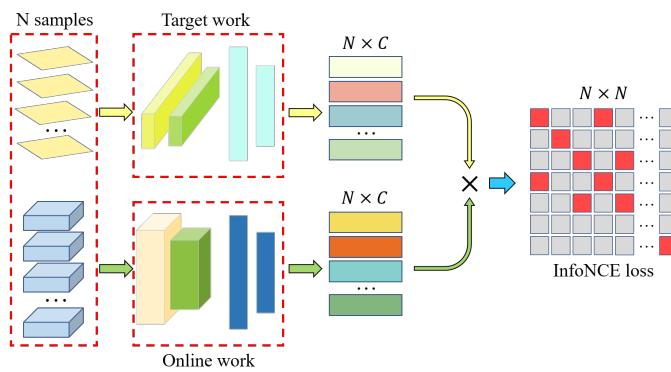


Fig. 3. The architecture of the heterogeneous network in HCLNet. It contains two networks with different architectures and is updated with InfoNCE loss. The output of the target network belonging to different superpixels in the same minibatch will be served as negative samples.

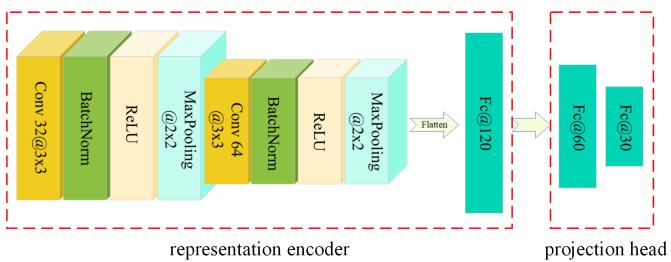


Fig. 4. The architecture of the online network in the heterogeneous network. It contains the representation encoder and the projection head; the former will be used for fine-tuning.

second component is Superpixel-based Instacne Discrimination, which improves the general Instance Discrimination selection of positive and negative samples. Moreover, the final component is the Heterogeneous Network, which is the

most important one and learns the high-level representation of PolSAR data. First, all target decomposition features are filtered using the beam search. Then the coherent matrix and the filtered target decomposition features are used as the input of the heterogeneous network. The superpixel-based instance discrimination is used for unsupervised training. The specific details of each component are as follows.

A. Beam Search

Here, we will introduce the beam search in detail. Like [46], we extract many features by the target decompositions mentioned above, N in total, and M groups; each target decomposition method generates a group of features. Then, we design a 1-D CNN model to evaluate the performance of different combinations of features. The network's input is an $N \times 1$ vector V , representing the N features of a pixel. During training, we use all N features to predict the label of each pixel for supervised learning with least-squares loss function [51].

After training the 1-D CNN as the classifier, we use beam search to select the appropriate combination of features and use classifier accuracy as the selection basis. Specifically, Step 1, we start with the initial M group of features and choose to remove the first k groups of features that cause the slightest reduction in classification accuracy to form k branches, where one group of features is removed from each branch. Step 2, in each branch, the above steps are repeated such that each branch forms k branches, and the total branches are $k \times k$. Step 3, we choose the first k branches, according to the classification accuracy from high to low. Repeating Step 2 and Step 3 until the feature groups number is reduced to the threshold θ . The process of selecting features by beam search is outlined in Algorithm 1.

To ensure the unity of the input dimensions of the classifier, we directly set the value of the features removed each time to

1 **Algorithm 1:** Beam Search for selecting the appropriate combination of feature groups

2 **Input:** feature groups set M_i , the number of feature

3 groups N , threshold θ , branch number k ,

4 classifier F

5 **Output:** selected feature groups set M_o

6 set $Q = \{M_i\}$

7 **while** $N > \theta$ **do**

8 $Q' \leftarrow \{\}$

9 **for** feature groups M in Q **do**

10 **for** feature f in M **do**

11 remove f from M

12 **if** $\text{len}(Q') == k$ **then**

13 | **if** $F(M) \geq \max(F(Q'))$ **then**

14 | | pop $Q'(\max(F(Q')))$

15 | | push M into Q'

16 | **end**

17 | **else**

18 | | push M into Q'

19 | **end**

20 **end**

21 **end**

22 $N \leftarrow N-1$

23 $Q \leftarrow Q'$

24 **end**

28 0. Finally, we obtain θ group, θ_N features as the input of the
29 target network, as described in Section III-C.

B. Superpixel-based Instance Discrimination

32 As mentioned in Section II-A, the general pretext task of
33 the optical images is Instance Discrimination with InfoNCE
34 loss function. The input of CL is usually the whole optical
35 image, with both pixel and semantic features. Significant
36 differences exist between different optical images, so the
37 common Instance Discrimination can perform well. However,
38 PolSAR mainly uses pixels as instances for CL training, which
39 has a considerable similarity between pixels. Therefore, we
40 can no longer treat all other pixels as negative samples of
41 the current pixel. To continue to take advantage of Instance
42 Discrimination, we improve the way of selecting negative
43 samples. The details are as follows:

44 We segment the PolSAR image into some superpixels,
45 as shown in Fig. 5. Specifically, we choose the classical
46 superpixel segmentation method: Simple Linear Iterative Clus-
47 tering (SLIC) [3], which is fast, memory efficient, boundary
48 adherence, and needs to set a few parameters. It used the idea
49 of clustering to cluster part of the pixels into a superpixel.
50 So the similarity of pixels within the same superpixel is
51 high and vice versa. So pixels within different superpixels
52 are defined as negative samples of each other and within the
53 same superpixel as positive samples. Assumed the size of the
54 PolSAR image is $H \times W$, and we obtain N_s superpixels, N_{s_i}
55 pixels in i th superpixel, the pixel in i th superpixel has at most
56 $N_{s_i} - 1$ positive sample, $H \times W - N_{s_i}$ negative samples.
57 Then, we can use traditional Instance Discrimination to train
58 the Heterogeneous Network described in Section III-C.

C. Heterogeneous Network

59 Inspired by SimCLR [37], PolSAR custom-made Hetero-
60 geneous Network is proposed in this article, shown in Fig.
61 3. It is similar to the traditional CL Network, which has
62 two networks; however, the architecture is different between
63 the two networks. The online network of the heterogeneous
64 network is a 2-D CNN, while the target network of the
65 heterogeneous network is a 1-D CNN. The details of the two
66 networks are as follows:

67 1) *Online network:* The online network consists of a 2-D
68 CNN with a representation encoder and a projection head. The
69 typical architecture of the representation encoder is shown in
70 Fig. 4. The encoder consists of four main parts: convolution
71 layer, pooling layer, linear embedding layer, and nonlinear
72 activation.

73 The formulation of convolution is follow

$$y_i^{(l+1)} = \sum_j^J \omega_{ij}^{(l)} \otimes x_j^{(l)} + b_i^{(l+1)} \quad (7)$$

74 where $x_j^{(l)}$ represents the j th input in the $l + 1$ th layer and
75 J represents the number of inputs, y_i^{l+1} represents the i th
76 output in the $l+1$ th layer, ω represents the kernel matrix and
77 b represents the bias, \otimes represents convolution operation.

78 The pooling operation is usually a subsampling operation
79 that reduces the input dimension and the computation amount.
80 Moreover, it facilitates the identification of displacement,
81 scaling, and other distortion invariants. It has two common
82 choices: max pooling and average pooling. In this article, we
83 choose max pooling, which selects the maximum value of the
84 specified region.

85 The linear embedding layer is a linear weighted transfor-
86 mation of the representation in a 1-D space. The formulation
87 is similar to convolution and is as follows

$$y_i^{(l+1)} = \sum_j^J \omega_{ij}^{(l)} * x_j^{(l)} + b_i^{(l+1)} \quad (8)$$

88 where $x_j^{(l)}$, J , $y_i^{(l+1)}$, ω , b is the same as convolution, $*$
89 represents vector multiplication.

90 The nonlinear activation is to improve the nonlinear ability
91 of the network. Avoid using softmax and tanh, which leads
92 to gradient vanishing. We select the ReLU, the nonlinear
93 activation function. The formula of it is as follows

$$f(x) = \max_{\text{kernel}} (0, x) \quad (9)$$

94 where x is the input, $\max_{\text{kernel}}(\cdot)$ does the kernel size specify
95 the maximum value operation in the region.

96 By combining the four operations defined above, we obtain
97 the representation encoder $f_e(\cdot)$. Then a projection head is
98 followed by $f_p(\cdot)$, which aims to embed the representation
99 into the more high semantic space and is defined as $g_p(\cdot)$. It
100 is also a linear embedding layer. So we obtain the final online
101 network $g_p(f_e(\cdot))$.

102 The input of the online network is the coherency matrix, as
103 mentioned in Section II-B. Specifically, for pixel a, we crop out

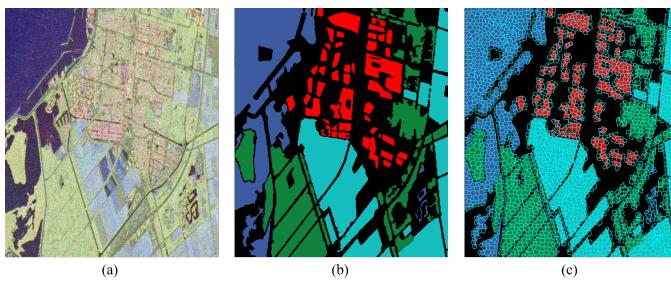


Fig. 5. RADARSAT-2 Flevoland data. (a) Pauli RGB image. (b) Ground-truth image. (c) Superpixel image.

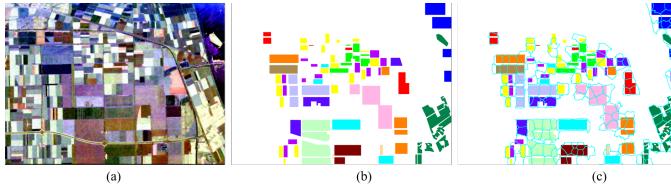


Fig. 6. AIRSAR Flevoland data. (a) Pauli RGB image. (b) Ground-truth image. (c) Superpixel image.

a pixel block of size $k \times k$ with a as the center, and the value of each pixel is represented by the straightened coherency matrix T . Finally, the input dimension is $k \times k \times 9$.

2) *Target network*: The target network is a 1-D CNN. The overall architecture and convolution, pooling, linear embedding operations, and nonlinear activation are similar to 2-D CNN, except they change from two dimensions to one. The input of the target network is the combination of the target decomposition features according to beam search as described in Section III-A. For pixel a , has θ_N features that represent a $\theta_N \times 1$ vector.

Finally, for pixel a , the online network inputs the matrix of dimension $k \times k \times 9$, and outputs the $m \times 1$ vector r_o as representation; the target network inputs the vector of dimension $\theta_N \times 1$, and outputs the vector r_{t+} whose dimension is the same as r_o . Then we use the InfoNCE loss function $L(r_o, r)$ to compute the similarity, where $r \in \{r_{t+}, r_{t-}\}$, r_{t-} represents the other target network outputs of the other pixel. After training, we can obtain the final online network to transform PolSAR data into a high-level representation. It can be used as the backbone network for downstream tasks and performs well with fine-tuning.

IV. EXPERIMENTAL RESULTS AND ANALYSIS

A. data sets Description

In this section, we employ three standard PolSAR data sets to verify the superiority of the HCLNet. They include RADARSAT-2 Flevoland, AIRSAR Flevoland, and ESAR Oberpfaffenhofen.

- RADARSAT-2 Flevoland: As shown in Fig. 5. a C-band, fully polarimetric image of the area of Netherland is obtained through the RADARSAT-2 system and was produced in April 2008. The size of the sub-image is 2375×1635 . It identifies four types of ground objects: forest, farmland, city, and water area.

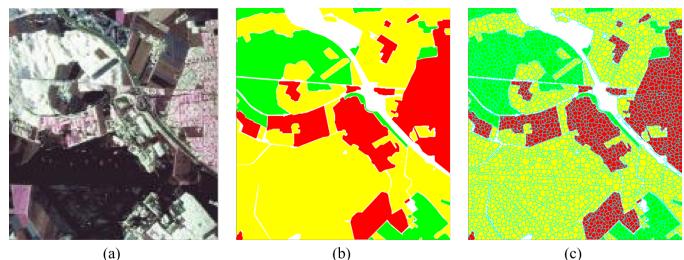


Fig. 7. ESAR Oberpfaffenhofen data. (a) Pauli RGB image. (b) Ground-truth image. (c) Superpixel image.

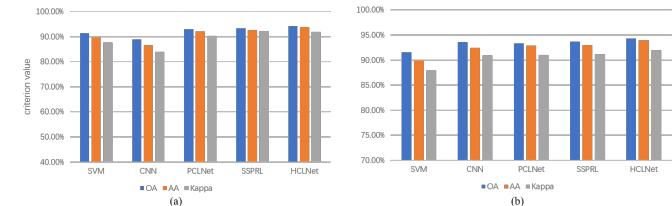


Fig. 8. Comparisons of different methods on the RADARSAT-2 Flevoland data set. (a) Few-shot result. (b) Full-sample result.

- AIRSAR Flevoland: An L-band, full polarimetric PolSAR image of the region of Flevoland, Netherlands, 750×1024 , is obtained through the NASA/Jet Propulsion Laboratory AIRSAR. There are 15 labeled objects, including forest, rapeseed, beet, bare soil, grasses, peas, lucerne, barley, buildings, potatoes, water, stembeans, and three kinds of wheat seen in Fig. 6.
- ESAR Oberpfaffenhofen: It covers Oberpfaffenhofen, Germany, which is an L-band, full polarimetric image and is obtained through ESAR airborne platform. The size of the image is 1200×1300 . Its ground-truth map can be seen in Fig. 7. It contains three classes: built-up areas, wood land, and open areas.

B. Experimental Settings

- Implement details: The online network, in which the input patch size is 15×15 , has two 2-D convolution layers in which the kernel size is 3×3 , the padding is 2×2 and 1×1 , and two linear embedding layers. The target network has two 1-D convolution layers in which the kernel size is the same as the online network, and the padding is 2×2 , and one linear embedding layer. After each convolution layer, the batch norm and max pooling layers are added. Furthermore, the same as MoCo, we use normalization in the outputs of both networks. The optimizer is the SGD which the learning rate is initialized to 0.01 and decreases with a cosine trend with training epoch, the momentum is 0.9, and weight decay is 0.0001. The τ is 0.07. The network is trained for 30 epochs and the minibatch size of 4096. All experiments were conducted independently on a single GeForce 3070 GPU with the PyTorch library.
- Multi-features: The Refined Lee filter with the window size 7×7 is used to preprocess the three PolSAR data sets to reduce the influence of speckles on the result of

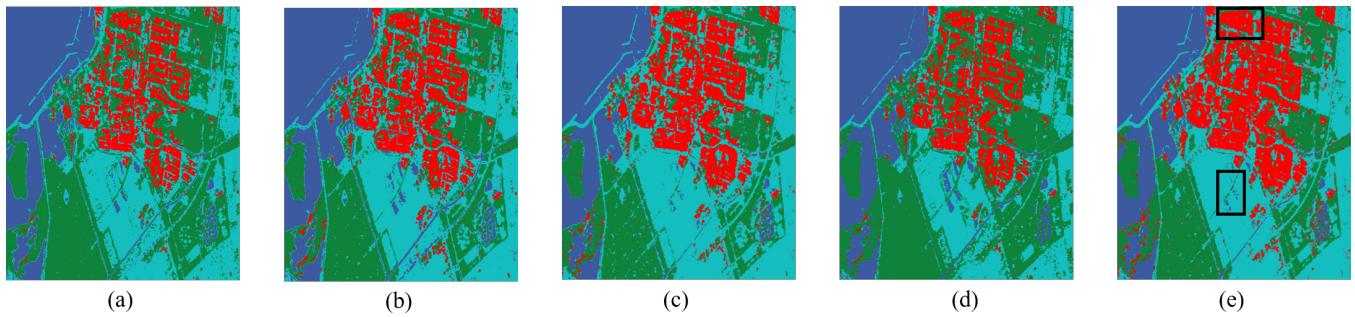


Fig. 9. Few-shot classification results with different methods on the RADARSAT-2 Flevoland data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

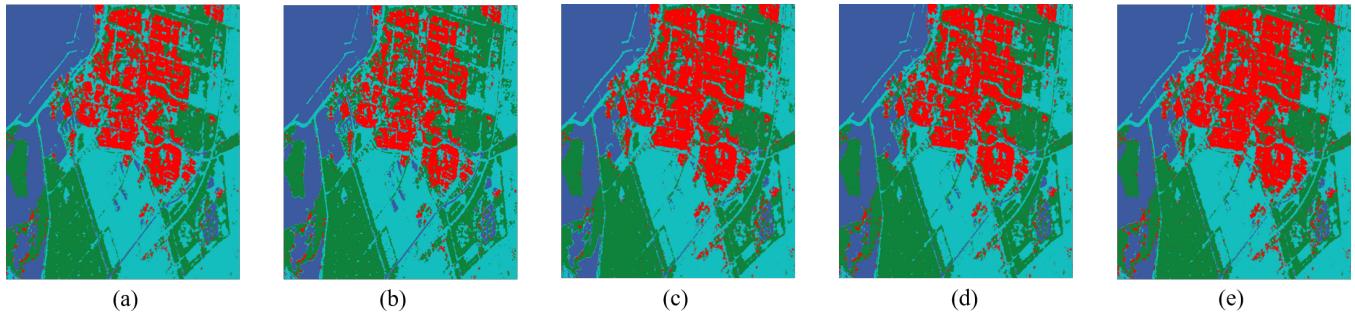


Fig. 10. Full-sample classification results with different methods on the RADARSAT-2 Flevoland data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

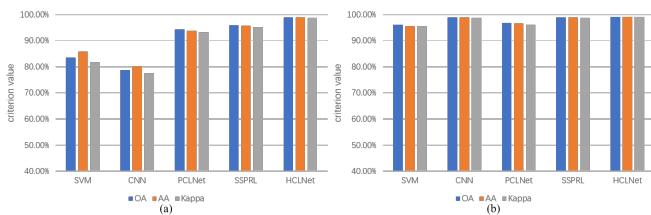


Fig. 11. Comparisons of different methods on the AIRSAR Flevoland data set. (a) Few-shot result. (b) Full-sample result.

classification. Then the same as [46], we use 14 groups of target decomposition features and obtain 70 decompositions features as the initial features of the Beam Search. The whole target decompositions are shown in Table I. These features are sufficient to represent PolSAR data.

- Compared method: To evaluate the superiority of the proposed method, we select several supervised and CL methods for comparison. Specifically, two traditional classification methods are chosen, including CNN with the same online network architecture, Support Vector Machine (SVM) [12], PCLNet [41], and SSPRL [42].

C. Experimental Results and Analysis

1) *Classification accuracy*:: In the experiment, our method is pre-trained with 10% unlabeled data in each data set. To verify the effectiveness of HCLNet, we choose 0.1% and 10% samples per category for training, denoted as few-shot and full-sample classifications. The rest of the samples are used as the test set for evaluation. The overall classification accuracy

TABLE I
ALLOF THE TARGET DECOMPOSITION FEATURES

Target Decomposition	Feature Name	Number
Krogager	sphere, diplane, helix	3
TSVM	alpha-s, phi-s, phi, tau-m	4
Neuman	delta, psi, tau	3
Huynen	(T11,T22,T33)dB	3
Holm	Holm1:(T11,T22,T33)dB, Holm2:(T11,T22,T33)dB	6
Freeman	Freeman2:(Vol,Ground)dB, Freeman3:(Odd,Dbl,Vol)dB	5
Cloude	(T11,T22,T33)dB	3
Barnes	Barnes1:(T11, T22, T33)dB, Barnes2:(T11, T22, T33)dB	6
ANNED	(Odd, Dbl, Vol)dB	3
AnYang	AnYang3:(Odd,Dbl,Vol)dB, AnYang4:(Odd, Dbl, Vol, Hlx)dB	7
H/A/ α	alpha, anisotropy, beta, delta, entropy, gamma, lambda, combination: HA, (1-H)A, H(1-A), (1-H)(1-A)	11
Yamaguchi	Yamaguchi3:(Odd, Dbl, Vol)dB, Yamaguchi4:(Odd, Dbl, Vol, Hlx)dB	7
Vanzyl	(Odd, Dbl, Vol)dB	3
MCSM	(Odd, Dbl, Vol, Hlx, Dbl-Hlx, Wire)dB	6
SUM		70

(OA), average accuracy (AA), and kappa coefficient (Kappa) are used as criteria to assess the performance of all methods.

The results of the three data sets are shown in Tables II to VI, respectively, demonstrating the superiority of HCLNet

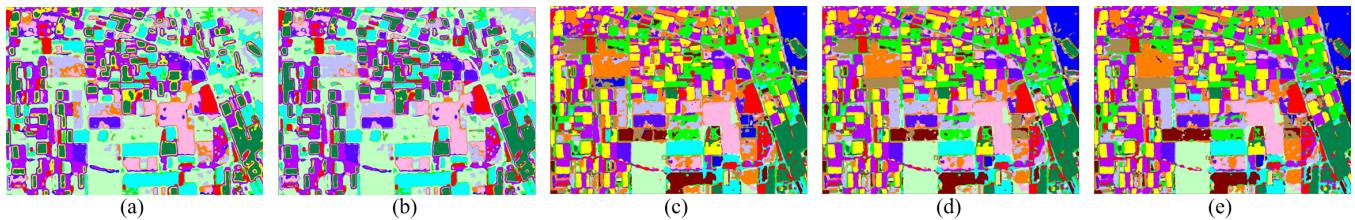


Fig. 12. Few-shot classification results with different methods on the AIRSAR Flevoland data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

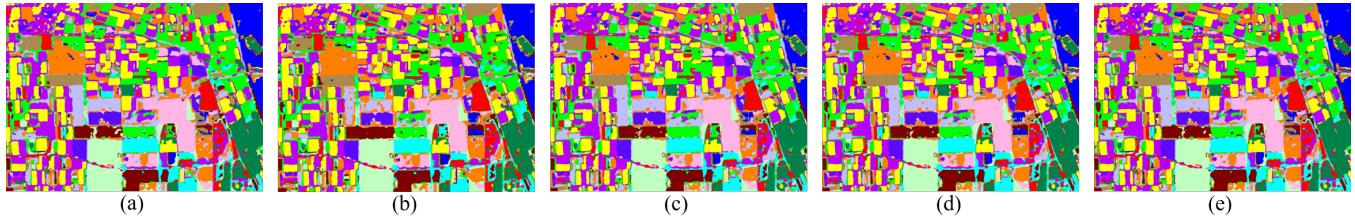


Fig. 13. Full-sample classification results with different methods on the AIRSAR Flevoland data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

TABLE II
FEW-SHOT CLASSIFICATION RESULTS (%) ON THE RADARSAT-2 FLEVOLAND WITH DIFFERENT METHODS

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Forest</i>	78.01	67.84	84.68	89.24	90.23
<i>Cropland</i>	98.40	97.73	98.26	98.19	98.55
<i>Water</i>	92.63	96.43	90.28	91.78	92.33
<i>Urban</i>	90.13	84.14	94.79	93.80	93.86
<i>OA</i>	91.31	88.89	93.01	93.25	94.15
<i>AA</i>	89.79	86.53	92.00	92.50	93.74
<i>Kappa</i>	87.61	84.00	90.20	92.15	91.84

TABLE III
FULL-SAMPLE CLASSIFICATION RESULTS (%) ON THE RADARSAT-2 FLEVOLAND WITH DIFFERENT METHODS

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Forest</i>	77.93	84.12	90.33	86.62	90.34
<i>Cropland</i>	98.13	98.84	97.90	98.68	98.36
<i>Water</i>	93.78	92.73	94.37	95.15	94.24
<i>Urban</i>	90.12	94.07	90.56	91.38	92.71
<i>OA</i>	91.55	93.55	93.29	93.70	94.30
<i>AA</i>	89.99	92.44	92.91	92.96	93.91
<i>Kappa</i>	87.95	90.92	91.06	91.17	92.04

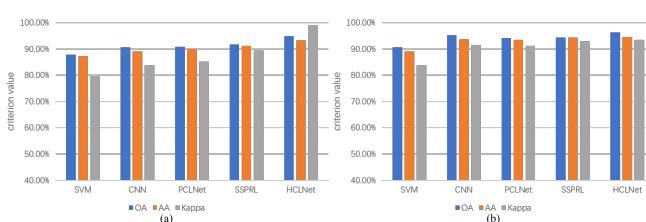


Fig. 14. Comparisons of different methods on the ESAR Oberpfaffenhofen data set. (a) Few-shot result. (b) Full-sample result.

in the small number of labeled data. Different cases will generally have different results, but the trend is consistent across different data sets.

Specifically, in Tables II and III, for RADARSAT-2

Flevoland, the traditional CNN model performs well in full-sample classification; however, its performance drops sharply, even lower than the SVM, when the number of data decreases, that is when few-shot classification. Numerically, from full-sample classification to few-shot classification, the OA, AA, and Kappa of the CNN decrease by 4.66%, 5.91%, and 6.92%. The SVM is relatively stable, dropping only 0.24%, 0.2%, and 0.34% from full-sample classification to few-shot classification, but its overall classification accuracy is not high and only 91.55% in full-sample classification. The overall number of network parameters of PCLNet is similar to that of HCLNet, which is 1.5 larger than HCLNet. Through the customized task and positive/negative sample selection of PCLNet, its accuracy is much improved compared with the traditional CNN, especially in few-shot classification, in which the over-

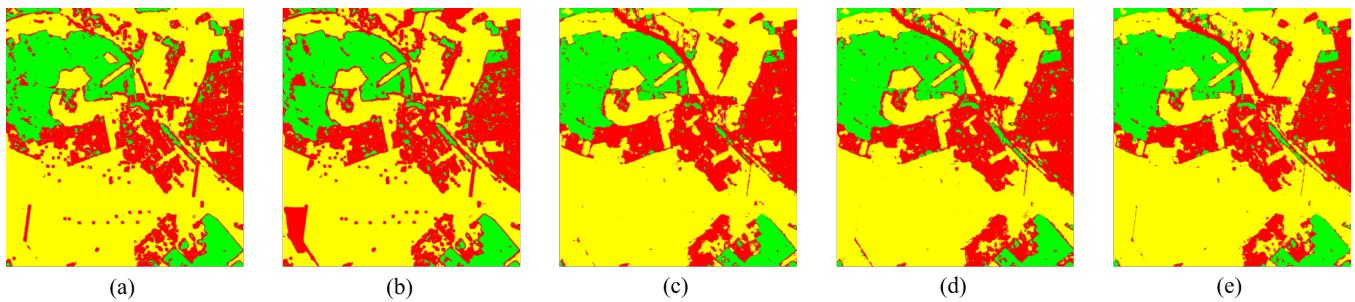


Fig. 15. Few-shot classification results with different methods on the ESAR Oberpfaffenhofen data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

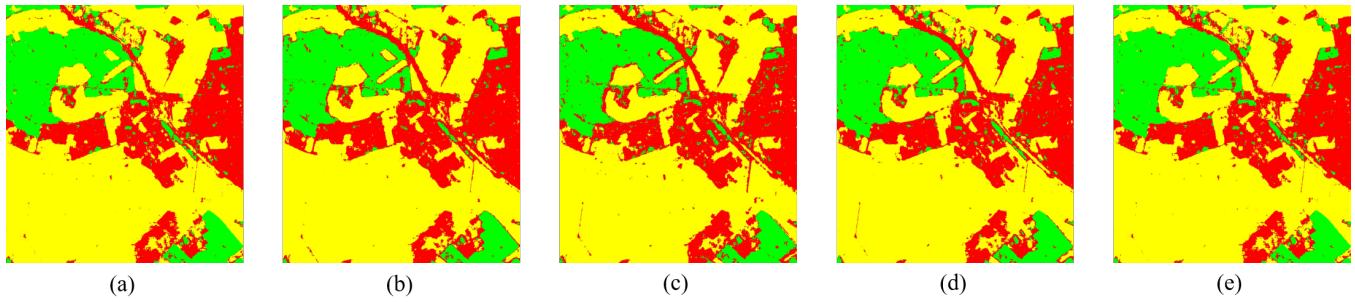


Fig. 16. Full-sample classification results with different methods on the ESAR Oberpfaffenhofen data set. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

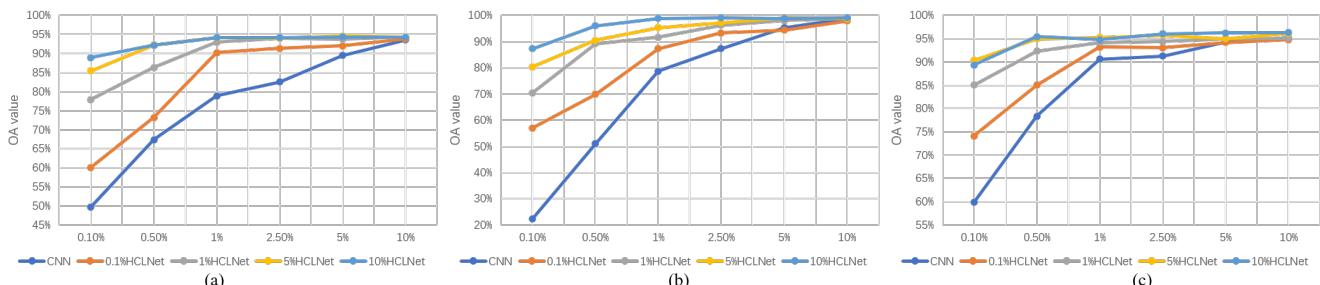


Fig. 17. Comparisons of the performance(OA) with different ratios of unlabeled and labeled samples between CNN and HCLNet on three data sets. (a) RADARSAT-2 Flevoland. (b) AIRSAR Flevoland. (c) ESAR Oberpfaffenhofen.

all improvements are 4.12%, 5.47%, and 6.20%. However, HCLNet outperforms it in both full-sample classification and few-shot classification. In addition, due to the dual contrastive learning architecture of SSPRL and the unique selection of positive samples, its performance exceeds that of PCLNet. Numerically, the OA, AA, and Kappa of SSPRL are 0.24%, 0.5%, and 1.95% higher than that of PCLNet, respectively. However, the network architecture of SSPRL is too complex, and the number of parameters is too large, which is more than ten times that of HCLNet. Furthermore, its performance fails to surpass HCLNet. The proposed HCLNet obtains the best results in full-sample and few-shots, in which the OA, AA, and Kappa are 94.15%, 93.74%, and 91.84%.

For a more straightforward comparison, Fig. 8 illustrates the results of different methods in few-shot classification and full-sample classification of RADARSAT-2 Flevoland. The result clearly shows that HCLNet comprehensively outperforms all other methods. Compared to the results of the traditional CNN

and HCLNet between few-shot and full-sample classification, we can find that the gap between the conventional CNN and HCLNet becomes larger with less labeled data, shown in Fig. 17(a) more clearly. Furthermore, the classification maps of few-shot and full-sample classification results for different methods are presented in Figs. 9 and 10. It can be observed that the HCLNet achieves the best result of the different landforms. Moreover, as indicated by the black box in Fig. 9, we find many scattered, isolated pixel in the four compared methods. In comparison, our approach can solve this problem well due to the introduction of superpixels to ensure better contextual consistency in the phase of CL.

To explore the importance of the number of unlabeled and labeled samples, we designed comparison experiments between HCLNet and its backbone model using OA as the metric. Specifically, we use the same online network of HCLNet without the projector head as the backbone model. Its parameters are initialized randomly. The backbone network

1
2 TABLE IV
3 FEW-SHOT CLASSIFICATION RESULTS (%) ON THE AIRSAR FLEVOLAND WITH DIFFERENT METHODS
4

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Buildings</i>	96.18	95.80	97.24	98.32	99.95
<i>Rapeseed</i>	42.72	35.99	87.32	90.60	96.91
<i>Beet</i>	67.34	50.34	98.33	99.02	99.77
<i>Stembeans</i>	89.29	87.91	98.18	98.89	99.93
<i>Peas</i>	86.05	86.04	95.24	96.32	99.21
<i>Forest</i>	84.60	67.14	96.38	97.93	95.11
<i>Lucerne</i>	92.06	91.84	95.02	97.28	99.13
<i>Potatoes</i>	93.37	92.35	92.37	95.00	99.84
<i>Bare soil</i>	93.95	90.23	89.88	92.13	98.52
<i>Grass</i>	90.32	84.26	85.42	89.26	98.47
<i>Barley</i>	91.81	74.58	95.77	96.23	97.75
<i>Water</i>	71.85	67.49	90.28	95.26	98.87
<i>Wheat one</i>	92.22	90.26	97.46	97.33	99.71
<i>Wheat two</i>	93.14	91.42	93.02	95.40	99.92
<i>Wheat three</i>	67.39	66.93	95.38	97.56	99.32
OA	83.45	78.72	94.30	95.83	98.80
AA	85.82	80.14	93.82	95.77	98.82
Kappa	81.79	77.43	93.27	95.23	98.73

21 TABLE V
22 FULL-SAMPLE CLASSIFICATION RESULTS (%) ON THE AIRSAR FLEVOLAND WITH DIFFERENT METHODS
23

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Buildings</i>	99.70	99.87	98.30	99.53	99.91
<i>Rapeseed</i>	94.15	94.54	93.24	96.31	97.14
<i>Beet</i>	96.77	99.69	99.01	99.02	99.78
<i>Stembeans</i>	99.92	99.98	99.26	99.55	99.96
<i>Peas</i>	97.74	98.93	98.93	99.08	99.12
<i>Forest</i>	70.66	95.78	97.41	94.37	96.04
<i>Lucerne</i>	98.18	99.54	96.02	98.99	99.81
<i>Potatoes</i>	99.77	99.76	95.44	99.82	99.74
<i>Bare soil</i>	96.64	98.80	98.21	99.05	99.26
<i>Grass</i>	97.53	97.91	90.43	98.17	98.49
<i>Barley</i>	98.74	99.30	98.09	99.56	98.77
<i>Water</i>	88.18	99.22	91.52	99.79	98.84
<i>Wheat one</i>	99.67	99.81	97.40	99.93	99.69
<i>Wheat two</i>	99.57	99.94	95.88	99.84	99.95
<i>Wheat three</i>	96.73	99.32	98.31	99.40	99.32
OA	96.10	98.81	96.77	98.84	99.05
AA	95.60	98.83	96.50	98.83	99.06
Kappa	95.57	98.75	96.14	98.79	98.99

41 and HCLNet are evaluated using 0.1%, 0.5%, 1%, 2.5%,
42 5%, and 10% samples per category. Moreover, for unlabeled
43 samples, we set different ratios of 0.1%, 1%, 5%, and 10%
44 to pre-training HCLNet. HCLNet performs similarly to the
45 backbone network with 1% labeled data when trained with
46 0.1% labeled data and 10% unlabeled data, as shown in
47 Fig. 17(a). When the ratio of labeled data is over 1%, the
48 performance of HCLNet is better than the backbone network in
49 any labeled data ratio. It fully demonstrates the effectiveness of
50 HCLNet, which learns robust high-level representations from
51 unlabeled data.

52 Similar experimental results for AIRSAR Flevoland are
53 shown in Tables IV and V; the predicted map is shown in
54 Figs. 11 and 12. Since this data set has more categories and
55 fewer data per class, the actual data amount of 1% of training
56 samples is small. It further widens the performance gap be-
57 tween HCLNet and other methods. Since there are only a few
58 superpixels in each class, the difference between samples is

59 tremendous, which makes the training of HCLNet more effec-
60 tive. The OA, AA, and Kappa of HCLNet are 2.97%, 3.05%,
61 and 3.5% higher than the best previous method, demonstrating
62 that HCLNet has a greater advantage over PCLNet and SSPRL
63 in this case. Visually, as shown in Fig. 13, SVM and CNN
64 misclassified many data due to the lack of training labeled data.
65 The performance is significantly degraded compared to the
66 other methods. As mentioned, HCLNet reduces the scattered,
67 isolated pixels to obtain better contextual consistency in the
68 maps than the other four methods.

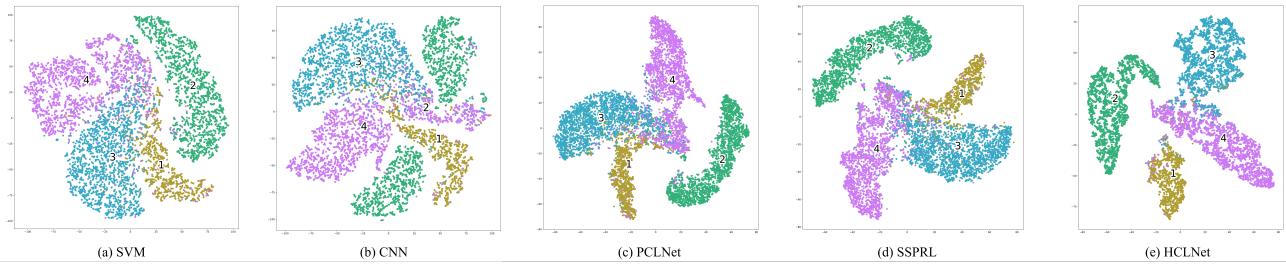
69 Furthermore, as shown in Fig. 17(b), the performance of
70 HCLNet with the ratio of labeled data is 0.5%, and the ratio
71 of unlabeled data is 10% is better than backbone network with
72 the ratio of labeled data is 5%, and achieves the highest result.
73 When only 0.1% of training samples per category are used,
74 the OA gap between the two methods is the largest, reaching
75 65.34%. In contrast, the backbone network requires at least
76 2.5% training samples per category to achieve the same result.

1
2 TABLE VI
3 FEW-SHOT CLASSIFICATION RESULTS (%) ON THE ESAR OBERPFAFFENHOFEN WITH DIFFERENT METHODS
4

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Built-up areas</i>	84.47	82.42	85.38	87.79	87.02
<i>Wood land</i>	85.54	90.52	89.70	90.88	94.50
<i>Open areas</i>	89.04	94.21	95.06	94.94	98.33
<i>OA</i>	87.81	90.59	90.83	91.78	94.80
<i>AA</i>	87.35	89.05	90.05	91.20	93.29
<i>Kappa</i>	79.87	83.82	85.21	89.32	92.99

12 TABLE VII
13 FULL-SAMPLE CLASSIFICATION RESULTS (%) ON THE ESAR OBERPFAFFENHOFEN WITH DIFFERENT METHODS
14

method	SVM	CNN	PCLNet	SSPRL	HCLNet
<i>Built-up areas</i>	82.42	87.34	89.08	90.18	90.34
<i>Wood land</i>	90.52	95.13	95.23	94.27	95.20
<i>Open areas</i>	94.21	98.68	96.30	98.35	97.92
<i>OA</i>	90.59	95.20	94.09	94.30	96.33
<i>AA</i>	89.05	93.72	93.54	94.27	94.49
<i>Kappa</i>	83.82	91.60	91.22	92.99	93.41



32 Fig. 18. T-SNE visualization of the representations learned on the RADARSAT-2 Flevoland data set with different methods. (a) SVM. (b) CNN. (c) PCLNet.
33 (d) SSPRL. (e) HCLNet.

36 It demonstrates the great generalization of the representation
37 learned by HCLNet.

38 Compared with the first two data sets, ESAR Oberpfaf-
39 fenhofen has fewer categories and a larger number of each
40 category. It may result in similar data, but our method still
41 works well in this case. As shown in Table VI, Table VII, and
42 Fig. 16, compared to the other two methods, which provide
43 a different strategy to reduce the similarity between data,
44 HCLNet is 3.97%, 3.24%, 7.78% higher than PCLNet and
45 3.02%, 2.09%, 3.76% higher than SSPRL in the few-shot.
46 It demonstrates the effectiveness of superpixel-based instance
47 discrimination to reduce the similarity between data. It still can
48 keep an excellent contextual consistency shown in Figs. 14 and
49 15. Moreover, due to the more labeled data, it can be observed
50 that when the ratio of labeled data is 0.5%, HCLNet shows
51 excellent performance, outperforming the backbone network
52 with any ratio of labeled data, as shown in Fig. 17(c).

53 To sum up, the experimental results on three data sets can
54 confirm the effectiveness of HCLNet in generalization and
55 contextual consistency.

56 2) *Representation visualization:* In the above experiments,
57 HCLNet has demonstrated the powerful ability of represen-
58 tation learning with unlabeled data using supervised-based

59 instance discrimination. In order to further explore the quality
60 of representation, 2-D t-stochastic neighbor embedding(t-SNE)
61 [52] is used to visualize the learned representation. Specifi-
62 cally, SVM, CNN, PCLNet, SSPRL, and HCLNet utilize 1%
63 labeled samples per category in the map. For SVM, we refer
64 to the coherent matrix as the representation. For CNN, we use
65 the output of the last hidden layer of the model trained in the
66 few-shot. For PCLNet, SSPRL, and HCLNet, the output of the
67 representation encoder is used as the representation without
68 any labeled samples. These visualizations are performed on
69 three benchmark data sets above. The results are shown
70 in Figs. 18, 19, and 20, different colors indicate different
71 categories.

72 From the results, it can be observed that the coherent matrix
73 cannot distinguish each category well, with existing multiple
74 categories highly overlapping. Moreover, the closeness is rela-
75 tively weak. Some features of the same category are distributed
76 in different locations and form multiple disconnected regions.

77 For CNN, under the guidance of only a few labeled data,
78 it leads to the poor improvement of closeness overlapping,
79 and some even deteriorate. Benefiting from training on large
80 amounts of unlabeled data, PCLNet and SSPRL drastically re-
81 duce the overlapping but still exist some disconnected regions.

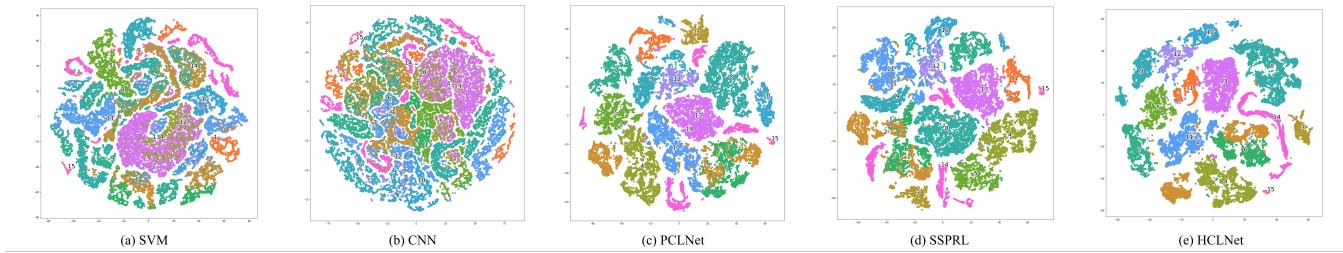


Fig. 19. T-SNE visualization of the representations learned on the AIRSAR Flevoland data set with different methods. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

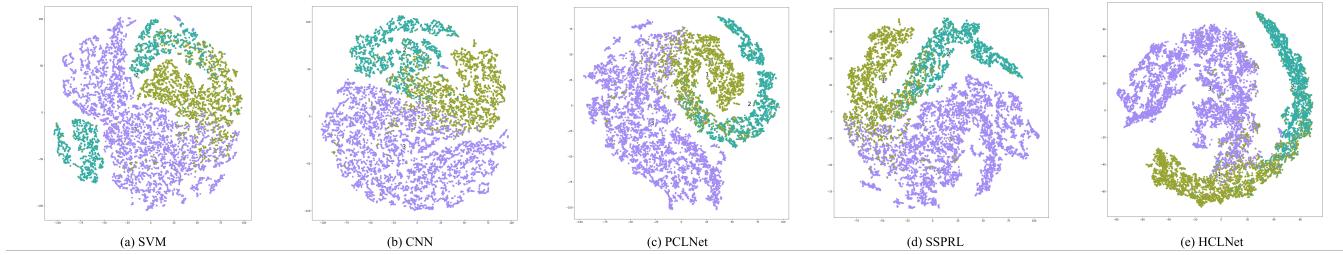


Fig. 20. T-SNE visualization of the representations learned on the ESAR Oberpfaffenhofen data set with different methods. (a) SVM. (b) CNN. (c) PCLNet. (d) SSPRL. (e) HCLNet.

On the contrary, HCLNet significantly alleviates disconnection and overlapping. It turns out that HCLNet provides a better representation by learning from different perspectives for the downstream task to improve the classification ability.

3) *Ablation study*: To better understand the effectiveness of the individual components of HCLNet, we experiment with combinations of different components on the three data sets above and design the three groups of ablation experiments. The ablation results are reported in Table VIII. The abbreviations in Table VIII represent the different components: CNN is the Siamese network in which the target network and the online network are the same network architecture, and both of them share parameters; Hnet is the heterogeneous network with the same architecture as HCLNet; SID is the superpixel-based Instance Discrimination; BS is the beam search. The results implicate that each component can improve the classification results. The first two experiments involve that compared with CNN as the architecture, Hnet exhibits the more powerful representation extraction ability and has fewer parameters. When BS is added to remove the redundancy between features, the model can learn better high-level representation, demonstrated in the experiment's third group.

TABLE VIII
OA (%) ON RADARSAT-2 FLEVOLAND (OA_{RF}), AIRSAR FLEVOLAND (OA_{AF}) AND ESAR OBERPFAFFENHOFEN (OA_{EO}) DATA SETS FOR ABLATION STUDY

method	OA_{RF}	OA_{AF}	OA_{EO}
CNN+SID+BS	90.22	95.35	91.37
Hnet+SID+BS	94.15	98.80	94.80
Hnet+SID	91.47	93.65	92.91
Hnet+BS	64.50	85.32	79.66

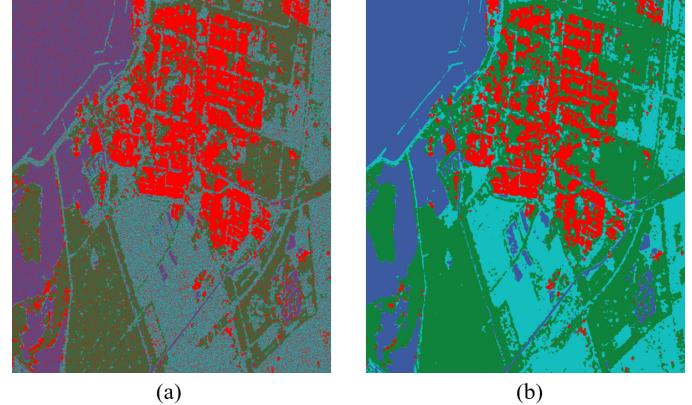


Fig. 21. (a) The salt and pepper noise problem and (b) the improved result.

Significantly, the fourth group of experiments demonstrates that the SID is essential for heterogeneous networks. As shown in Fig. 21(a), without SID, the training of HCLNet likely falls into the salt and pepper noise problem. We analyze that the network is challenging to learn the high-level representation when CL without SID to distinguish the differences between pixels and even confuse the representations between pixels of different categories. To verify our conjecture, we replace the SID and use pixels of different categories as negative samples and pixels of the same category as positive samples to increase the dissimilarity between negative samples. Finally, it successfully solves the salt and pepper noise problem, as shown in Fig. 21(b). So starting from the direction of SSL, Hnet, and SID may be a match made in heaven.

V. CONCLUSION

This article proposes a self-supervised learning method based on CL with a heterogeneous network for the first time. We use HCLNet to extract the high-level representation of PolSAR data from the physical and statistical properties and propose two plugins. The beam search is introduced to select the appropriate combination of features to reduce the redundancy of multi-target decomposition features. The superpixel-based Instance Discrimination is proposed to reduce the similarity between pixels and learn better representation. Therefore, with the help of unsupervised pre-training to learn representation, the online network can achieve high results of few-shot PolSAR classification by fine-tuning. Experiments are conducted on three widely used benchmark data sets, and the experimental results demonstrate the performance of HCLNet compared with several mainstream methods in both few-shot and full-sample classification.

Compared with the CL of optical images, PolSAR has more valuable features under different properties, while the heterogeneous network has the natural advantage of entirely using these features. This work creates a precedent for future research on heterogeneous network learning. And we believe that more in-depth and comprehensive research about the heterogeneous network may further improve the PolSAR classifier performance. Our future interest is to explore the problem of positive and negative selection for heterogeneous networks in depth.

REFERENCES

- [1] W. Nie, K. Huang, J. Yang, and P. Li, "A deep reinforcement learning-based framework for polarimetric imagery classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022.
- [2] H. Bi, F. Xu, Z. Wei, Y. Xue, and Z. Xu, "An active deep learning approach for minimally supervised polarimetric image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 9378–9395, 2019.
- [3] G. Hong, S. Wang, J. Li, and J. Huang, "Fully polarimetric synthetic aperture radar (sar) processing for crop type identification," *Photogramm. Eng. Remote Sens.*, vol. 81, no. 2, pp. 109–117, 2015.
- [4] F. T. Ulaby and C. Elachi, "Radar polarimetry for geoscience applications," 1990.
- [5] R. Shirvany, M. Chabert, and J.-Y. Tourneret, "Ship and oil-spill detection using the degree of polarization in linear and hybrid/compact dual-pol sar," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 5, no. 3, pp. 885–892, 2012.
- [6] A. Freeman and S. L. Durden, "A three-component scattering model for polarimetric sar data," *IEEE Trans. Geosci. Remote Sens.*, vol. 36, no. 3, pp. 963–973, 1998.
- [7] S. R. Cloude and E. Pottier, "An entropy based classification scheme for land applications of polarimetric sar," *IEEE Trans. Geosci. Remote Sens.*, vol. 35, no. 1, pp. 68–78, 1997.
- [8] E. Pottier, "Dr. jr huynen's main contributions in the development of polarimetric radar techniques and how the radar targets phenomenological concept becomes a theory," in *Proc. SPIE*, vol. 1748. SPIE, 1993, pp. 72–85.
- [9] J. R. Huynen, "Physical reality of radar targets," in *Proc. SPIE*, vol. 1748. SPIE, 1993, pp. 86–96.
- [10] W. L. Cameron and L. K. Leung, "Feature motivated polarization scattering matrix decomposition," in *Proc. IEEE Int. Conf. Radar*. IEEE, 1990, pp. 549–557.
- [11] E. Krogager, "New decomposition of the radar target scattering matrix," *Electron. Lett.*, vol. 18, no. 26, pp. 1525–1527, 1990.
- [12] C. Lardeux, P.-L. Frison, C. Tison, J.-C. Souyris, B. Stoll, B. Fruneau, and J.-P. Rudant, "Support vector machine for multifrequency sar polarimetric data classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 47, no. 12, pp. 4143–4152, 2009.
- [13] B. Zou, H. Li, and L. Zhang, "Polarimetric image classification using bp neural network based on quantum clonal evolutionary algorithm," in *2010 IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*. IEEE, 2010, pp. 1573–1576.
- [14] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, vol. 60, no. 6, pp. 84–90, 2017.
- [15] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2015, pp. 1–9.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.
- [17] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.
- [18] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, and J. Schmidhuber, "Lstm: A search space odyssey," *IEEE Trans. Neur. Net. Lear.*, vol. 28, no. 10, pp. 2222–2232, 2016.
- [19] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [20] X. X. Zhu, D. Tuia, L. Mou, G.-S. Xia, L. Zhang, F. Xu, and F. Fraundorfer, "Deep learning in remote sensing: A comprehensive review and list of resources," *IEEE Geosc. Rem. Sen. M.*, vol. 5, no. 4, pp. 8–36, 2017.
- [21] L. Zhang, L. Zhang, and B. Du, "Deep learning for remote sensing data: A technical tutorial on the state of the art," *IEEE Geosci. Remote Sens. Mag.*, vol. 4, no. 2, pp. 22–40, 2016.
- [22] L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson, "Deep learning in remote sensing applications: A meta-analysis and review," *ISPRS J. Photogramm.*, vol. 152, pp. 166–177, 2019.
- [23] Y. Zhou, H. Wang, F. Xu, and Y.-Q. Jin, "Polarimetric sar image classification using deep convolutional neural networks," *IEEE Geosci. Remote Sens.*, vol. 13, no. 12, pp. 1935–1939, 2016.
- [24] Z. Zhang, H. Wang, F. Xu, and Y.-Q. Jin, "Complex-valued convolutional neural network and its application in polarimetric sar image classification," *IEEE Trans. Geosci. Remote Sens.*, vol. 55, no. 12, pp. 7177–7188, 2017.
- [25] F. Liu, L. Jiao, B. Hou, and S. Yang, "Pol-sar image classification based on wishart dbn and local spatial information," *IEEE Trans. Geosci. Remote Sens.*, vol. 54, no. 6, pp. 3292–3308, 2016.
- [26] A. G. Mullissa, C. Persello, and A. Stein, "Polsarnet: A deep fully convolutional network for polarimetric sar image classification," *IEEE J. Stars.*, vol. 12, no. 12, pp. 5300–5309, 2019.
- [27] X. Tan, M. Li, P. Zhang, Y. Wu, and W. Song, "Complex-valued 3-d convolutional neural network for polarimetric sar image classification," *IEEE Geosci. Remote Sens.*, vol. 17, no. 6, pp. 1022–1026, 2019.
- [28] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *J. Mach. Learn. Res.*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [29] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proc. 27th Int. Conf. Mach. Learn. (ICML)*, 2010, pp. 807–814.
- [30] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Highway networks," *arXiv preprint arXiv:1505.00387*, 2015.
- [31] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 11, pp. 4037–4058, 2020.
- [32] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 3733–3742.
- [33] M. Ye, X. Zhang, P. C. Yuen, and S.-F. Chang, "Unsupervised embedding learning via invariant and spreading instance feature," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 6210–6219.
- [34] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [35] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*. Springer, 2020, pp. 776–794.
- [36] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 9729–9738.
- [37] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *Int. Conf. Mach. Learn. (ICML)*. PMLR, 2020, pp. 1597–1607.

- 1 [38] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya,
2 C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*,
3 “Bootstrap your own latent-a new approach to self-supervised learning,”
4 *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, vol. 33, pp. 21 271–21 284,
5 2020.
6 [39] X. Chen and K. He, “Exploring simple siamese representation learning,”
7 in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp.
8 15 750–15 758.
9 [40] B. Ren, Y. Zhao, B. Hou, J. Chanussot, and L. Jiao, “A mutual
10 information-based self-supervised learning model for polsar land cover
11 classification,” *IEEE Trans. Geosci. Remote Sens.*, vol. 59, no. 11, pp.
12 9224–9237, 2021.
13 [41] L. Zhang, S. Zhang, B. Zou, and H. Dong, “Unsupervised deep repre-
14 sentation learning and few-shot classification of polsar images,” *IEEE*
15 *Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–16, 2020.
16 [42] W. Zhang, Z. Pan, and Y. Hu, “Exploring polsar images representation
17 via self-supervised learning and its application on few-shot classifica-
18 tion,” *IEEE Geosci. Remote Sens.*, vol. 19, pp. 1–5, 2022.
19 [43] Y. Cui, F. Liu, X. Liu, L. Li, and X. Qian, “Tcspanet: two-staged
20 contrastive learning and sub-patch attention based network for polsar
21 image classification,” *ISPRS J. Photogramm. Remote Sens.*, vol. 14,
22 no. 10, p. 2451, 2022.
23 [44] R. Zhang, P. Isola, and A. A. Efros, “Split-brain autoencoders: Un-
24 supervised learning by cross-channel prediction,” in *Proc. IEEE Conf.*
25 *Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 1058–1067.
26 [45] D. Eigen and R. Fergus, “Predicting depth, surface normals and semantic
27 labels with a common multi-scale convolutional architecture,” in *Proc.*
28 *IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 2650–2658.
29 [46] C. Yang, B. Hou, B. Ren, Y. Hu, and L. Jiao, “Cnn-based polarimetric
30 decomposition feature selection for polsar image classification,” *IEEE*
31 *Trans. Geosci. Remote Sens.*, vol. 57, no. 11, pp. 8796–8812, 2019.
32 [47] S. R. Cloude and E. Pottier, “A review of target decomposition theorems
33 in radar polarimetry,” *IEEE Trans. Geosci. Remote Sens.*, vol. 34, no. 2,
34 pp. 498–518, 1996.
35 [48] Y. Yamaguchi, T. Moriyama, M. Ishido, and H. Yamada, “Four-
36 component scattering model for polarimetric sar image decomposition,”
37 *IEEE Trans. Geosci. Remote Sens.*, vol. 43, no. 8, pp. 1699–1706, 2005.
38 [49] J. J. van Zyl, “Application of cloude’s target decomposition theorem
39 to polarimetric imaging radar data,” in *Radar polarimetry*, vol. 1748.
40 SPIE, 1993, pp. 184–191.
41 [50] L. Zhang, B. Zou, H. Cai, and Y. Zhang, “Multiple-component scattering
42 model for polarimetric sar image decomposition,” *IEEE Geosci. Remote*
43 *Sens.*, vol. 5, no. 4, pp. 603–607, 2008.
44 [51] W. A. Holm and R. M. Barnes, “On radar polarization mixed target state
45 decomposition techniques,” in *Proc. IEEE Int. Conf. Radar.* IEEE,
46 1988, pp. 249–254.
47 [52] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne.” *J.*
48 *Mach. Learn. Res.*, vol. 9, no. 11, 2008.
49
50
51
52
53
54
55
56
57
58
59
60