# A Deep-Tree-Model-Based Radio Resource Distribution for 5G Networks

M. Shamim Hossain and Ghulam Muhammad

## Abstract

Deep learning is a branch of machine learning that learns the high-level abstraction of data in a layered structure. Since its invention, it has been successfully applied in many image and speech processing applications. The success of deep learning depends on how big the data size is. Recently, the number of smart sensors and the Internet of Things have increased exponentially. This, in turn, has created huge traffic congestion in mobile and wireless communication networks. The available network resources need to be carefully utilized for seamless transmission of this large amount of data. Fortunately, deep learning performs very well with the big size of data. Therefore, the gap between machine learning research and advanced communication research should be narrowed down. In this article, we target an intelligent allocation of radio resources for 5G networks using deep learning. A framework consisting of a deep tree model and a long short-term memory network is proposed to predict future traffic congestion. Based on the prediction, the uplink and downlink ratio is adapted to utilize the resources optimally. Experimental results demonstrate that the proposed framework can achieve a low packet loss ratio and high throughput.

## Introduction

Deep learning is a powerful technique of machine learning that has brought numerous advantages over the last decade. Many models of deep learning have been proposed in various applications. The problems of conventional neural networks such as the weight vanishing problem and regularization can be overcome using a very deep and narrow model, or a shallow and broad model. Both problems cannot be solved simultaneously so far using these two types of models. Therefore, a compromise is often realized to design a new deep model. Some of the state-of-the-art deep learning models are restricted Boltzmann machine (RBM), auto-encoder (AE), sparse auto-encoder (SAE), convolutional neural network (CNN), recurrent neural network (RNN), and generative adversarial network (GAN). These models are mostly used in various image processing and speech processing applications.

Internet-based devices are increasing rapidly in terms of both device types and applications. The Internet of Things (IoT) is now realized in many aspects of our daily life such as medical, entertainment, and household items. In recent years, the use of IoT has increased so much that it is predicted that the annual usage of IP mobile traffic may reach up to 1015 MB by 2021. People are more inclined to use wireless connectivity than before. Smartphones and IoT are now part and parcel of our lives. Demands for capacity, seamless transmission, and real-time execution have increased manifold. Various IoT devices supply big data, which are heterogeneous and sparse in nature. To cope with this type of data, intelligent heterogeneous networks are needed that can meet strict end-user requirements.

Big data helps deep learning to be realized in an excellent way because it eliminates domain expertise and extracts hierarchical features. An abstract correlation can also be established that reduces the effort of traditional feature extraction methods. Now, deep learning can be run using a graphical processing unit (GPU) in parallel, which enables execution in a very short time. Deep learning can facilitate a network manager in analyzing and distributing data in a timely manner and with high accuracy, and can replace traditional algorithms such as game theory and metaheuristics.

So far, deep learning has been used sparsely in some communication applications. Using big data, deep-learning-driven mobile and wireless networking can be categorized into seven groups:
- Network-level mobile data analysis (e.g., network prediction and traffic classification)
- Application-level mobile data analysis, such as healthcare and speech recognition applications
- Mobility analysis
- User localization
- Network control, such as optimization, routing, scheduling, resource allocation, and radio control
- Network security (e.g., user privacy and infrastructure)
- Signal processing, for example, modulation and multiple-input multiple-output (MIMO) systems.

In addition, lifelong learning and transfer learning are used to cooperate with the changing mobile environments, and parallel deep models are used in mobile systems and distributed data centers.

A specific type of deep learning can be utilized for a particular type of mobile network solution [1]. For instance, RBM is used in model weight initialization and network flow prediction, AE and SAE are used in mobile data reduction and mobile data anomaly prediction, CNN is widely used for spatial

The authors are with King Saud University. The corresponding author is M. Shamim Hossain.

mobile data analysis, while RNN is mainly used for mobile spatio-temporal data analysis. GAN is mostly used for virtual mobile data generation.

The main ingredients of a 5G mobile network are massive MIMO, beamforming, ultra-dense networks (UDNs), and software-defined networking (SDN). 5G has distinctive and dynamic network features. Early radio systems used non-coherent modulation techniques such as amplitude modulation and frequency modulation. Now, new systems need dynamic channel signal interference (CSI) because of the high data rate and unavailability of the wideband spectrum. 5G systems require spatial signal processing of massive MIMO. As the traffic is continuously growing over time, the 5G network operators need to use dynamic traffic control policies instead of traditional policies to prevent frequent packet loss and poor user quality of experience (QoE).

If there are no intelligent ways to assign radio resources to user equipment (UE) by base stations such as the evolved node B (eNB), there is a high possibility that regular or irregular traffic congestion will occur at UDNs. A predictive congestion technique can be deployed to change the uplink and downlink arrangements before actual congestion occurs, and some of the eNB occupancies can be reduced for smooth traffic.

In this article, we propose a deep-model-based radio resource distribution framework for 5G UDNs. The main objective of the framework is to predict traffic congestion and the occupancy state of the eNBs so that an adaptive uplink and downlink ratio can be applied to avoid traffic congestion. To this end, the proposed framework introduces a tree-based deep model, which has high information density and less model complexity. Following the tree-based model, a long short-term memory (LSTM) structure is applied to predict traffic congestion. The tree-based model uses convolutional layers, which are useful when dealing with spatial data generated by the UE. LSTM is useful to predict an output based on current and past data. Thus, the proposed distribution framework can effectively use past and present spatial data to predict eNB occupancy so that appropriate resource management can be deployed in 5G UDNs.

The main contributions of this article are as follows:
- Propose a framework with a new tree-based deep model together with an LSTM for radio resource management purposes in 5G UDNs.
- Perform experiments with the proposed framework and provide comparisons with conventional approaches.
- Achieve less computational complexity with the proposed deep model.

To the best of our knowledge, the combination of the tree-based deep model and LSTM has never been used in any previous works whether they are related to mobile and wireless communication applications or not.

The rest of the article is organized as follows. The next section briefly mentions some previous works related to deep learning and wireless communication. Then the proposed framework is elaborated. After that, some experimental results and comparisons with conventional methods are presented. Finally, a conclusion with some future research directions is given at the end of the article.

# Related Previous Work

The authors of [2] proposed an LSTM-based deep learning approach to assigning resources to UEs and predicting future network characteristics for 5G UDNs. The future characteristics were obtained from the previous and current data. The approach cannot handle spatial data, which is a concern of this approach.

CSI structures were reviewed in [3], and a deep neural network (DNN)-based nonlinear model was used to solve the problems of inadequate channel modeling in the CSI structure. CSI is essential for a high beamforming gain, which is an ingredient of the transmission of millimeter-wave signals. CSI is also responsible for ultra-reliable and low-latency communication. Good management of CSI can contribute to user scheduling and buffer design for high-dimensional data transfer. The DNN used in [3] was not described in full detail, although it achieved good results.

Several deep learning models were investigated for RF fingerprinting and identification in [4]. As the data rate has been increased manifold, and a single channel can be shared by many transmitters, RF waveforms need to be identified and fingerprinted for a smooth allocation of the channel. Another deep-learning-based radio fingerprinting approach was proposed in [5]. In this approach, a CNN architecture was used, which was trained using I/Q sequence data. Five devices were used for fingerprinting in the experiments.

The authors of [6] examined the role of artificial intelligence (AI) on the key 5G enablers and proposed an AI-based 5G radio access structure. The AI-based system was deployed for the physical and medium access control (MAC) layer processing. A spectrum monitoring system using CNN was proposed in [7]. The system could identify various wireless signals such as the raw I/Q temporal wireless data, time-domain amplitude and phase data, and frequency-domain magnitude data. In the system, the CNN model was computationally expensive.

A good survey on deep-learning-based applications in mobile and wireless communications is given in [1]. Deep learning has been applied in the physical layer, data link layer, and routing layer in the form of resource allocation, modulation and coding, and traffic balancing [8]. Network security and compressive sensing of mobile data are two emerging domains where deep learning is being utilized. Not all the existing deep learning models can be directly used in mobile and wireless communication; therefore, some tailor-made models are necessary. The survey in [1] also focused on this important aspect.

A CNN-based automatic modulation recognition system was developed in [9]. The modulation recognition is an important attribute of cognitive radio. The CNN was trained using in-phase and quadrature component signals to identify eight modulation modes. An inception module was incorporated in the CNN model.

An audio-visual emotion recognition system using fog computing was proposed in [10]. The system utilized edge caches to store the parameters of the CNN model, which was the core ingredient in the system. Similarly, an environment classification system using the CNN model

5G systems require spatial signal processing of massive MIMO. As the traffic is continuously growing over time, 5G network operators need to use dynamic traffic control policies instead of traditional policies to prevent frequent packet loss and poor user quality of experience.
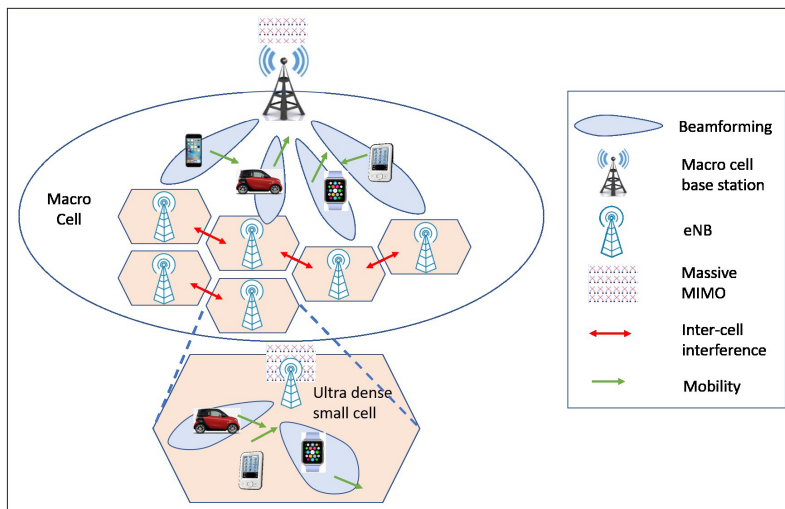
**FIGURE 1.** The 5G network environment considered in this work.

was developed using urban big data in [11]. An extreme learning machine together with CNN was used in the system.

Heterogeneous networks are necessary to deal with an unstructured and huge volume of data to be distributed over the network. AI-enabled cognitive computation for resource distribution in heterogeneous networks was proposed in [12]. A game-based technique was used for dynamic uplink-to-downlink ratio allocation for resource management in [13]. The main objective of this work was to minimize interference and maximize throughput. The work mainly used the currently available dataset, and there was no use of previous datasets. Time-division duplex (TDD) was proved superior to frequency-division duplex (FDD). A solution of dealing with healthcare big data for 5G communication using edge and cloud computing was proposed in [14].

Although there are several existing works using deep learning that deal with different aspects of mobile and wireless communications, there is little prior work that can predict future traffic congestion using previous and current data for optimal resource management.

## Network Environment

Figure 1 shows a simplified 5G network environment that is considered in this work. The environment is a part of a heterogeneous network. There is a macrocell consisting of several ultra-dense small cells. The base station of a UDN is referred to as an eNB. Each UDN can have several UEs that exchange data with the eNB. The data exchange is done using uplink and downlink directions with dynamic TDD. The macrocell eNB receives aggregated data from all the UDNs. The data rate from the UE is not constant; sometimes the traffic exceeds the capacity limit.

If a UE moves around the edge of a UDN cell, there is interference between this cell and the neighboring cell. In this case, there is a small signal-to-interference-plus-noise ratio. The eNB uses beamforming with different bandwidths for some specific UEs.

In FDD transmission, different bandwidths for uplink and downlink are deployed, while in TDD transmission, the same bandwidth for uplink and downlink is used. It has already been mentioned that TDD is better than FDD in 5G networks where there are UDNs involving massive MIMO and beamforming. Allocating different bandwidths will leave some resources unutilized while there is huge traffic. Therefore, a dynamic uplink and downlink ratio is selected using TDD transmission. In a conventional way, there is no intelligent technique to predict future traffic; it simply waits for traffic congestion to occur. Hence, we cannot avoid traffic congestion in a conventional way.

In the proposed framework, future traffic congestion is predicted using previous and current data. The data are gathered in the spatial domain from many UEs. Depending on the traffic prediction, the uplink-to-downlink ratio is adjusted to use all the resources properly.

## Background of the Proposed Deep Model

Deep models that were used in the domain of mobile and wireless communication are all either very deep and narrow, or shallow and broad. This is because all the famous deep models fall into these categories and are available with their pretrained versions. Very deep and thin models such as AlexNet and VGG Net are winners of various image processing applications. These models provide high accuracy; however, they have many weights to train. Residual networks consider depth as the most dominant factor to solve the weight vanishing problem. On the contrary, wide residual networks place emphasis on filters in each convolutional layer. Both very deep and shallow networks suffer from problems. For example, very deep and narrow models have the problem of redundant learning of weights, while shallow and broad models suffer from the lack of regularization. If there is a shallow model that can enhance the regularization, the model complexity can be reduced. Such models are crucial in mobile and wireless communication where the real-time result is one of the most important criteria for a user's QoE.

An inception model such as GoogLeNet uses a branching technique to enhance regularization. The problem of the inception model is that it is complex and difficult to tailor according to the applications. A root-tree-like model is another type of model that uses channel shuffle and pointwise operators to enhance regularization and reduce complexity. However, this type of model loses the advantages of convolutional operations.

Based on the above discussion, a deep model having a tree structure is proposed in this article. The model was initially realized in [15]. The proposed model has the following main characteristics:
- The model is structured as a tree, where there are two important parameters: branching factor and tree height.
- It has the benefit of multiple activation functions and batch normalization layers from the root of the tree to the leaves.
- The model uses a split technique to reduce the number of parameters in each branch.
- The model performs a group convolution in each tree module in a hierarchical way.

## The Proposed Deep Model with Tree Structure

In the proposed model, the input is divided into multiple branches with residual functions. Then the model performs some convolutional oper-
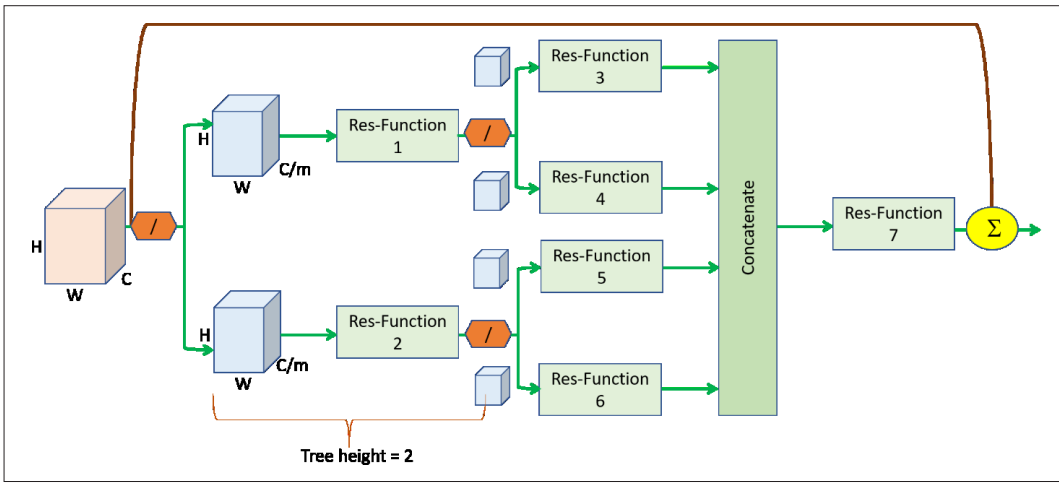
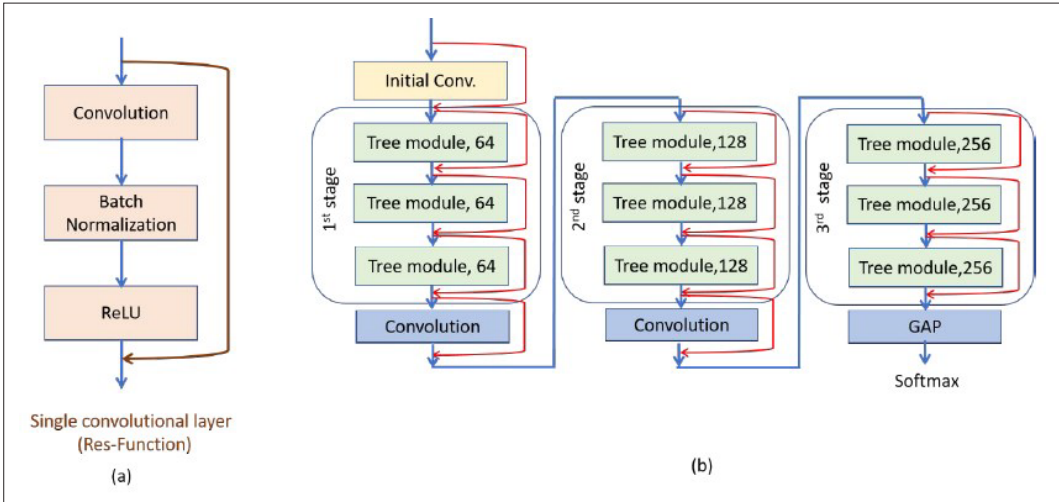FIGURE 2. Architecture of the proposed deep model with tree structure.

FIGURE 3. a) Residual function; b) the three stages of the proposed model.

ations and concatenates them to form a single output. To reduce complexity, the input volume is divided into two sub-volumes. The residual function is applied to each sub-volume, and the output is again divided into sub-volumes. This division continues until a prescribed tree height is reached. In this way, convolutional layers are disseminated in a tree-like structure. Figure 2 shows the structure of the proposed tree-based deep model. In the figure, the tree height is two, and the branching factor is also two. The split operation is shown as a division (/) sign.

The basic building block of the proposed model is a residual function (in the figure, it is marked as "Res-Function"). The residual function has three operations in succession: convolution, batch normalization, and nonlinear activation. In the convolution, the filter size is $3 \times 3$. Rectified linear unit (ReLU) is used as the nonlinear activation as it is quicker than the sigmoid function for convergence. The residual function is shown in Fig. 3a.

The proposed model works as follows. Suppose the input volume has the dimension $W \times H \times C$, where $W \times H$ is the spatial size of the input and $C$ is the number of channels. If the branching factor is $m$, the input volume is split into sub-volumes of size $W \times H \times C/m$. The residual function is applied to each sub-volume. The output of the residual function is again split into sub-volumes. This process is

repeated until the tree height is reached to build a complete tree. The outputs of all the residual functions at the last layer are concatenated. The last residual function is employed to the concatenated output and is augmented by the tree input. This last residual function can be used to reduce or enhance the dimensionality of the output.

There are three stages in the proposed model, as shown in Fig. 3. Before the first stage, there is an early convolutional layer that extends the input channels to a predetermined number of maps. The three stages bring in flexibility in the model in the sense that one stage structure does not depend on other stages when hyperparameters vary. Increasing the number of maps will increase the number of parameters, which may cause overfitting. Therefore, we need to downsample. Between the stages, there are convolutional layers whose job is to reduce the number of parameters. This is done by choosing a stride equal to two, which is equivalent to a downsample by a factor of two. This type of downsampling is better than a pooling layer.

A global average pooling (GAP) layer is used after the final stage to reduce the spatial dimensionality. Therefore, the output of the GAP layer has a dimension of $1 \times 1 \times C$, by taking the average per map. There is a softmax layer after the GAP layer. The filter size of all the convolutional layers except the last one is $3 \times 3$. The filter size of
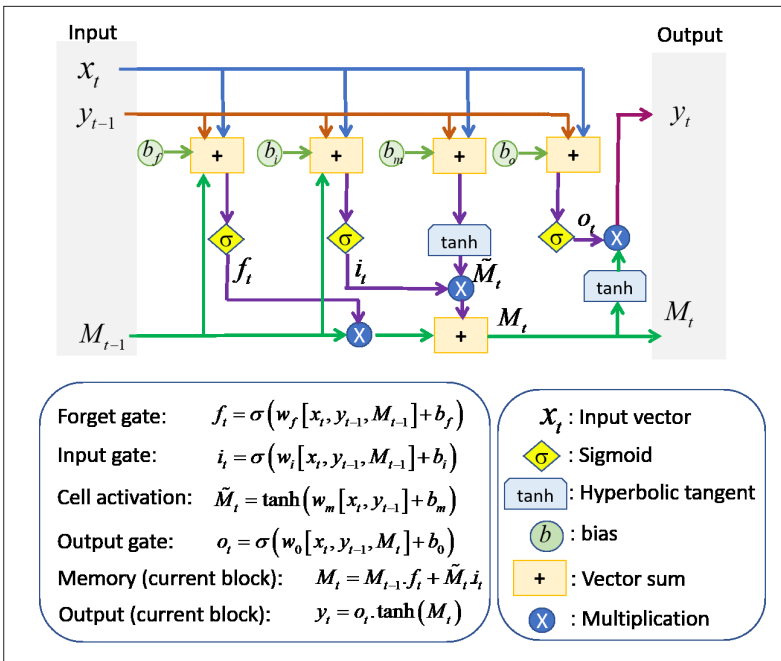
**FIGURE 4.** A simplified building block of LSTM.

the last convolutional layer is 1 × 1. The proposed three-stage model has the advantage of running in parallel without affecting the parameters' variation in each stage. This is particularly important when different UEs are available, and they can run in parallel to reduce the total execution time.

The traffic congestion is predicted using the information of previous traffic conditions. LSTM is a perfect match to use temporal information for a future prediction. In the proposed framework, the LSTM succeeds the tree-based deep model. The tree-based deep model extracts deep-learned features of traffic congestion. These features are fed into the LSTM. The LSTM is an adaptation of the RNN where long-term contextual dependencies are learned from the temporal feature sequence. Figure 4 shows a simple building architecture of the LSTM network. A unit of LSTM is characterized by three gates, which are the input gate, forget gate, and output gate. These gates control the operation in the unit with the help of three inputs: the input feature vector, the memory of the previous timestamp, and the output of the previous timestamp. The nonlinearity is obtained by the sigmoid function and the hyperbolic tangent function. The memory and the output are updated at each timestamp. The transition equations of LSTM at time instant $t$ are given in Fig. 4.

The output of the softmax layer of the tree-based deep model is fed to the LSTM. A timestamp is linked with the softmax layer so that the output of the deep model can be fed to the appropriate temporal input of the LSTM.

## EXPERIMENTS

The whole model was trained with fixed uplink and downlink ratios. The data were resized according to the input structure of the model, where the input dimension was 42 × 42 × 3. The third value, which is 3, was obtained using the temporal first-order and second-order derivatives. For a specific uplink and downlink ratio, all the rel-

evant data were used to update the weights. For the next specific uplink and downlink ratio, all the relevant data were used to update the previously updated weights. This updating process was continued until all the fixed uplink and downlink ratio related data were used.

A pretrained deep model with tree structure was the starting point of the experiments [15]. The weights were updated using the ratio data. A stochastic gradient descent algorithm was used to optimize the parameters of the model. The learning rate was set to be 0.01, and the batch size was 100. All the experiments were conducted on a workstation having a specification of 3.1 GHz, 32 GB RAM, and 8 GB GPU. The simulation was done in the TensorFlow environment.

In the configuration, we initially set the number of UEs per UDN small cell to be 15, where five of them used beamforming with various frequencies and the rest are non-beamforming. Each small cell had only one eNB. The number of UEs were gradually increased with an increment of 10. The packet loss rate (percent) and throughput for each case were measured. A packet loss was considered if the number of packets in the eNB exceeded 80 percent of the maximum capacity of the buffer size. A mean opinion score (MoS) was also obtained from 20 participants. The MoS had a scale between 0 and 5, where 5 means fully satisfied and 0 means not at all satisfied. The proposed model was compared with the conventional method in terms of all three measurements.

Figure 5 demonstrates the evaluation performance of the proposed model and the conventional one. Packet loss gradually increased with the increase of UEs. The proposed model had significantly less packet loss than did the conventional one. The throughput of the proposed model is much higher than that of the conventional one. Users were satisfied with the QoE of the proposed model, as can be seen from the MoS.

Based on the performance, we can state that the proposed model has less packet loss and high throughput, and is attractive to users. As the proposed model can intelligently predict the traffic congestion and thereby adjust the uplink and downlink ratio, it has high performance.

## CONCLUSION AND FUTURE WORK

An intelligent radio resource management system is proposed using deep learning. As the traffic of mobile networks is increasing exponentially, conventional resource management will not utilize the resource properly; the conventional method just waits for congestion to occur, and then tries to make a solution. The proposed model can predict traffic congestion using the previous and current data, and adjust the uplink and downlink ratio to utilize the resource in an efficient way. The proposed model uses a new tree-based deep model followed by LSTM. The deep model is computationally less expensive because the convolution operations are performed in parallel in the tree-like structure. The parallel processing is done using GPU, so the resource management was real-time.

The future of deep learning in 5G is very promising. Some of the directions can be as follows:
• The real-time processing of video data from various UEs using deep learning
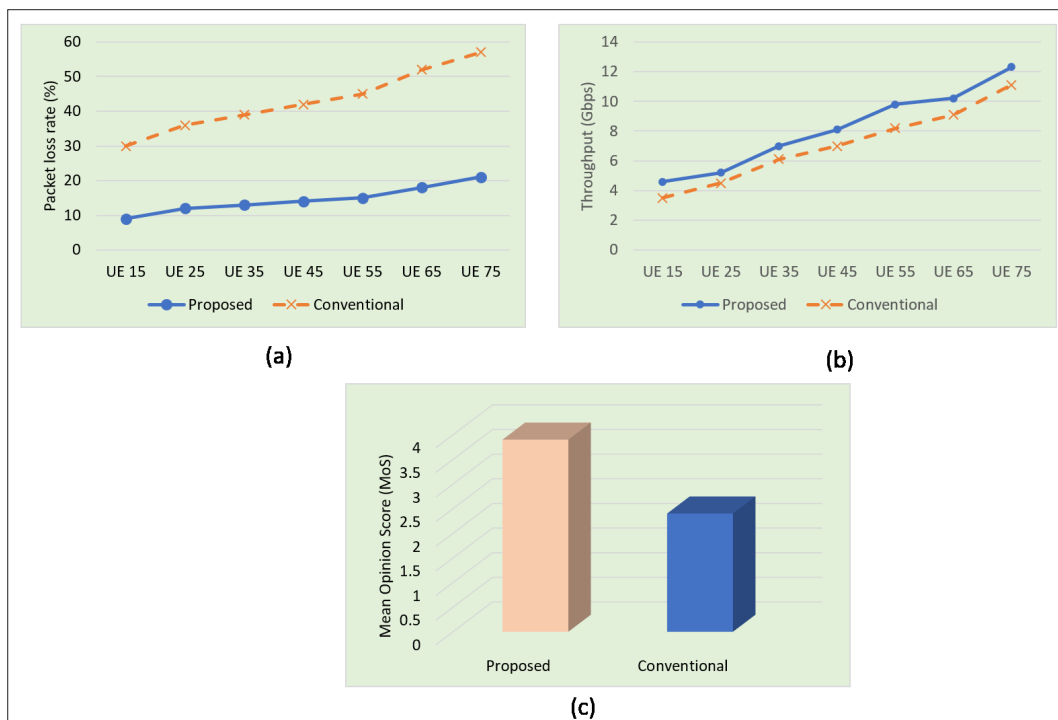
FIGURE 5. Performance evaluation of the proposed model and the conventional one: a) packet loss rate (percent); b) throughput; c) MoS.

The proposed model used a new tree-based deep model followed by the LSTM. The deep model was computationally less expensive because the convolution operations were performed in parallel in the tree-like structure. The parallel processing was done using GPU, so the resource management was real-time.

- A performance investigation of different tree heights and tree branch factors of the tree-based deep model
- An investigation of shallow models followed by LSTM for traffic congestion prediction

## Acknowledgement

## References

[1] C. Zhang, P. Patras, and H. Haddadi, "Deep Learning in Mobile and Wireless Networking: A Survey," *IEEE Commun. Surveys & Tutorials*, vol. 21, no. 3, 3rd qtr. 2019, pp. 2224–87.

[2] Y. Zhou *et al.*, "A Deep-Learning-Based Radio Resource Assignment Technique for 5G Ultra Dense Networks," *IEEE Network*, vol. 32, no. 6, Nov./Dec. 2018, pp. 28–34.

[3] Z. Jiang *et al.*, "Exploiting Wireless Channel State Information Structures Beyond Linear Correlations: A Deep Learning Approach," *IEEE Commun. Mag.*, vol. 57, no. 3, Mar. 2019, pp. 28–34.

[4] K. Youssef *et al.*, "Machine Learning Approach to RF Transmitter Identification," *IEEE J. Radio Frequency Identification*, vol. 2, no. 4, Dec. 2018, pp. 197–205.

[5] S. Riyaz *et al.*, "Deep Learning Convolutional Neural Networks for Radio Identification," *IEEE Commun. Mag.*, vol. 56, no. 9, Sept. 2018, pp. 146–52.

[6] M. Yao *et al.*, "Artificial Intelligence Defined 5G Radio Access Networks," *IEEE Commun. Mag.*, vol. 57, no. 3, Mar. 2019, pp. 14–20.

[7] M. Kulin *et al.*, "End-to-End Learning from Spectrum Data: A Deep Learning Approach for Wireless Signal Identification in Spectrum Monitoring Applications," *IEEE Access*, vol. 6, 2018, pp. 18484–18501.

[8] J. Wang *et al.*, "A Software Defined Network Routing in Wireless Multihop Network," *J. Net. Comp. Appl.*, vol. 85, no. 2017, May 2017, pp. 76–83.

[9] Y. Wang *et al.*, "Data-Driven Deep Learning for Automatic Modulation Recognition in Cognitive Radios," *IEEE Trans. Vehic. Tech.*, vol. 68, no. 4, Apr. 2019, pp. 4074–77.

[10] M. S. Hossain, and G. Muhammad, "Emotion Recognition Using Deep Learning Approach from Audio-Visual Emotional Big Data," *Info. Fusion*, vol. 49, no. 2019, Sept. 2019, pp. 69–78.

[11] M. S. Hossain and G. Muhammad, "Environment Classification for Urban Big Data Using Deep Learning," *IEEE Commun. Mag.*, vol. 56, no.11, Nov. 2018, pp. 44–50.

[12] K. Lin *et al.*, "Artificial-Intelligence-Based Data Analytics for Cognitive Communication in Heterogeneous Wireless Networks," *IEEE Wireless Commun.*, vol. 26, no. 3, June 2019.

[13] C. C. Chao *et al.*, "Distributed Dynamic-TDD Resource Allocation in Femtocell Networks Using Evolutionary Game," *Proc. 2015 IEEE 26th Annual PIMRC*, Hong Kong, China, Aug. 2015, pp. 1157–62.

[14] G. Muhammad *et al.*, "Edge Computing with Cloud for Voice Disorders Assessment and Treatment," *IEEE Commun. Mag.*, vol. 56, no. 4, Apr. 2018, pp. 60–65.

[15] A. Amory, G. Muhammad, and H. Mathkour, "Deep Convolutional Tree Networks," *Future Generation Computer Systems*, vol. 101, Dec. 2019, pp. 152–68.

## Biographies

M. Shamim Hossain [SM'09] (mshossain@ksu.edu.sa; corresponding author) is a professor in the Department of Software Engineering, College of Computer and Information Sciences King Saud University, Riyadh, Saudi Arabia. He is on the Editorial Boards of *IEEE Transactions on Multimedia, IEEE Multimedia, IEEE Network, IEEE Wireless Communications*, the *Journal of Network and Computer Applications*, and the *Journal of Multimedia Tools and Applications*. His research interests include cloud networking, smart environment (smart city, smart health), IoT, edge computing, multimedia for health care, the deep learning approach to multimedia processing, and multimedia big data. He is a Senior Member of ACM.

Ghulam Muhammad [M'10, SM'19] (ghulam@ksu.edu.sa) is a professor in the Department of Computer Engineering, College of Computer and Information Sciences, King Saud University. He received his Ph.D. degree in 2006 from Toyohashi University of Technology, Japan. His research interests include signal processing, image and speech signal processing, multimedia forensics, and healthcare. He has authored more than 200 journal and conference papers. He has supervised more than 15 Ph.D. and Master's theses. He owns two U.S. patents.