# Learning to Denoise and Decode: A Novel Residual Neural Network Decoder for Polar Codes

Zhiwei Cao, Hongfei Zhu, Yuping Zhao, Dou Li
School of Electronics Engineering and Computer Science
Peking University, Beijing, 100871, China
Email:{cao_zhiwei, zhuhongfei, yuping.zhao, lidou}@pku.edu.cn

*Abstract*—**Polar codes have been adopted as the control channel coding scheme in the fifth generation new radio (5G NR) standard due to its capacity-achievable property. Traditional polar decoding algorithms such as successive cancellation (SC) suffer from high latency problem because of their sequential decoding nature. Neural network decoder (NND) has been proved to be a candidate for polar decoder since it is capable of *one-shot* decoding and parallel computing. Whereas, the bit-error-rate (BER) performance of NND is still inferior to that of SC algorithm. In this paper, we propose a residual neural network decoder (RNND) for polar codes. Different from previous works which directly use neural network for decoding symbols received from the channel, the proposed RNND introduces a denoising module based on residual learning before NND. The proposed residual learning denoiser is able to remove remarkable amount of noise from received signals. Numerical results show that our proposed RNND outperforms traditional NND with regard to the BER performance under comparable latency.**

## I. INTRODUCTION

Polar codes proposed by Erdal Arikan are the first provable capacity achieving codes for symmetric binary-input discrete memoryless channels (B-DMCs) [1]. Although successive cancellation (SC) decoding algorithm proposed by Arikan initially has a lower complexity $O(NlogN)$ in terms of the code length $N$, its performance for finite block lengths is not satisfactory. Later successive cancellation list [2], [3], successive cancellation stack [4] and CRC-aided SCL/SCS decoders [5] were introduced to improve the decoding performance of polar codes. However, they still suffer from high latency as well as limited throughput due to their sequential decoding property, which is disadvantageous to the practical wireless communication systems requiring high reliability and low latency.

Recently deep learning (DL) has attracted worldwide attentions because of its powerful capabilities to solve complicated tasks. With the help of deep learning, significant improvements have been achieved in many fields, such as computer vision [6], natural language processing [7], autonomous vehicles [8] and many other areas. In the field of communication, the general channel decoding problems can also be solved with deep learning, since they can be regarded as a type of classification problem. The channel decoder based on deep neural network is called the neural network decoder (NND)

consisting of two stages: NND training and NND testing. Non-iterative and consequently low-latency decoding are two main advantages of NND, because NND calculates the estimated value of information bits by passing each layer only once with the pre-trained neural network, which is referred to as *one-shot* decoding. What's more, the current deep learning platforms, such as Pytorch [9], and the powerful hardwares like graphical processing units (GPUs) can enable the efficient implementation of NND.

NND for polar codes was proposed in [10], where the authors found that structured codes are indeed easier to learn than random codes, and the bit-error-rate (BER) performance of NND is close to that of maximum a posteriori (MAP) decoding. In [11] the authors considered a deep feed-forward neural network for polar codes and investigated its decoding performances with respect to numerous configurations: the number of hidden layers, the number of nodes for each layer, and activation functions. Later, more discussion on the activation function of the neural network for decoding polar codes and decoding under Reighlay fading channels using NND can be found in [12] and [13], respectively. In [14], the authors compared the decoding performance of three types of neural network, i.e., multi-layer perceptron (MLP), convolutional neural network (CNN) and recurrent neural network (RNN) with the same parameter magnitude. The authors found that the neural network is capable of learning the complete encoding structure with noiseless codewords as training data. The comparison of the three types of neural network was also discussed in [15], where the authors proposed a unified polar-LDPC NND by concatenating an indicator section.

All results shown in the above papers suggest that higher test-SNR (Signal-to-Noise Ratio) during the NND testing stage facilitates decoding with NND. Dramatically thinking, if we can place a denoising neural network before NND to improve the test-SNR, the modified NND will certainly show a better decoding performance than the original NND. It is well known that residual learning has been widely used in image denoising [16]. Thus, it may shed some lights on the codeword denoising.

Motivated by this, in this paper we propose a residual neural network decoder (RNND) for polar codes. The proposed RNND introduces a denoiser based on residual learning before NND to improve the SNR of received symbols. Simulation results demonstrate that the proposed RNND has a considerable advantage over existing NNDs w.r.t. bit-error-rate (BER)
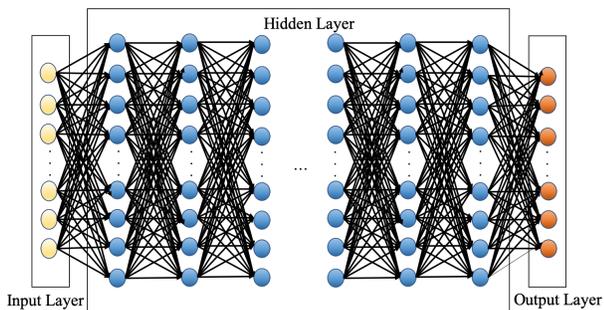
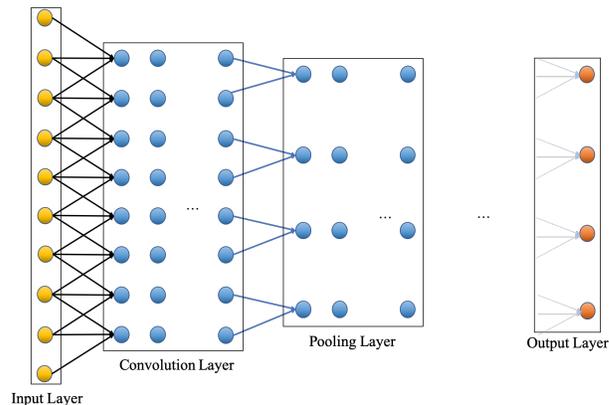Fig. 1. Diagram of multi-layer perceptron (MLP).



Fig. 2. Diagram of convolutional neural network (CNN). The black and blue arrows represent affine transformation and pooling operators, respectively. Better view in color.
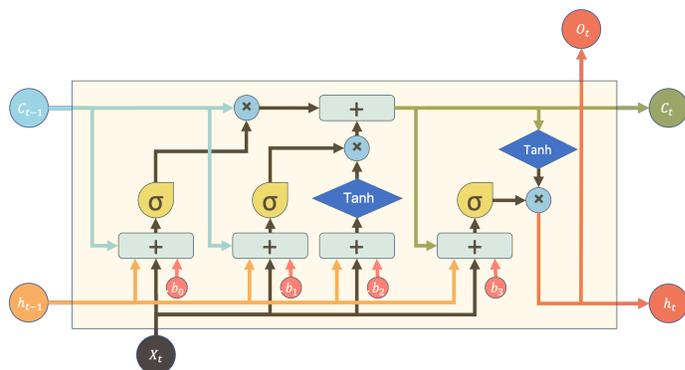
performance under low latency.

The rest of this paper is organized as follows. Section II provides some preliminary knowledge of polar codes, three types of neural network (MLP, CNN, RNN) and the residual learning. The system model is introduced in Section III. The architecture of the proposed RNND with a novel multi-task learning objective and the training and testing strategy are discussed in Section IV. We provide the experiment configuration and corresponding numerical results in Section V. Eventually, concluding remarks are given in Section VI.

## II. PRELIMINARIES

### A. Polar Codes

To construct an $(N, K)$ polar codes, the $K$ information bits $u_{\mathcal{A}}$ and the other $N - K$ frozen bits $u_{\mathcal{A}^c}$ are first assigned to the reliable and unreliable positions of the $N$-bit message $u_1^N$, respectively. $\mathcal{A}$ is the information set, while $\mathcal{A}^c$ is the frozen set which is the complementary set of $\mathcal{A}$. The $N - K$ frozen bits with indices in $\mathcal{A}^c$ are always fixed to zeros. Then the $N$-bit transmitted codeword $x_1^N$ can be obtained by multiplying $u_1^N$ with the generator matrix $G_N$ as follows

$$x_1^N = u_1^N G_N = u_1^N B_N F_2^{\otimes n}, n = log_2 N \qquad (1)$$

where $\otimes$ denotes the Kronecker product, $F_2 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ and $B_N$ is the bit-reversal permutation matrix.

### B. Neural Network

In this paper, we utilize three popular and effective neural network architectures: 1) MLP; 2) CNN; 3) RNN to construct RNNDs and compare their BER performances. Here, we briefly describe the general architecture of these neural networks.

*1) MLP:* MLP is a kind of feedforward and densely connected neural network. There is no feedback or loop in the neural network and each node is connected with all its ancestors. The input data is processed layer by layer with affine transformation and nonlinear activation function between layers and finally reaches the output layer. Fig. 1 illustrates the general architecture of MLP.



Fig. 3. Diagram of long short-time memory (LSTM). $h_t, C_t, X_t$ and $O_t$ represent the hidden vector, cell state, input vector and output vector at time $t$, respectively.

*2) CNN:* The hidden layers of classical CNN are either convolutional layers or pooling layers. It is universally known that CNN has strong ability of extracting features automatically. Fig. 2 shows the classical diagram for CNN.

*3) RNN:* RNN is a class of neural network that has recurrent structure, i.e., previous states will influence current outputs. Therefore, RNN is very powerful for time series modeling. Traditional RNN suffers from severe vanishing and exploding gradient problem [17], making it difficult to train. Hence, in practice, people usually utilize its variants like LSTM [18] and GRU [19]. In this paper, we use LSTM as our practical RNN implementation. Fig. 3 presents the classical structure of the LSTM.

### C. Residual Learing

Residual learning of deep neural network was first proposed in [20] to solve the performance degradation problem, i.e., the training accuracy begins to decrease as the depth of the neural network increases. By assuming the residual mapping is much easier to learn, the authors in [20] stacked a few nonlinear layers for explicitly learning the residual mapping. Residual learning of neural networks is achieved through
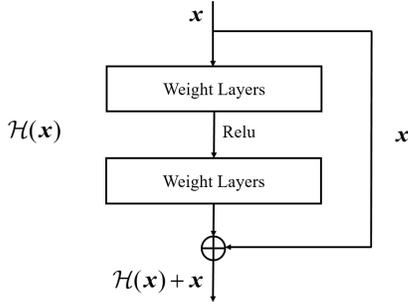
Fig. 4. Diagram of the residual learning block. $\mathcal{H}(\boldsymbol{x})$ corresponds to the stacked weight layers with Relu non-linearity.

shortcut connections. Fig. 4 illustrates the general architecture of the residual learning block.

## III. SYSTEM MODEL

The system model in this paper is shown in Fig. 5. At the transmitter, $K$-bit information $u_{\mathcal{A}}$ is first encoded into a $N$-bit codeword $x_1^N$ with a polar encoder. Then $x_1^N$ is mapped to $N$ modulated symbols $s_1^N$ through binary phase shift keying (BPSK) modulation by $s_i = 1 - 2x_i$ for $i = 1, 2, ..., N$. The modulated symbol $s_1^N$ is subsequently trainsmitted over an additive white Gaussian noise (AWGN) channel.

The signal model for received symbols can be formulated as follows

$$y_1^N = s_1^N + n_1^N \tag{2}$$

where $n_1^N \sim \mathcal{N}(\boldsymbol{0}, \sigma^2 \mathbf{I}_N)$ represents the i.i.d. Gaussian noise vector.

At the receiver, the received symbol vector $y_1^N$ first enters the residual learning denoiser to reduce the noise. After the denoising process, the filtered signal vector $\hat{s}_1^N$ is then sent to the neural network decoder. Finally, we obtain the estimated value of $u_{\mathcal{A}}$, denoted as $\hat{u}_{\mathcal{A}}$. The two modules, residual learning denoiser and neural network decoder, make up the architecture of the proposed RNND which will be addressed in Section IV. Addtionally, the three loss function calculation modules which are denoising loss, decoding loss and multi-task loss shown in Fig. 5 will be discussed in detail in Section IV.

## IV. PROPOSED RESIDUAL NEURAL NETWORK DECODER

In this section, we first delineate the architecture of the proposed RNND which contains the residual learning denoiser and the neural network decoder. The denoiser and decoder will be jointly trained through a novel multi-task learning objective which will be discussed in the following. Ultimately, we depict the training and testing strategy for jointly learning the decoder and the denoiser.

### A. Residual Learning Denoiser

Residual learning has been shown to be a powerful tool for image denoising [16] and the codeword denoising can

be regarded as the special case of image denoising in one-dimension input. Residual learning module, i.e., $\mathcal{H}(\boldsymbol{x})$ in Fig. 4, can be directly viewed as a denoising module as we optimize its parameters, so that the output of it, i.e., $\mathcal{H}(\boldsymbol{x}) + \boldsymbol{x}$ in Fig. 4, can be as close as possible to the transmitted symbols. The residual learning denoiser consists of some stacked layers and a shortcut connection, as shown in Fig. 4. The weight layers in Fig. 4 can be any type of neural network. In this work, we investigate three different types of neural network as the denoiser, which are MLP, CNN and RNN, respectively.

As shown in Fig. 5, denote the function corresponding to the stacked weight layers in the denoiser as $\mathcal{H}(y_1^N)$ and the output of the denoiser is the element-wise summation over the received signal vector $y_1^N$ and $\mathcal{H}(y_1^N)$

$$\hat{s}_1^N = y_1^N + \mathcal{H}(y_1^N) \tag{3}$$

For the denoising task, we aim to minimize the difference of $\hat{s}_1^N$ and the transmitted symbol $s_1^N$. We take mean squared error (MSE) to measure the discrepency between $\hat{s}_1^N$ and $s_1^N$ and it leads to our loss function for the denoising task

$$\mathcal{L}_{denoise} = \frac{1}{N} ||\hat{s}_1^N - s_1^N||_2^2 \tag{4}$$

### B. Neural Network Decoder

The neural network decoder is illustrated in Fig. 5, following the residual learning denoiser. In our work, we regard channel decoding as a binary classification problem: the decoder attempts to categorize each received symbol into 0 or 1. For the classification problem, there are some commonly used loss functions like binary cross entropy (BCE) and MSE. In [10] the author has shown via experiments that MSE and BCE have little difference w.r.t. the BER performance of NNDs. Here, we utilize MSE as our loss function for decoding task

$$\mathcal{L}_{decode}(\hat{u}_{\mathcal{A}}, u_{\mathcal{A}}) = \frac{1}{K} ||\hat{u}_{\mathcal{A}} - u_{\mathcal{A}}||_2^2 \tag{5}$$
$$= \frac{1}{K} ||\mathcal{G}(\hat{s}_1^N) - u_{\mathcal{A}}||_2^2$$

where, $\mathcal{G}$ is the function corresponds to the neural network decoder.

### C. Multi-task Learning

As mentioned above, our goal is to train the entire neural network so that it can accomplish both denoising and decoding tasks. One potential way for achieving such goal is multi-task learning. In our work, we explicitly sum the decoding loss and denoising loss together as our final loss function

$$\mathcal{L} = \mathcal{L}_{denoise} + \mathcal{L}_{decode} \tag{6}$$

With multi-task learning, we hope the learning process of the denoiser and that of the decoder reinforce mutually. The gradient signal from the decoder may lead to better learning of the denoiser. Meanwhile the proposed denoiser will reduce noise level of the received signal, which facilitates the decoder to learn decoding rules. Since $\mathcal{L}_{denoise}$ and $\mathcal{L}_{denoise}$ are both continous and differentiable everywhere, it can be efficiently optimized with some famous optimizers like SGD [21] and Adam [22].
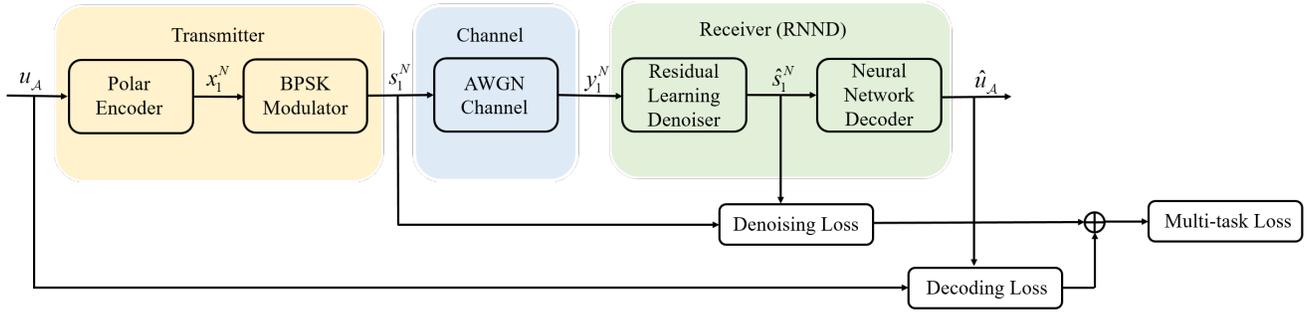
Fig. 5. The system model in this paper. Residual learning denoiser and neural network decoder make up the architecture of the proposed RNND. The denoising loss and decoding loss correspond to (4) and (5), respectively. The multi-task loss is the proposed joint optimization objective (6).

### D. Training and Testing Strategy

*1) Training data generation:* To avoid *curse of dimensionality* [23], we only take polar codes with code length $N = 16$ and rate $R = 1/2$. With such code setting, there are $2^8 = 256$ distinct codewords with length $N$. We use all of them for training NNDs. During training, these codewords will be modulated with BPSK and contaminated by the AWGN channel. We set the train-SNR = 0 dB as [10] suggests.

*2) Training procedure:* We use batch-based training method. The batch size is set to be 64 in our experiments. At each iteration, we first select a batch of codewords from the codeword set in order. Then we push these codewords into the model and evaluate the loss function with (6). Finally, we calculate gradients with backpropagation and update model parameters with Adam [22] in an end-to-end fashion. Leaning rate is set to be 0.001 and the momentum is 0.99. After we traverse the whole codeword set, we call one training epoch is over. We jointly train our denoiser and decoder with $2^{16}$ epochs.

*3) Testing procedure:* After the training stages of different RNNDs, we start the NND testing stage. In contrast to the training stage, the modulated symbols are transmitted over the AWGN channel with different noise levels at the test stage. All results presented in this paper are obtained with the trained models during the testing stage.

## V. NUMERICAL RESULTS

In this section, we first describe the design of the neural network architecture. We subsequently calculate the SNR of $y_1^N$ and $\hat{s}_1^N$ with different neural network structures as denoiser for verifying their strong denoising capacity. The probability density function (PDF) of $y_1^N$ and $\hat{s}_1^N$ are further calculated to validate the design. Finally, decoding performance and computation time comparison between NNDs, proposed RNNDs and the SC algorithm will be discussed.

### A. Design of Neural Network Architecture

In this paper, we name the RNNDs with MLP, CNN and RNN architecture MLP-RNND, CNN-RNND, RNN-RNND, respectively. The configuration of each RNND is portrayed in the following.

*1) MLP-RNND:* The proposed architecture of MLP has 3 layers for denoiser and each layer contains 128, 64, 32 nodes, respectively. Following the denoiser, there are another 3 layers MLP for the decoder with 128, 64, 32 nodes in each layer. The input and output size of the entire system are $N$ and $K$ with no doubt.

*2) CNN-RNND:* As CNN is widely used in image processing task, most CNN architectures deal with 2-D input data. In our channel decoding task, we modify the CNN architecture and set the input of CNN as 1-D received signal vectors instead of 2-D images. Meanwhile, we also revise the convolutional layer with 1-D convolution instead of classical 2-D convolution. The denoiser consists of 3 convolutional layers with 64, 48, 32 channels, while the decoder contains 3 convolutional layers with 64, 32, 32 channels. We set the convolutional kernel size to 3. In addition, MaxPooling [24] with kernel size 2 and stride 2 is utilized between the adjacent convolutional layers.

*3) RNN-RNND:* In order to make RNN suitable for channel decoding, we regard the channel decoding task as a time series classification problem. The RNN runs through the input vector with one symbol as input and produces an output vector at each timestep. After the RNN reads the entire input vector, it produces a sequence of output vectors. Due to the recurrent nature of RNN, the last vector in the output sequence encodes the input vector. Thus we select it as the feature vector. It will be further processed with one affine transformation layer to produce the denoising or decoding output. We leverage one layer RNN with output dimension 64 for the denoiser and another one layer RNN with output dimension 48 for the decoder.

Each RNND has its NND counterpart named MLP-NND, CNN-NND and RNN-NND for performance comparison. For NNDs, we just cancel the shortcut connection and keep the number of layers the same as their counterpart. All NNDs and RNNDs have similar amount of parameters to avoid the performance difference coming from the difference of parameter number. The total number of parameters of six types of neural network is described in detail in Table I.

TABLE I
TOTAL NUMBER OF PARAMETERS OF SIX TYPES OF NEURAL NETWORK

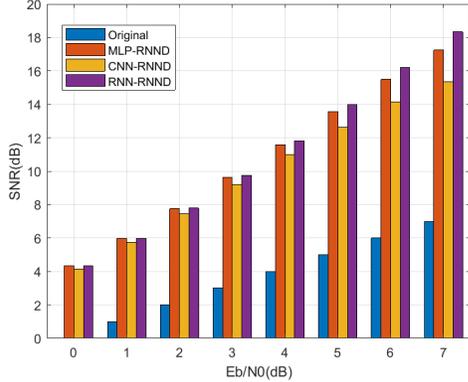| Type of Neural Network | Total Number of Parameters |
| --- | --- |
| MLP-NND | 27336 |
| MLP-RNND | 25816 |
| CNN-NND | 25576 |
| CNN-RNND | 25256 |
| RNN-NND | 27208 |
| RNN-RNND | 28376 |



Fig. 6. The SNR comparison of signals before and after denoising with different RNNDs. Original refers to the SNR of original received signals.



Fig. 7. The PDF of signals before and after denoising of the MLP-RNND under different test-Eb/N0. Blue and orange region represent the PDF of received signals and denoised signals, respectively. Better view in color.



Fig. 8. The BER performance comparison of NNDs, proposed RNNDs and the SC algorithm.

### B. Denoising Capacity Comparison

We investigate the denoising performance of different denoisers in this section. Fig. 6 shows the denoising performance of the denoiser of different neural networks with test-SNR ranging from 0 dB to 7 dB. From Fig. 6, we can see the great SNR improvement achieved by the denoisers. For example, MLP-RNND improve the SNR by about 4 dB at test-SNR = dB and more gain can be achieved when the test-SNR increases. It can also be observed that MLP-RNND has similar denoising performance compared with RNN-RNND and both of them substantially outperform the CNN-RNND. It is reasonable to conjecture that the BER performance of MLP-RNND and RNN-RNND will also significantly outperform CNN-RNND since the MLP and RNN denoiser erase more noise than the CNN denoiser.

Fig. 7 studies the PDF of $y_1^N$ and $\hat{s}_1^N$ of the MLP denoiser under different test-Eb/N0. It is worth noting that the PDF of $\hat{s}_1^N$ is still similar to Gaussian distribution even though the input Gaussian distributed signals $y_1^N$ have gone through nonlinear operations. What's more, it can be found that the variance of the $\hat{s}_1^N$ is much smaller than that of $y_1^N$, which implies considerably amount of Gaussian noise has been removed from received signals.

### C. BER Performance Comparison

After comparing denoising capacity of different denoisers, we study the BER performance of NNDs, proposed RNNDs and the SC algorithm, which is shown in Fig. 8. We can observe that each RNND remarkably outperforms its NND
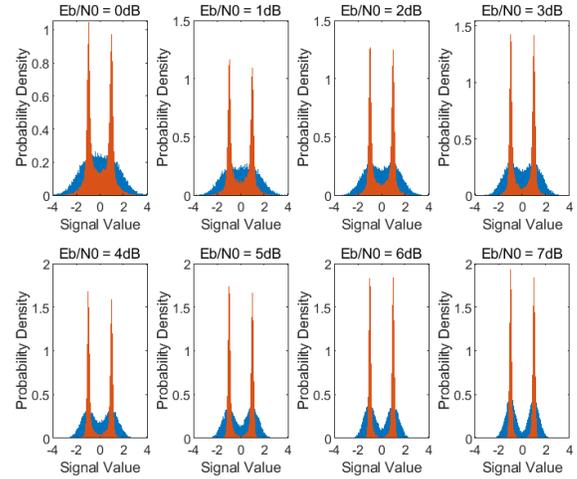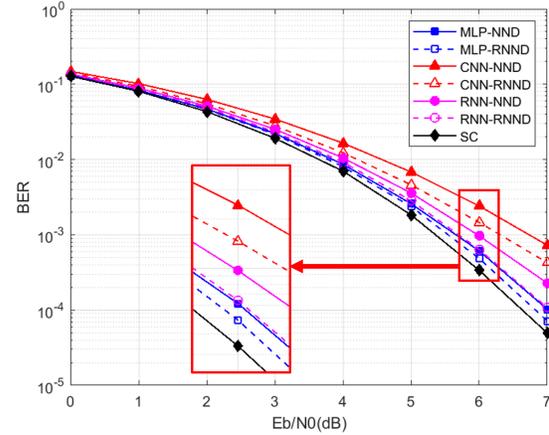
counterpart. For instance, MLP-RNND obtains a gain of roughly 0.2dB over MLP-NND at BER $10^{-4}$. Meanwhile, the BER performance of MLP-RNND is very close to the SC algorithm which is near optimal in our code setting. It should be pointed out that the BER performances of MLP-RNND and RNN-RNND exceed those of CNN-RNND by a large margin. This is consistent with the conjecture we made in V-B. It is also important to highlight that MLP-RNND outperforms both CNN-RNND and RNN-RNND, which conforms to previous research [15]. It implies that MLP may be the optimal structure for polar NND.

### D. Computation Time Comparison

We further study the computational time of NNDs, proposed RNNDs and the SC algorithm. The results are illustrated in Fig. 9. It is worthwhile mentioning that although there exists small gaps between the BER performance of MLP-RNND
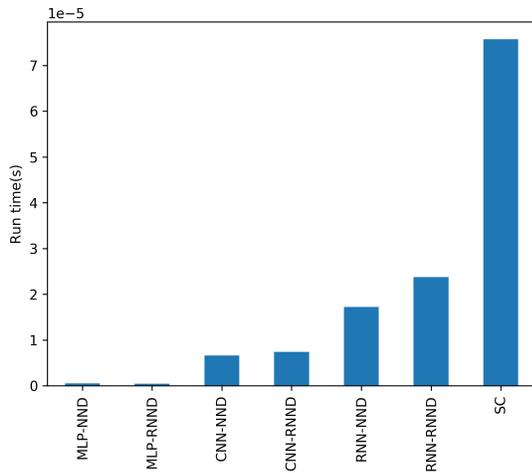
Fig. 9. The computation time comparison of NNDs, proposed RNNDs and the SC algorithm.

and that of the traditional SC algorithm, MLP-RNND runs more than 100 times faster than the SC algorithm. It implies that the proposed MLP-RNND may be a strong alternative for traditional SC algorithm and its variants. It must also be mentioned that RNNDs run slightly slower than NNDs whereas such differences are neglectable in practice.

## VI. CONCLUTION

In this paper, we propose a novel residual neural network decoder (RNND), which jointly learns to denoise and decode with a multi-task training strategy. We present the PDF and SNR of the signals before and after the trained denoiser, which demonstrates the significant noise level supression of our proposed denoiser. Simulation results show that with the aid of the denoiser, one can obtain considerable BER performance gain since less noisy signals are advantageous to learning decoding rules. The computation time results demonstrate that our proposed MLP-RNND may be a strong competitor with classical SC algorithm due to its near optimal BER performance and ultra low latency. Notably, although our research focuses on polar codes with short code length, it is also promising to use the proposed RNND to decode polar codes with longer code length since the residual structure facilitates training deep neural network [20].

## REFERENCES

[1] E. Arikan, "Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels," *IEEE Transactions on Information Theory*, vol. 55, no. 7, pp. 3051–3073, July 2009.

[2] I. Tal and A. Vardy, "List decoding of polar codes," *IEEE Transactions on Information Theory*, vol. 61, no. 5, pp. 2213–2226, May 2015.

[3] K. Chen, K. Niu, and J. R. Lin, "List successive cancellation decoding of polar codes," *Electronics Letters*, vol. 48, no. 9, pp. 500–501, April 2012.

[4] K. Niu and K. Chen, "Stack decoding of polar codes," *Electronics Letters*, vol. 48, no. 12, pp. 695 –697, June 2012.

[5] ——, "Crc-aided decoding of polar codes," *IEEE Communications Letters*, vol. 16, no. 10, pp. 1668–1671, October 2012.

[6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 770–778.

[7] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to Sequence Learning with Neural Networks," *arXiv e-prints*, p. arXiv:1409.3215, Sep 2014.

[8] C. Chen, A. Seff, A. Kornhauser, and J. Xiao, "Deepdriving: Learning affordance for direct perception in autonomous driving," in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec 2015, pp. 2722–2730.

[9] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *NIPS Autodiff Workshop*, 2017.

[10] T. Gruber, S. Cammerer, J. Hoydis, and S. t. Brink, "On deep learning-based channel decoding," in *2017 51st Annual Conference on Information Sciences and Systems (CISS)*, March 2017, pp. 1–6.

[11] J. Seo, J. Lee, and K. Kim, "Decoding of polar code by using deep feed-forward neural networks," in *2018 International Conference on Computing, Networking and Communications (ICNC)*, March 2018, pp. 238–242.

[12] ——, "Activation functions of deep neural networks for polar decoding applications," in *2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*, Oct 2017, pp. 1–5.

[13] A. Irawan, G. Witjaksono, and W. K. Wibowo, "Deep learning for polar codes over flat fading channels," in *2019 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*, Feb 2019, pp. 488–491.

[14] W. Lyu, Z. Zhang, C. Jiao, K. Qin, and H. Zhang, "Performance evaluation of channel decoding with deep neural networks," in *2018 IEEE International Conference on Communications (ICC)*, May 2018, pp. 1–6.

[15] Y. Wang, Z. Zhang, S. Zhang, S. Cao, and S. Xu, "A unified deep learning based polar-ldpc decoder for 5g communication systems," in *2018 10th International Conference on Wireless Communications and Signal Processing (WCSP)*, Oct 2018, pp. 1–6.

[16] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, July 2017.

[17] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *International conference on machine learning*, 2013, pp. 1310–1318.

[18] J. Schmidhuber and S. Hochreiter, "Long short-term memory," *Neural Computing*, vol. 9, no. 8, pp. 1735–1780, 1997.

[19] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Gated feedback recurrent neural networks," in *International Conference on Machine Learning*, 2015, pp. 2067–2075.

[20] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[21] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*. Springer, 2010, pp. 177–186.

[22] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[23] Xiao-An Wang and S. B. Wicker, "An artificial neural net viterbi decoder," *IEEE Transactions on Communications*, vol. 44, no. 2, pp. 165–171, Feb 1996.

[24] Y.-L. Boureau, J. Ponce, and Y. LeCun, "A theoretical analysis of feature pooling in visual recognition," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 111–118.