Check for updates

#### REVIEW





# Deep learning in physical layer communications: Evolution and prospects in 5G and 6G networks

Chengchen Mao<sup>1</sup> | Zongwen Mu<sup>2</sup> | Qilian Liang<sup>1</sup> | Ioannis Schizas<sup>1</sup> | Chenyun Pan<sup>1</sup>

#### Correspondence

Qilian Liang, Department of Electrical Engineering, The University of Texas at Arlington, Texas, USA. Email: liang@uta.edu

**Present address**: 701 S. Nedderman Drive, Arlington, Texas 76019, USA.

#### Funding information

U.S. National Science Foundation (NSF), Grant/Award Number: CCF-2219753

### Abstract

With the rapid development of the communication industry in the fifth generation and the advance towards the intelligent society of the sixth generation wireless networks, traditional methods are unable to meet the ever-growing demands for higher data rates and improved quality of service. Deep learning (DL) has achieved unprecedented success in various fields such as computer vision, large language model processing, and speech recognition due to its powerful representation capabilities and computational convenience. It has also made significant progress in the communication field in meeting stringent demands and overcoming deficiencies in existing technologies. The main purpose of this article is to uncover the latest advancements in the field of DL-based algorithm methods in the physical layer of wireless communication, introduce their potential applications in the next generation of communication mechanisms, and finally summarize the open research questions.

### 1 | INTRODUCTION

The fifth generation (5G) technology marks the beginning of a new era in wireless communication, providing unprecedented transmission speeds, reducing latency, and supporting concurrent connections from a multitude of devices [1]. Its emergence has created possibilities for the development of emerging fields such as the Internet of Things (IoT) [2], remote healthcare, and vehicle-to-everything (V2X) [3]. More importantly, 5G is the key to driving innovation and productivity in other sectors, injecting new vitality into broader economic growth.

Despite 5G introducing more flexibility and efficiency to wireless networks through the use of new technologies such as massive multiple-input and multiple-output (MIMO) and millimeter waves, numerous severe challenges still remain [4]. Even though cellular communication systems have advanced to a new level with the development of 5G, they still cannot meet all future requirements by 2030, and hence, researchers have now started to focus on the sixth generation (6G) wireless communication networks [5].

Deep learning (DL) has demonstrated its remarkable power across various domains. For instance, in the field of computer vision, deep learning models facilitate deep analysis of images, making tasks like facial recognition, object detection,

and semantic segmentation possible [6]. In the realm of speech recognition, deep learning has taken the lead, enabling applications such as speech-to-text and voice assistants, playing a vital role in our daily lives [7]. Large language models (LLM), like Generative Pre-training Transformer (GPT), leverage the capabilities of deep learning to deeply understand natural language [8], providing powerful tools for tasks like chatbots, machine translation, and automated text generation. Deep learning data analysis is a powerful technique that can be applied to various domains, including signal propagation in complex networks [9]. It helps extract valuable insights and patterns from large data sets, enabling a better understanding of how signals propagate and interact within these networks. Building on these successes, the applications of deep learning are continuously expanding into more specialized fields. For example, it has been used in the analysis of criminal networks [10], helping in understanding their structure and operations. It also plays a crucial role in EEG signal analysis [11], helping to decode complex brain signals for healthcare and research purposes

DL is also gradually entering the sophisticated physical layer (PHY) of wireless communication, to handle the optimization of multiple performance objectives while adhering to a range of intricate restrictions [12]. Traditionally, signals flow from optimally designed transmitters with modulation, coding, and

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. IET Communications published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology.

IET Commun. 2023;17:1863—1876. wileyonlinelibrary.com/iet-com

<sup>&</sup>lt;sup>1</sup>Department of Electrical Engineering, The University of Texas at Arlington, Texas, USA

<sup>&</sup>lt;sup>2</sup>AutoX, Inc., San Jose, California, USA

signalling schemes, through a series of mathematically defined channel models, to be reliably detected in the receivers, with each block being separately optimized. Existing results in the PHY suggest that DL can aid in understanding wireless content, identifying undiscovered patterns, reducing complexity, and generating results equivalent to conventional methods, even surpassing them in some cases [13].

The objective of this paper is to present a thorough summary of current research on PHY, with an emphasis on the potential benefits and challenges of DL-based wireless communication systems. This paper provides a blueprint for future inquiries by elaborating on the driving forces, proposed tactics, achieved results, and limitations of these studies. Figure 1 illustrates the structure of this article. The review makes important contributions in these key aspects:

- 1. This article initially highlights DL methodologies that tackle the obstacles in MIMO detection, channel estimation, channel coding/decoding, and resource allocation aiming to counter the limitations of conventional techniques. As shown in Figure 1, the application of DL is discussed within the context of two types of receiver frameworks, namely, joint symbol detection and channel estimation, and joint equalizing and decoding.
- 2. Besides the above existing block in wireless communication, we also focus on two recently developed wireless communication technologies: reconfigurable intelligent surfaces (RISs) and MIMO-based index modulation (MIMO-IM). We believe that applying DL techniques to these areas could lead to significant future advancements in the field.
- 3. Contemporary research has zeroed in on transceiver designs, optimizing the entire transmitter and receiver pipeline by leveraging the end-to-end (E2E) approach, which is being championed as a promising trajectory. We focus on the application of an E2E approach based on DL to overcome the weaknesses of orthogonal frequency-division multiplexing (OFDM), as well as its role in semantic communications.

The rest of this paper is organized as follows: Section 2 provides an overview and analysis of prior studies conducted by researchers in the field. In Section 3, there is a brief discussion on the structure of DL, along with an introduction to basic DL techniques. Section 4 presents several examples of using DL as alternatives for wireless communication systems. In Section 5, two applications of DL in emerging wireless communication technologies are present. Section 6 summarizes DL-based E2E communications. Section 7 introduces DL solutions for wireless platforms. Section 8 focuses on potential areas for future research. Section 9 provides a conclusion to the paper.

#### 2 | RELATED WORK

Here, the researchers have conducted a comprehensive review of survey articles that discuss the application of DL in PHY of wireless communication.

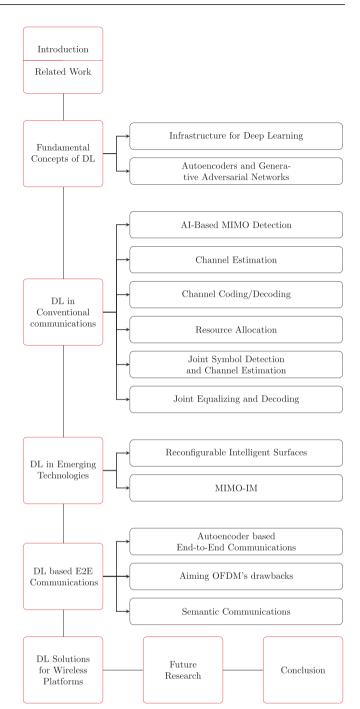


FIGURE 1 Organization of the article.

Wang et al.[14] have presented an overview of development and investigation in DL-based PHY processing. This includes re-engineering conventional system modules (e.g. modulation recognition, channel decoding, and detection) using DL, and going as far as supplanting the traditional communication system with a revolutionary autoencoder-based architecture. However, this paper was published in 2017 and cannot address rapid changes since then.

Similarly, Qin et al.[12] classified the applications of DL in PHY communications into two categories: systems with block structures and systems without block structures. Nevertheless, it

is worth noting that this paper was published in 2019. Additionally, concerning systems with block structures, they only selected a few blocks and did not provide a comprehensive overview.

Mao et al.[15] provided a thorough examination of the applications of DL algorithms across various network layers. It covered PHY modulation/coding, data link layer access control/resource allocation, routing layer path search, and traffic balancing. However, this paper only discussed a few application scenarios for PHY, including interference alignment, jamming resistance, modulation classification, and physical coding. It did not delve into in-depth discussions on 5G.

In this insightful review [16], it is noted that the recent surge of research in AI-driven communication technologies holds great potential for enhancing data rates and elevating QoS while keeping the implementation costs manageable. This paper astutely encapsulates the cutting-edge advancements in AI-integrated 5G and beyond 5G (B5G) techniques, examining them at the algorithmic, implementation, and optimization stages. But this article restricted the topic to the existing 5G framework and did not discuss possible 6G technologies.

Ozpoyraz et al.[17] shined a spotlight on the latest developments in the realm of DL-based PHY methods, aiming to catalyze the extraordinary potentials of 6G applications. Specifically, four avant-garde PHY concepts, poised to revolutionize next-generation communications, are expertly dissected: massive MIMO systems, intricate multi-carrier waveform designs, RIS-empowered communications, and PHY security. Through this exploration, Ozpoyraz et al. navigated the reader towards understanding the future trajectory of this rapidly evolving field. However, that article fell short in effectively summarizing and consolidating existing papers. It could be seen more as an aggregation of information rather than a comprehensive analysis.

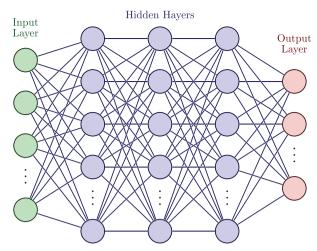
In essence, the prior literature reviews are lacking in three significant areas:

- Many have overlooked the inclusion of newly published research.
- 2. A substantial number of reviews have neglected the consideration of advancements in 5G or 6G technologies.
- 3. There is a noticeable deficiency in the way some articles summarize the information accurately and comprehensively.

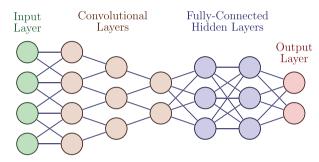
The driving force behind this research hinges on these observed gaps. It seeks to address these insufficiencies to provide a more thorough and up-to-date examination of the subject matter.

# 3 | FUNDAMENTAL CONCEPTS OF DEEP LEARNING

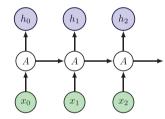
Here, we arere going to provide a brief overview the general structure of three foundational deep learning models: deep neural network (DNN), convolutional neural network (CNN), and recurrent neural network (RNN). Following this, we will explore autoencoders and generative adversarial networks (GANs).



(a) A DNN utilizes hidden layers to extract meaningful features between the input and output layers.



(b) A CNN architecture begins by incorporating convolution layers prior to dense layers.



(c) An RNN architecture incorporates the extracted features from the previous state as part of the current input information.

FIGURE 2 Three common deep learning architectures.

# 3.1 | Infrastructure for deep learning

DNNs are sophisticated machine learning structures composed of interconnected layers. Every individual neuron, or computational unit, within a particular layer is intricately interlinked with all the neurons of its adjacent layers, be it the preceding one or the subsequent one (Figure 2a). This comprehensive interconnectivity establishes an extensive and fully integrated structure, a distinctive characteristic of DNNs that sets them apart from simpler models.

Each layer plays a unique role within the overall network: the input layer, intermediary hidden layers, and the output layer. The input layer corresponds directly to the raw input data,

TABLE 1 List of activation functions.

Name	[σ(u)] <sub>i</sub>	Range
Linear	$u_i$	$(-\infty,\infty)$
ReLU	$\max(0, u_i)$	$[0, \infty)$
Tanh	$tanh(u_i)$	(-1, 1)
Sigmoid	1 1+a=ni	(0,1)
Softmax	$\frac{1+e^{-n_i}}{\sum_j e^{n_j}}$	(0,1)

with the number of neurons being equivalent to the number of features presented in the data set. Similarly, the output layer's size matches the number of categories, classes, or other targets depending on the specific task the DNN is designed to perform. Thus, the structure of a DNN is inherently adaptive and can flexibly accommodate various application needs.

One of the crucial components that enable DNNs to learn and adapt is the presence of real-valued weights. Each connection, or link, between two neurons within the network is represented by such a weight. These weights encapsulate the accumulated knowledge of the network, adjusting as the model learns from the data it is exposed to. Therefore, the learning mechanism of DNNs is fundamentally rooted in the continual adaptation and refinement of these connection weights.

Activation functions are crucial to DNNs. They transform the weighted input into the neuron's output, adding non-linearity to the network which helps it learn complex patterns. Common activation functions include rectified linear unit (ReLU), tanh, sigmoid, and softmax. Each has specific benefits and is suitable for different scenarios. Choosing the right function depends on the problem and data characteristics. Table 1 summarizes the commonly used activation functions.

CNNs are a special type of DNN designed to counteract the explosion of parameters that can occur in traditional fully connected networks, particularly when applied to complex tasks like image recognition. CNNs are customized architectures designed to suit the specific needs of different scenarios. The underlying principle of a CNN involves introducing convolutional and pooling layers before the data reaches a fully connected network (Figure 2b). In a convolutional layer, the neural connections are localized, which means that each neuron is connected only to a subset of neurons from the preceding layer. Neurons are arranged in a grid-like formation, creating feature maps, each identifying different features within the data by applying a shared set of weights across the entire input. Pooling layers follow the convolutional layers in a CNN. The role of these layers is to reduce the dimensionality of each feature map while preserving the most important information. This is achieved by grouping neurons in each feature map and then calculating either the average value (average pooling) or the maximum value (max pooling) of each group. This action drastically reduces the number of parameters, making the neural network more manageable.

RNNs have been designed to imbue neural networks with the capability of memory. This property is key for handling sequential data such as in natural language processing (NLP) where the

context plays a significant role. In more traditional, memoryless neural networks, the neurons in each layer are connected only to those in the preceding and subsequent layers, with no intralayer connections. However, this architecture does not provide the network with the ability to maintain and utilize any contextual or sequential information from prior states, which could limit its performance on tasks that inherently require knowledge about prior inputs. RNNs address this limitation by incorporating feedback connections in the hidden layers (Figure 2c). This means that the neurons in a given layer receive not just inputs from the preceding layer, but also the outputs of their own layer from previous steps. In essence, a recurrent neuron maintains a kind of memory by using its output from the previous step as part of its input for the current step. This allows the network to 'remember' and use information from the past, effectively enabling it to handle data where temporal dynamics and dependencies matter. Various types of RNNs have been proposed to address different challenges and usecases. Bidirectional RNNs process data from both ends to the middle, providing more context by considering both past and future data, which can enhance performance in tasks like NLP. Long short-term memory (LSTM) networks incorporate 'gates' in their structure, effectively controlling the flow of information and making them proficient in learning long-range dependencies in sequential data. Gated recurrent units (GRUs) simplify the LSTM structure, maintaining its ability to mitigate the vanishing gradient problem but with fewer parameters, improving computational efficiency.

# 3.2 | Autoencoders and generative adversarial networks

An autoencoder (AE) is an unsupervised learning algorithm [18]. It is a neural network that learns to generate its output which is almost close to its input.

As shown in Figure 3, an AE consists of two parts, the encoder, which computes a latent representation of the input, and the decoder, which reproduces the original input from the latent representation. Define the encoder parameters as  $\phi$  and the decoder parameters as  $\psi$  [19], that is,

$$\phi: \mathcal{X} \to \mathcal{F} \tag{1}$$

$$\psi: \mathcal{F} \to \mathcal{X} \tag{2}$$

$$\phi, \psi = \underset{\phi, \psi}{\arg \min} \|X - (\psi \circ \phi)X\|^2$$
 (3)

where  $\mathcal{X}$  is the data space,  $\mathcal{F}$  is the latent space and  $X \in \mathcal{X}$ .

Without loss of generality, we assume that the data and the latent spaces are real valued with dimension d and p, respectively. The encoder takes the input  $\mathbf{x} \in \mathbb{R}^d = \mathcal{X}$  and maps it to  $\mathbf{h} \in \mathbb{R}^p = \mathcal{F}$ :

$$\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \tag{4}$$

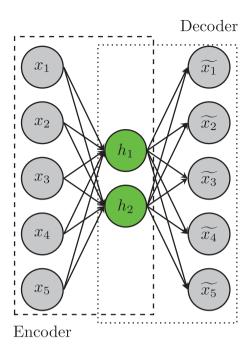


FIGURE 3 An autoencoder is a neural network architecture that learns to encode and decode data.

where  $\sigma$  is an activation function, **W** is a weight matrix, and **b** is a bias vector. The decoder maps **h** to the reconstruction  $\mathbf{x}'$ :

$$\mathbf{x}' = \sigma'(\mathbf{W}'\mathbf{h} + \mathbf{b}') \tag{5}$$

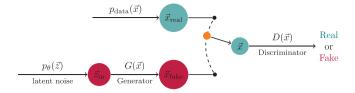
where  $\sigma'$  is an activation function,  $\mathbf{W}'$  is a weight matrix, and  $\mathbf{b}'$  is a bias vector.

When  $\mathbf{x} \approx \mathbf{x}'$ , it is considered that the trained AE reconstructs the input. The cost function could be defined as follows:

$$\mathcal{L}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2 + \sum_{i,j,k} \left(\omega_{jk}^{(i)}\right)^2$$
 (6)

where  $\omega_{jk}^{(i)}$  is the connection weights between the *j*th neurons of layer *i* and the *k*th neurons of layer i + 1. The first term (mean squared errors) is the reconstruction error, and the second term (weight decay) is a regularizing penalty which is typically included to prevent overfitting.

A (GAN can be compared to an inverted autoencoder in terms of its structure [20]. However, in contrast to autoencoders which shrink input data, GANs transform low-dimensional input into complex, high-dimensional data using its internal network. GANs, utilizing two rival neural networks (hence the term adversarial), are engineered systems capable of generating novel, synthetic data that can often be mistaken for authentic data. As depicted in Figure 4, a GAN comprises of two components. The discriminator, denoted as D, functions as a classifier, distinguishing between genuine data and synthetic data created by the generator, denoted as G. If D identifies any unreal outcomes, G receives a penalty. This signal is then channeled back through the generator, modifying its weights such that G



**FIGURE 4** A generative adversarial network consists of two neural networks, a generator and a discriminator, competing with each other to produce and evaluate realistic synthetic data.

gradually learns to create increasingly realistic samples. If the training is successful, *G* eventually manages to deceive *D*.

# 4 | DEEP LEARNING IN CONVENTIONAL COMMUNICATIONS

As shown in Figure 5, a conventional wireless communication system typically features various blocks such as source encoding and decoding, channel encoding and decoding, modulation and demodulation, channel estimation, equalization and detection, along with RF transceiving. Each of these signal processing blocks is fine-tuned individually to ensure secure and reliable communication from the originating source to the intended target destination.

# 4.1 | AI-based MIMO detection

In large-scale MIMO, employing the maximum a posteriori (MAP) detector, which offers the best detection performance, is infeasible due to its exponential computational complexity. Therefore, linear detectors such as the matched filter (MF), zero forcing (ZF), and linear minimum mean square error (LMMSE) have been developed. These have lower complexity but display inferior performance when compared to the MAP detector. Additionally, there are iterative detection algorithms like approximate message passing (AMP), sphere decoding (SD), and soft interference cancellation (SIC), which can achieve good performance under certain conditions and have moderate complexity. All these detectors require a full understanding of the channel state information (CSI). If the system model does not match the actual transmission model or imperfect CSI is present, the performance will significantly degrade [21].

Data-driven DL detectors using DNN architectures can recover transmitted symbols in various scenarios with high precision [22, 23], albeit at the cost of a large amount of trainable parameters and training samples. Reference [21] proves that the data-driven DL detector with a ReLU DNN can effectively approximate the MAP detector. The rate of convergence of the DL detector to the MAP detector scales at least polynomially fast with the size of the training samples. Furthermore, these detectors are robust to CSI uncertainty.

Model-driven DL detectors, which have evolved from traditional iterative detection algorithms such as DetNet [24] and OAMPNet [25], often yield detectors with superior

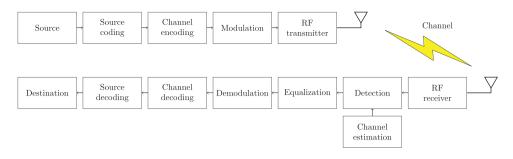


FIGURE 5 A conventional wireless communications system diagram consists of multiple blocks. Each of these signal processing blocks is individually optimized to ensure dependable communication between the source and the intended destination.

performance and faster convergence. Each layer of their networks represents an iteration and incorporates some trainable variables. However, current model-driven DL detectors are premised on the availability of linear channel models and CSI, which limits their application in complex environments.

### 4.2 | Channel estimation

Accurate channel estimation is essential for optimizing the quality and reliability of a wireless communication link. Typically the receiver uses the known pilots to estimate the channel. Conventional pilot-based channel estimation methods include least squares (LS) and LMMSE techniques. LS is simple and efficient but sensitive to noise, while LMMSE provides improved performance by considering statistical properties at the cost of increased computational complexity.

In reference [26], channel matrix is treated as a 2D image. The authors utilize a denoising CNN (DnCNN) as the denoiser within the learned denoising-based approximated message passing (LDAMP) network. With a single pilot, the LDAMP network learns the channel's structure. The LDAMP network surpasses its conventional cousins. Similar approaches have also been employed in reference [27].

A deep learning-based joint pilot design and channel estimation scheme is presented in reference [28]. The proposed channel estimator is structured into two stages. In the first stage, two DNNs collaborate to handle pilot design and pilot-aided channel estimation. The second stage employs another DNN to iteratively improve the channel estimation and symbol detection, referred to as data-aided channel estimation, thereby enhancing the overall estimation performance.

# 4.3 | Channel coding/decoding

In the standard of 5G, both low-density parity-check (LDPC) codes and polar codes play significant roles. The LDPC codes are decoded by the belief-propagation (BP) algorithm, a robust message passing method that offers near-optimal error-rate performance through iterative decoding [29]. On the other hand, polar codes have the unique ability to asymptotically achieve the Shannon capacity, utilizing the successive cancellation (SC)

decoding algorithm, especially as the code length approaches infinity [30]. Current research in both SC and BP decoding is focused on striking a balance between error-rate performance and complexity, thereby advancing the efficiency and effectiveness of these communication systems.

Exploring the enhancement of BP/SC decoders through the application of neural networks is one direction. In reference [31], the authors improved the BP decoding algorithm for LDPC codes through deep learning. They used a Tanner graph and trained these weights using stochastic gradient descent, leading to a significant reduction in the bit error rate (BER). In reference [32], a DNN is utilized to estimate the least number of iterations necessary. The study's simulations suggest that, especially in high SNR scenarios, the DNN can precisely predict the required iteration count. In reference [33], the introduction of an iterative BP-CNN decoder for tackling correlated noise is suggested. This involves the integration of the BP decoder's output into a feed-forward CNN to ascertain the correlated noise across different channels.

Another exploration involves utilizing neural networks as a direct substitution for the decoder. One such study [34], employs a neural network decoder (NND) to decode both unstructured (LDPC and HDPC) and structured codes (polar codes). Simulations revealed that the NND is capable of achieving MAP performance for short block lengths, whether for structured or unstructured codes. Building upon the findings of reference[34], another study, referenced as [35], took the research a step further. In this study, the researchers developed a method to break down a long block length polar code into manageable sub-blocks. Each of these smaller pieces is then decoded using a compact NND, making the training process more feasible. The results obtained from the NND are subsequently propagated through a traditional BP structure.

In addition to enhancing BP/SC decoders through neural networks and substituting decoders directly with neural networks, another direction to consider is the realm of code construction. Grounded in genetic algorithm, studies [36] and [37] have optimized the design of LDPC and polar codes, respectively. They considered factors such as channel conditions, code length, and the number of iterations in their optimization process. A more detailed discussion on the integration of AI-aided encoding and decoding is set to be presented in the subsequent E2E section.

### 4.4 | Resource allocation

In 5G and future wireless communications, a key task affecting network performance is how to improve spectral efficiency and system energy efficiency through optimized resource allocation strategies, while ensuring quality of service (QoS) [38]. Although the resource allocation problem can obtain closed-form expressions in some cases [39], in most situations, the optimization problem is non-convex. Therefore, deep learning has become a promising approach.

In reference [40], a proposal is presented that focuses on the core idea of viewing the input and output of a resource allocation algorithm as an unknown non-linear mapping. This method employs a DNN as a means for its approximation. In reference [41], the deep power control (DPC) is introduced as the first-ever transmit power control framework based on a CNN. Within this framework, the CNN is trained to learn a transmit power control strategy aiming to maximize either spectral efficiency (SE) or energy efficiency (EE).

With its ability to control non-orthogonality, non-orthogonal multiple access (NOMA) has been shown to employ resources more efficiently than traditional methods [42]. Furthermore, there is an emerging trend of deep learning-enhanced NOMA systems with resource allocation. In reference [43], deep belief network (DBN) is used in simultaneous wireless information and power transfer (SWIPT) and multi-carrier NOMA (MC-NOMA). The objective is to minimize total transmit power while meeting each user's QoS requirements. In reference [44], constrained deep reinforcement learning (CDRL) is utilized to investigate the complex issue of simultaneous multi-UAV altitude control and random channel access management within a multi-cell UAV-based wireless network employing NOMA.

# 4.5 | Joint symbol detection and channel estimation

Sections 4.1 and 4.2 discuss symbol detection and channel estimation, respectively. Some attempts have already begun to explore the joint channel estimation and signal detection (JCESD).

The authors of reference [45] introduce the DeepSM network, comprising two synchronized DNNs, and compare it with traditional receivers in diverse channel contexts. Using the LS method, DeepSM updates the CSI and detects transmitted symbols. Initially, a conventional DNN is used for channel estimation and signal detection in time-invariant spatial modulation (SM) systems. This evolves into the DeepSM structure for systems in time-varying fading channels. Both model-based and DNN-based receivers apply LS channel estimation, but use different methods for detecting transmitted symbols.

In reference [46], two novel DL-based receiver structures, FullCon and MdNet, are introduced for uplink multi-user MIMO systems. FullCon, a data-driven algorithm employing a DNN, detects information bits directly from the received

signal, bypassing explicit channel estimation. It acts as a comprehensive MIMO receiver, amalgamating a channel estimator, signal detector, and demodulator in a fully connected deep learning network. Contrastingly, MdNet employs a model-driven approach, fusing traditional communication knowledge with DL to perform channel estimation and symbol detection in distinct phases.

Reference [47] introduces a novel approach to modelling environmental noise using the Student's t-distribution rather than a Gaussian distribution, enhancing robustness against outliers. This complex model is made tractable through a generalized EM (GEM) algorithm which jointly and robustly estimates the channel matrix and transmitted signals, showing superior performance compared to conventional independent methods. In order to decrease the computational load of the GEM algorithm, it is unfolded into a DNN for fluctuating channels, utilizing a modified trainable projected gradient (TPG) detector in the M-step. The TPG detector, initially crafted for Quadrastic Phase Shift Keying (QPSK) modulation, is further adapted for high-order modulations. The unfolded GEM network, requiring fewer iterations, surpasses the original GEM algorithm's performance.

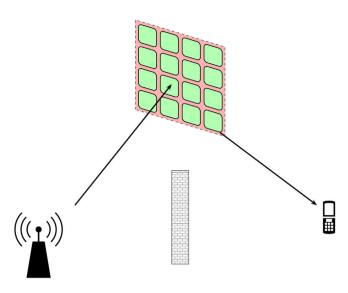
# 4.6 | Joint equalizing and decoding

Channel equalization is used to address the problem of potential performance degradation in wireless communication systems due to factors such as wireless channel impairments. These impairments include inter-symbol interference (ISI), Doppler frequency shift, and fading effects. A variety of neural network variants are utilized to enhance the performance of equalization, including the multi-layer perceptrons (MLP) [48], function-link artificial neural networks (FLANN) [49], radial-based function (RBF) neural networks [50], and (RNNs [51].

Considering both equalization and decoding jointly has emerged in addressing a notable issue. This issue pertains to the fact that after applying neural network-based equalization, the noise present in the system tends to deviate from Gaussian characteristics. Consequently, this non-Gaussian noise can lead to a decline in the performance of the channel decoder. Reference[52] employs a DNN for joint channel equalization and decoding by the received symbols as input and the estimated codeword as output. This method outperforms traditional methods for polar code (16,8). To address the issues of flexibility and complexity codes by DNN, several other neural networks have also been introduced. A joint training of CNN equalizer and DNN decoder is discussed in reference [53], while reference [54] utilizes RNN-based joint equalization and decoding.

### 5 | DL IN EMERGING TECHNOLOGIES

Recently, several new wireless communication technologies have appeared. This section focuses on two notable ones, RISs and MIMO-IM, exploring how DL can be applied



**FIGURE 6** Illustrations of an reconfigurable intelligent surface (RIS)-aided downlink transmission. RIS enhance wireless communication by manipulating electromagnetic waves through programmable surfaces.

within these technologies, potentially opening avenues for future advancements.

# 5.1 Reconfigurable intelligent surfaces

Driven by limited communication spectrum, a shift towards higher frequency bands occurs. This, however, raises issues with electromagnetic wave obstructions, particularly in urban areas. Adding more relays and base stations increases energy use, making traditional cellular methods challenging. To mitigate this, RISs have emerged as a key solution, improving both spectrum and energy efficiency in wireless networks [55]. An RIS-aided downlink transmission is illustrated with no direct connectivity as depicted in Figure 6.

Primarily built using metamaterials, which consist of periodically aligned subwavelength elements, RIS models provide full control over the electromagnetic actions of the metasurface, and they interact intelligently with incoming signals to enhance energy efficiency and coverage in radio communication systems. Applications of DL in RIS almost cover all modules of wireless communication. Here are some examples on the beamforming design and resource allocation.

In reference [56], the authors propose a dual-phase neural network featuring an unsupervised learning approach aimed at addressing the joint passive and active beamforming design in RIS-aided multi-user MISO downlink platforms. The objective is to effectively resolve the complex problem of joint optimization.

Reference[57] introduces a DL-assisted RIS scheme utilizing the deep deterministic policy gradient (DDPG) method, which is a fusion of the deep-Q network and policy gradient (PG). The scheme leverages a continuous action space, which hastens the training phase. This model, thanks to the DDPG algorithm's continuous action space, exhibits a remarkable ability to swiftly

adapt to fluctuations in channel data and environmental conditions. In reference [58], a joint optimization problem in a NOMA downlink network using RIS is addressed. The development of RIS phase shifts is guided by the application of a DDPG algorithm. This model uses a reward function based on the sum rate of mobile clients, aiding the agent to identify the optimal path. The system's performance can be enhanced by adjusting the number and complexity of the RIS's reflecting elements. In reference [59], the authors addressed the joint vehicle scheduling and passive beamforming in RIS-empowered vehicular communication by employing a DRL framework with a multi-binary action space, a strategy proposed to maximize the minimum average bit rate for vehicles using wireless scheduling.

#### **5.2** | **MIMO-IM**

Index modulation (IM) pioneers a unique approach to next-generation communication, diverging from standard amplitude-phase modulation [60]. By engaging certain transmit antennas or time slots to form unique activation patterns, it enables a high-dimensional modulation scheme. This approach boosts spectral efficiency under suitable configurations. Additionally, the index of activation patterns can carry extra bit streams, supplementing those modulated by constellation symbols, improving overall data transmission efficiency.

The advancement of deep learning has opened up opportunities for various applications in the field of index modulation. Specifically, three main areas of focus have emerged: transmit antenna selection (TAS) and power allocation [61], modulation and coding scheme selection [62], and detection. While the latter two approaches share similarities with those discussed in previous sections, the following section will specifically concentrate on the application of deep learning in TAS and power allocation.

In reference [61], the researchers initially tackle the challenges of TAS and power allocation (PA) in SM-MIMO by transforming them into data-driven prediction problems instead of relying on traditional optimization-driven decisions. To achieve this, they develop supervised-learning classifiers (SLC) like the Knearest neighbours (KNN) and support vector machine (SVM) algorithms, which provide statistically consistent solutions. Additionally, they explore the integration of DNNs with these adaptive SM-MIMO techniques and propose a novel DNN-based multi-label classifier to evaluate TAS and PA parameters. Moreover, they investigate the creation of feature vectors for the SLC and DNN approaches and introduce a unique feature vector generator specifically tailored to the transmission mode of SM.

In reference [63], the generalized TAS pipeline is formulated in both neural networks (NN) and gradient boosting decision trees (GBDT). In this formulation, the importance of different features that reflect the different elements from CSI is analysed, taking into consideration the empirical data as well. The results confirm that both GBDT and NN are capable of achieving

a near-optimal BER curve, with the former demonstrating an even better balance between efficiency and performance. A similar TAS-GSM scheme based on DNN and decision tree could be found in reference [64].

# 6 | DEEP LEARNING-BASED END-TO-END COMMUNICATIONS

In the previous section, we discussed several deep learning-based approaches that can replace one or two processing blocks in the conventional communication systems. Optimizing each processing block individually does not guarantee the optimal solution for the entire communication problem. This is because the performance of a communication system is influenced by multiple factors, including channel characteristics, noise, interference, and so on. Therefore, optimizing each block alone does not ensure the best overall optimization for the communication problem.

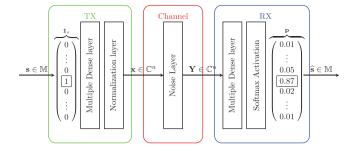
Recently, a new concept based on deep learning has been proposed, which makes significant improvements over existing ideas. This concept redefines the communication task as an E2E reconstruction optimization task, eliminating the modular structure that needs to be manually constructed in traditional communication systems. This novel concept is based on using an autoencoder system to achieve an E2E communication system. Initial studies have shown that this approach is comparable to traditional systems in terms of performance and demonstrates the enormous potential of E2E methods in becoming a general solution for various channel models.

Here, we will explore the emerging concept of E2E communications and provide an overview of some recent studies in this field.

# 6.1 | Autoencoder-based end-to-end communications

Autoencoder (AE) is a neural network that learns to generate its output which is almost close to its input. The concept of utilizing AE in the field of communication systems was initially introduced in reference [13]. In the communications systems, it is possible to view it as an AE that aims to minimize the error by reconstructing the transmitted message at the receiver. The encoder and decoder of the AE can be analogously considered as the transmitter and receiver blocks of the system, respectively. Figure 7 illustrates an AE architecture that is used in E2E learning of a communications system.

As shown in Figure 7, at the transmitter, the incoming message s is mapped to an  $M \times 1$  one-hot vector, which is used as an input vector of the encoder, where S is the set of all  $M = 2^m$  possible messages, each having m data bits. Then, this M-dimensional vector is passed through a feed-forward neural network (NN) with multiple dense layers, followed by a normalization layer to satisfy the physical constraints of the transmit vector  $\mathbf{x}$ . The output of the transmitter is a complex vector of dimension m. The channel is represented by a noise layer.



**FIGURE 7** An autoencoder-based end-to-end communications. The input signal  $\mathbf{s}$  is encoded into a one-hot representation, resulting in transmitted signal  $\mathbf{x}$ . After introducing noise to this encoded signal, it is decoded. Finally, the original signal  $\hat{\mathbf{s}}$  is determined.

The receiver also consists of a feed-forward neural network (NN) with one or multiple dense layers, followed by an output layer with a softmax activation function. The output layer yields a probability vector, denoted as  $\mathbf{p} \in (0,1)^M$ , representing the probabilities of all possible messages. The decoded information,  $\hat{\mathbf{s}}$ , corresponds to the index of the element in  $\mathbf{p}$  with the highest probability.

As for the training aspect, the autoencoder utilizes stochastic gradient descent or any other suitable optimization method to perform E2E training on the set of all possible information  $s \in M$ . The reconstruction quality of the autoencoder is measured using the categorical cross-entropy loss function.

Based on the foundation of reference [13], a series of studies have emerged. For instance, in the field of constellation design, reference [65] utilizes an AE for the design of optimal constellations and receiver architectures in AWGN channels affected by additive radar interference. Reference [66] discusses the application of AE-based constellation design in a multi-user interference channel to tackle issues related to dynamic interference.

Others have delved further into the research on whether the channel model is known, dividing into model assumed channel AE and model free channel AE. Moving into the domain of model assumed channel AE, there appears to be a couple of relevant studies. In reference [67], a model based on variational autoencoder (VAE) was introduced. This model incorporated prior knowledge about the channels into its cost function, which effectively mitigated the impact of noisy latent codes. Furthermore, this innovative approach to noise reduction significantly enhanced the speed of training. In reference [68], interference in multi-path settings is investigated using random channel parameters. The study further determines the fewest channel samples required for optimal performance via the confidence interval method, paving the way for a robust encoding and decoding scheme. For model free application, an attempt to resolve this particular problem is documented in reference [69] via the implementation of a two-phase training strategy. Initially, the entire system undergoes training based on a preconceived channel model. Subsequently, fine-tuning of the receiver occurs over the actual channel, which serves to correct any discrepancies that may arise.

In addition to the AE-based E2E learning systems, there have been compelling advancements in the implementation of RL and conditional GAN-based methodologies in E2E communication systems as well. Similar to AE, the essence of both RL and conditional GANs lies in training a model to produce an output that aligns as closely as possible with the desired result [12]. For example, in reference [70], the authors suggest employing a conditional GAN as a tool to mimic channel effects. This approach serves as a conduit linking the transmitter DNN and the receiver DNN. The primary advantage of this methodology lies in enabling the backpropagation of the transmitter DNN's gradient from the receiver DNN.

# 6.2 | Aiming OFDM's drawbacks

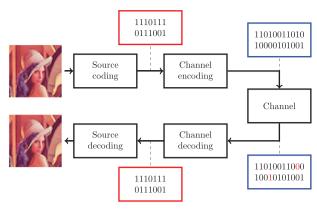
OFDM emerges as the most prevalent multi-carrier waveform [71], finding widespread use in many standards like the IEEE 802.11 family and 5G. Its adoption is largely due to its straightforward and efficient structure, boosting the performance of wireless communication. However, OFDM has drawbacks. It suffers from high peak-to-average power ratio (PAPR), cyclic prefix (CP) overhead, and pilot overhead. The following will provide a summary of the research literature focusing on AE-based techniques aimed at addressing the primary limitations of OFDM.

In reference [72], a DL-oriented scheme for PAPR reduction in an OFDM system, called PRNet, is introduced. Utilizing AE, this proposed method demonstrates the ability to reduce PAPR, simultaneously preserving the BER. Similar structures could be found in references [73] and [74].

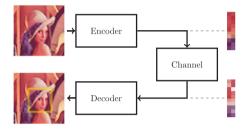
In reference [75], the authors first demonstrate a partial reduction in the number of pilot symbols without a decrease in error performance. Then, they utilize an AE-based Neural Network to completely eliminate pilots, while learning a constellation or superimposed pilots simultaneously. They also demonstrated that it is entirely possible to eliminate both CP and pilot by leveraging E2E learning in reference [76]. In reference [77], a pilot-free E2E paradigm, developed for various wireless channels including frequency-selective and flat-fading MIMO channels, utilizes a model wherein the wireless channels are seen as a stochastic convolutional layer. This system has shown effectiveness under diverse channel conditions, particularly demonstrating its capacity to conserve pilot resources and leverage the correlation in wireless channels and source data.

### 6.3 | Semantic communications

Semantic communications is a concept focused on understanding meaning in data transmission, rather than the literal data itself. It applies to interactions among humans, and between humans and computers, ensuring accurate interpretation of the sender's intentions despite potential inaccuracies or ambiguities [78]. Although semantic communication is not a new



(a) Traditional Communications Systems. After passing through a channel, the signal will experience errors, but they can be corrected through error correction, allowing for normal decoding.



(b) Semantic Communications Systems. After passing through the channel, the signal will experience errors, but facial recognition functionality can be achieved through semantic communications.

FIGURE 8 Comparison of two communications systems.

research topic [79], with the advancements in micro-electronics and AI technologies, coupled with the evolution of DL and E2E approaches, it has once again become one of the emerging communication paradigms [80]. Figure 8 demonstrates the concept of semantic communications, specifically focusing on face recognition.

In text-based semantic communication, DL has been proposed for use in joint source-channel coding (JSCC) due to the powerful representation capabilities of DNNs. Notably, research in reference [81] have advanced a JSCC scheme specifically for text-based semantic communications. In this scheme, the encoder and decoder are implemented by two RNNs, with a dropout layer representing the channel. This was largely inspired by the successful application of DNNs in NLP. In reference [82], enabled by intelligent E2E communications, a novel framework for semantic communication systems has been proposed, aiming to design a JSCC scheme to maximize system capacity. A significant feature of this system is its utilization of a transformer and self-attention mechanism, which notably simplifies the comprehension of long sentences at the destination.

In image-based semantic communication, CNNs offer significant advantages. The first CNN-based E2E JSCC scheme was introduced in reference [83], demonstrating its potential for transmitting high-resolution images over both AWGN and Rayleigh channels. This scheme notably outperforms traditional compression algorithms like JPEG and JPEG2000, showcasing

superior performance without succumbing to the "cliff effect" often seen in these conventional methods.

In other scenarios of semantic communication, such as those based on sound and multimedia, the exploration is still in the initial stage due to factors such as complex dimensions and difficulty in evaluation.

# 7 | DEEP LEARNING SOLUTIONS FOR WIRELESS PLATFORMS

There are numerous approaches for implementing DL, and some of these methods have been applied in PHY of wireless networks. In the following, we provide a summary of DL implementations that have been employed in wireless communication.

- MATLAB: The Neural Network Toolbox of MATLAB is a versatile software package for designing, training, and analysing neural networks. It offers a user-friendly graphical interface, allowing interactive network design and configuration. The toolbox includes a variety of training algorithms and supports different types of neural networks for tasks such as classification, regression, and clustering. It also provides tools for data preprocessing, feature selection, and performance evaluation. The integration with MATLAB enables seamless data manipulation and statistical analysis. With its extensive capabilities and ease of use, the MATLAB Neural Network Toolbox is a popular choice among researchers, engineers, and data scientists working on machine learning applications.
- TensorFlow [84]: It is an open-source library for numerical computation and machine learning. Developed by Google, it offers a flexible and scalable framework for building and deploying machine learning models. With its high-level application programming interface (API), Keras, users can easily construct and train complex neural networks. Tensor-Flow also provides a low-level API for fine-grained control over model architecture and training. It excels in leveraging hardware acceleration, running computations efficiently on central processing units (CPUs), graphics processing units (GPUs), and specialized hardware like tensor processing units (TPUs). TensorFlow's extensive ecosystem includes pre-built models, tools, and integration with other popular libraries. Its distributed computing capabilities enable training large-scale models across multiple machines, making it a powerful and widely used tool in the field of machine learning.
- PyTorch [85]: It is a popular open-source machine learning library for Python, developed by Facebook's AI Research lab. Its main features include tensor computation with strong GPU acceleration support and DNNs built on a tape-based autograd system, which allows flexible creation and modification of computational graphs. PyTorch's imperative programming model offers an intuitive and interactive experience, making it an excellent tool for both research and application development. Furthermore, its extensive ecosystem, inclusive of tools and libraries such as TorchServe and

- TorchVision, facilitates model serving and computer vision tasks, enhancing the overall ML workflow.
- Sionna [86]: It is a TensorFlow-based open-source Python library for wireless communication simulations and system design. It provides comprehensive toolsets for defining, optimizing, and executing complex simulations such as multi-dimensional tensor operations, forward error correction, and E2E system models. Users can use Sionna to implement and simulate neural network models on high-performance GPU architectures. The system model is built using various components that represent mathematical operations and data tensors. The library includes a wide range of tutorials for both beginners and experts, making it a versatile tool for research and development in wireless communication systems.

The current literature on the application of deep learning in the PHY of wireless communications is mostly based on the aforementioned architecture. However, there have also been attempts that explore architectures such as Caffe [87] and Theano [88] in this field.

### 8 | FUTURE RESEARCH TRENDS

As previously discussed, the global implementation of 5G technology has achieved substantial progress, and rigorous exploration into 6G technology has been greatly stimulated. Undoubtedly, artificial intelligence will play a pivotal role in the inevitable revolution spurred by 6G. Even though past research has presented promising outcomes, extensive challenges still merit further investigation in the future.

- 1. Performance metric: Currently, in DL-based wireless communication systems, especially in E2E scenarios, the key performance indicators primarily focus on block error rate (BLER) and BER. The goal is to recover as much accurate information from the transmitted signals as possible. Additional performance metrics like latency and power may also be considered. To effectively handle multiple metrics, a constraint-based training strategy can be employed. This allows for dynamic management of the trade-offs among these metrics. In certain cases, such as semantic communications, all data transmission is not considered equally important. The recovered data may contain transmission errors, yet the semantic information within the data must remain intact. This necessitates the development of a new performance metric for evaluation.
- 2. Real-world data sets: Just as massive real-world data has promoted the rapid development of (LLM, real-world data sets are equally indispensable for further development of DL in the field of wireless communication. At present, most research is still only using data generated from their own simulations. Moreover, regulations concerning data protection and privacy pose additional constraints on the open access to real-world data.

- 3. **E2E concerns**: The preliminary findings [69, 70] indicate that the efficacy of DL-based E2E communications aligns with that of traditional methods. Yet, it remains uncertain if DL-based E2E communication systems will eventually surpass their conventional counterparts in aspects such as performance and complexity, or to what extent they can offer any enhancements. Given the ability of DL methodologies to holistically optimize a communication system in an E2E manner, it raises a query about the potential for next-generation wireless communication technologies to evolve beyond the confines of a stringent standardization process that mandates ongoing regulations.
- 4. Cross layer: This article concentrates on the application of DL in PHY, with numerous studies discussing its application within individual layers [4]. Some researchers argue against confining the application of artificial intelligence solely to PHY and advocate for considering cross-layer applications. For instance, current researchers have embarked on AI-based network optimization. Learning that spans two or three layers holds the promise of greater advantages.
- 5. Hardware learning: Due to the convenience and adaptability of software, learning via software will always be the preferred choice. However, for complex learning tasks, software implementation can introduce high time and spatial complexity. As such, learning through hardware is also being considered. Even though hardware learning is a trending direction in other DL domains, effective implementation is equally needed for the various DL modules within wireless communication systems, with corresponding algorithms that should be hardware friendly.

# 9 | CONCLUSION

This paper provides a thorough review of the methodologies for applying DL schemes to enhance the performance of PHY in wireless networks, particularly within the context of 5G and 6G environments. The focus of these methodologies spans three main aspects: (1) Replacing parts of the conventional wireless communication system. (2) Integrating with emerging wireless communication technologies. (3) Constructing E2E wireless communications. Moreover, the paper outlines several vital research challenges in this field that need to be addressed in the near future. Its goal is to guide readers to understand the state-of-the-art of DL-based wireless networks, and to identify intriguing and challenging research subjects in this critical domain.

### **AUTHOR CONTRIBUTIONS**

Chengchen Mao: Conceptualization, formal analysis, methodology, software, writing - original draft, writing - review and editing. Zongwen Mu: Conceptualization, formal analysis, software, validation, writing - review and editing. Qilian Liang: Conceptualization, formal analysis, funding acquisition, investigation, methodology, project administration, supervision, writing - review and editing. Ioannis Schizas: Conceptualization, formal analysis, visualization, writing - review and editing.

Chenyun Pan: Conceptualization, formal analysis, funding acquisition, writing - review and editing.

#### **ACKNOWLEDGEMENTS**

This work was supported by U.S. National Science Foundation (NSF) under Grant CCF-2219753.

#### CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

#### DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no datasets were generated or analysed during the current study.

#### ORCID

*Oilian Liang* https://orcid.org/0000-0002-3630-8010

#### REFERENCES

- Wang, C.X., Haider, F., Gao, X., You, X.H., Yang, Y., Yuan, D., et al.: Cellular architecture and key technologies for 5G wireless communication networks. IEEE Commun. Mag. 52(2), 122–130 (2014)
- Ghasempour, A.: Internet of things in smart grid: Architecture, applications, services, key technologies, and challenges. Inventions 4(1), 22 (2019)
- Ullah, H., Nair, N.G., Moore, A., Nugent, C., Muschamp, P., Cuevas, M.: 5G communication: An overview of vehicle-to-everything, drones, and healthcare use-cases. IEEE Access 7, 37251–37268 (2019)
- You, X., Wang, C.X., Huang, J., Gao, X., Zhang, Z., Wang, M., et al.: Towards 6G wireless communication networks: Vision, enabling technologies, and new paradigm shifts. Sci. China Inf. Sci. 64, 1–74 (2021)
- Zong, B., Fan, C., Wang, X., Duan, X., Wang, B., Wang, J.: 6G technologies: Key drivers, core requirements, system architectures, and enabling technologies. IEEE Veh. Technol. Mag. 14(3), 18–27 (2019)
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., et al.: Segment anything. (2023) arXiv:2304.02643
- Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I.: Robust speech recognition via large-scale weak supervision. In: International Conference on Machine Learning. PMLR, pp. 28492–28518 (2023)
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., et al.: Language models are few-shot learners. Adv. Neural Inf. Process. 33, 1877–1901 (2020)
- 9. Ji, P., Ye, J., Mu, Y., Lin, W., Tian, Y., Hens, C., et al.: Signal propagation in complex networks. Phys. Rep. 1017, 1–96 (2023)
- Ribeiro, H.V., Lopes, D.D., Pessa, A.A., Martins, A.F., da Cunha, B.R., Gonçalves, S., et al.: Deep learning criminal networks. Chaos Solitons Fractals 172, 113579 (2023)
- Gao, Z., Dang, W., Wang, X., Hong, X., Hou, L., Ma, K., et al.: Complex networks and deep learning for EEG signal analysis. Cogn. Neurodynamics 15, 369–388 (2021)
- Qin, Z., Ye, H., Li, G.Y., Juang, B.H.F.: Deep learning in physical layer communications. IEEE Wireless Commun. 26(2), 93–99 (2019)
- 13. O'shea, T., Hoydis, J.: An introduction to deep learning for the physical layer. IEEE Trans. Cognit. Commun. Netw. 3(4), 563–575 (2017)
- Wang, T., Wen, C.K., Wang, H., Gao, F., Jiang, T., Jin, S.: Deep learning for wireless physical layer: Opportunities and challenges. China Commun. 14(11), 92–111 (2017)
- Mao, Q., Hu, F., Hao, Q.: Deep learning for intelligent wireless networks: A comprehensive survey. IEEE Commun. Surv. Tutorials 20(4), 2595–2621 (2018)
- Zhang, C., Ueng, Y.L., Studer, C., Burg, A.: Artificial intelligence for 5G and beyond 5G: Implementations, algorithms, and optimizations. IEEE J. Emerging Sel. Top. Circuits Syst. 10(2), 149–163 (2020)

 Ozpoyraz, B., Dogukan, A.T., Gevez, Y., Altun, U., Basar, E.: Deep learning-aided 6G wireless networks: A comprehensive survey of revolutionary PHY architectures. IEEE Open J. Commun. Soc. (2022)

- Kramer, M.A.: Nonlinear principal component analysis using autoassociative neural networks. AlChE J. 37(2), 233–243 (1991)
- Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. MIT Press, Massachusetts (2016)
- Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., Bharath, A.A.: Generative adversarial networks: An overview. IEEE Signal Process Mag. 35(1), 53–65 (2018)
- Hu, Q., Gao, F., Zhang, H., Li, G.Y., Xu, Z.: Understanding deep MIMO detection. IEEE Trans. Wireless Commun. (2023). https://doi.org/10.1109/TWC.2023.3272525, early access.
- Farsad, N., Goldsmith, A.: Neural network detection of data sequences in communication systems. IEEE Trans. Signal Process. 66(21), 5663–5678 (2018)
- Shlezinger, N., Fu, R., Eldar, Y.C.: Deepsic: Deep soft interference cancellation for multiuser MIMO detection. IEEE Trans. Wireless Commun. 20(2), 1349–1362 (2020)
- Samuel, N., Diskin, T., Wiesel, A.: Learning to detect. IEEE Trans. Signal Process. 67(10), 2554–2564 (2019)
- He, H., Wen, C.K., Jin, S., Li, G.Y.: Model-driven deep learning for mimo detection. IEEE Trans. Signal Process. 68, 1702–1715 (2020)
- He, H., Wen, C.K., Jin, S., Li, G.Y.: Deep learning-based channel estimation for beamspace mmwave massive MIMO systems. IEEE Wireless Commun. Lett. 7(5), 852–855 (2018)
- Soltani, M., Pourahmadi, V., Mirzaei, A., Sheikhzadeh, H.: Deep learningbased channel estimation. IEEE Commun. Lett. 23(4), 652–655 (2019)
- Chun, C.J., Kang, J.M., Kim, I.M.: Deep learning-based channel estimation for massive MIMO systems. IEEE Wireless Commun. Lett. 8(4), 1228– 1231 (2019)
- Kschischang, F.R., Frey, B.J., Loeliger, H.A.: Factor graphs and the sumproduct algorithm. IEEE Trans. Inf. Theory 47(2), 498–519 (2001)
- Arikan, E.: Channel polarization: A method for constructing capacityachieving codes for symmetric binary-input memoryless channels. IEEE Trans. Inf. Theory 55(7), 3051–3073 (2009)
- Nachmani, E., Marciano, E., Burshtein, D., Be'ery, Y.: RNN decoding of linear block codes. arXiv:170207560 (2017)
- Wang, Y., Zhang, S., Zhang, C., Chen, X., Xu, S.: A low-complexity belief propagation based decoding scheme for polar codes-decodability detection and early stopping prediction. IEEE Access 7, 159808–159820 (2019)
- Liang, F., Shen, C., Wu, F.: An iterative BP-CNN architecture for channel decoding. IEEE J. Sel. Top. Signal Process. 12(1), 144–159 (2018)
- Gruber, T., Cammerer, S., Hoydis, J., ten Brink, S.: On deep learningbased channel decoding. In: 2017 51st annual conference on information sciences and systems (CISS), pp. 1–6. IEEE, New York (2017)
- Cammerer, S., Gruber, T., Hoydis, J., ten Brink, S.: Scaling deep learningbased decoding of polar codes via partitioning. In: 2017 IEEE Global Communications Conference, pp. 1–6. IEEE, New York (2017)
- Elkelesh, A., Ebada, M., Cammerer, S., Schmalen, L., ten Brink, S.: Decoder-in-the-loop: Genetic optimization-based LDPC code design. IEEE Access 7, 141161–141170 (2019)
- Elkelesh, A., Ebada, M., Cammerer, S., ten Brink, S.: Decoder-tailored polar code design using the genetic algorithm. IEEE Trans. Commun. 67(7), 4521–4534 (2019)
- Buzzi, S., Chih-Lin, I., Klein, T.E., Poor, H.V., Yang, C., Zappone, A.: A survey of energy-efficient techniques for 5G networks and challenges ahead. IEEE J. Sel. Areas Commun. 34(4), 697–709 (2016)
- Sun, C., She, C., Yang, C., Quek, T.Q., Li, Y., Vucetic, B.: Optimizing resource allocation in the short blocklength regime for ultra-reliable and low-latency communications. IEEE Trans. Wireless Commun. 18(1), 402–415 (2018)
- Sun, H., Chen, X., Shi, Q., Hong, M., Fu, X., Sidiropoulos, N.D.: Learning to optimize: Training deep neural networks for interference management. IEEE Trans. Signal Process. 66(20), 5438–5453 (2018)

Lee, W., Kim, M., Cho, D.H.: Deep power control: Transmit power control scheme based on convolutional neural network. IEEE Commun. Lett. 22(6), 1276–1279 (2018)

- Chen, X., Liu, G., Ma, Z., Zhang, X., Xu, W., Fan, P.: Optimal power allocations for non-orthogonal multiple access over 5G full/half-duplex relaying mobile wireless networks. IEEE Trans. Wireless Commun. 18(1), 77–92 (2018)
- Luo, J., Tang, J., So, D.K., Chen, G., Cumanan, K., Chambers, J.A.: A deep learning-based approach to power minimization in multi-carrier NOMA with SWIPT. IEEE Access 7, 17450–17460 (2019)
- Khairy, S., Balaprakash, P., Cai, L.X., Cheng, Y.: Constrained deep reinforcement learning for energy sustainable multi-UAV based random access IOT networks with NOMA. IEEE J. Sel. Areas Commun. 39(4), 1101–1115 (2020)
- Xiang, L., Liu, Y., Van Luong, T., Maunder, R.G., Yang, L.L., Hanzo, L.: Deep-learning-aided joint channel estimation and data detection for spatial modulation. IEEE Access 8, 191910–191919 (2020)
- Wang, X., Hua, H., Xu, Y.: Pilot-assisted channel estimation and signal detection in uplink multi-user MIMO systems with deep learning. IEEE Access 8, 44936

  –44946 (2020)
- Zhang, Y., Sun, J., Xue, J., Li, G.Y., Xu, Z.: Deep expectation-maximization for joint MIMO channel estimation and signal detection. IEEE Trans. Signal Process. 70, 4483–4497 (2022)
- Chen, S., Gibson, G., Cowan, C., Grant, P.: Adaptive equalization of finite non-linear channels using multilayer perceptrons. Signal Process. 20(2), 107–119 (1990)
- Patra, J.C., Pal, R.N.: A functional link artificial neural network for adaptive channel equalization. Signal Process. 43(2), 181–195 (1995)
- Lee, J., Beach, C., Tepedelenlioglu, N.: A practical radial basis function equalizer. IEEE Trans. Neural Networks 10(2), 450–455 (1999)
- Kechriotis, G., Zervas, E., Manolakos, E.S.: Using recurrent neural networks for adaptive communication channel equalization. IEEE Trans. Neural Networks 5(2), 267–278 (1994)
- Ye, H., Li, G.Y.: Initial results on deep learning for joint channel equalization and decoding. In: 2017 IEEE 86th Vehicular Technology Conference (VTC-Fall), pp. 1–5. IEEE, New York (2017)
- Xu, W., Zhong, Z., Be'ery, Y., You, X., Zhang, C.: Joint neural network equalizer and decoder. In: 2018 15th International Symposium on Wireless Communication Systems (ISWCS), pp. 1–5. IEEE, New York (2018)
- 54. Hu, Y., Zhao, L., Hu, Y.: Joint channel equalization and decoding with one recurrent neural network. In: 2019 IEEE International Symposium on Broadband Multimedia Systems and Broadcasting (BMSB), pp. 1–4. IEEE, New York (2019)
- Hassouna, S., Jamshed, M.A., Rains, J., Kazim, J.U.R., Rehman, M.U., Abualhayja, M., et al.: A survey on reconfigurable intelligent surfaces: Wireless communication perspective. IET Commun. 17(5), 497–537 (2023)
- Song, H., Zhang, M., Gao, J., Zhong, C.: Unsupervised learning-based joint active and passive beamforming design for reconfigurable intelligent surfaces aided wireless networks. IEEE Commun. Lett. 25(3), 892–896 (2020)
- Feng, K., Wang, Q., Li, X., Wen, C.K.: Deep reinforcement learning based intelligent reflecting surface optimization for MISO communication systems. IEEE Wireless Commun. Lett. 9(5), 745–749 (2020)
- Yang, Z., Liu, Y., Chen, Y., Al-Dhahir, N.: Machine learning for user partitioning and phase shifters design in RIS-aided noma networks. IEEE Trans. Commun. 69(11), 7414

  –7428 (2021)
- Al-Hilo, A., Samir, M., Elhattab, M., Assi, C., Sharafeddine, S.: Reconfigurable intelligent surface enabled vehicular communication: Joint user scheduling and passive beamforming. IEEE Trans. Veh. Technol. 71(3), 2333–2345 (2022)
- Basar, E.: Index modulation techniques for 5G wireless networks. IEEE Commun. Mag. 54(7), 168–175 (2016)
- Yang, P., Xiao, Y., Xiao, M., Guan, Y.L., Li, S., Xiang, W.: Adaptive spatial modulation MIMO based on machine learning. IEEE J. Sel. Areas Commun. 37(9), 2117–2131 (2019)

62. Tato, A., Mosquera, C., Henarejos, P., Pérez-Neira, A.: Neural network aided computation of mutual information for adaptation of spatial modulation. IEEE Trans. Commun. 68(5), 2809–2822 (2020)

- Zhang, Y., Wang, J., Wang, X., Xue, Y., Song, J.: Efficient selection on spatial modulation antennas: Learning or boosting. IEEE Wireless Commun. Lett. 9(8), 1249–1252 (2020)
- Gecgel, S., Goztepe, C., Kurt, G.K.: Transmit antenna selection for largescale MIMO GSM with machine learning. IEEE Wireless Commun. Lett. 9(1), 113–116 (2019)
- Alberge, F.: Deep learning constellation design for the AWGN channel with additive radar interference. IEEE Trans. Commun. 67(2), 1413–1423 (2018)
- Wu, D., Nekovee, M., Wang, Y.: Deep learning-based autoencoder for muser wireless interference channel physical layer design. IEEE Access 8, 174679–174691 (2020)
- Raj, V., Kalyani, S.: Design of communication systems using deep learning: A variational inference perspective. IEEE Trans. Cognit. Commun. Networking 6(4), 1320–1334 (2020)
- Zhang, H., Lan, M., Huang, J., Huang, C., Cui, S.: Noncoherent energy-modulated massive SIMO in multipath channels: A machine learning approach. IEEE Internet Things J. 7(9), 8263–8270 (2020)
- Dörner, S., Cammerer, S., Hoydis, J., Ten.Brink, S.: Deep learning based communication over the air. IEEE J. Sel. Top. Signal Process. 12(1), 132– 143 (2017)
- Ye, H., Liang, L., Li, G.Y., Juang, B.H.: Deep learning-based end-to-end wireless communication systems with conditional GANs as unknown channels. IEEE Trans. Wireless Commun. 19(5), 3133–3143 (2020)
- Liu, Y., Chen, X., Zhong, Z., Ai, B., Miao, D., Zhao, Z., et al.: Waveform design for 5G networks: Analysis and comparison. IEEE Access 5, 19282– 19292 (2017)
- Kim, M., Lee, W., Cho, D.H.: A novel PAPR reduction scheme for OFDM system based on deep learning. IEEE Commun. Lett. 22(3), 510–513 (2017)
- Liu, Z., Hu, X., Han, K., Zhang, S., Sun, L., Xu, L., et al.: Low-complexity PAPR reduction method for OFDM systems based on real-valued neural networks. IEEE Wireless Commun. Lett. 9(11), 1840–1844 (2020)
- Goutay, M., Aoudia, F.A., Hoydis, J., Gorce, J.M.: End-to-end learning of ofdm waveforms with PAPR and ACLR constraints. In: 2021 IEEE Globecom Workshops (GC Workshops), pp. 1–6. IEEE, New York (2021)
- Aoudia, F.A., Hoydis, J.: End-to-end learning for OFDM: From neural receivers to pilotless communication. IEEE Trans. Wireless Commun. 21(2), 1049–1063 (2021)
- Aoudia, F.A., Hoydis, J.: Trimming the fat from OFDM: Pilot-and cp-less communication with end-to-end learning. In: 2021 IEEE International Conference on Communications Workshops (ICC Workshops), pp. 1–6. IEEE, New York (2021)
- 77. Ye, H., Li, G.Y., Juang, B.H.: Deep learning based end-to-end wireless communication systems without pilots. IEEE Trans. Cognit. Commun. Networking 7(3), 702–714 (2021)

- Luo, X., Chen, H.H., Guo, Q.: Semantic communications: Overview, open issues, and future research directions. IEEE Wireless Commun. 29(1), 210–219 (2022)
- Carnap, R., Bar-Hillel, Y., et al.: An outline of a theory of semantic information. (1952) https://dspace.mit.edu/bitstream/handle/1721.1/4821/RLE-TR-247-03150899.pdf?sequence=1MIT RLE Technical Report
- Bao, J., Basu, P., Dean, M., Partridge, C., Swami, A., Leland, W., et al.: Towards a theory of semantic communication. In: 2011 IEEE Network Science Workshop, pp. 110–117. IEEE, New York (2011)
- Farsad, N., Rao, M., Goldsmith, A.: Deep learning for joint source-channel coding of text. In: 2018 IEEE international conference on acoustics, speech and signal processing (ICASSP), pp. 2326–2330. IEEE, New York (2018)
- Xie, H., Qin, Z., Li, G.Y., Juang, B.H.: Deep learning based semantic communications: An initial investigation. In: GLOBECOM 2020-2020 IEEE Global Communications Conference, pp. 1–6. IEEE, New York (2020)
- Bourtsoulatze, E., Kurka, D.B., Gündüz, D.: Deep joint source-channel coding for wireless image transmission. IEEE Trans. Cognit. Commun. Networking 5(3), 567–579 (2019)
- 84. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., et al.: TensorFlow: Large-scale machine learning on heterogeneous systems. (2015). https://www.tensorflow.org/
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., et al.: Pytorch: An imperative style, high-performance deep learning library. Adv. Neural Inf. Process. 32, 1-12 (2019). https://papers.nips.cc/paper\_files/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdfpp
- Hoydis, J., Cammerer, S., Ait Aoudia, F., Vem, A., Binder, N., Marcus, G., et al.: Sionna: An open-source library for next-generation physical layer research. arXiv (2022). https://arxiv.org/abs/2203.11854
- Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., et al.: Caffe: Convolutional architecture for fast feature embedding. In: Proceedings of the 22nd ACM international conference on Multimedia, pp. 675–678 (2014)
- Al-Rfou, R., Alain, G., Almahairi, A., Angermueller, C., Bahdanau, D., Ballas, N., et al.: Theano: A python framework for fast computation of mathematical expressions. arXiv–1605 (2016). https://arxiv.org/abs/ 1605.02688

How to cite this article: Mao, C., Mu, Z., Liang, Q., Schizas, I., Pan, C.: Deep learning in physical layer communications: Evolution and prospects in 5G and 6G networks. IET Commun. 17, 1863–1876 (2023). https://doi.org/10.1049/cmu2.12669