

Progressive Data Augmentation Method for Remote Sensing Ship Image Classification based on Simulation Samples and Neural Style Transfer

Qi Xiao, Bo Liu, Zengyi Li, Wei Ni, Zhen Yang, and Ligang Li*

Abstract—Deep learning has shown great power in processing remote sensing data, especially for fine-grained remote sensing ship image classification. However, the lack of a large amount of effective training data greatly limits the performance of neural networks. Based on current data augmentation methods, images of ships on the sea generated for remote sensing have the problem of distortion, blurring, and poor diversity. To tackle this problem, we propose a novel progressive remote sensing ship image data augmentation method that combines ship simulation samples and a neural style transfer-based network to generate a large amount of transferred remote sensing ship images. Our method consists of two stages. The first stage uses a visible light imaging simulation system to generate ship simulation samples through 3D models of real images. This stage can significantly increase the diversity of the training dataset. For the second stage, to eliminate the domain gap between real ship images and ship simulation samples, a few real images and a newly designed neural style transfer-based network called *Sim2RealNet* are employed to realize style transfer from simulation samples to real images. The proposed method was applied to a variety of ship targets to verify its effectiveness compared to other data augmentation methods on remote sensing image classification tasks. The experimental results demonstrate the effectiveness of the proposed method.

Index Terms—Domain gap, image classification, neural style transfer, remote sensing, ship simulation samples.

I. INTRODUCTION

WITH the rapid improvement of computer technology and GPU computing power, given sufficient data, deep learning has shown a strong dominance in the field of computer vision, such as image classification [[1],[2],[3]], object detection [[4],[5],[6]], and semantic segmentation [[7],[8]]. Following the success of deep learning and the increasing availability of remote sensing data, deep learning has been playing an increasingly important role in the field of remote sensing. With sufficient remote sensing data for training, researchers have focused on designing convolutional neural

networks (CNNs) to perform feature selection [[9],[10],[11],[12]], extraction [[13],[14]], and coding [[15]] on high-resolution remote sensing images, thereby improving network performance. Meanwhile, remote sensing data also bring unprecedented challenges to deep learning. For example, many application domains do not have access to big data, such as images of targeted ships, which are difficult to capture. At the same time, the labeling of remote sensing images requires considerable human and material resources. A lack of sufficient data will lead to the CNN overfitting problem.

On the one hand, from the model's architecture itself, many standard techniques have been proposed to alleviate the overfitting problem in the case of insufficient data. Nitish *et al.* [[17]] proposed the dropout method, which is a regularization technique that allows the network to learn more robust features by randomly zeroing some of the neuron values during the CNN training process. Another regularization technique is batch normalization [[18]], which normalizes the activation value of each layer into a vector with mean 0 and variance 1 to improve the stability and training speed of the neural network. This technique has now become a standard data preprocessing method for neural networks. Transfer learning [[19],[20]], pretraining [[21]], and one-shot [[22]] and zero-shot learning [[23]] algorithms can also effectively relieve the overfitting problem of deep neural networks.

On the other hand, in terms of training data to solve the overfitting problem, researchers have proposed the concept of data augmentation, a data-space solution to the problem of limited data that aims to intelligently augment the quality and quantity of the original small amount of data through a suite of techniques to improve deep neural network performance. Data augmentation can generate the most possible effective data from any amount of available data. Currently, the data augmentation methods most commonly used and considered most effective include traditional transformation, deep learning-based, and imaging simulation system-based methods.

Q. Xiao and B. Liu are with the Key Laboratory of Electronics and Information Technology for Space System, National Space Science Center, Chinese Academy of Sciences, Beijing 100190, China (e-mail (Q. X.): xiaoqi19@mails.ucas.ac.cn) and the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 100049, China (e-mail (B. L.): liubo183@mails.ucas.ac.cn).

Z. Li, W. Ni, Z. Yang, and L. Li (*corresponding author) are with the Key Laboratory of Electronics and Information Technology for Space System, National Space Science Center, Chinese Academy of Sciences, Beijing 100190, China (*e-mail: liligang@nssc.ac.cn).

The first type mainly generates new samples by, e.g., rotating, cropping, and adding random noise to the original images [[24]]. The second type mainly encode and decode the input images to generate new images using neural networks, such as generative adversarial networks (GANs) first introduced by Goodfellow *et al.* [[25]]. The third type uses imaging simulation systems to generate a large number of simulation samples and combine these samples with real images via CNNs [[26],[27],[28]].

In this paper, we propose a novel progressive data augmentation method (Fig. 1) consisting of two stages. The first stage uses a small number of targeted ship images from the original training data to implement visible light imaging simulation system (VLISS) to generate a large number of high-quality simulation samples of the targeted ship through 3D models of the real images. In the second stage, to eliminate the domain gap between the simulated samples and real images, which mainly manifests in the background of the two domains, a few real images and a neural style transfer (NST)-based network (*Sim2RealNet*) designed by us using DenseNet-121 [[29]] are employed to realize style transfer from the simulation samples to real images. We conducted extensive experimental verifications of the proposed method on remote sensing image classification tasks. The main contributions of this study can be summarized as follows:

- 1) We built a VLISS and utilized the system to generate a large number of simulation samples of a targeted ship for style transfer, and we designed *Sim2RealNet* to realize style transfer from the simulation samples to real images.
- 2) We introduced the structural similarity index measure (SSIM) [[30]] into our framework to achieve high-quality transferred image selection and *Sim2RealNet* model optimization.
- 3) We applied the proposed method to a variety of ship targets and conducted extensive experiments to demonstrate that our proposed method can achieve excellent performance and efficiency on fine-grained remote sensing ship image classification tasks.

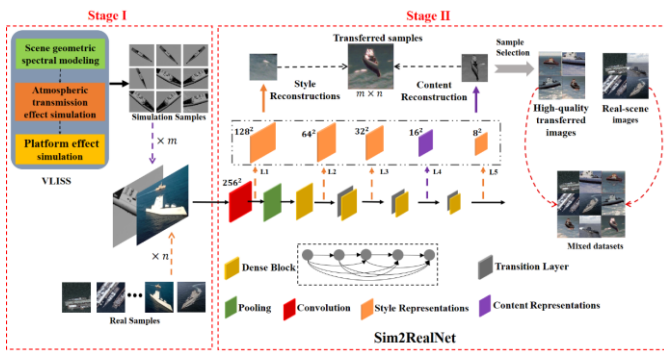


Fig. 1. Framework of our data augmentation method. First, for a small number of real images of the targeted ship, we use VLISS to generate m ship simulation samples through a 3D model. The original data set contains n ship images. After that, we make each real image a style image and each ship simulation sample a content image, and we input them into *Sim2RealNet* to obtain $m \times n$ transferred samples. We select the k highest-quality transferred images through the SSIM algorithm and add them to the original dataset to constitute a mixed dataset. The mixed dataset is used to train the CNN for image classification.

II. RELATED WORK

A. Data Augmentation Methods

Currently, there are three main categories of image augmentation method. The first is the traditional white-box methods based on space and color manipulations, which can be divided into single-image and mixed-image methods. The second is black-box methods based on deep learning proposed in recent years, which include variational autoencoders (VAEs) [[31]] and GANs. The third has emerged over the past few years and consists of imaging simulation systems to generate simulation samples for data augmentation. This subsection describes each of these three categories in detail.

1) Traditional White-box Data Augmentation Methods

Single-image augmentation methods perform all manipulations around a single image itself, which mainly includes geometric and color transformations, such as zoom in/out, flipping, cropping, rotation, noise injection, and grayscale. These methods are very easy to implement compared to other methods and were originally mainly applied on ImageNet [[32]] to increase the number of images for training CNNs.

In contrast, mixed-image augmentation refers to the simultaneous use of multiple images to generate one image. The synthetic minority oversampling technique (SMOTE) [[33]] uses interpolation to generate new samples by sampling small samples in the feature space of their neighboring samples to deal with the problem of sample imbalance, thereby improving classifier performance. However, mixed-image augmentation has problems such as marginalization and blindness when selecting neighbors.

Mixup [[34]] is a data augmentation method based on the principle of neighborhood risk minimization proposed by the Facebook Artificial Intelligence Research Institute and MIT in the paper “Beyond Empirical Risk Minimization”. It uses linear interpolation to obtain new sample data. Lusa *et al.* [[35]]. proved that for high-dimensional data, the above two data augmentation algorithms have difficulty playing an effective role. Ioune *et al.* [[36]] proposed a method called sample pairing to realize data augmentation, which randomly selects two images and calculates the value of each pixel on each RGB channel of the two images as the pixel value of the new image. This method reduced the error rate from 8.22% to 6.93% on the CIFAR-10 dataset [[37]].

2) Deep Learning Black-box Data Augmentation Methods

Among the data augmentation methods based on deep learning, CNNs have demonstrated powerful feature extraction ability. A VAE maps an input image to a low-dimensional space through an encoder, learns the inherent encoding method of the image, and simultaneously uses a decoder to recover the image to the original size. Through the encoder and decoder, a VAE can generate new samples. The traditional VAE has a problem in that the generated images are relatively fuzzy. Although many studies have continued to improve its structure regarding this problem [[38]], because the VAE focuses on minimizing the interpolation between image pixels, there is still a lack of vivid image performance.

In 2014, Goodfellow [[25]] proposed GANs. By building a GAN, the features of the original data can be flexibly supplemented, modified, and reconstructed to form a new feature vector with attributes similar to the original sample and generate a new sample. A GAN is a two-player game between two networks called the discriminator and generator. The generator generates fake images from the given noise, and the discriminator effectively distinguishes whether the image is a real image or one generated by the generator. Owing to the constant competition between the discriminator and generator, the images generated by a GAN constantly approach the real image. Finally, the discriminator can no longer distinguish whether an image comes from the real data or the generator. In a deep convolutional GAN [[39]], both the discriminator and generator of use a CNN to replace the multi-layer perceptron in the traditional GAN. At the same time, to make the entire network differentiable, the pooling layer in the CNN is removed, and the fully connected layer is replaced with a global pooling layer to reduce the amount of calculation. Karras *et al.* [[40]] proposed the progressive growing GAN to generate high-resolution images of size 1024×1024 for the first time. This method uses a progressive approach to generate images. First, a 4×4 image is generated. When the network is trained to a certain level, the resolution is gradually increased to 8×8 , 16×16 , and finally up to 1024×1024 . In the field of remote sensing image generation, Lin *et al.* [[41]] designed a multiple-layer feature-matching generative adversarial networks to generate images containing simple objects, but the images were still blurred and distorted.

3) Imaging Simulation Data Augmentation Methods

With the development of imaging simulation systems and virtual engines, some attempts have been made in recent years to use simulation or virtual samples to train network models instead of real scene samples.

Li *et al.* [[42]] released a virtual image dataset. They used CityEngine and Unity3D to construct a large-scale urban street scene with the aim of applying the virtual images to traffic vision research. By using transfer learning, Wang *et al.* [[43]] applied the simulation samples in an infrared band to the object detection task, which effectively improved the performance of the detector under the condition of extreme scarcity of real scene samples. In semantic segmentation research, Richter *et al.* [[44]] presented a method for building virtual datasets via modern video games and obtained the corresponding annotations using outside graphics hardware without access to the source code of the game. The above works applied virtual images directly as the training set to increase the feature diversity. However, ignoring the domain gap between simulation images and real scene images, when the number of simulation images is much larger than those of real images, the network will appear over-fitted to the simulated features. Nowadays, GANs are widely used to produce photorealistic synthetic images [[45]]; however, these images lack the corresponding annotations.

All of the above data augmentation methods have many limitations for a small number of remote sensing ship images. The traditional data augmentation methods can only simply

increase the quantities of the dataset, and the increased data will easily cause the neural network to fall into overfitting. For deep learning data augmentation methods, the training of GANs relies on a large amount of data, and the generated images will have poor diversity and blurring problems. For imaging simulation data augmentation methods, the generated simulation samples and real images will often have large domain gaps, so the improvement of neural network performance is very limited. At the same time, the addition of simulation samples will reduce the accuracy of other types of target recognition. The proposed progressive data augmentation method can effectively solve these abovementioned problems as follows. In its first stage, the VLISS can generate a large number of ship simulation grayscale samples from different angles and under different environmental conditions through several 3D models of real images. The structures and sizes of the simulated samples of these ships are strictly consistent with the real samples. At this stage, the simulation system can generate a large number of new simulation samples with diversity, which solves the problems of poor diversity, blurring, and distortion in current data augmentation methods. In the second stage, because there is a huge domain gap between the simulated samples and real images in their backgrounds, colors, and textures, the designed *Sim2RealNet* is employed to align the style-level features between the simulation samples and real images to realize style transfer from the simulation domain to the real domain, thereby improving the generalization ability of the network. We use style-aligned images to train the network and increase its robustness.

B. Neural Style Transfer

Due to the powerful feature extraction capabilities of CNNs, a style transfer network uses a CNN to render the content of an image into images of different styles. Style transfer networks have a wide range of applications, such as the generation of famous paintings and migration of starry sky backgrounds. A style transfer network has two inputs: content image and style image. A style transfer network extracts the features of the style image at its shallow level as the style representation and extracts the features of the content image at its deep level as the content representation. Through the training of the neural network, the generated image is close to the style image in style and content image in content.

It is not appropriate to directly add the simulation samples generated by VLISS to the original dataset to train the neural network because the simulation samples and real images have two different domains. From the perspective of computer graphics, the information of an image can be divided into style and content information. The style information mainly includes the texture, background, and color of the image, and the content information includes the shapes, structures, and positions of the objects in the image. The shape, structure, and position of the targeted ship in a simulation sample and real image remain highly similar, so the content information of the two is essentially the same. However, for style information, the simulation samples generated by VLISS are single-channel grayscale images, which lack the rich style information of real

remote sensing images. Therefore, the domain gap between the simulation samples and real images is mainly due to the large difference in style information. Inspired by the NST network, we use a CNN to extract the content features of the simulation samples and the style features of the real images and use a style transfer network to complete the style transfer from the simulated samples to the real images.

The closest method to ours is the NST method first proposed by Gatys *et al.* [[46]] in 2016. We made two improvements to this method. First, we used DenseNet-121 to replace VGG-19 [[2]] in the original method as the backbone. Second, we added structural similarity loss to improve the loss function. We will introduce our designed *Sim2RealNet* in detail in Section 3.2.

C. Structural Similarity

Structural similarity is an evaluation index used to measure the similarity between two images. It was first proposed by the Laboratory for Image and Video Engineering of the University of Texas at Austin in 2004 [[30]]. The SSIM algorithm evaluates the similarity of two images according to the following equation:

$$SSIM_{(x,y)} = \frac{(2\mu_x\mu_y+c_1)(2\sigma_{xy}+c_2)}{(\mu_x^2+\mu_y^2+c_1)(\sigma_x^2+\sigma_y^2+c_2)} \quad (1)$$

where x and y are input images, μ_x and μ_y are the mean values of x and y , σ_x^2 and σ_y^2 are variances of x and y , respectively, and σ_{xy} is the covariance of x and y . $SSIM_{(x,y)}$ takes values in the range of -1 to 1. When images x and y are highly similar, $SSIM_{(x,y)} = 1$; otherwise, $SSIM_{(x,y)} = -1$.

SSIM measures the similarity between two images in terms of image brightness, contrast, and covariance. Therefore, using SSIM to measure the similarity between the image *Gen* generated by the style transfer network and the input simulation sample *Sim* is very effective. We use the SSIM value to achieve both high-quality generated sample selection and loss function improvement of *Sim2RealNet*.

III. METHODS

In this section, we introduce our progressive data augmentation method in detail. In Section 3.1, we explain the basic principles of the VLISS and how we can generate a large number of simulation samples. In Section 3.2, we introduce the architecture of *Sim2RealNet*.

A. Visible Light Imaging Simulation System

Based on the analysis of the visible light imaging radiation process and the mechanism of radiation transmission, we arrange the parameters that affect the characteristics of radiation transmission, including spatial scale, observation orientation, and atmospheric conditions, in order of importance. These parameters are organized into a discrete observation space as the input of the imaging simulation framework. At the same time, combined with the simulation principle of visible light imaging, the simulation of samples is conducted via the following steps (Fig. 2).

Scene diversity construction. This step mainly includes ship target geometry modeling, ship target spectral characteristic modeling, sea surface dynamic geometric

modeling, and sea surface spectral characteristic modeling. Among them, ship target geometric structure modeling is achieved according to the typical real ship size and structure based on 3ds Max. Spectral feature modeling of ship targets can generate multispectral data curves by collecting the spectral data of typical materials. Sea surface dynamic geometry modeling generates a dynamic sea surface level using a geometry grid modeling method based on the Phillips sea wave spectrum. Spectral characteristic modeling of the sea surface is achieved using the model from [[47]] to generate a multispectral curve of the sea surface.

Atmospheric data off-line calculation. We apply MODTRAN to calculate the solar and sky irradiance and atmospheric downward transmittance of the ocean scene under specific observed conditions. At the same time, the upward path radiation and atmospheric upward transmittance of the detection platform and the moving ship target are calculated under the observation conditions. Furthermore, we organize and manage these parameters using look-up tables (LUTs).

Incident radiation calculation. According to the simulation conditions, the atmospheric LUT is used to obtain the solar and sky irradiance, atmospheric downward transmittance, and other parameters. The incident radiation energy is calculated one-by-one for each patch unit of the ship target and sea background.

Outward radiation calculation. After calculating the incident radiation, we obtain the corresponding reflectance spectral curve data from the material data table based on the observed azimuth information and related attributes of the scene geometry and material. Based on this, the scene radiation energy at the entrance pupil is calculated by adding the effects

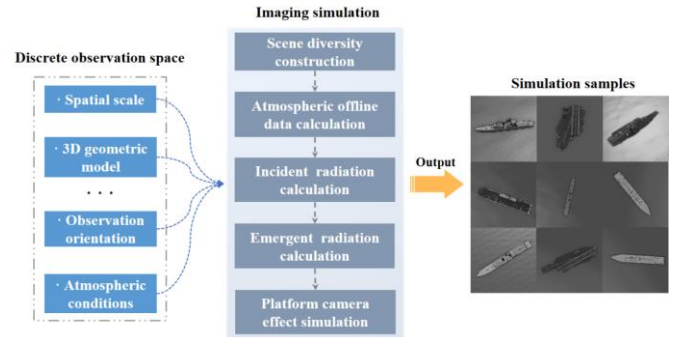


Fig. 2. Framework of VLISS. We organized imaging parameters (e.g., spatial size, 3D geometric model, observation orientation, atmospheric conditions, etc.) into a discrete observation space (left). These imaging parameters are used as the input of the VLISS. The imaging simulation module (middle) contains the main steps to complete the full link simulation modeling and generates a large number of high-quality simulation images (right).

of atmospheric range radiation and atmospheric upward transmittance.

Platform camera effect simulation. After the outward radiation calculation, the geometric relationship between the scene and camera focal plane is mapped by considering the platform position, attitude, and camera imaging geometry. At the same time, according to the photoelectric conversion calibration coefficient, we transform the radiation value before the pupil into a quantitative gray value of the camera and complete the simulation image generation under the current platform and observation conditions.

Based on the VLISS described above, we performed imaging simulations of six types of ships. The number of real images of these ships was very low. Fig. 3 shows several of the simulation results of the visible light band under the conditions of 15° solar altitude, mid-latitude summer atmospheric model, and marine aerosol type. The squint distance is the distance between the camera and ship target. The observation angle is determined by the observation altitude and azimuth angles. The observation altitude angle is the angle between the optical axis of the starting camera and the vertical axis at sea level. The observation azimuth angle is the angle between the projection at sea level and the direction of the bow (positive in the clockwise direction and negative in the anticlockwise direction).

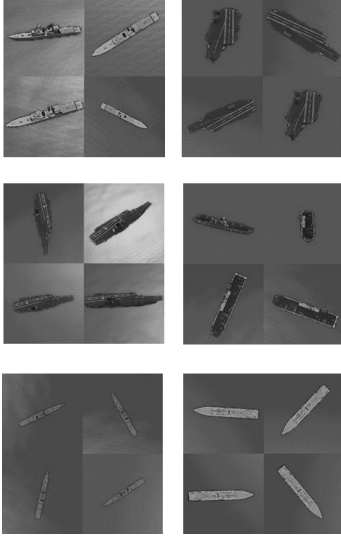


Fig. 3. Some simulation results for six target ships using our VLISS. From top left to bottom right: Asagri, Kittyhawk, Kuznetsov, Lzumo, Murasame, Sacramento.

Table I lists the detailed simulation parameters of all the simulation samples shown in Fig. 3.

TABLE I

DETAILED SIMULATION PARAMETERS CORRESPONDING TO FIG. 3. THERE ARE 4 SIMULATION SAMPLES FOR EACH TYPE OF SHIP: 1 (TOP LEFT), 2 (TOP RIGHT), 3 (BOTTOM LEFT), AND 4 (BOTTOM RIGHT).

Ship	Sample no.	Simulation Parameters		
		Squint distance (km)	Observation azimuth ($^\circ$)	Observation altitude angle ($^\circ$)
Asagri	1	1.53	87.3	88.2
	2	1.37	62.5	90.0
	3	1.10	67.8	89.1
	4	2.05	-42.4	90
Kittyhawk	1	1.12	15.7	83.1
	2	1.66	-33.7	90
	3	1.53	58.9	90
	4	1.37	31.5	72.6
Kuznetsov	1	1.66	-3.4	90.0
	2	1.12	-127.8	88.2
	3	1.12	-99.1	90.0

Lzumo	4	1.34	-90.0	90.0
	1	1.56	-66.2	80.2
	2	1.42	0.00	82.2
	3	1.66	37.8	90.0
Murasame	4	1.68	137.8	90.0
	1	1.13	-138.4	90.0
	2	1.51	166.2	90.0
	3	1.24	-159.2	90.0
Sacramento	4	1.33	72.1	90.0
	1	1.22	92.1	90.0
	2	1.17	46.5	90.0
	3	1.19	90.0	90.0
	4	1.22	-45.2	90.0

B. Sim2RealNet

This section describes the architecture of the proposed network in detail. The design of *Sim2RealNet* strongly follows the rules of NST and borrows insights from Gatys *et al.* [[46]].

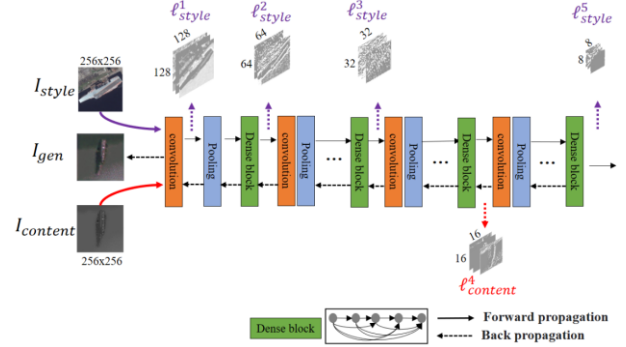


Fig. 4. Architecture of *Sim2RealNet*. I_{style} represent a real image, $I_{content}$ represents a simulation sample, and I_{gen} represents a transferred image generated by *Sim2RealNet*. ℓ_{style}^i represents the style loss of the i -th layer, and $\ell_{content}^i$ represents the content loss of the i -th layer.

In the following, we compare each step of our network to those of another NST network. The full architecture of *Sim2RealNet* is illustrated in Fig. 4.

Background. To easily understand our work, we briefly summarize the NST algorithm by Gatys *et al.* [[46]]. Given style image S and content image C , the NST network produces an output image O , which appears as S in style and C in content. The algorithm minimizes the objective function

$$\mathcal{L}_{total} = \sum_{l=1}^L \alpha_l \mathcal{L}_{content}^l + K \sum_{l=1}^L \beta_l \mathcal{L}_{style}^l \quad (2a)$$

with

$$\mathcal{L}_{content}^l = \frac{1}{2N_l D_l} \sum_{ij} (F_l[O] - F_l[C])_{ij}^2 \quad (2b)$$

$$\mathcal{L}_{style}^l = \frac{1}{2N_l^2} \sum_{ij} (G_l[O] - G_l[S])_{ij}^2 \quad (2c)$$

where L represents the number of convolutional layers used to extract image features, and l represents the l -th convolutional layer. N_l is the number of filters in each convolutional layer, and D_l represents the size of the feature map in each layer. The Gram matrix $G_l[\cdot] = F_l[\cdot] F_l[\cdot]^T$ is defined as the inner product of the feature maps. α_l and β_l are the weights of each layer to

configure the layer preferences. K is a weight that balances *content loss* (2b) and *style loss* (2c).

Network architecture. Among the previous NST networks, most of them [[48], [49], [50], [51]] continue to employ the pre-trained VGG-19 network [[2]] proposed in 2014 as the feature extractor, while we employ the pre-trained DenseNet-121 [[29]] proposed in 2018 to replace VGG-19 in the original method as our backbone network. DenseNet is a CNN with dense connections. In this network, there is a direct connection between any two layers; that is, the input of each layer of the network is the union of the outputs of all previous layers, and

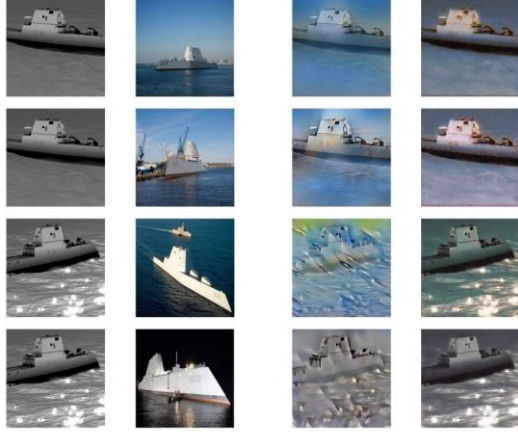


Fig. 5. Comparison of DenseNet against VGG-19 as our feature extractor.

the feature map learned by this layer will also be directly passed to all the following layers as input. In Fig. 5, the results show that employing DenseNet as our feature extractor can achieve better style transfer than VGG-19.

DenseNet-121 consists of 121 layers, including four dense blocks. We refer the reader to [26] for more details on DenseNet-121; this paper will not describe it in detail. As shown in Fig.4, we selected the output features of five layers in DenseNet to build the feature extractor, which is the output after the first convolutional layer $L1$ and the output of the four dense blocks $L2$, $L3$, $L4$, and $L5$.

Loss Function. In *Sim2RealNet*, for the network to realize style transfer from a simulation sample to a real image faster, the output image O is initialized to content image C , which is the simulation sample. At the same time, we also introduced a new *ssim loss* as follows:

$$\mathcal{L}_{ssim} = SSIM(O, C) \quad (3)$$

Because O is initialized to C , at the beginning of the network, we have $SSIM(O, C) \approx 1$, and we must assign the network a larger penalty loss. With the training of the network, $SSIM(O, C)$ decreases alongside \mathcal{L}_{ssim} when O gradually appears to be image I in style, which is in line with our expectations.

We formulate the *Sim2RealNet* objective function by combining all three components as follows:

$$\mathcal{L}_{total} = \sum_{l=1}^L \alpha_l \mathcal{L}_c^l + K \sum_{l=1}^L \beta_l \mathcal{L}_s^l + W \mathcal{L}_{ssim} \quad (4)$$

where W is a weight that balances *ssim loss*.

Implementation details. We now describe the implementation details of *Sim2RealNet*. We employed pre-

trained DenseNet-121 as the feature extractor. As mentioned above, we selected the output of five layers ($L1$ – $L5$) in DenseNet-121 for content and style feature representations. The effects of using different layers to build the content and style representations are illustrated in Fig. 6. We chose the highest $SSIM_{(O,G)}$ value, which was $L4$ as the content representation ($\alpha_l = 1$ for this layer and $\alpha_l = 0$ for other layers) and $L1$, $L2$, $L3$, $L4$, and $L5$ as the style representations ($\beta_l = 1$ for these layers and $\beta_l = 0$ for other layers). Another important parameter is W ; the effect of W is illustrated in Fig. 7. We chose the best result, $W = 25$, which has the highest $SSIM_{(O,G)}$ value. The other parameter $K = 1000$ for all the results. Our work is based on the

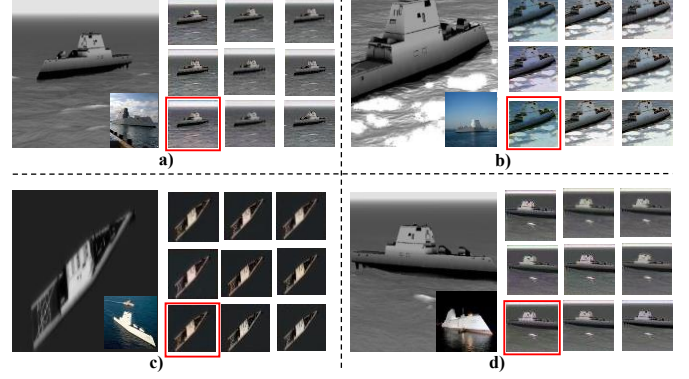


Fig. 6. Using different layers to build content representation and style representation. In each row, from left to right, using ($L4$), ($L5$), and ($L4$, $L5$) to build the content representation. In each column, from top to bottom, using ($L1$, $L2$, $L3$), ($L1$, $L2$, $L3$, $L4$), and ($L1$, $L2$, $L3$, $L4$, $L5$) to build style representation. The highest average $SSIM_{(O,G)}$ of the transferred images marked with a red box are equal to **0.909**.

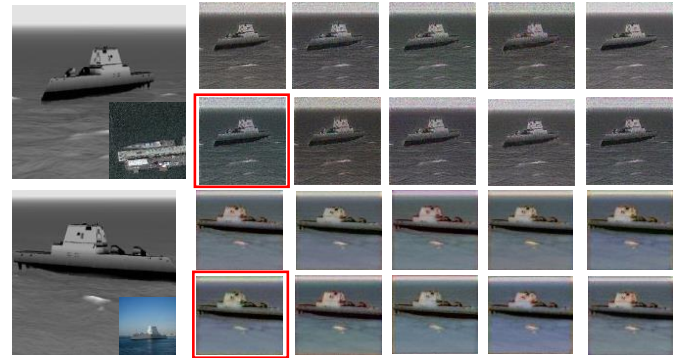


Fig. 7. Effect of using different W values. We set 10 groups of W parameters, from 0 to 45 with increments of 5. From left to right and top to bottom. The highest $SSIM_{(O,G)}$ of the transferred images marked with a red box are equal to **0.942**.

implementation of Gatys *et al.* [[46]]. Our code is available at <https://github.com/xiaoqi25478/Progressive-Data-Augmentation-Method>.

IV. EXPERIMENTS

In this section, we report our extensive experiments conducted to verify the effectiveness of our method compared with traditional and imaging simulation data augmentation methods on remote sensing ship image classification tasks. We used three datasets, each containing 12 types of ships: *Asagri* (AS), *Lzumo* (LZ), *Kuznetsov* (KU), *Sacramento* (SA), *Kittyhawk* (KI), *Murasame* (MU), *ArleighBruke* (AR), *Nimitz*

(NI), Wasp (WA), Whitby (WH), Midway (MI), and Independence (IN).

A. Datasets

TABLE II
SETTINGS OF ORIGINAL DATASET.

Ship	Training set	Test set	Total
Asagri (AS)	35	35	70
Lzumo (LZ)	32	31	63
Kuznetsov (KU)	34	34	68
Sacramento (SA)	25	25	50
Kittyhawk (KI)	34	34	68
Murasame (MU)	32	31	63
ArleighBruke (AR)	407	174	581
Nimitz (NI)	388	165	553
Wasp (WA)	318	135	453
Whitby (WH)	195	83	278
Midway (MI)	146	62	208
Independence (IN)	148	62	210

1) Real Image Dataset

To evaluate our method, we used FGSCR-42 [[52]], which is a public dataset for fine-grained ship classification in remote sensing images. The entire dataset contains approximately 9320 images, which are divided into 42 categories. The number of ships in each category varies significantly. To fully examine the performance of the proposed method, we selected six types of ships that are rare in number and applied our proposed data augmentation method to them. Considering the distribution of the numbers of each ship in the real dataset, we also selected six types of ship targets with a larger number in the FGSCR-42 dataset. For these six types of ship targets, we did not apply any data augmentation methods. Overall, the 12 types of ships made up the original dataset for our experiments, as shown in Fig. 8.

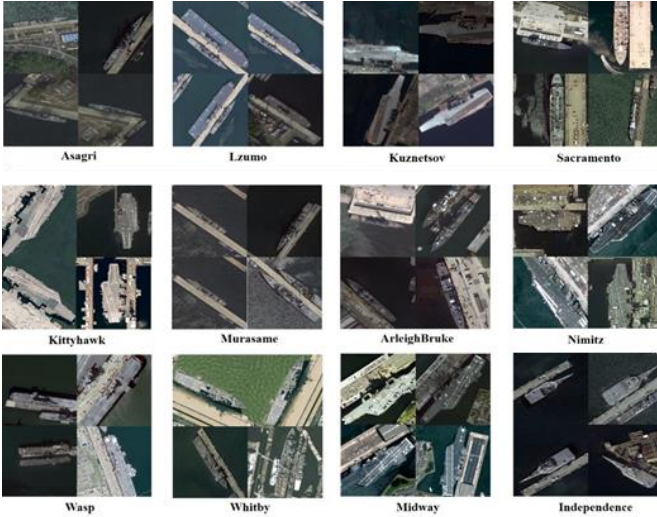


Fig. 8. Our original dataset consisting of 12 types of ships, of which six types (Asagri, Lzumo, Kuznetsov, Sacramento, Kittyhawk, and Murasame) are scarce, and the remaining six types have an appropriate number.

As shown in Table II, we split the original dataset into a training set and test set with scientific proportions.

2) Simulation Sample Dataset

To compare our method with the data augmentation method based on imaging simulation, we used VLISS to obtain simulation samples of six types of ships: Asagri, Lzumo, Kuznetsov, Sacramento, Kittyhawk, and Murasame. The

number of simulation samples for each ship was ten times its number in the training set of the original dataset, as listed in Table III

TABLE III
SETTINGS OF SIMULATION SAMPLE DATASET.

Ship	Training set	Test set	Total
Asagiri	350	0	350
Lzumo	320	0	320
Kuznetsov	340	0	340
Sacramento	250	0	250
Kittyhawk	340	0	340
Murasame	320	0	320

3) Transferred Sample Dataset

Similarly, for our progressive data augmentation method, we strictly took the following steps to generate transferred samples. First, for each of the six types of ships that require data augmentation, we used all the simulation samples in the simulation sample dataset and real images in the original dataset to perform one-to-one style transfer. Second, we used the $SSIM_{(sim, transferred)}$ value to measure the quality of each transferred sample. The larger the $SSIM$ value, the better the transfer effect of the transferred sample is. Finally, we selected a certain number of transferred samples for each ship in the order of $SSIM_{(sim, transferred)}$ from high to low, and the number of transferred samples for each ship was 10 times that in the

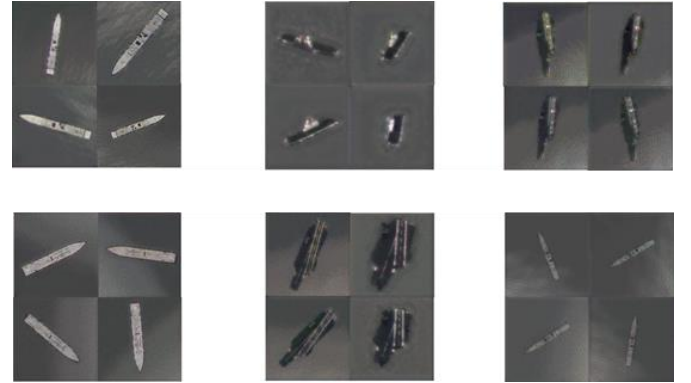


Fig. 9. Some transferred results for six target ships by our *Sim2RealNet*. From top left to bottom right: Asagri, Lzumo, Kuznetsov, Sacramento, Kittyhawk, and Murasame.

training set of the original dataset, as listed in Table IV

Fig. 9 shows some of the transferred samples for each ship by *Sim2RealNet*.

TABLE IV
SETTINGS OF TRANSFERRED SAMPLE DATASET.

Transferred samples based on <i>Sim2RealNet</i>			
Ship	Training set	Test set	Total
Asagiri	350	0	350
Lzumo	320	0	320
Kuznetsov	340	0	340
Sacramento	250	0	250

Kittyhawk	340	0	340
Murasame	320	0	320

B. Experimental Settings

To fully test the performance of our proposed method, we conducted 46 sets of experiments, divided into two types: independent and mutual experiments. In the independent experiments, we only used one data augmentation method for each experiment and compared it with other methods to illustrate the excellent performance of our method (ID 0 to 29). In the mutual experiments, each experiment used other data augmentation methods and our method at the same time to verify the compatibility of our method in cooperating with other data augmentation methods for CNNs (ID 30 to 45). As listed in Table V, for traditional data amplification methods (*Trad Augs*), we used a variety of augmentation methods that are currently commonly used. These included *random crop*, *random horizontal flip*, *cutout*, *random erasing*, *dual cutout*, *mix up*, *ricap*, *cutmix*, and *label smoothing*. We set up a series of trapezoidal experiments using imaging simulation methods (*Sim Augs*) and our method (*Ours*). In the table, 1× mixed means adding one times images, which are simulation sample (imaging simulation method) or transferred sample (our method) to the training set in the original dataset, etc. It should be noted that only *Asagri (AS)*, *Lzumo (LZ)*, *Kuznetsov (KU)*, *Sacramento (SA)*, *Kittyhawk (KI)*, and *Murasame (MU)* have simulation samples and transferred samples, and the number of each type of ship in the test set was constant throughout.

TABLE V
FULL LISTINGS OF ALL EXPERIMENTS.

Method	ID	Operations	ID	Trad Augs	Sim Augs	Ours
Only Trad Augs	0	<i>no Augs</i>	30	✓	×	×
	1	<i>random crop</i>	31	✓	×	2× mixed
	2	<i>random horizontal flip</i>	32	✓	×	4× mixed
	3	<i>cutout</i>	33	✓	×	6× mixed
	4	<i>random erasing</i>	34	✓	×	8× mixed
	5	<i>dual cutout</i>	35	✓	×	10× mixed
	6	<i>mix up</i>	36	×	2×	×
	7	<i>ricap</i>	37	×	mixed	2× mixed
	8	<i>cutmix</i>	38	×	4×	×
Only Sim Augs	9	<i>label smoothing</i>	39	×	mixed	4× mixed
	10	1× mixed	40	×	6×	×
	11	2× mixed	41	×	mixed	6× mixed
	12	3× mixed	42	×	8×	×
	13	4× mixed	43	×	mixed	8× mixed
	14	5× mixed	44	×		×

15	6× mixed	45	×	10× mixed	10× mixed
16	7× mixed				
17	8× mixed				
18	9× mixed				
19	10× mixed				
Only Ours	20	1× mixed			
	21	2× mixed			
	22	3× mixed			
	23	4× mixed			
	24	5× mixed			
	25	6× mixed			
	26	7× mixed			
	27	8× mixed			
	28	9× mixed			
	29	10× mixed			

C. Implementation Details

We conducted all experiments in PyTorch 1.7.0, CUDA 11.1, and Python 3.8.7, and we used an NVIDIA RTX 3090 GPU for training and testing. We used the average precision (AP) of each ship to measure the performances of all methods in this study.

For remote sensing ship image classification tasks, we used ResNet-34[[53]] as our classification network, where we set the input image size to 256×256, learning rate to 0.003, batch size to 16, and we adopted a mini-batch SGD optimizer [[54]] with a momentum of 0.9. We used the regularization method, and the weight decay was set to 1e-4. Each network was trained for 100 epochs.

For *Sim2RealNet*, we set the input image size to 256×256, learning rate to 0.01, batch size to 64, and adopted the Adam optimizer [[55]] with a momentum of 0.3. We also used the regularization method, and the weight decay was set to 1e-3.

D. Experimental Results

We used mean AP (mAP) as the evaluation method for all data amplification methods in our experiments. The reason for this is that in the field of computer vision, mAP is a general evaluation method for computer vision tasks.

Table VI lists the results of the independent experiments. The experimental groups with Exp ID 7, 12, and 16 represent the best results of the traditional methods, imaging simulation methods, and our method, respectively. From these three results, our method improves the mAP value of remote sensing ship classification tasks by almost 4%. The experimental results show that our method exhibits excellent performance as an independent data augmentation method.

Table VII lists the results of the mutual experiments. From the comparative experimental results of ID 31–35 to 30, 37 to 36, etc. Our method has excellent performance in improving the performance of CNNs based on other data augmentation methods.

V. CONCLUSION

In this paper, we proposed a progressive data augmentation method that combines simulation samples and neural style transfer for remote sensing ship image classification. Compared to other data augmentation methods, our method can effectively solve the problems of distortion, blurring, and poor diversity of the generated images. Extensive experimental results prove that, despite the scarcity of real images, our method can still effectively improve the accuracy of remote sensing ship classification tasks without reducing the recognition accuracy for other types of ships. We further expanded the application

prospects of simulated images in the field of deep learning. However, our method still has many disadvantages, such as the generated image resolution being too low, the time cost of style transfer being large, and the amount of augmented data generated still inevitably causing the CNN overfitting problem. In future work, we will focus on solving these problems and aim to apply this method to remote sensing image target detection tasks.

TABLE VI
RESULTS OF INDEPENDENT EXPERIMENT (ID 0 TO 29).

Method	ID	Operations	AP (%)												mAP (%)
			AS	LZ	KU	SA	KI	MU	AR	NI	WA	WH	MI	IN	
Base	0	<i>no Augs</i>	71.4 3	78.4 3	67.6 5	72.0 0	63.4 4	58.0 6	82.8 8	90.9 4	92.04	71.0 8	92.68	85.4 8	77.18
Only Trad Augs	1	<i>random crop</i>	75.8 6	74.1 9	67.6 5	76.0 0	79.4 1	64.5 2	85.6 3	96.3 6	100.0 0	85.5 4	100.0 0	88.7 1	82.82
	2	<i>random horizontal flip</i>	74.2 9	80.3 2	70.5 9	68.0 0	69.4 1	61.2 9	87.3 6	96.9 7	99.26	81.9 3	100.0 0	88.7 1	81.51
	3	<i>cutout</i>	74.2 9	74.1 9	73.5 3	64.0 0	64.7 1	59.6 1	89.6 6	94.5 5	100.0 0	74.7 0	96.77	93.5 5	79.96
	4	<i>random erasing</i>	80.0 0	82.3 3	73.5 3	72.0 0	76.4 7	64.1 9	86.7 8	94.5 5	99.26	81.9 3	100.0 0	88.7 1	83.31
	5	<i>dual cutout</i>	77.1 4	80.6 5	64.7 1	52.0 0	67.6 5	54.8 4	88.1 3	90.3 0	97.78	74.7 0	93.55	93.5 5	77.92
	6	<i>mix up</i>	62.8 6	70.9 7	66.9 8	44.0 0	68.2 4	58.0 6	86.7 8	86.0 6	94.81	66.2 7	87.10	85.4 8	73.13
	7	<i>ricap</i>	73.1 3	80.6 5	76.4 7	72.0 0	79.4 1	67.7 4	89.5 3	91.5 2	100.0 0	81.9 3	100.0 0	91.9 4	83.69
	8	<i>cutmix</i>	79.8 6	83.8 7	73.5 3	56.0 0	73.5 3	59.6 1	89.6 6	93.3 3	99.26	78.3 1	93.55	91.9 4	81.04
	9	<i>label smoothing</i>	80.0 0	83.8 7	61.7 6	56.0 0	70.5 9	51.6 1	88.5 1	91.5 2	100.0 0	77.1 1	96.77	91.9 4	79.14
Only Sim Augs	10	1× mixed	74.2 9	80.6 5	64.7 1	52.0 0	67.6 5	51.6 1	90.8 0	90.9 1	100.0 0	77.1 1	98.39	90.3 2	78.20
	11	2× mixed	77.1 4	83.5 5	70.5 9	52.0 0	79.4 1	67.7 4	91.9 5	92.1 2	100.0 0	75.9 0	100.0 0	91.9 4	81.86
	12	3× mixed	80.0 0	80.3 2	76.4 7	76.0 0	72.3 5	64.5 2	94.2 5	92.7 3	100.0 0	80.7 2	100.0 0	88.7 1	83.84
	13	4× mixed	82.8 6	81.2 9	78.2 4	60.0 0	68.3 5	75.9 7	91.3 8	92.7 3	100.0 0	78.3 1	98.39	91.9 4	83.29
	14	5× mixed	83.8 8	77.4 2	73.5 3	60.0 0	76.4 7	54.8 4	93.1 0	92.1 2	100.0 0	81.9 3	96.77	87.1 0	81.43
	15	6× mixed	82.8 6	80.6 5	64.7 1	60.0 0	67.3 5	61.2 9	93.1 0	93.3 3	100.0 0	83.1 3	100.0 0	88.7 1	81.26
	16	7× mixed	82.5 1	74.1 9	73.5 3	56.0 0	65.2 9	61.2 9	91.9 5	93.3 3	100.0 0	81.9 3	100.0 0	88.7 1	80.73
	17	8× mixed	81.3 3	83.8 7	72.3 5	72.0 0	70.5 9	67.7 4	87.3 6	96.9 7	100.0 0	79.5 2	100.0 0	90.3 2	83.50
	18	9× mixed	84.1 5	83.8 7	73.5 3	56.0 0	79.4 1	61.2 9	90.8 0	93.9 4	100.0 0	80.7 2	100.0 0	87.1 0	82.57
	19	10× mixed	80.0 0	77.4 2	64.7 1	56.0 0	76.4 7	54.8 4	92.5 3	92.7 3	98.52	81.9 3	96.77	85.4 8	79.78
Only Ours	20	1× mixed	88.0 3	87.1 0	83.5 3	56.0 0	70.5 9	67.5 4	90.2 3	89.7 0	100.0 0	78.3 1	96.77	88.7 1	83.04
	21	2× mixed	88.5 7	82.1 1	79.4 1	48.0 0	76.3 3	71.2 9	89.0 8	93.3 3	100.0 0	84.3 4	100.0 0	87.1 0	83.30
	22	3× mixed	85.0 0	85.6 3	73.5 3	52.0 0	77.0 1	76.5 5	89.6 6	88.4 8	99.26	79.5 2	98.39	90.3 2	82.95
	23	4× mixed	82.8 6	87.1 0	91.1 8	64.0 0	79.2 1	67.7 4	90.8 0	94.5 5	100.0 0	81.9 3	98.39	93.5 3	85.94
	24	5× mixed	87.1 4	84.1 9	70.5 9	60.0 0	73.5 3	75.5 4	87.9 3	94.5 5	99.26	81.9 3	98.39	91.9 4	83.75
	25	6× mixed	87.1 4	84.1 9	70.5 9	52.0 0	82.3 5	70.9 7	90.2 3	92.7 3	100.0 0	74.7 0	100.0 0	91.9 4	83.07

26	7× mixed	82.8 6	87.1 0	84.7 1	78.0 0	82.3 5	74.6 6	87.3 6	94.5 5	100.0 0	84.3 4	100.0 0	91.9 4	87.32
27	8× mixed	86.0 0	83.8 7	85.2 9	82.0 0	79.4 1	72.3 3	91.3 8	92.7 3	100.0 0	79.5 2	100.0 0	88.7 1	86.77
28	9× mixed	87.1 4	90.3 2	82.5 3	52.0 0	79.4 1	72.1 2	89.6 6	93.3 3	99.26	78.3 1	98.39	90.3 2	84.40
29	10× mixed	89.4 3	83.8 7	83.2 0	78.0 0	70.5 9	77.4 2	89.6 6	90.9 1	99.26	80.7 2	100.0 0	87.1 0	85.85

TABLE VII
RESULTS OF MUTUAL EXPERIMENT (ID 30 TO 45).

ID	Trad Augs	Sim Augs	Ours	AP (%)												mAP (%)
				AS	LZ	KU	SA	KI	MU	AR	NI	WA	WH	MI	IN	
30	√	×	×	78.6 5	87.8 7	74.0 2	72.0 0	77.65	64.8 8	88.3 3	93.9 4	94.81	79.2 1	94.11	87.2 1	82.72
31	√	×	2× mixed	85.7 1	87.1 0	77.1 8	88.0 0	87.12	67.7 4	89.0 8	98.7 9	94.81	89.1 6	100.0 0	93.5 5	88.19
32	√	×	4× mixed	88.5 7	90.0 6	79.4 1	84.0 0	100.0 0	74.1 9	81.6 1	99.3 9	94.07	83.1 3	98.39	90.3 2	88.60
33	√	×	6× mixed	79.5 7	91.3 3	79.4 1	76.0 0	97.06	74.1 9	84.4 8	96.9 7	97.78	87.9 5	98.39	98.3 9	88.46
34	√	×	8× mixed	80.0 0	89.1 0	77.4 3	80.0 0	94.12	70.9 7	83.3 3	95.1 5	97.78	84.3 4	98.39	95.1 6	87.15
35	√	×	10× mixed	83.1 8	89.8 7	85.4 8	74.0 0	94.12	64.5 2	88.5 1	96.9 7	92.59	84.3 4	100.0 0	83.8 7	86.45
36	×	2×	×	77.1 4	83.5 5	70.5 9	52.0 0	79.41	67.7 4	91.9 5	92.1 2	100.0 0	75.9 0	100.0 0	91.9 4	81.86
37	×	mixed	2× mixed	82.8 6	87.1 0	79.4 1	76.0 0	82.35	77.7 4	92.5 3	93.9 4	100.0 0	81.9 3	98.39	91.9 4	87.01
38	×	4×	×	82.8 6	81.2 9	78.2 4	60.0 0	68.35	75.9 7	91.3 8	92.7 3	100.0 0	78.3 1	98.39	91.9 4	83.29
39	×	mixed	4× mixed	85.7 1	90.3 2	70.5 9	80.0 0	82.35	79.5 2	89.6 6	95.7 6	100.0 0	79.5 2	100.0 0	90.3 2	86.98
40	×	6×	×	82.8 6	80.6 5	64.7 1	60.0 0	67.35	61.2 9	93.1 0	93.3 3	100.0 0	83.1 3	100.0 0	88.7 1	81.26
41	×	mixed	6× mixed	84.2 9	87.4 2	82.3 5	72.0 0	76.47	74.1 9	95.2 3	96.9 7	100.0 0	80.7 2	100.0 0	90.3 2	86.66
42	×	8×	×	81.3 3	83.8 7	72.3 5	72.0 0	70.59	67.7 4	87.3 6	96.9 7	100.0 0	79.5 2	100.0 0	90.3 2	83.50
43	×	mixed	8× mixed	84.1 2	86.4 4	79.4 1	78.0 0	79.41	67.7 4	91.9 5	96.9 1	100.0 0	80.7 2	100.0 0	93.5 5	86.52
44	×	10×	×	80.0 0	77.4 2	64.7 1	56.0 0	76.47	54.8 4	92.5 3	92.7 3	98.52	81.9 3	96.77	85.4 8	79.78
45	×	mixed	10× mixed	80.0 0	83.8 7	70.5 9	58.0 0	70.59	70.9 7	89.0 8	93.9 4	100.0 0	83.1 3	100.0 0	88.7 1	82.41

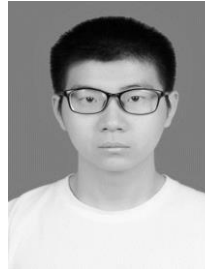
REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2012, pp. 1097–1105.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *IEEE International Conference on Learning Representation (ICLR)*, 2015.
- [3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [4] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
- [5] X. Sun, P. Wang, C. Wang, Y. Liu, and K. Fu, "PBNet: Part-based convolutional neural network for complex composite object detection in remote sensing imagery," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 173, pp. 50–65, 2021. DOI: <https://doi.org/10.1016/j.isprsjprs.2020.12.015>.
- [6] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [7] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [8] H. Noh, S. Hong, and B. Han, "Learning deconvolutional network for semantic segmentation," in *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [9] H. Gu, H. Li, L. Yan, Z. Liu, T. Blaschke, and U. Soergel, "An Object-Based Semantic Classification Method for High Resolution Remote Sensing Imagery Using Ontology," *Remote Sens.*, vol. 9, pp. 329, 2017. DOI: <https://doi.org/10.3390/rs9040329>.
- [10] X. Sun, B. Wang, Z. Wang, H. Li, H. Li, and K. Fu, "Research Progress on Few-Shot Learning for Remote Sensing Image Interpretation," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 2387–2402, 2021. DOI: 10.1109/JSTARS.2021.3052869.
- [11] Liu, G.; Gao, L.; Qi, L. "Hyperspectral Image Classification via Multi-Feature-Based Correlation Adaptive Representation," *Remote Sens.*, vol. 13, pp. 1253, 2021. DOI: <https://doi.org/10.3390/rs13071253>.
- [12] S. Xiang, Q. Xie, and M. Wang, "Semantic Segmentation for Remote Sensing Images Based on Adaptive Feature Selection Network," in *IEEE Geoscience and Remote Sensing Letters*. DOI: 10.1109/LGRS.2021.3049125.
- [13] W. Hu, Y. Huang, L. Wei, F. Zhang, and H. Li, "Deep convolutional neural networks for hyperspectral image classification," *Journal of Sensors*, vol. 2015, no. 258619, 2015.
- [14] W. Liu, M. Yang, M. Xie, Z. Guo, E. Li, L. Zhang, T. Pei, and D. Wang, "Accurate Building Extraction from Fused DSM and UAV Images Using a Chain Fully Convolutional Neural Network," *Remote Sens.*, vol. 11, pp. 2912, 2019. <https://doi.org/10.3390/rs11242912>.
- [15] Q. He, X. Sun, Z. Yan, and K. Fu, "DABNet: Deformable Contextual and Boundary-Weighted Network for Cloud Detection in Remote Sensing Images," *IEEE Transactions on Geoscience and Remote Sensing*. DOI: 10.1109/TGRS.2020.3045474.
- [16] Q. Chen, L. Wang, Y. Wu, G. Wu, Z. Guo, and S. L. Waslander, "Aerial imagery for roof segmentation: A large-scale dataset towards automatic mapping of buildings," *arXiv preprint*, 2018. DOI: arXiv:1807.09532.
- [17] 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.
- [18] S. Ioffe and C. Szegedy, "Batch normalization: accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning (ICML)*, 2015.
- [19] K. Weiss, T. M. Khoshgoftaar, and D. Wang, "A survey of transfer learning," *J. Big Data*, vol. 3, pp. 9, 2016.
- [20] L. Shao, "Transfer learning for visual categorization: a survey," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 5, pp. 1019–1034, 2015.
- [21] D. Erhan, Y. Bengio, A. Courville, P. A. Manzagol, P. Vincent, "Why does unsupervised pre-training help deep learning?" *J. Mach. Learn. Res.*, vol. 11, pp. 625–660, 2010.
- [22] J. Pennington, R. Socher, C.D. Manning, "GloVe: global vectors for word representation," in *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 2014, pp. 12.
- [23] A. Halevy, P. Norvig, and F. Pereira, "The unreasonable effectiveness of data," *IEEE Intell. Syst.*, vol. 24, pp. 8–12, 2009.
- [24] J. Wang and L. Perez, "The effectiveness of data augmentation in image classification using deep learning," 2017. DOI: arXiv:1712.04621.
- [25] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, 2014.
- [26] A.G. Howard, M. Zhu, B. Chen, and J. Sun, "Mobile-nets: Efficient convolutional neural networks for mobile vision applications," 2017. DOI: arxiv.org/abs/1704.04861.
- [27] M. Sandler, A. Howard, M. L. Zhu, A. Zhmoginov, and L. C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, June 18–23, 2018, pp. 4510–4520.
- [28] X. Y. Zhang, X. Y. Zhou, M. X. Lin, and J. Sun, "Shuffle-net: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Salt Lake City, UT, USA, June 18–23, 2018, pp. 6848–6856.
- [29] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4700–4708.
- [30] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004. DOI: 10.1109/TIP.2003.819861.
- [31] D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," *arXiv preprint*, 2013. DOI: arXiv:1312.6114.
- [32] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *Int. J. Comput. Vis.*, vol. 115, pp. 211–252, 2015. DOI: <https://doi.org/10.1007/s11263-015-0816-y>.
- [33] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, no. 1, pp. 321–357.
- [34] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond Empirical Risk Minimization," *arXiv preprint*, 2017. DOI: arXiv:1710.09412.
- [35] L. Lusa, "SMOTE for high-dimensional class-imbalanced data," *BMC Bioinformatics*, vol. 14, no. 1, pp. 106, 2013.
- [36] I. Hiroshi, "Data augmentation by pairing samples for images classification," *ArXiv e-prints*, 2018.
- [37] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Tech Report, 2009. 5.
- [38] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey, "Adversarial Autoencoders," *arXiv preprint*, 2016. DOI: arXiv:1511.05644.
- [39] A. Radford, L. Metz, and S. Chintala, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks," *Computer Science*, 2015.

- [40] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," 2017. DOI: arXiv:1710.10196.
- [41] D. Lin, K. Fu, Y. Wang, G. Xu, and X. Sun, "MARTA GANs: Unsupervised representation learning for remote sensing image classification," *IEEE Geosci. Remote Sens. Lett.*, vol. 14, pp. 2092–2096, 2017.
- [42] X. Li, K. Wang, Y. Tian, L. Yan, and F. -Y. Wang, "The ParallelEye dataset: Constructing large-scale artificial scenes for traffic vision research," in *IEEE 20th International Conference on Intelligent Transportation Systems*, 2017.
- [43] Z. B. Wang, "Feasibility Study on Application of Simulated Images in Deep Learning," M.S. thesis, Xidian University, Xi'an, 2019..
- [44] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *European Conference on Computer Vision*, 2016, pp. 102–118.
- [45] K. Wang, C. Gou, Y. Duan, Y. Lin, X. Zheng, and F. -Y. Wang, "Generative adversarial networks: introduction and outlook," *IEEE/CAA Journal of Automatica Sinica*, vol. 4, no. 4, pp. 588–598, 2017.
- [46] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2414–2423.
- [47] J. L. Mitchell, "Real-Time Synthesis and Rendering of Ocean water," ATI, Marlboro, Tech. Rep. 445 (121-126P), 2005.
- [48] F. Luan, S. Paris, E. Shechtman, K. Bala, in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4990–4998.
- [49] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual Losses for Real-Time Style Transfer and Super-Resolution." 2016. DOI: arXiv:1603.08155.
- [50] H. Zhang and K. Dana, "Multi-style Generative network for Real-time Transfer" in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 1029–1037.
- [51] Y. Li, M. -Y. Liu, X. Li, M. -H. Yang, and J. Kautz, "A Closed-form Solution to Photorealistic Image Stylization," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 453–468.
- [52] Y. Di, Z. Jiang, and H. Zhang, "A Public Dataset for Fine-Grained Ship Classification in Optical Remote Sensing Images," *Remote Sens.*, vol. 13, pp. 747, 2021. DOI: <https://doi.org/10.3390/rs13040747>.
- [53] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [54] P. Goyal, P. Dollar, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch sgd: training imagenet in 1 hour," *arXiv preprint*, 2017. DOI: arXiv:1706.02677.
- [55] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *IEEE International Conference on Learning Representation (ICLR)*, 2015.



QI XIAO received the B.E. degree in communication engineering from Xiangtan University, Xiangtan, China, in 2019. He is currently pursuing the M.S. degree with the University of Chinese Academy of Sciences. His research interests include domain adaptation and neural style transfer.



BO LIU received the B.E. degree from the North University of China, Taiyuan, China, 2018. He is pursuing the Ph.D. degree at the University of Chinese Academy of Science, majoring in computer applications. His research

interests include domain adaptation and object detection based on simulation-scene images.



ZENGYI LI received the B.E. degree in automation from Henan University, Kaifeng, China, in 2021. He is currently pursuing the M.S. degree with the University of Chinese Academy of Sciences. His research interests include deep learning and simulation system.



WEI NI received the B.E. degree in automation from Northwest Normal University, Lanzhou, China, in 2011. Senior Engineer, National Space Science Center, Chinese Academy of Sciences.



ZHEN YANG received the B.S. and M.S. degrees in communication and electronic engineering from National Defense University, Changsha, China, in 1994 and 1997, respectively, and the Ph.D. degree in computer application from the University of Chinese Academy of Sciences, Beijing, China, in 2014. He is currently a Professor with the National

Space Science Center, Chinese Academy of Sciences. His research interests include complex system simulation, spatial task collaborative design and demonstration, spatial information services, and distributed space systems.



LIGANG LI received the Ph.D. in optical engineering from the Chinese Academy of Sciences in 2006 and is currently with the professor at the Chinese Academy of Sciences Mainly engaged in optical imaging simulation system technology, spatial information

processing and other aspects of research.