

Automatic Modulation Classification Using Involution Enabled Residual Networks

Hao Zhang, *Student Member, IEEE*, Lu Yuan, Guangyu Wu,
Fuhui Zhou, *Senior Member, IEEE*, and Qihui Wu, *Senior Member, IEEE*

Abstract

Automatic modulation classification (AMC) is of crucial importance for realizing wireless intelligence communications. Many deep learning based models especially convolution neural networks (CNNs) have been proposed for AMC. However, the computation cost is very high, which makes them inappropriate for beyond the fifth generation wireless communication networks that have stringent requirements on the classification accuracy and computing time. In order to tackle those challenges, a novel involution enabled AMC scheme is proposed by using the bottleneck structure of the residual networks. Involution is utilized instead of convolution to enhance the discrimination capability and expressiveness of the model by incorporating a self-attention mechanism. Simulation results demonstrate that our proposed scheme achieves superior classification performance and faster convergence speed comparing with other benchmark schemes.

Index Terms

Automatic modulation classification, deep learning, involution, residual networks.

I. INTRODUCTION

AUTOMATIC modulation classification (AMC) is a vital technique to identify the modulation formats under noise and interference [1]. AMC has been widely used in military and

H. Zhang, L. Yuan, F. Zhou, and Qihui Wu are with College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 211106 China. They are also with Key Laboratory of Dynamic Cognitive System of Electromagnetic Spectrum Space (Nanjing University of Aeronautics and Astronautics), and with Ministry of Industry and Information Technology, Nanjing, 211106, China (email: haozhangcn@nuaa.edu.cn, yuanlu@nuaa.edu.cn, zhoufuhui@ieee.org, wuqihui2014@sina.com)

Guangyu Wu is with Department of Computer Science and Technology, University of Science and Technology of China, Hefei 230026 China (email: gywu9908@163.com)

civilian applications, such as spectrum management, electronic warfare, interference identification [2], [3]. The AMC schemes mainly have two categories, namely, model-driven AMC and data-driven AMC. The model-driven schemes can be mainly include likelihood-based (LB) schemes and feature-based (FB) schemes [1]. The LB schemes can obtain the optimal solution from the Bayes' sense by calculating the likelihood function under the modulation hypothesis. The FB schemes aim to find better features and have lower complexity with robust performances.

Recent studies have demonstrated that data-driven AMC can achieve superior classification performances compared to the model-driven schemes by learning effective representations from data [4]–[6]. Machine learning (ML) based schemes such as support vector machine (SVM), K-nearest neighbor (KNN), and logistic regression can recognize the modulation formats by using a large number of data from the received signals. However, ML-based models still rely on the features generated by FB-based schemes. Thus, deep learning (DL) based algorithms [4]–[6] were proposed to extract features automatically from the original data, such as I/Q samples, cyclic spectrum, constellation diagrams.

Recently, many novel neural networks were designed to extract discriminative representations for AMC in order to improve the classification performance. A long short-term memory (LSTM) based AMC algorithm was proposed in [7] to identify modulation formats, but its recurrent structure results in high computational complexity. Inspired by the residual learning for image classification, a modified residual network (ResNet) was applied to extract features from the received I/Q symbols for AMC [4]. However, the classification performance is limit due to the over-fitting problem caused by a large number of network parameters [6]. Besides those networks, other promising techniques were also used for advancing the performance. For example, a graph convolutional network was investigated for AMC algorithm in [8]. These novel neural networks have improved the classification performance significantly. However, the performance is still limit due to the complex environments, and deep networks need a longer time to converge. In the beyond fifth generation (5G wireless communication networks), intelligent communications with high reliability and low latency are the main characteristics. However, the traditional AMC schemes based on convolution cannot satisfy the classification performance and low computing cost requirements of the beyond 5G wireless communication networks.

In this letter, we propose a residual network (ResNet) based AMC scheme using involution [9]. The contributions are summarized as follows. Firstly, involution is utilized instead of convolution to enhance the discrimination capability and expressiveness of the model by incorporating a self-

attention mechanism. Secondly, a novel network that inherits the advantages of residual learning and involution is designed to learn high-dimensional representations of different modulations from I/Q signals and construct the classifier. Thirdly, simulation results demonstrate the effectiveness of the proposed involution based ResNet for AMC. It is shown that our proposed scheme has a better classification performance and a faster convergence speed.

The remainder of this paper is organized as follows. Section II presents the problem of AMC. Section III presents our proposed AMC scheme. Simulation results are given in Section IV and Section V concludes this letter.

II. PROBLEM STATEMENT

According to the classical modulation classification problem statement [4], [5], the received signal can be given as

$$x(n) = s(n) + \omega(n), n = 1, 2, \dots, N, \quad (1)$$

where N is the total number of signal symbols, $s(n)$ denotes the n th (complex) symbol, and $\omega(n)$ is the additive white Gaussian noise (AWGN) with zero mean and variance δ_ω^2 .

The real and imaginary parts of the received signal from the In-phase and Quadrature (I/Q) parts are utilized, which can be expressed as a vector, given as

$$\begin{aligned} \mathbf{x} &= \mathbf{I}_x + \mathbf{Q}_x \\ &= \Re(\mathbf{x}) + j\Im(\mathbf{x}), \end{aligned} \quad (2)$$

where \mathbf{I}_x and \mathbf{Q}_x denote the real and imaginary parts of the received signal, respectively, and $j = \sqrt{-1}$. $\Re(\cdot)$ and $\Im(\cdot)$ represent the operators of the real and imaginary parts of the signal, respectively. \mathbf{x} can be specifically expressed as

$$\mathbf{x} = \begin{pmatrix} \Re[x(1), x(2), \dots, x(N)] \\ \Im[x(1), x(2), \dots, x(N)] \end{pmatrix}. \quad (3)$$

The average probability of correct classification (\Pr_{cc}) is utilized as the performance metric, which is defined as $\Pr_{cc} = \sum_{s=1}^{|\mathcal{S}|} \Pr(\hat{H} = H_s | H_s) \Pr(H_s)$, $H_s \in \mathcal{S}$, where \mathcal{S} denotes the candidate modulation formats. $\Pr(H_s)$ is the prior probability of modulation format H_s , which is equal for each format. $\Pr(\hat{H} = H_s | H_s) \Pr(H_s)$ represents the probability that the modulation format is correctly determined as H_s .

III. INVOLUTION ENABLED RESNET FOR AMC

In this part, we first introduce the standard convolution operation to make the definition of the proposed involution clearly. Then, involution which inverses characteristics of convolution in the spatial and channel domain with low complexity is presented. Finally, a residual network based on involution for AMC is proposed.

Let $\mathbf{X} \in \mathbb{R}^{H \times W \times C_i}$ denote the input feature map, where H , W , and C_i represent its height, width and input channels, respectively. A series of convolution filters C_o with the fixed kernel size of $K \times K$ are expressed as $\mathcal{F} \in \mathbb{R}^{C_o \times C_i \times K \times K}$, where each filter $\mathcal{F}_k \in \mathbb{R}^{C_i \times K \times K}$, $k = 1, 2, \dots, C_o$, contains C_i convolution kernels $\mathcal{F}_{k,c} \in \mathbb{R}^{K \times K}$, $c = 1, 2, \dots, C_i$. The filter executes multiply-add operations on the input feature map using a sliding window to generate the output feature map $\mathbf{Y} \in \mathbb{R}^{H \times W \times C_o}$, given as

$$\mathbf{Y}_{i,j,k} = \sum_{c=1}^{C_i} \sum_{(u,v) \in \Delta_K} \mathcal{F}_{k,c,u+[K/2],v+[K/2]} \mathbf{X}_{i+u,j+v,c}, \quad (4)$$

where $\Delta_K \in \mathbb{Z}^2$ denotes the set of offsets in the neighborhood considering convolution conducted on the center pixel, given as

$$\Delta_K = [-\lfloor K/2 \rfloor, \dots, \lfloor K/2 \rfloor] \times [-\lfloor K/2 \rfloor, \dots, \lfloor K/2 \rfloor]. \quad (5)$$

It is well known that there is inter-channel redundancy inside convolution filters, which results in the flexibility problem in convolution operation [9]. Compared to the standard convolution, the **involution** kernel $\mathcal{H} \in \mathbb{R}^{H \times W \times K \times K \times G}$ is designed to realize transforms with inverse characteristics in the spatial and channel domain. Specifically, an involution kernel $\mathcal{H}_{i,j,\cdot,\cdot,g} \in \mathbb{R}^{K \times K}$, $g = 1, 2, \dots, G$, is specially adapted for the pixel $\mathbf{X}_{i,j} \in \mathcal{R}^C$ located at the corresponding coordinate (i, j) , but shared over the channels. G denotes the number of groups and each group shares the same involution kernel. The output feature map of involution is derived by executing multiply-add operations on the input with the involution kernels, given as

$$\mathbf{Y}_{i,j,k} = \sum_{(u,v) \in \Delta_K} \mathcal{H}_{i,j,u+[K/2],v+[K/2],\lfloor kG/C \rfloor} \mathbf{X}_{i+u,j+v,k}. \quad (6)$$

Different from convolution kernels, the shape of involution kernels \mathcal{H} depends on the input feature map \mathbf{X} . The output kernels are aligned to the input by generating the involution kernels based on the original input tensor. Thus, the kernel generation function ϕ and mapping function

at each location (i, j) can be expressed as

$$\mathcal{H}_{i,j} = \phi(\mathbf{X}_{\Psi_{i,j}}), \quad (7)$$

where $\Psi_{i,j}$ indexes the set of pixels $\mathcal{H}_{i,j}$ is conditioned on. The kernel generation function $\phi: \mathbb{R}^C \rightarrow \mathbb{R}^{K \times K \times G}$ with $\Psi_{i,j} = \{(i, j)\}$ is given as

$$\mathcal{H}_{i,j} = \phi(\mathbf{X}_{i,j}) = \mathbf{W}_1 \sigma(\mathbf{W}_0 \mathbf{X}_{i,j}), \quad (8)$$

where $\mathbf{W}_0 \in \mathbb{R}^{\frac{C}{r} \times C}$ and $\mathbf{W}_1 \in \mathbb{R}^{(K \times K \times G) \times \frac{C}{r}}$ denote two linear transformations that collectively constitute a bottleneck structure. The channel reduction operation under a ratio r is used for efficient processing, and σ represents Batch Normalization and non-linear activation functions between two linear projections. The pseudo code of Algorithm 1 shows the computation flow of involution, which is visualized in Fig. 1.

Algorithm 1 Pseudo code of involution.

Input: Batch size B , height H , width W , channel C , group number G , kernel size K , stride s , and reduction ratio r ;

Output: The involution kernel out .

1: **Initialize the network operations;**

2: Define the operation o as average pooling with kernel s if $s > 1$, otherwise o is the identity mapping.

3: Define the operation $reduce$ as convolution with kernel of $C \times C // r \times 1$.

4: Define the operation $span$ as convolution with kernel of $C // r \times K * K * G \times 1$.

5: Define the operation $unfold$ as unfold with kernel of $K \times dilation \times padding \times s$.

6: **Forward pass;**

7: Calculate x_u using $x_u = unfold(x)$, and reshape it into $(B, G, C // G, K * K, H, W)$;

8: Generate the involution kernel using eq. 8 as $kernel = span(reduce(o(x)))$, and reshape it into $(B, G, K * K, H, W)$;

9: Execute Multiply-Add operation according to eq. 6 as $out = mul(kernel, x_u).sum(dim = 3)$, and reshape it into (B, C, H, W) .

The feature generation process in eq. 6 can be considered a generalized version of self-attention [10]. The self-attention pools $values$ V depending on the affinities obtained by computing similarity between the $query$ Q and key K , are formulated as

$$Y_{i,j,k} = \sum_{(p,q) \in \Omega} (QK^T)_{i,j,p,q, \lceil kH/C \rceil} V_{p,q,k}, \quad (9)$$

where Q , K , and V are linearly transformed from the input X , and H is the number of heads in multi-head self-attention [10]. The similarity is related to that both operators collect pixels

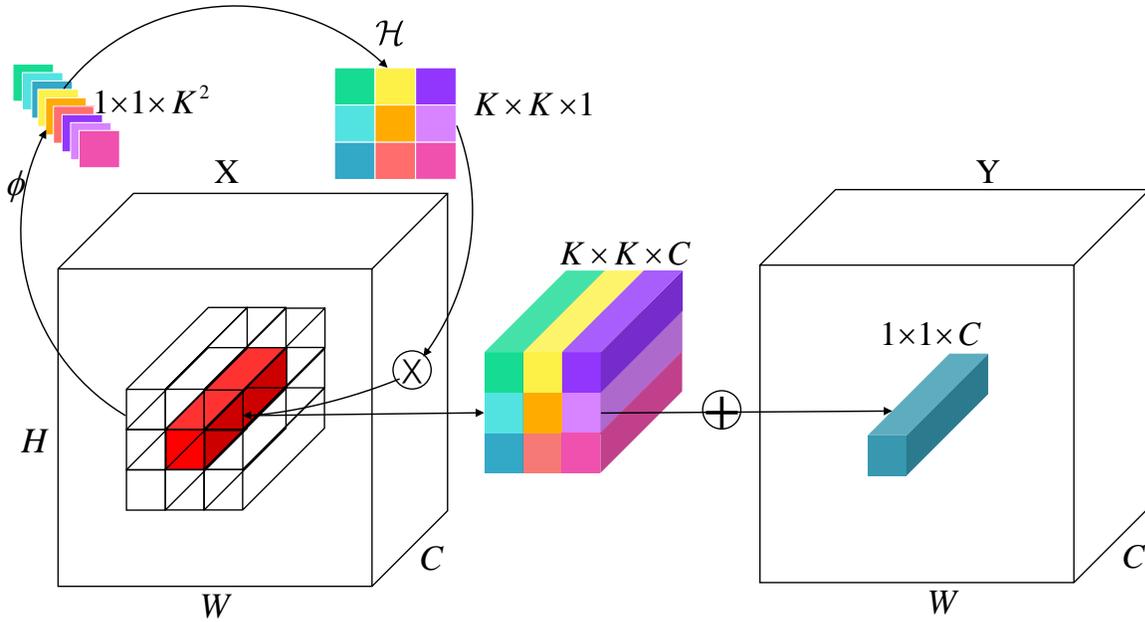


Fig. 1. Schematic illustration of the proposed involution. The involution kernel $\mathcal{H}_{i,j} \in \mathbb{R}^{K \times K \times 1}$ ($G = 1$ for ease of demonstration) is produced from the function ϕ based on a single pixel at (i, j) , followed by a channel-to-space rearrangement. The multiply-add operation of involution consists of two steps, where N indicates multiplication broadcast across C channels and L represents summation aggregated within the $K \times K$ spatial neighborhood.

in the neighborhood δ or a less bounded range Ω through a weighted sum. On one hand, the computation of involution can be viewed as a spatial attentive aggregation. On the other hand, the attention map (also called affinity or similarity matrix) QK^\top in the self-attention mechanism can be considered as a kind of involution kernel \mathcal{H} .

To build an entire network with involution for AMC, we design a novel lightweight network named Invo-ResNet by stacking bottleneck blocks and using involution kernels since the elegant architecture of ResNet makes it successful for many applications [11]. As shown in Fig. 2, the proposed Invo-ResNet consists of two modules, namely, the feature extraction module and classification module. The feature extraction module aims at extracting the underlying high-level representations from the input signals. Then, these representations are transformed into fixed-dimensional feature vectors by the global average pooling (GAP) layer. Modulation classification is subsequently conducted using these vectors in the fully connected (FC) layer of the classification module.

To balance the tradeoff between accuracy and efficiency, the feature extraction module consists of one convolutional layer and two bottlenecks. A convolutional layer with a kernel size of 3×1 is used to extract low-level information and execute the feature aggregation. Then, the

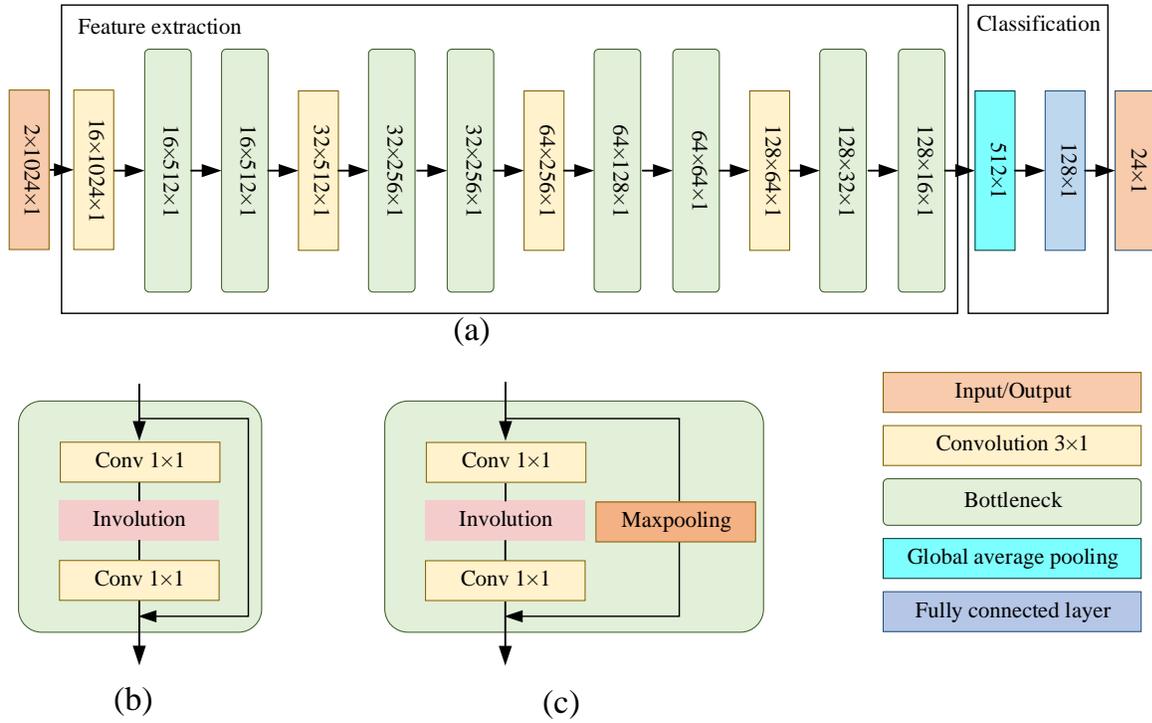


Fig. 2. Structure of the proposed Invo-ResNet, (a) network structure, (b) bottleneck, (c) bottleneck with downsampling using maxpooling.

bottleneck is utilized to learn high-level features by using involution with fewer parameters compared to convolution. As shown in Fig. 2 (b), the bottleneck is constructed by replacing the 3×3 convolutional layer in the original ResNet [4] with involution and retaining all the 1×1 convolution for channel projection and fusion [9]. Instead of using a constant channel number as in [4], a pyramid architecture with increasing channel numbers is utilized in the Invo-ResNet. Therefore, maxpooling is adopted to decrease the shape of the feature map in order to reduce the computation cost, as shown in Fig. 2(c).

For the last part, a GAP layer and an FC layer are served as the classification module. GAP layer aggregates information from the feature extraction module and enables the network to take samples with an arbitrary length. *Softmax* is used as the activation function in the last FC layer to normalize the output of each neuron, indicating the probability that the target signal belongs to the corresponding modulation format. Moreover, we adopt rectified linear units (ReLU) as the activation function in the convolutional layers to introduce nonlinearity and sparsity.

IV. SIMULATION

The public RadiomL 2018.01A that contains 24 kinds of modulations under an SNR range from -20dB to 30dB with a step of 2dB is adopted. There are over 2 million samples with 1024 points. 80% of these data are utilized for training and the rest 20% are used for testing. The models are trained by using SGD with an initial learning rate of 0.01, a weight decay of 5×10^{-4} , and a momentum of 0.9 for 50 epochs.

Fig. 3(a) shows the classification performance comparison of our proposed proposed Invo-ResNet with those achieved by three representative DL-based models for AMC including VGG [4], ResNet [4] and MCNet [5] on RadiomL 2018.01A dataset. It is evident that the proposed Invo-ResNet is superior to other traditional models, and it can provide 2% gains over our previous work MSNet, 2 dB gains over ResNet, 3 dB gains over MCNet and 4 dB gains over VGG. Moreover, Invo-ResNet can reach over 90% accuracy when the SNR is larger than 6dB, and it can achieve about 95% accuracy at 10dB, while the best performance of ResNet is 91.47% at 10dB. To further show the superiority of the proposed Invo-ResNet, Fig. 3(b) illustrates the training loss during the training process of the compared models. The proposed Invo-ResNet can achieve a lower loss in the training set compared to the other models and obtains a faster convergence speed than other models. To show the effectiveness of the involution for AMC, a fully convolutional network with the same structure is trained and tested under the same dataset. As shown in Fig. 3, compared to our proposed scheme with involution, the convolutional counterpart (Conv-ResNet) only reaches a best performance of 94.5% at high SNR condition, which is about 1% lower than involution based model, and it performs worse than our previous work MSNet. For the analysis of the computational complexity, we calculate the network parameters of the compared models, as shown in Fig. 3(a). Our proposed scheme achieves a better performance with less parameters than these conventional schemes. Moreover, compared to the convolution based scheme under the same structure, our proposed scheme with involution has about 40% reduced parameters, which can be implemented in resource limited devices.

To further demonstrate the effectiveness of the proposed Invo-ResNet, the comparison of the confusion matrices at 10dB of ResNet [4], MSNet [1] with those of our proposed scheme are shown in Fig. 4. It is seen that the proposed Invo-ResNet has less confusion compared to other traditional models. Specifically, only AM-DSB-WC receives the worst performance in our proposed scheme. On the contrary, two modulation formats perform worse in ResNet and

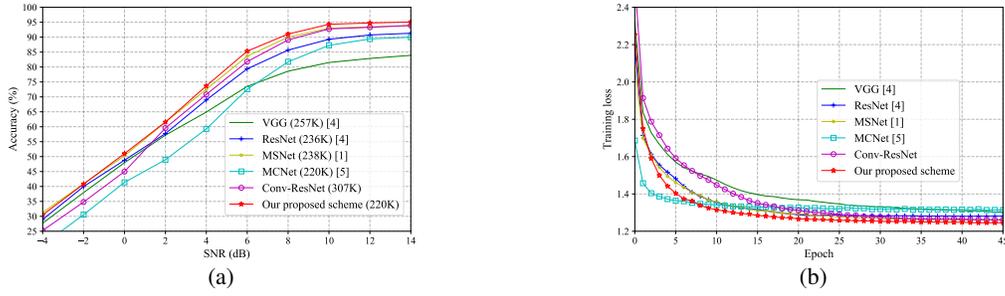


Fig. 3. Comparison of (a) the classification performance and (b) the training loss among VGG, ResNet, MSNet, MCNet, our proposed scheme with convolution layers and our proposed scheme with involution layers.

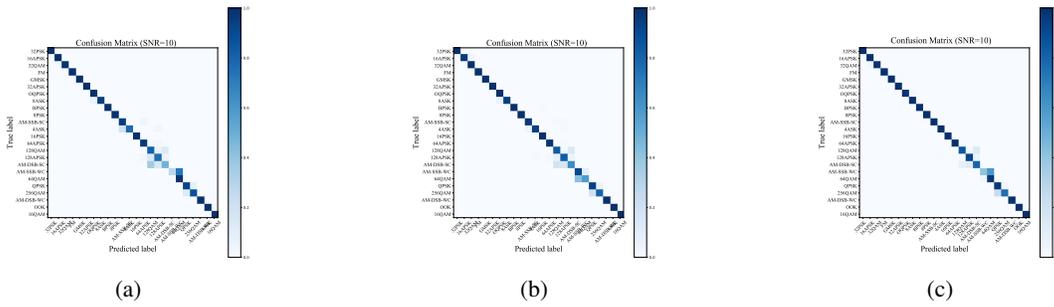


Fig. 4. Comparison of confusion matrix at 10dB, (a) ResNet [4], (b) MSNet [1], and (c) Our proposed scheme.

MSNet, which are AM-DSB-SC and 64QAM, and AM-DSB-SC and AM-SSB-WC as shown in Fig. 4(a) and Fig. 4(b), respectively.

V. CONCLUSION

A novel AMC scheme was proposed by designing a novel network using involution. In order to improve the classification accuracy and decrease the computation cost of convolution based AMC schemes, involution was utilized to enhance the discrimination capability and expressiveness of the model and reduce the training time by incorporating a self-attention mechanism. Simulation results demonstrated the superiority of our proposed scheme in terms of classification accuracy and the training time. In this case, our proposed scheme is more appropriate in beyond 5G wireless communication networks

REFERENCES

- [1] H. Zhang *et al.*, "A novel automatic modulation classification scheme based on multi-scale networks," *IEEE Trans. Cogn. Commun. Netw.*, 2021.

- [2] Q. Wu *et al.*, “Spatial-temporal opportunity detection for spectrum-heterogeneous cognitive radio networks: Two-dimensional sensing,” *IEEE Trans. Wireless Commun.*, vol. 12, no. 2, pp. 516–526, 2013.
- [3] Q. Wu *et al.*, “Cognitive internet of things: a new paradigm beyond connection,” *IEEE Internet Things J.*, vol. 1, no. 2, pp. 129–143, 2014.
- [4] T. J. O’Shea *et al.*, “Over-the-air deep learning based radio signal classification,” *IEEE J. Sel. Top. Sign. Proces.*, vol. 12, no. 1, pp. 168–179, 2018.
- [5] H.T. Thien *et al.*, “MCNet: An efficient CNN architecture for robust automatic modulation classification,” *IEEE Commun. Lett.*, vol. 24, no. 4, pp. 811–815, 2020.
- [6] S. Huang *et al.*, “Automatic modulation classification using contrastive fully convolutional network,” *IEEE Wireless Commun. Lett.*, vol. 8, no. 4, pp. 1044–1047, 2019.
- [7] S. Rajendran *et al.*, “Deep learning models for wireless signal classification with distributed low-cost spectrum sensors,” *IEEE Trans. Cogn. Commun. Netw.*, vol. 4, no. 3, pp. 433–445, 2018.
- [8] Y. Liu *et al.*, “Modulation recognition with graph convolutional network,” *IEEE Wireless Commun. Lett.*, vol. 9, no. 5, pp. 624–627, 2020.
- [9] D. Li *et al.*, “Involution: Inverting the inherence of convolution for visual recognition,” *arXiv preprint arXiv:2103.06255*, 2021.
- [10] A. Vaswani *et al.*, “Attention is all you need,” in *Adv. neural inf. proces. syst.*, 2017, pp. 5999–6009.
- [11] P. Qi *et al.*, “Automatic modulation classification based on deep residual networks with multimodal information,” *IEEE Trans. Cogn. Commun. Netw.*, 2020.