# Edge Learning Based Collaborative Automatic Modulation Classification for Hierarchical Cognitive Radio Networks

Peihao Dong, *Member, IEEE*, Chaowei He, Shen Gao, Fuhui Zhou, *Senior Member, IEEE*, and Qihui Wu, *Fellow, IEEE*

*Abstract*—In hierarchical cognitive radio networks, edge or cloud servers utilize the data collected by edge devices for modulation classification, which, however, is faced with problems of the computation load, transmission overhead, and data privacy. In this article, an edge learning (EL) based framework jointly mobilizing the edge device and the edge server for intelligent co-inference is proposed to realize the collaborative automatic modulation classification (C-AMC) between them. A spectrum semantic compression neural network is designed for the edge device to compress the collected raw data into a compact semantic embedding that is then sent to the edge server via the wireless channel. On the edge server side, a modulation classification neural network combining the bidirectional long-short term memory and attention structures is elaborated to determine the modulation type from the noisy semantic embedding. The C-AMC framework decently balances the computation resources of both sides while avoiding the high transmission overhead and data privacy leakage. Both the offline and online training procedures of the C-AMC framework are elaborated. The compression strategy of the C-AMC framework is also developed to further facilitate the deployment, especially for the resource-constrained edge device. Simulation results show the superiority of the EL-based C-AMC framework in terms of the classification accuracy, computational complexity, and the data compression rate as well as reveal useful insights paving the practical implementation.

*Index Terms*—Edge learning, cognitive radio, automatic modulation classification, spectrum semantics, model compression

## I. Introduction

**D**ESPITE the skyrocketing development, wireless networks are faced with unprecedented challenges in terms of the spectrum scarcity and transmission security. Automatic modulation classification (AMC) is anticipated to play a crucial role in constructing effective solutions since it endows the intended receiver or monitor with the ability of recognizing the signal modulation type with the minimal prior information [1]. AMC was born for the military purpose, and was then widespreadly applied to the civilian wireless networks in past decades, such as the spectrum surveillance, intelligent modem design, and malicious attack identification [2]–[4].

P. Dong is with the College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China, and also with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 211111, China (e-mail: phdong@nuaa.edu.cn).

C. He, F. Zhou, and Q. Wu are with the College of Electronic and Information Engineering, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China (e-mail: hcw0110@nuaa.edu.cn; zhoufuhui@ieee.org; wuqihui2014@sina.com).

S. Gao was with the National Mobile Communications Research Laboratory, Southeast University, Nanjing 211111, China (e-mail: gaoshen@seu.edu.cn).

The AMC algorithms are mainly categorized into two branches: likelihood-based (LB) and feature-based (FB) approaches. LB approaches calculate the likelihoods of all candidate modulation schemes with respect to the received signal, and select that with the maximal likelihood [2], [5]. Despite the optimality in the Bayesian sense, the LB approaches require the accurate knowledge of channel characteristics and high computational complexity hindering the practical implementation. This problem can be well treated by the FB approaches, which extract the received signal features, e.g., cyclic moments, wavelet-based features, high-order cumulants, in a low-complexity manner for classification decision while possessing the near-optimal performance [6]–[8]. However, due to the nature of relying on manually extracted features, traditional FB approaches may be overwhelmed by the increasingly complicated wireless environments as well as the soaring number of signal emitters, which invokes the more scalable tool for feature extraction.

The revival of deep learning (DL) invigorates the FB approaches and the marriage of them tends to become the mainstream solution for the modern AMC since the deep neural network (DNN) can act as a versatile feature extractor [9]. In [10], one of the works spearheading this direction, the DL-based AMC approach was proved feasible and superior under the real propagation environment inclusive of the effects of some key system parameters. In [11], convolutional neural network (CNN) was exploited for AMC and the data format of the received signals matching best with the proposed CNN structure was revealed. A modified generative adversarial network was designed to improve the classification accuracy by augmenting the training set composed of contour stellar images in [12]. In [13], an end-to-end CNN architecture using the two-step training was developed for AMC with the enhanced generalization ability. Considering different recognition difficulties of the modulation types, a DL structure consisted of two concatenate CNNs with respective recognition objects was developed in [14]. The direction of DL-based AMC became further prosperous as various subsequent studies mushroomed focusing on feature extraction enhancement [15]–[20], lightweight design [21]–[23], few-shot learning [24]–[26], distributed framework [4], [27], [28], adversarial defense [29], and so on.

As the wireless networks are increasingly expanded and complicated, the wide-area spectrum surveillance and management become pretty necessary considering both the spectrum scarcity and transmission security. As a result, the traditional

cognitive radio (CR) system will evolve into a hierarchical network consists of the edge device, edge server and/or cloud server [30]. In this architecture, the spectrum semantics need to be known by the edge/cloud server for the global spectrum management or decision-making while the edge device mainly takes on the task of sensing data. In other words, the data sensing and spectrum cognition are conducted at different locations in the network. The signal modulation type can be regarded as a kind of spectrum semantics extracted from the spectrum data while the existing DL models for AMC cannot be mechanically applied in the hierarchical CR network. In [31], a long-short term memory (LSTM) network was proposed to enable the AMC at distributed low-cost sensors. However, the manner of directly reporting the classification result brudens the resource-constrained edge device with the whole computation load and may leak the secret information. On the other hand, if the DNNs for classification are deployed at the edge server, the edge device has to deliver the sensed data via the wireless channel, incurring the high communication overhead and the risk of exposing the raw data privacy. Edge learning (EL), a paradigm enabling the more flexible deployment of DL models at the edge [32], provides the potential solution for this problem. That is, the DNN can be split into two parts in order to be respectively deployed at the edge device and edge server, aiming to balance the computation load, reduce the communication overhead, and improve the safety simultaneously [30]. In [33], the idea of DNN model splitting was adopted for AMC by partitioning a residual network, where the in-phase/quadrature (I/Q) samples are first processed by the model segment at the device and then the cut-layer representation is passed to the remaining model segment at the server for modulation classification. Although the cut-layer representation hides the private information, its high-dimension requires a considerable transmission overhead. In addition, the performance relies on a large set of I/Q samples, burdening the edge device in the data sensing and processing phases.

To address the mentioned-above problems, in this article, we propose an EL-based collaborative automatic modulation classification (C-AMC) framework consisted of a spectrum semantics compression neural network (SSCNet) and a modulation classification neural network (MCNet) deployed at the edge device and edge server, respectively, by treating the signal modulation type as a kind of spectrum semantics. The main novelty and contribution can be summarized as follows:

1) SSCNet and MCNet are well designed to achieve the goal of the C-AMC framework in terms of balancing the computation load, reducing the transmission overhead, and guaranteeing the information security simultaneously. Specifically, a quite lightweight structure accommodating the resource-constrained edge device is designed for SS-CNet to compress the collected raw data into a compact semantic embedding. Thanks to the low dimension of the semantic embedding and the more powerful computation capability at the edge server, MCNet incorporates the bidirectional long-short term memory (Bi-LSTM) and multi-head attention structures to achieve the high clas-

sification accuracy via the sufficient feature extraction.
2) The offline training procedure of the C-AMC framework is elaborated along with the insight on the generalization capability. A simple online training procedure is proposed by considering the update of SSCNet and MCNet over the air. The combination of offline and online training enables the C-AMC framework to adapt to a new scenario fast.
3) The compression strategy of the C-AMC framework is developed to further facilitate the deployment, especially for the resource-constrained edge device. The magnitude-based importance for a weight is analyzed, based on which the weight pruning procedure is elaborated. Then the post-training quantization is applied to further accelerate the model inference and reduce the model size. A layer-by-layer complexity analysis of the C-AMC framework is also provided.
4) Extensive simulation results are provided to show the superiority of the EL-based C-AMC framework over baseline schemes in terms of the classification accuracy, computational complexity, and data compression rate. Useful insights are extracted from the results to shed light on the practical implementation of the C-AMC framework.

By using the C-AMC framework, the following benefits can be gained for the hierarchical CR system: 1) The computation load for classification is allocated between the resource-constrained edge device and the edge server more properly so that the endurance of the edge device can be improved. 2) By transmitting the compact and intricate semantic embedding instead of the raw data or the classification result from the edge device to the edge server, the transmission overhead is reduced while the system safety is enhanced. 3) The framework is scalable and provides a paradigm for the hierarchical CR system to realize the cognition of various spectrum semantics besides the signal modulation type considered in this paper.

The rest of this paper is organized as follows. Section II introduces the basic signal model, based on which the EL-based C-AMC framework including SSCNet and MCNet is developed in Section III. Section IV elaborates the offline and online training procedures of the C-AMC framework. Section V developes the weight pruning and quantization based compression strategy for the C-AMC framework along with the complexity analysis. Simulation results are presented in Section VI, followed by Section VII giving concluding remarks.

## II. System Model

In this section, the hierarchical CR network model is elaborated by considering the modulation classification as the cognitive task.

As illustrated in Fig. 1, consider a hierarchical CR network consisted of an edge device and an edge server aiming to recognize the modulation type of a source signal in the spec-
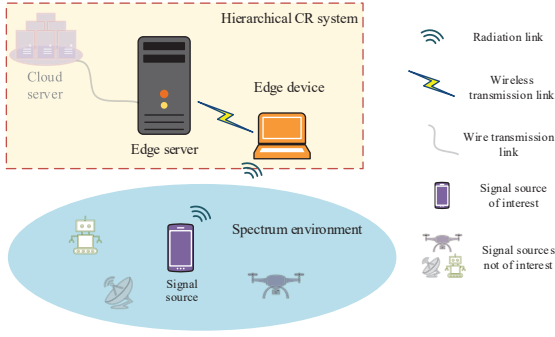
Fig. 1. Hierarchical CR network model.

trum environment.[1] The cognitive process includes two phases, data sensing at the edge device and information transmission from the edge device to the edge server, along with the corresponding signal processing procedures.

*Phase I (Data Sensing at the Edge Device):* The edge device keeps sounding the source signal $s(t)$ within $L$ sampling time instants and the received signal is given by

$$x[l] = h[l]e^{-j2\pi(\nu l T_s + \vartheta)}s[l] + z[l], \quad l = 1, \ldots, L, \quad (1)$$

where $h[l]$, $s[l]$, and $z[l]$ respectively denote the channel gain, modulated source signal, and additive white Gaussian noise (AWGN) at the $l$th sampling time instant, and $T_s$, $\nu$, and $\vartheta$ represent the sampling period, frequency shift, and phase shift, respectively. The modulation type of the source signal, $C(s)$, is taken from the candidate set $\mathcal{M}$ including $M$ elements. Denote $\mathbf{x} = [x[1], \ldots, x[L]]$ as the vector form of the received signal. Then $\mathbf{x}$ is processed at the edge device to yield an $N$-dimensional vector expressed as

$$\mathbf{x}_s = \phi(\mathbf{x}), \quad (2)$$

where $\phi(\cdot)$ represents the general mapping function of the signal processing. Specifically, $\mathbf{x}_s$ will give the classification result if $\phi(\cdot)$ is the modulation classification mapping while $\mathbf{x}_s = \mathbf{x}$ holds if $\phi(\cdot)$ is the identity mapping.

*Phase II (Information Transmission from the Edge Device to the Edge Server):* In this phase, the edge device transmits $\mathbf{x}_s$ to the edge server via a certain air interface. The received signal after equalization at the edge server is expressed as

$$y_i = x_{s,i} + w_i, \quad i = 1, \ldots, N, \quad (3)$$

where $x_{s,i}$ is the $i$th element of $\mathbf{x}_s$ and $w_i$ denotes the corresponding effective noise compounding the interference and AWGN. Then the edge server conduct the processing, $\varphi(\cdot)$, on $\mathbf{y} = [y_1, \ldots, y_N]$ to recognize the modulation type of the source signal. Without loss of generality, consider the noise-free case, yielding

$$\hat{\mathbf{p}} = \varphi(\mathbf{y}) = \varphi(\phi(\mathbf{x})), \quad (4)$$

[1]Since the limited capacity of the wireless link from the edge device to the edge server hinders the accurate global spectrum cognition for the hierarchical CR network, we focus on addressing this crux along with the computation load and transmission security problems. The transmission from the edge server to the cloud server can be carried out via the wire link and thus is not considered.

where $\hat{\mathbf{p}} \in \mathbb{R}^M$ is a probability vector indicating the most likely modulation type of $s(t)$ according to the index of its largest entry. From (4), $\phi(\cdot)$ and $\varphi(\cdot)$ collaborate to finish the modulation classification task in a complementary manner. For the specific cases mentioned above, $\varphi(\cdot)$ can be the identity mapping if $\phi(\cdot)$ is the modulation classification mapping and vice versa, which, however, aggravates the computation load for the edge device or incurs the high transmission overhead, in addition to the risk of exposing the secure information. Therefore, we aim to address these problems by elaborating $\phi(\cdot)$ and $\varphi(\cdot)$ in the following sections.

## III. EL-BASED C-AMC FRAMEWORK

In this section, the EL-empowered C-AMC framework is constructed. The framework is overviewed first, followed by detailing its two component DL models, i.e., SSCNet and MCNet.

### A. Overview of Framework

Inspired by EL, the C-AMC framework is able to balance the computation load of modulation classification between the edge device and edge server by respectively deploying SSCNet and MCNet for them, as illustrated in Fig. 2. SSCNet compresses the sensed data into the low-dimensional spectrum semantic embedding for transmission by extracting features related to the modulation type therein. At the other end of the wireless channel, MCNet utilizes the noisy semantic embedding to predict the modulation type of the source signal. By doing this, a part of computation task of modulation classification can be offloaded from the edge device with the limited computing resource to the edge server. Besides, the low-dimensional semantic embedding reduces the transmission overhead and is difficult to decode for the potential eavesdroppers in the wireless channel.

Before being processed by SSCNet at the edge device, the sensed data $\mathbf{x}$ in the I/Q form is converted to the amplitude/phase (A/P) form as

$$\mathbf{X}_{AP} = \mathcal{P}(\mathbf{x}) = \left[ \begin{array}{c} |x[1]|, \ldots, |x[L]| \\ \theta(x[1]), \ldots, \theta(x[L]) \end{array} \right]^T, \quad (5)$$

where $|x|$ and $\theta(x) = \arctan\frac{\text{Im}(x)}{\text{Re}(x)}$ respectively denote the amplitude and phase of $x$. It is noted that this preprocessing step helps SSCNet better extract the useful features. Then $\mathbf{x}_{AP}$ is processed by SSCNet to generate the spectrum semantic embedding with the compression rate $r = \frac{2L}{N}$, that is,

$$\mathbf{x}_s = f(\mathbf{X}_{AP}; \mathbf{\Theta}), \quad (6)$$

where $f(\cdot)$ denotes the mapping function of SSCNet parameterized by the weight set $\mathbf{\Theta}$. So the general mapping function $\phi(\cdot)$ can be represented as the composite of $\mathcal{P}(\cdot)$ and $f(\cdot)$, i.e., $\phi = f \circ \mathcal{P}$.

At the edge server, the noisy semantic embedding $\mathbf{y} = \mathbf{x}_s + \mathbf{w}$ with $\mathbf{w} = [w_1, \ldots, w_N]$ is processed by MCNet to predict the modulation type as

$$\hat{\mathbf{p}} = g(\mathbf{y}; \mathbf{\Phi}) = g(f(\mathcal{P}(\mathbf{x}); \mathbf{\Theta}) + \mathbf{w}; \mathbf{\Phi}), \quad (7)$$
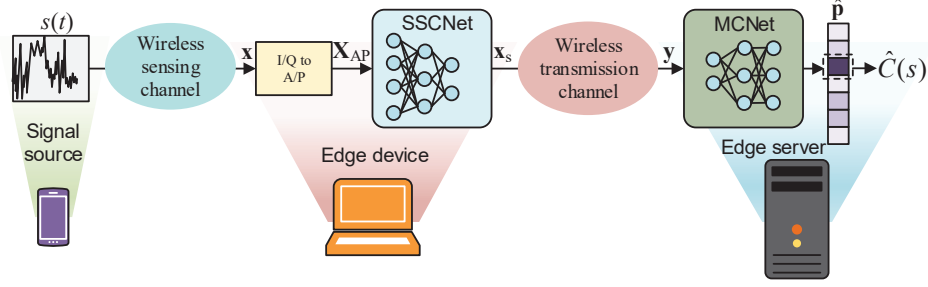
Fig. 2. EL-based C-AMC framework.

where $g(\cdot)$ denotes the mapping function of MCNet parameterized by the weight set $\mathbf{\Phi}$. Then the general mapping function $\varphi(\cdot)$ is instantiated as $g(\cdot)$, i.e., $\varphi = g$. Since $\mathcal{P}(\cdot)$ is a fixed operation, the design of $\phi(\cdot)$ and $\varphi(\cdot)$ becomes the design of SSCNet and MCNet represented by $f(\cdot)$ and $g(\cdot)$, respectively.

### B. SSCNet

In the C-AMC framework, the role of SSCNet is to compress the sensed data with a high dimension into the compact semantic embedding using a lightweight architecture to accommodate the resource-limited edge device. To sufficiently compress the sensed data, the temporal correlation therein can be exploited by the powerful LSTM structure [31]. Generally, more than one LSTM layer is needed to fully extract the temporal correlation, which, however, causes the high computation load for the edge device. To address this problem, we use the one-dimensional convolution (Conv1D) as a lightweight alternative to tentatively handle the temporal correlation instead of stacking multiple LSTM layers.

Fig. 3(a) shows the detailed architecture of SSCNet. A Conv1D layer using $64$ kernels with length $8$ and rectified linear unit (ReLU) activation function is utilized first to filter the input $\mathbf{x}_{\mathrm{AP}} \in \mathbb{R}^{L \times 2}$, in order to extract the local temporal correlation therein, after which a dropout layer with the dropout rate of $0.5$ is added to prevent overfitting and to improve the robustness of neurons. Then the $L \times 64$ feature map is processed by another Conv1D layer using $32$ kernels with length $8$ and ReLU to further distill the desired high-order spectrum feature, after which the size of the feature map is compressed from the $L \times 64$ to $L \times 32$. However, since the sequence length of the sampling data, $L$, usually ranges from several hundred to over a thousand, the feature map with size $L \times 32$ still consumes too much transmission overhead and thus needs to be further compressed. To this end, the column-sum pooling operation is applied to sum the feature map in a column-wise manner so that the feature map size can be dramatically reduced to $1 \times 32$.[2] Then a batch normalization (BN) layer is appended. Dense1 layer transforms the feature map to an $N$-dimensional vector. The scaled exponential linear unit (SELU) activation function is applied to Dense1 layer to

[2]As the $L \times 32$ feature map before pooling actually represents 32 $L$-dimensional feature vectors, the column-sum pooling operation can reserve the feature diversity and indeed achieves the better performance compared with other pooling ways according to simulation trials.
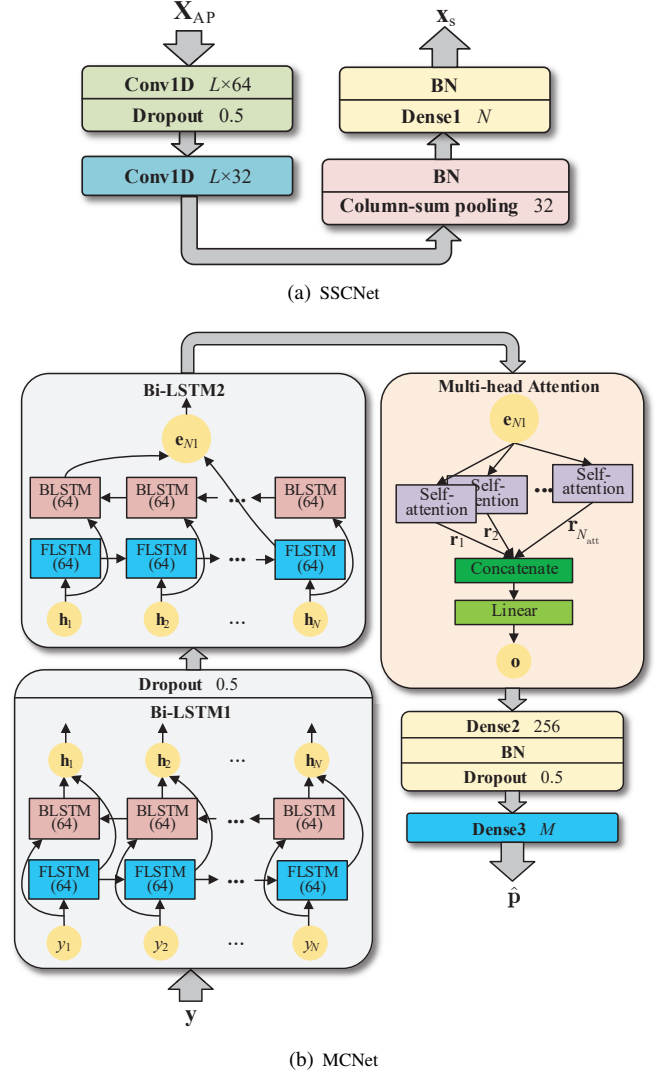


(a) SSCNet



(b) MCNet

Fig. 3. Architecture of SSCNet and MCNet.

ensure the neuron activity regardless of the sign of the input, that is,

$$\mathrm{SELU}(x) = \delta \begin{cases} x, & x > 0 \\ \alpha\left(e^x - 1\right), & x \leq 0 \end{cases} \tag{8}$$

with $\alpha > 0$ and the scaling factor $\delta > 1$. After BN, the $N$-dimensional real-valued spectrum semantic embedding, $\mathbf{x}_{\mathrm{s}}$, with $N \ll 2L$ is obtained for low-overhead transmission.

## C. MCNet

MCNet should be designed to coordinate with SSCNet so that they can be integrated across the wireless channel to output the predicted modulation type at the edge server. Thanks to the much lower dimension of the noisy semantic embedding, $\mathbf{y}$, the powerful LSTM and attention layers can be applied in MCNet to fully extract the desired feature from $\mathbf{y}$ as well as to remove the noise imposed by the wireless channel by exploiting the global correlation therein while avoiding the high computational complexity.

Fig. 3(b) shows the detailed architecture of MCNet. Considering the mutual dependence of the elements in $\mathbf{y}$, the Bi-LSTM structure is applied to extract more useful information. Bi-LSTM1 layer includes $N$ forward LSTM units and $N$ backward LSTM units, and the processing procedure can be expressed as

$$\mathbf{h}_{\mathrm{F},i} = \mathrm{FLSTM}(y_i; \mathbf{h}_{\mathrm{F},i-1}, \mathbf{c}_{\mathrm{F},i-1}, \boldsymbol{\Xi}_{\mathrm{F}}), \tag{9}$$

$$\mathbf{h}_{\mathrm{B},i} = \mathrm{BLSTM}(y_i; \mathbf{h}_{\mathrm{B},i+1}, \mathbf{c}_{\mathrm{B},i+1}, \boldsymbol{\Xi}_{\mathrm{B}}), \tag{10}$$

$$\mathbf{h}_i = [\mathbf{h}_{\mathrm{F},i}, \mathbf{h}_{\mathrm{B},i}], \quad i = 1, \ldots, N, \tag{11}$$

$$\mathbf{H} = [\mathbf{h}_1^T, \ldots, \mathbf{h}_N^T]^T, \tag{12}$$

where $\mathbf{h}_{\mathrm{F},i-1}$ and $\mathbf{c}_{\mathrm{F},i-1}$ denote vectors passed from the previous forward LSTM unit with the parameter set $\boldsymbol{\Xi}_{\mathrm{F}}$ including all the weight matrices and bias vectors of the unit, and $\mathbf{h}_{\mathrm{B},i+1}$ and $\mathbf{c}_{\mathrm{B},i+1}$ denote vectors passed from the previous backward LSTM unit with the parameter set $\boldsymbol{\Xi}_{\mathrm{B}}$ including all the weight matrices and bias vectors of the unit. The $i$th pair of forward and backward LSTM units respectively transform $y_i$ into the vectors $\mathbf{h}_{\mathrm{F},i} \in \mathbb{R}^{64}$ and $\mathbf{h}_{\mathrm{B},i} \in \mathbb{R}^{64}$, which are then concatenated to yield $\mathbf{h}_i$ as the output of this pair. Stacking $\mathbf{h}_1, \ldots, \mathbf{h}_N$ gives the output of of Bi-LSTM1 layer, i.e., $\mathbf{H} \in \mathbb{R}^{N \times 128}$. After Bi-LSTM1 layer, a dropout layer with the dropout rate of $0.5$ is appended. Bi-LSTM2 layer also consists of $N$ forward LSTM units and $N$ backward LSTM units while it only passes the output of the final units of the two directions to the following layer for processing, that is,

$$\mathbf{e}_{\mathrm{F},N} = \mathrm{FLSTM}(\mathbf{h}_N; \mathbf{e}_{\mathrm{F},N-1}, \mathbf{d}_{\mathrm{F},N-1}, \boldsymbol{\Pi}_{\mathrm{F}}), \tag{13}$$

$$\mathbf{e}_{\mathrm{B},1} = \mathrm{BLSTM}(\mathbf{h}_1; \mathbf{e}_{\mathrm{B},2}, \mathbf{d}_{\mathrm{B},2}, \boldsymbol{\Pi}_{\mathrm{B}}), \tag{14}$$

$$\mathbf{e}_{N1} = [\mathbf{e}_{\mathrm{F},N}, \mathbf{e}_{\mathrm{B},1}], \tag{15}$$

where the notations are similar to Bi-LSTM1 layer. Since the units of two directions in Bi-LSTM2 layer respectively transform $\mathbf{h}_i \in \mathbb{R}^{128}$ to $\mathbf{e}_{\mathrm{F},i} \in \mathbb{R}^{64}$ and $\mathbf{e}_{\mathrm{B},i} \in \mathbb{R}^{64}$, the size of the feature map output by this layer, $\mathbf{e}_{N1}$, is $1 \times 128$. Then $\mathbf{e}_{N1}$ is processed by an $N_{\mathrm{A}}$-head self-attention layer. The output of the $n$th self-attention head is given by

$$\mathbf{r}_n = \mathrm{S\text{-}Att}(\mathbf{e}_{N1}\mathbf{W}_{\mathrm{Q}}^{(n)}, \mathbf{e}_{N1}\mathbf{W}_{\mathrm{K}}^{(n)}, \mathbf{e}_{N1}\mathbf{W}_{\mathrm{V}}^{(n)}), \tag{16}$$

where $\mathbf{W}_{\mathrm{Q}}^{(n)}, \mathbf{W}_{\mathrm{K}}^{(n)}, \mathbf{W}_{\mathrm{V}}^{(n)} \in \mathbb{R}^{128 \times d_{\mathrm{A}}}$ denote the weight matrices used to generate the query, key, and value matrices
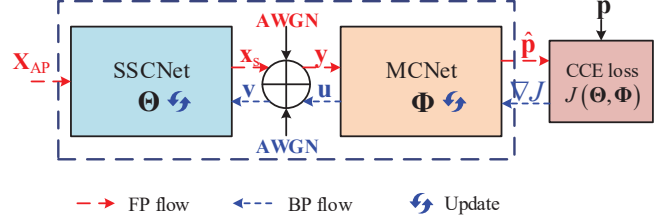


Fig. 4. Offline training procedure of C-AMC framework.

with $n = 1, \ldots, N_{\mathrm{A}}$. The outputs of $N_{\mathrm{A}}$ heads are vertically concatenated as $\mathbf{R} = [\mathbf{r}_1^T, \ldots, \mathbf{r}_{N_{\mathrm{A}}}^T]^T$. Then $\mathbf{R}$ is transformed linearly to yield $\mathbf{o} \in \mathbb{R}^{128}$ as the output of the multi-head attention layer. The feature map $\mathbf{o}$ is processed by Dense2 layer with the dimension increased to $256$, followed by the BN and dropout operation. Finally, Dense3 layer outputs the probability vector $\hat{\mathbf{p}}$ by using Softmax activation function. SELU is applied for Bi-LSTM1, Bi-LSTM2, and Dense2 layers.

## IV. TRAINING PROCEDURES OF C-AMC FRAMEWORK

In this section, the offline and online training procedures of the C-AMC framework are developed for different application scenarios. Since SSCNet and MCNet are integrated to fulfill the modulation classification task, they are trained jointly in an end-to-end manner for both of two modes.

### A. Offline Training

The offline training means to train SSCNet and MCNet centralizedly in the simulation environment and then deploy them separately at the edge device and server. Resorting to the low-cost offline computing resource, this mode can provide a decent base model provided that a certain amount of representative training data are available.

Fig. 4 illustrates the offline training procedure, where SSCNet and MCNet are connected by the AWGN. Note that the Gaussian distributed noise can be safely used to model the impact of the effective noise $\mathbf{w}$ imposed by the wireless transmission channel since it corresponds to the lower bound of the transmission capacity. Each training sample is paired by the input $\mathbf{X}_{\mathrm{AP}}$ and the label $\mathbf{p}$, where $\mathbf{p}$ is a one-hot vector indicating the true modulation type of the source signal. Each time of weight update includes two phases detailed as follows.

**1) Forward Propagation (FP):** In the FP flow, SSCNet receives the input $\mathbf{X}_{\mathrm{AP}}$ and outputs the semantic embedding $\mathbf{x}_{\mathrm{s}}$. By adding the generated AWGN $\bar{\mathbf{w}}$ on $\mathbf{x}_{\mathrm{s}}$, the noisy semantic embedding can be written as

$$\mathbf{y} = \mathbf{x}_{\mathrm{s}} + \bar{\mathbf{w}}, \tag{17}$$

which acts as the input of MCNet. Then MCNet outputs the predicted probability vector $\hat{\mathbf{p}}$ to approximate $\mathbf{p}$. The categorical cross-entropy (CCE) loss function is used to measure the discrepancy between the prediction and label vectors on a batch of training samples, that is,

$$J(\boldsymbol{\Theta}, \boldsymbol{\Phi}) = -\frac{1}{N_{\mathrm{bat}}} \sum_{n=1}^{N_{\mathrm{tr}}} \sum_{m=1}^{M} p_m^{(n)} \log \hat{p}_m^{(n)}, \tag{18}$$
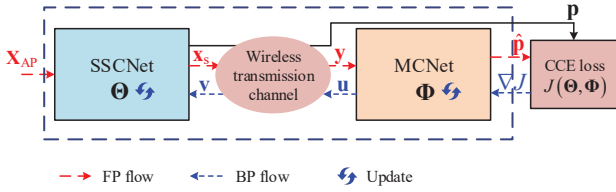
Fig. 5. Online training procedure of C-AMC framework.

where $N_{\text{bat}}$ denotes the batch size, $p_m^{(n)}$ and $\hat{p}_m^{(n)}$ respectively represent the $m$th element of $\mathbf{p}$ and $\hat{\mathbf{p}}$ corresponding to the $n$th sample.

**2) Back Propagation (BP) & Update:** The error gradient $\nabla J$ is back propagated to SSCNet and MCNet to update $\mathbf{\Theta}$ and $\mathbf{\Phi}$. In the BP flow, $\mathbf{\Phi}$ is directly updated as

$$\mathbf{\Phi} \leftarrow \mathbf{\Phi} - \eta \nabla_{\mathbf{\Phi}} J, \tag{19}$$

where $\eta$ denotes learning rate for the weight update. Based on $\nabla_{\mathbf{\Phi}} J$, the error gradient with respect to $\mathbf{x}_{\text{s}}$

$$\mathbf{u} = \nabla_{\mathbf{x}_{\text{s}}} J = \left[ \frac{\partial J}{\partial x_{\text{s},1}}, \dots, \frac{\partial J}{\partial x_{\text{s},N}} \right] \tag{20}$$

is calculated. Adding the AWGN $\bar{\mathbf{w}}'$ on $\mathbf{u}$ yields the noisy gradient expressed as

$$\mathbf{v} = \mathbf{u} + \bar{\mathbf{w}}', \tag{21}$$

which is used to calculate the approximated error gradient with respect to the weights in $\mathbf{\Theta}$ as per the chain rule. For an arbitrary weight $\theta$ in $\mathbf{\Theta}$, its approximated error gradient is given by

$$\hat{\nabla}_{\theta} J = \sum_{n=1}^{N} v_n \frac{\partial x_{\text{s},n}}{\partial \theta} = \sum_{n=1}^{N} \left( \frac{\partial J}{\partial x_{\text{s},n}} \frac{\partial x_{\text{s},n}}{\partial \theta} + \bar{w}_n' \frac{\partial x_{\text{s},n}}{\partial \theta} \right)$$
$$= \nabla_{\theta} J + \sum_{n=1}^{N} \bar{w}_n' \frac{\partial x_{\text{s},n}}{\partial \theta}, \tag{22}$$

where $v_n$ and $\bar{w}_n'$ denote the $n$th elements of $\mathbf{v}$ and $\bar{\mathbf{w}}'$, respectively. Denoting $\hat{\nabla}_{\mathbf{\Theta}} J$ as the set including all the approximated error gradients, $\mathbf{\Theta}$ can be updated as

$$\mathbf{\Theta} \leftarrow \mathbf{\Theta} - \eta \hat{\nabla}_{\mathbf{\Theta}} J. \tag{23}$$

The addition of $\bar{\mathbf{w}}$ and $\bar{\mathbf{w}}'$ introduces the perturbation in the training stage and thus improves the generalization capability of the C-AMC framework.

*B. Online Training*

In some practical scenarios, there are not enough training data available prior to model deployment or the sensed data used for modulation classification exhibit the non-stationary property, in which cases the online training needs to be invoked.

Fig. 5 illustrates the online training procedure. Different from the offline training mode inserting the generated AWGN between SSCNet and MCNet, SSCNet and MCNet are separated by the wireless transmission channel in the online

training mode as they have been deployed at the edge device and server, respectively. For the online mode, another main difference from the offline mode is that the training data with form of $\langle \mathbf{X}_{\text{AP}}, \mathbf{p} \rangle$ are collected by the edge device sequentially. In the following, the FP and BP & update phases are briefly described by indicating the detailed differences therein.

**1) FP:** In the online mode, the noisy semantic embedding acting as the input of MCNet experiences the interface processing and becomes

$$\mathbf{y} = \mathbf{x}_{\text{s}} + \mathbf{w}, \tag{24}$$

where $\mathbf{w}$ is the the effective noise compounding AWGN and the interference after the channel equalization. In addition to $\mathbf{x}_{\text{s}}$, the edge device also needs to transmit the label $\mathbf{p}$ to the edge server for loss calculation. After accumulating a batch of training samples, the edge server can calculate the CCE loss as per (18).

**2) BP & Update:** Most steps of this phase are same as the counterparts in the offline mode while the difference lies in the noisy gradient. Similar to the form of $\mathbf{y}$, the noisy gradient undergoes the inverse transmission processing and becomes the superposition of $\mathbf{u}$ and the effective noise $\mathbf{w}'$, that is,

$$\mathbf{v} = \mathbf{u} + \mathbf{w}'. \tag{25}$$

Then SSCNet and MCNet can be jointly trained by using the sequentially arrived training samples even though they are deployed at two ends of the wireless channel.

It is noted that the offline and online training modes also can be adopted together in practical application. Specifically, the offline training provides the decent base models of SSCNet and MCNet for deployment. After that the online training is executed to adapt the models by using only a few training samples so that the C-AMC framework is able to cope with the dynamics of the spectrum environment and system state.

## V. C-AMC MODEL COMPRESSION

In this section, the model compression strategy of the C-AMC framework is studied by resorting to weight pruning and quantization, in order to lower the requirements of computation and storage, especially for the resource-constrained edge device. The computational complexity of the C-AMC models is finally analyzed to shed light on the effect of model compression.

*A. Weight Pruning*

Weight pruning can reduce the quantity of parameters in a DL model significantly by cutting off neuronal connections with the trivial contribution to the prediction accuracy and then fine-tuning the model, incurring the limited performance loss. The weight magnitude is a metric to measure the weight contribution and those weights with relatively small absolute values will be removed. Specifically, consider a weight $w$ with the small absolute value connecting neuron $i$ and $j$ in a DL model. The output of neuron $i$ and the input of neuron $j$ are denoted by $\alpha$ and $\beta$, respectively. Then we have

$$\beta = w\alpha + B, \tag{26}$$

where $B$ denotes the part of input from other weights for neuron $j$. From the FP flow perspective, small $|w|$ weakens the impact of $\alpha$ on $\beta$, leading to the limited contribution of $w$. From the BP flow perspective, to update the weight $v$ of a connection contributing to the input of neuron $i$, its gradient is calculated as

$$\frac{\partial \mathcal{L}}{\partial v} = \frac{\partial \mathcal{L}}{\partial \beta}\frac{\partial \beta}{\partial \alpha}\frac{\partial \alpha}{\partial v} = \frac{\partial \mathcal{L}}{\partial \beta}w\frac{\partial \alpha}{\partial v}, \qquad (27)$$

where $\mathcal{L}$ denotes the loss function. It can be seen that $w$ with the small absolute value also weakens the gradients of weights related to $w$ in previous layers and thus is less important to the model training. Therefore, the weights with small absolute values contribute trivially to both the model training and inference and the magnitude-based criterion is applied to prune the C-AMC framework.

For the magnitude-based pruning, the pruning threshold is usually determined as per the pruning ratio. Specifically, denote $\mathcal{W}_i$ and $\rho_i$ as the weight set and pruning ratio of the $i$th layer in the C-AMC framework with $i = 1, \ldots, L_{\text{SSC}} + L_{\text{MC}}$, where $L_{\text{SSC}}$ and $L_{\text{MC}}$ represent the numbers of layers in SSCNet and MCNet, respectively. Sort all $N_{\mathcal{W}_i}$ weights in $\mathcal{W}_i$ in ascending order as per the absolute values to yield the ordered vector $\zeta_i$, whose $\lfloor\rho_i N_{\mathcal{W}_i}\rfloor$th element, $\zeta_{i,\lfloor\rho_i N_{\mathcal{W}_i}\rfloor}$, is selected as the pruning threshold. Then the $j$th weight in $\mathcal{W}_i$ is filtered as

$$\mathcal{W}_i(j) = \begin{cases} \mathcal{W}_i(j), & |\mathcal{W}_i(j)| \geq \zeta_{i,\lfloor\rho_i N_{\mathcal{W}_i}\rfloor} \\ 0, & \text{otherwise} \end{cases}, \qquad (28)$$

for $j = 1, \ldots, N_{\mathcal{W}_i}$. After pruning all layers of interest in the C-AMC framework, both SSCNet and MCNet are fine-tuned based on the training set to compensate the performance loss.

### B. Weight Quantization

Weight quantization can further reduce the model size and accelerate the model inference, which is a good partner of weight pruning to facilitate the model deployment, especially for the edge device with limited resources.

To reduce the complexity, the post-training quantization is adopted since it is a push-button strategy decoupling the model training and quantization. Consider the $b$-bit uniform affine quantization mapping each reserved weight after pruning to an unsigned integer from the set $\{0, \ldots, 2^b - 1\}$ [34], which is expressed as

$$\mathbb{Q}\left(\mathcal{W}_i(j)\right) = \text{clamp}\left(\text{round}\left(\frac{\mathcal{W}_i(j)}{S}\right) + Z; 0, 2^b - 1\right) \quad (29)$$

where round$(\cdot)$ represents the round-to-nearest operation with clamp$(\cdot; \cdot, \cdot)$ given by

$$\text{clamp}(x; \mu_1, \mu_2) = \begin{cases} \mu_1, & x < \mu_1, \\ x, & \mu_1 \leq x \leq \mu_2, \\ \mu_2, & x > \mu_2, \end{cases}. \qquad (30)$$

In (29), $S$ and $Z$ denote a scaling factor and an integer zero point quantized from the real zero, which are respectively calculated as

$$S = \frac{\max(\mathcal{W}_i) - \min(\mathcal{W}_i)}{2^b - 1}, \qquad (31)$$

### TABLE I
### COMPUTATIONAL COMPLEXITY OF SSCNET AND MCNET

| | Layer type | Layer index | Complexity |
|---|---|---|---|
| SSCNet | Conv1D-1 | $l = 1$ | $\mathcal{O}(2\rho_l D_{l,1} D_{l,2} N_{\mathcal{W}_l})$ |
| | Conv1D-2 | $l = 2$ | $\mathcal{O}(2\rho_l D_{l,1} D_{l,2} N_{\mathcal{W}_l})$ |
| | Dense1 | $l = 3$ | $\mathcal{O}(2\rho_l N_{\mathcal{W}_l})$ |
| MCNet | Bi-LSTM1 | $l = 4$ | $\mathcal{O}(4\rho_l N(K_{\text{Bin},l} K_{\text{Bout},l} + K_{\text{Bout},l}^2))$ |
| | Bi-LSTM2 | $l = 5$ | $\mathcal{O}(4\rho_l N(K_{\text{Bin},l} K_{\text{Bout},l} + K_{\text{Bout},l}^2))$ |
| | Multi-head attention | $l = 6$ | $\mathcal{O}(\rho_l N_A d_A(3K_{\text{Ain}} + K_{\text{Aout}}) + 2N_A d_A^2)$ |
| | Dense2 | $l = 7$ | $\mathcal{O}(2\rho_l N_{\mathcal{W}_l})$ |
| | Dense3 | $l = 8$ | $\mathcal{O}(2\rho_l N_{\mathcal{W}_l})$ |

$$Z = -\text{round}\left(\frac{(2^b - 1)\min(\mathcal{W}_i)}{\max(\mathcal{W}_i) - \min(\mathcal{W}_i)}\right), \qquad (32)$$

where $\max(\mathcal{W}_i)$ and $\min(\mathcal{W}_i)$ denote the maximum and minimum values of the weights in $\mathcal{W}_i$. After weight pruning and quantization, the model compression ratios of SSCNet, MCNet, and the C-AMC framework are respectively given by

$$\gamma_{\text{SSC}} = \frac{32}{b}\sum_{l=1}^{L_{\text{SSC}}}\frac{N_{\mathcal{W}_l}}{N_{\text{SSC}}}\frac{1}{1 - \rho_l}, \qquad (33)$$

$$\gamma_{\text{MC}} = \frac{32}{b}\sum_{l=L_{\text{SSC}}+1}^{L_{\text{SSC}}+L_{\text{MC}}}\frac{N_{\mathcal{W}_l}}{N_{\text{MC}}}\frac{1}{1 - \rho_l}, \qquad (34)$$

$$\gamma = \frac{N_{\text{SSC}}\gamma_{\text{SSC}} + N_{\text{MC}}\gamma_{\text{MC}}}{N_{\text{SSC}} + N_{\text{MC}}}, \qquad (35)$$

where $N_{\text{SSC}}$ and $N_{\text{MC}}$ denote the numbers of weights of SSCNet and MCNet, respectively.

### C. Complexity Analysis

In this subsection, the model inference complexity of the C-AMC framework is analyzed by using the complexity of floating-point operations (FLOPs) as the metric.

The computational complexities of SSCNet and MCNet are listed in Table I layer by layer. Specifically, $D_l$ denotes the length of output feature maps for the Conv1D layer with the index $l$. $K_{\text{Bin},l}$ and $K_{\text{Bout},l}$ denote the input and output dimensions of the unit for the Bi-LSTM layer with the index $l$. $K_{\text{Ain}}$ and $K_{\text{Aout}}$ denote the input and output dimensions for the multi-head attention layer. The number of weights, $N_{\mathcal{W}_i}$, is introduced for Conv1D and dense layers to simplify the expression. From Table I, the theoretical complexity is reduced significantly thanks to the weight pruning. If further considering the weight quantization, the model can be more computationally effective with the smaller size.

Then the complexity of the C-AMC framework is given by

$$\mathcal{C}_{\text{C-AMC}} = \mathcal{C}_{\text{SSC}} + \lambda\mathcal{C}_{\text{MC}}, \qquad (36)$$

where $\lambda$ denotes the computation cost factor. Since the computation cost at the edge server is much lower than that in the edge device, $\lambda \ll 1$ usually holds true, leading to a much lower effective complexity for the C-AMC framework.
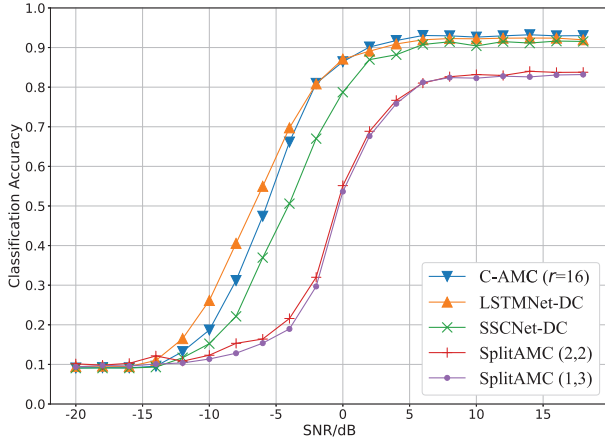
Fig. 6. Performance comparison between the C-AMC framework and baseline schemes.



Fig. 7. The confusion matrices of the C-AMC ($r = 16$) at different SNRs: (a) 0 dB, (b) 16 dB.

## VI. SIMULATION RESULTS

In this section, simulation results are provided to demonstrate the superiority of the proposed C-AMC framework as well as to reveal some insights for the practical design.

### A. Simulation Settings

The widely recognized RadioML2016.10A dataset is used, which includes 11 modulation types at different signal-to-noise ratios (SNRs). The sampling length of the source signal is set as $L = 512$. For the regular training of the C-AMC framework, the training, validation, and testing set contain $132,000$, $44,000$, and $44,000$ samples, respectively. The initial learning rate is $0.001$ and the batch size is 200. Three baseline schemes listed below are used for performance comparison.

- **LSTMNet-DC [31]**. This scheme deploys the LSTM network proposed in [31] at the edge device to classify the modulation types directly without the collaboration of the edge server.
- **SSCNet-DC**. This scheme deploys SSCNet at the edge device to classify the modulation types directly without the collaboration of the edge server.
- **SplitAMC [33]**. This scheme splits ResNet-18 with four residual blocks into two parts and deploys them at the edge device and server, respectively, for collaborative AMC.

Unless otherwise stated, the SNR in simulation results represents the wireless sensing channel SNR.

### B. Results

Fig. 6 shows the prediction accuracy versus SNR for the C-AMC framework with $r = 16$ and baseline schemes. From Fig. 6, the C-AMC framework outperforms SplitAMC significantly since the latter relies on a mass of I/Q samples. Compared to SSCNet-DC, the C-AMC framework improves the accuracy substantially thanks to the collaboration of MC-Net. LSTMNet-DC achieves the better performance than the C-AMC framework at the low SNR regime, which is not a regime of interest sin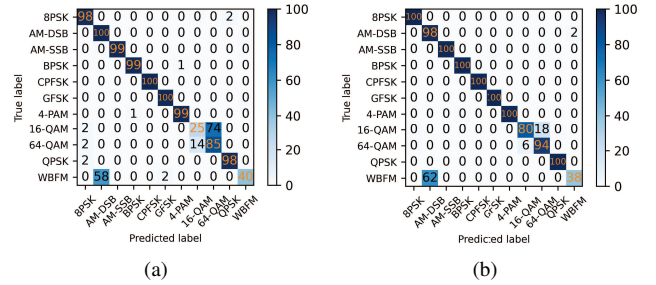ce the prediction accuracy is not high. As the SNR increases, the C-AMC framework gradually surpasses LSTMNet-DC and finally converges at the accuracy over 93% with the gain almost 1% compared to the latter. To provide a comprehensive comparison between the C-AMC framework and baseline schemes, Table II lists the prediction accuracy, the number of weights, FLOPs, the model inference time, the compression rate $r$, and the security for schemes in Fig. 6, where the second through forth metrics related to the computational complexity are considered for both the edge device and the edge server. In addition to the accuracy, the C-AMC framework also outperforms SplitAMC in terms of other four metrics, especially the compression rate. Instead of compressing the sampled data, SplitAMC enlarges the dimension of the feature vector and incurs the prohibitively high transmission overhead. The C-AMC framework and SSCNet-DC have the almost same number of weights, FLOPs, and inference time at the edge device, indicating that the performance gain achieved by C-AMC does not incur the additional cost for the edge device. Although the C-AMC framework and LSTMNet-DC have the almost same number of weights, the former incurs fewer FLOPs and thus reduces the inference time by $3\times$ for the edge device yet achieving the higher accuracy. Moreover, the C-AMC framework can protect the secure information by transmitting the intricate semantic embedding instead of the raw data or the classification result. Therefore, the C-AMC framework possesses the unique advantage under the overall consideration on these key metrics.

Fig. 7 shows the classification confusion matrices of the C-AMC framework with $r = 16$ at SNR = 0 dB and 16 dB, respectively. From Fig. 7, the accuracy loss of the C-AMC framework mainly comes from 16-QAM, 64-QAM, and WBFM. With the SNR increasing to 16 dB, the classification accuracy of 16-QAM and 64-QAM improves considerably even with a little confusion between them, leaving WBFM as the dominated destroyer of the overall accuracy.

To reveal the compression effectiveness of SSCNet, Fig. 8 shows the feature visualization of the semantic embedding $\mathbf{x}_s$ with $r = 16$. From the figure, the feature points of most modulation types can be clustered and separated from each other despite a few outliers while the feature points of WBFM and AM-DSB are severely overlapped. This phenomenon is coincide with the result in Fig. 7(b). Therefore, SSCNet can reserve most useful information when compressing the raw

TABLE II
COMPARISON OF PREDICTION ACCURACY, WEIGHTS, FLOPS, AND INFERENCE TIME

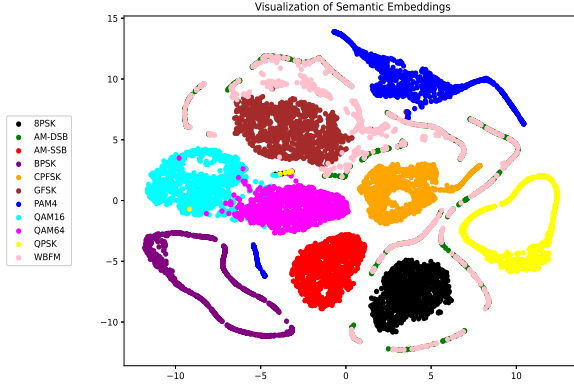| Model | $P_{\text{acc}}$ | Weights (Device) | Weights (Server) | FLOPs (Device) | FLOPs (Server) | Inference Time (Device) | Inference Time (Server) | Compression Rate | Security |
|---|---|---|---|---|---|---|---|---|---|
| C-AMC ($r=16$) | **0.932** | 20.0 K | **697.0 K** | **17.9 M** | **27.4 M** | **0.022 ms** | **0.034 ms** | **16** | ✓ |
| LSTMNet-DC | 0.924 | 20.5 K | - | 21.3 M | - | 0.065 ms | - | - | ✗ |
| SSCNet-DC | 0.916 | **18.0 K** | - | **17.9 M** | - | **0.022 ms** | - | - | ✗ |
| SplitAMC(2,2) | 0.839 | 679.7 K | 10.5 M | 150 M | 134 M | 0.133 ms | 0.107 ms | 0.125 | ✓ |
| SplitAMC(1,3) | 0.834 | 152.2 K | 11.0 M | 82.5 M | 202.5 M | 0.069 ms | 0.158 ms | 0.0625 | ✓ |



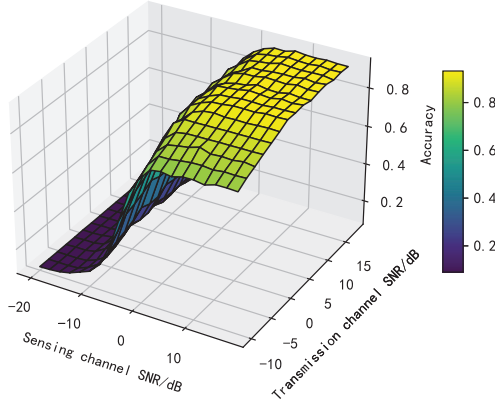Fig. 8. Feature visualization of the semantic embedding with $r = 16$.



Fig. 10. Classification accuracy of C-AMC with different compression rates.



Fig. 9. Classification accuracy of C-AMC ($r = 16$) versus the sensing and transmission channel SNRs.



Fig. 11. Classification accuracy of C-AMC with different compression rates in terms of the modulation type.

data and thus collaborate with MCNet to yield the superior classification accuracy.

Fig. 9 shows the classification accuracy of C-AMC versus the sensing and transmission channel SNRs with $r = 16$. From Fig. 9, the accuracy increases with the sensing channel SNR dramatically while with the transmission channel SNR moderately, revealing that the former dominates the classification performance.

Fig. 10 shows the classification accuracy of the C-AMC framework with different compression rates. Increasing $r$ from 8 to 16 does almost not degrade the accuracy. If the spectrum semantic embedding is compressed with $r = 32$, the accuracy decreases nontrivially. The accuracy can be maintained over 80% even if $r$ reaches to 64, demonstrating the robustness of the C-AMC framework to the high compression rate.
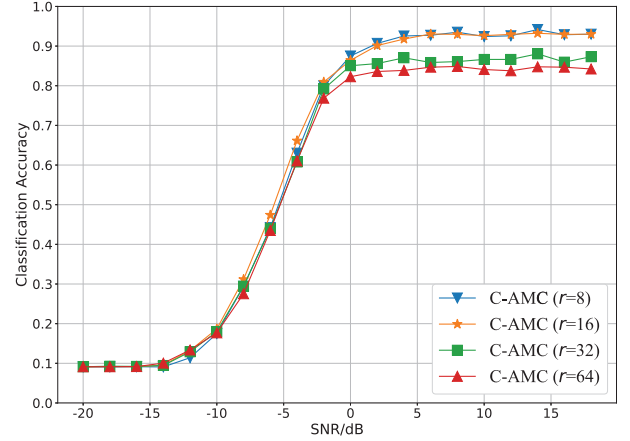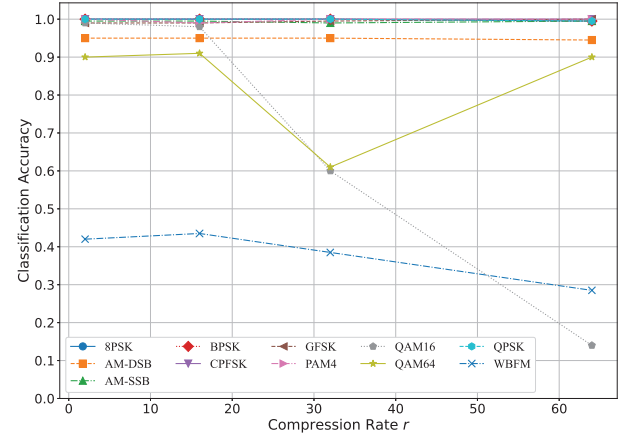
To further analyze the reason of performance change versus the compression rate, Fig. 11 shows the accuracy curve of each considered modulation type. From the figure, 16-QAM and 64-QAM are sensitive to the compression along with the confusion between them while WBFM keeps its role of dominated destroyer regardless of $r$. Therefore, if these three modulation types are not considered in some practical applications, the compression rate can reach to at least 64 without the performance loss.

Fig. 12 shows the classification accuracy of C-AMC with the mismatched length of sampling data, i.e., $L$ is halved. To match the input dimension of SSCNet, $L/2$ zeros are
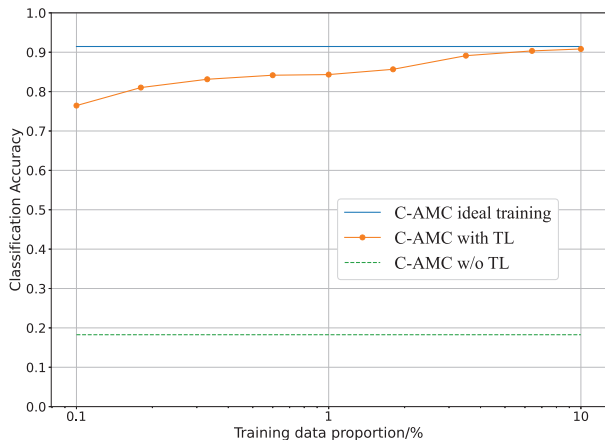
Fig. 12. Classification accuracy of C-AMC with the mismatched length of sampling data.
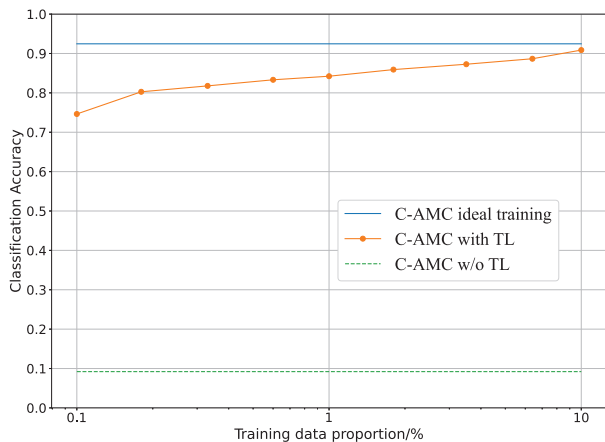


Fig. 13. Classification accuracy of C-AMC with the mismatched dataset.



Fig. 14. Classification accuracy of C-AMC versus the pruning ratio.

TABLE III
COMPARISON OF WEIGHTS, FLOPS, AND INFERENCE TIME FOR C-AMC

| | Model | Weights | FLOPs | Inference Time |
|---|---|---|---|---|
| Original | SSCNet | 20.0 K | 17.9 M | 0.022 ms |
| | MCNet | 697.0 K | 27.4 M | 0.034 ms |
| | C-AMC | 717.0K | 45.3 M | 0.056 ms |
| Pruned | SSCNet($\rho = 0.7$) | 6.4 K | 5.7 M | 0.019 ms |
| | MCNet($\rho = 0.8$) | 147.1 K | 5.7 M | 0.033 ms |
| | C-AMC | 153.4 K | 11.4 M | 0.052 ms |
| Pruned and Quantized | SSCNet ($\rho = 0.7, b = 8$) | 6.4 K | 5.7 M | 0.019 ms |
| | MCNet ($\rho = 0.8, b = 8$) | 147.1 K | 5.7 M | 0.032 ms |
| | C-AMC | 153.4 K | 11.4 M | 0.051 ms |

uniformly inserted into the sampled data sequence to yield a new sequence with length $L$. "C-AMC w/o TL" means directly inputting the new sequence into the C-AMC framework for classification and achieves the very poor performance. "C-AMC ideal training" means training the C-AMC framework from scratch using sufficient training data with the same format as the new sequence and acts as an upper bound. "C-AMC with TL" means adapting the original C-AMC framework with a few adaptation data, i.e., transfer learning (TL). The x-axis represents the ratio of the adaptation data quantity to the training data quantity. By resorting to TL, the C-AMC framework can rapidly adapt to the new scenario even with the proportion $0.1\%$ and achieve the satisfactory performance with the proportion $1\%$. If the proportion is increased to $10\%$, the accuracy of C-AMC becomes very close to that of ideal training. Furthermore, Fig. 13 shows the classification accuracy of C-AMC with the mismatched dataset, where 11 different modulation types in the testing set are chosen from RadioML2018 dataset. This figure reveals the similar phenomenon, further validating the robustness of C-AMC when faced with various mismatched cases.

Finally, Fig. 14 shows the classification accuracy of C-AMC versus the pruning ratio. Specifically, SSCNet pruning, MCNet
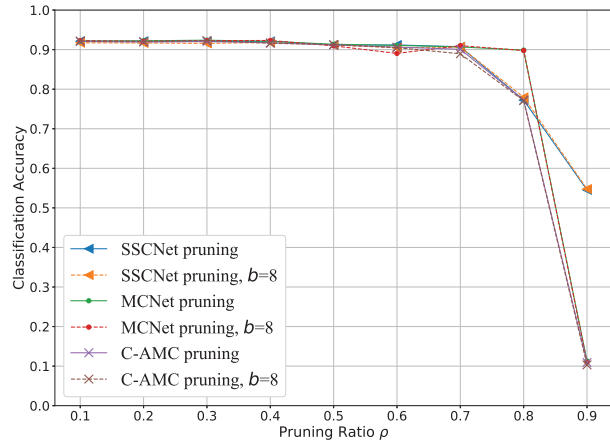
pruning, and C-AMC pruning correspond to the following three settings: 1) $\rho_1 = \cdots = \rho_{L_{\text{SSC}}} = \rho$, $\rho_{L_{\text{SSC}}+1} = \cdots = \rho_{L_{\text{SSC}}+L_{\text{MC}}} = 0$; 2) $\rho_1 = \cdots = \rho_{L_{\text{SSC}}} = 0$, $\rho_{L_{\text{SSC}}+1} = \cdots = \rho_{L_{\text{SSC}}+L_{\text{MC}}} = \rho$; 3) $\rho_1 = \cdots = \rho_{L_{\text{SSC}}+L_{\text{MC}}} = \rho$. From this figure, SSCNet can be pruned up to $\rho = 0.7$ with the almost unchanged performance while the threshold for MCNet is $\rho = 0.8$. The performance of C-AMC pruning is dominated by the one of SSCNet pruning and MCNet pruning that has the lower accuracy. So the curve of C-AMC pruning coincides with both SSCNet pruning and MCNet pruning with $\rho \leq 0.7$. As $\rho$ increases to $0.8$ and $0.9$, C-AMC pruning respectively coincides with SSCNet pruning and MCNet pruning. In addition, the 8-bit quantization considered for each pruning strategy does almost not cause the performance loss. Table III lists the number of weights, FLOPs, and the model inference time for the original, pruned, and pruned and quantized C-AMC frameworks. Through the appropriate weight pruning and quantization, i.e., $\rho = 0.7$ and $b = 8$, SSCNet can be compressed by $\frac{20}{6.4} \times 4 \approx 12\times$ with the inference accelerated by $14\%$, further lightening the computation load for the edge device.

## VII. CONCLUSION

In this article, an EL based C-AMC framework is proposed to meet the requirements of classification accuracy, computation load, transmission overhead, and data privacy for hierarchical CR networks. The C-AMC framework consists

of a lightweight SCCNet deployed at the edge device for spectrum semantic compression and an elaborated MCNet deployed at the edge server to predict the modulation type based on the noisy semantic embedding. Both the offline and online training procedures of the C-AMC framework are elaborated. The model pruning and quantization strategy for the C-AMC framework is developed to further facilitate the computation and storage, followed by the comprehensive complexity analysis. Simulation results demonstrate the superiority of the C-AMC framework and reveal useful insights paving the practical implementation.

Compared with the current semantic communication system focusing on semantics of the image, text or speech [35], [36], the spectrum semantics include various concrete types, e.g., the signal modulation type, the emitter type and position, the channel characteristics and the corresponding environment situation. In the future research, it is anticipated to construct a unified framework for the spectrum semantic cognition based on the C-AMC framework.

## REFERENCES

[1] O. A. Dobre, A. Abdi, Y. Bar-Ness, and W. Su, "Survey of automatic modulation classification techniques: Classical approaches and new trends," *IET Commun.*, vol. 1, no. 2, pp. 137–156, 2007.

[2] P. Panagiotou, A. Anastasopoulos, and A. Polydoros, "Likelihood ratio tests for modulation classification," in *Proc. 21st Century Mil. Commun. Archit. Technol. Inf. Superiority (MILCOM)*, vol. 2, 2000, pp. 670–674.

[3] Y. A. Eldemerdash, O. A. Dobre, and M. Öner, "Signal identification for multiple-antenna wireless systems: Achievements and challenges," *IEEE Commun. Surv. Tutor.*, vol. 18, no. 3, pp. 1524–1551, Jan. 2016.

[4] B. Dong *et al.*, "A lightweight decentralized learning-based automatic modulation classification method for resource-constrained edge devices," *IEEE Internet Things J.*, vol. 9, no. 24, pp. 24708–24720, Dec. 2022.

[5] J. L. Xu, W. Su, and M. Zhou, "Likelihood ratio approaches to automatic modulation classification," *IEEE Trans. Syst. Man. Cybern. PartC*, vol. 41, no. 4, pp. 455–469, Jul. 2011.

[6] W. A. Gardner, "Signal interception: A unifying theoretical framework for feature detection," *IEEE Trans. Commun.*, vol. 36, no. 8, pp. 897–906, Aug. 1988.

[7] D. Boutte and B. Santhanam, "A hybrid ICA-SVM approach to continuous phase modulation recognition," *IEEE Signal Proc. Lett.*, vol. 16, no. 5, pp. 402–405, May 2009.

[8] L. Han, F. Gao, Z. Li, and O. A. Dobre, "Low complexity automatic modulation classification based on order-statistics," *IEEE Trans. Wireless Commun.*, vol. 16, no. 1, pp. 400–411, Jan. 2017.

[9] T. J. O'Shea, J. Corgan, and T. C. Clancy, "Convolutional radio modulation recognition networks," in *Proc. Int. Conf. Eng. Appl. Neural Netw.*, 2016, pp. 213–226.

[10] T. J. O'Shea, T. Roy, and T. C. Clancy, "Over-the-air deep learning based radio signal classification," *IEEE J. Sel. Topics Signal Process.*, vol. 12, no. 1, pp. 168–179, Jan. 2018.

[11] S. Peng *et al.*, "Modulation classification based on signal constellation diagrams and deep learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 99, pp. 718–727, Mar. 2018.

[12] B. Tang, Y. Tu, Z. Zhang, and Y. Lin, "Digital signal modulation classification with data augmentation using generative adversarial nets in cognitive radio networks," *IEEE Access*, vol. 6, pp. 15713–15722, 2018.

[13] F. Meng, P. Chen, L. Wu, and X. Wang, "Automatic modulation classification: A deep learning enabled approach," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 10760–10772, Nov. 2018.

[14] Y. Wang, M. Liu, J. Yang, and G. Gui, "Data-driven deep learning for automatic modulation recognition in cognitive radios," *IEEE Trans. Veh. Technol.*, vol. 68, no. 4, pp. 4074–4077, Apr. 2019.

[15] T. Huynh-The, C.-H. Hua, Q.-V. Pham, and D.-S. Kim, "MCNet: An efficient CNN architecture for robust automatic modulation classification," *IEEE Commun. Lett.*, vol. 24, no. 4, pp. 811–815, Apr. 2020.

[16] Z. Zhang, H. Luo, C. Wang, C. Gan, and Y. Xiang, "Automatic modulation classification using CNN-LSTM based dual-stream structure," *IEEE Trans. Veh. Technol.*, vol. 69, no. 11, pp. 13521–13531, Nov. 2020.

[17] S. Chang, S. Huang, R. Zhang, Z. Feng, and L. Liu, "Multitask-learning based deep neural network for automatic modulation classification," *IEEE Internet Things J.*, vol. 9, no. 3, pp. 2192–2206, Feb. 2022.

[18] Z. Ke and H. Vikalo, "Real-time radio technology and modulation classification via an LSTM auto-encoder," *IEEE Trans. Wireless Commun.*, vol. 21, no. 1, pp. 370–382, Jan. 2022.

[19] H. Zhang, F. Zhou, Q. Wu, W. Wu, and R. Q. Hu, "A novel automatic modulation classification scheme based on multi-scale networks," *IEEE Trans. Cognit. Commun. Netw.*, vol. 8, no. 1, pp. 97—110, 2021.

[20] Y. Wang, G. Gui, T. Ohtsuki, and F. Adachi, "Multi-task learning for generalized automatic modulation classification under non-Gaussian noise with varying SNR conditions," *IEEE Trans. Wireless Commu.*, vol. 20, no. 6, pp. 3587–3596, June 2021.

[21] Y. Lin, Y. Tu, and Z. Dou, "An improved neural network pruning technology for automatic modulation classification in edge devices," *IEEE Trans. Veh. Technol.*, vol. 69, no. 5, pp. 5703–5706, May 2020.

[22] Z. Chen *et al.*, "Channel pruning method for signal modulation recognition deep learning models," *IEEE Trans. Cognit. Commun. Netw.*, to be published, 2023.

[23] L. Guo, Y. Wang, Y. Liu, Y. Lin, H. Zhao, and G. Gui, "Ultralight convolutional neural network for automatic modulation classification in internet of unmanned aerial vehicles," *IEEE Internet Things J.*, vol. 11, no. 11, pp. 20831–20839, June 2024.

[24] L. Li, J. Huang, Q. Cheng, H. Meng, and Z. Han, "Automatic modulation recognition: A few-shot learning method based on the capsule network," *IEEE Commun. Lett.*, vol. 10, no. 3, pp. 474–477, Mar. 2021.

[25] W. Lin, D. Hou, J. Huang, L. Li and Z. Han, "Transfer learning for automatic modulation recognition using a few modulated signal samples," *IEEE Trans. Veh. Technol.*, vol. 72, no. 9, pp. 12391–12395, Sep. 2023.

[26] W. Deng, X. Wang, Z. Huang, and Q. Xu, "Modulation classifier: A few-shot learning semi-supervised method based on multi-modal information and domain adversarial network," *IEEE Commun. Lett.*, vol. 27, no. 2, pp. 576—580, Feb. 2023.

[27] P. Qi, X. Zhou, Y. Ding, Z. Zhang, S. Zheng, and Z. Li, "FedBKD: Heterogenous federated learning via bidirectional knowledge distillation for modulation classification in IoT-edge system," *IEEE J. Sel. Topics Signal Process.*, vol. 17, no. 1, pp. 189–204, Jan. 2023.

[28] Y. Wang, G. Gui, H. Gacanin, B. Adebisi, H. Sari, and F. Adachi, "Federated learning for automatic modulation classification under class imbalance and varying noise condition," *IEEE Trans. Cogn. Commun. Netw.*, vol. 8, no. 1, pp. 86–96, March 2022.

[29] S. Zhang, Y. Yang, Z. Zhou, Z. Sun, and Y. Lin, "DIBAD: A Disentangled Information Bottleneck Adversarial Defense Method Using Hilbert-Schmidt Independence Criterion for Spectrum Security," *IEEE Trans. Inf. Forensics Secur.*, vol. 19, pp. 3879–3891, 2024.

[30] P. Dong, Q. Wu, X. Zhang and G. Ding, "Edge semantic cognitive intelligence for 6G networks: Novel theoretical models, enabling framework, and typical applications," *China Commun.*, vol. 19, no. 8, pp. 1–14, Aug. 2022.

[31] S. Rajendran, W. Meert, D. Giustiniano, V. Lenders, and S. Pollin, "Deep learning models for wireless signal classification with distributed low-cost spectrum sensors," *IEEE Trans. Cognit. Commun. Netw.*, vol. 4, no. 3, pp. 433–445, May 2018.

[32] W. Xu *et al.*, "Edge learning for B5G networks with distributed signal processing: Semantic communication, edge computing, and wireless sensing," *IEEE J. Sel. Topics Signal Process.*, vol. 17, no. 1, pp. 9–39, Jan. 2023..

[33] J. Park, S. Oh and S. -L. Kim, "SplitAMC: Split learning for robust automatic modulation classification," in *Proc. IEEE 97th Veh. Technol. Conf. (VTC-Spring)*, Florence, Italy, 2023, pp. 1–6.

[34] M. Nagel *et al.*, "A white paper on neural network quantization," *arXiv preprint arXiv: 2106.0829*, 2021.

[35] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE Trans. Signal Process.*, vol. 69, pp. 2663–2675, Apr. 2021.

[36] Z. Weng and Z. Qin, "Semantic communication systems for speech transmission," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2434–2444, Aug. 2021.