

# Performance Evaluation of Channel Decoder based on Recurrent Neural Network

Saisi Meng<sup>1,a</sup>, Xue-Qin Jiang<sup>1,b</sup>, Yongbin Gao<sup>2,c</sup>, Han Hai<sup>1,d</sup> and Jia Hou<sup>3,e</sup>

<sup>1</sup>School of Information Science and Technology, Donghua University, Shanghai, China

<sup>2</sup>School of Engineering Science, Shanghai University, Shanghai, China

<sup>3</sup>School of Electronics and Information, Soochow University, Suzhou, China

Email: <sup>a</sup>saisi\_meng@163.com, <sup>b</sup>xqjiang@dhu.edu.cn, <sup>c</sup>gaoyongbin@sues.ehu.cn, <sup>d</sup>hhai@dhu.edu.cn, <sup>e</sup>houjia@suda.edu.cn

**Abstract.** Recently, recurrent neural network (RNN) has demonstrated superior performance for channel decoder, which motivates us to explore which RNN decoder can be more efficient. In this paper, we propose three kinds of RNN decoders, which are built upon long short term memory (LSTM), gated recurrent unit (GRU) and bidirectional gated recurrent units (Bi-GRU), respectively. The performance of these three RNN decoders are evaluated through lots of simulations, which indicate that the GRU decoder with simplest structure and least computational time, has similar bit error rate (BER) performance as that of the LSTM decoder. The Bi-GRU decoder has the best BER performance at the expenses of more computational time. However, it is prone to overfitting. Furthermore, we find that the BER performance of RNN decoders without dropout is better than that of the decoders with dropout when decoding models are underfitting, while the RNN decoders are better to dropout when decoding models are overfitting.

## 1. Introduction

With the improvement of fifth generation (5G) network technology, a new generation of communication mode is coming soon. The theoretical transmission speed of 5G networks, which is hundreds of times faster than 4G, can reach tens of Gbs per second. Obviously, higher transmission rate requires lower decoding latency. Since most of the decoding algorithms involves large numbers of iterative calculations, we need to design a decoder that should satisfy the demand of high-speed and low-latency. It happens that deep learning can help us solve this problem.

Deep learning has shown great potential and advantage in some complex tasks. It has achieved remarkable results in image processing [1], machine translation [2], speech recognition [3] and many other areas. Inspired by powerful learning ability of deep learning, [4] provided some evidence that structured codes were easier to learn than random codes, and pointed out that the neural networks could learn a structure of decoding algorithm, rather than a simple classifier. [5] compared three types of neural network decoder (NND) with the same parameter magnitude and drew a conclusion that recurrent neural network (RNN) decoder had the best decoding BER performance but the highest computational time. In recent years, RNN has derived many variants, which demonstrate better performance in many tasks. According to the conclusion of [5], we want to explore the performance of

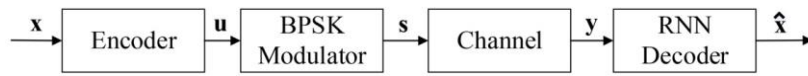


the decoders built upon RNN variants. In the following works, when the decoders of these RNN variants appear at the same time, we collectively refer to them as RNN decoders.

In this paper, we propose three kinds of RNN decoders, which are built upon long short-term memory (LSTM), gated recurrent unit (GRU) and bidirectional gated recurrent unit (Bi-GRU). We discuss them separately in the case of decoders with dropout and decoders without dropout. In this work we focus on decoding polar codes, which are mathematically proven to achieve channel capacity. Our goal for this paper is to find which RNN decoder can be more efficient, which is measured by two indicators, decoding BER and computational time.

## 2. System design

### 2.1. Channel model



**Figure 1.** The architecture of channel model with RNN decoder.

The architecture of channel model with RNN decoder is illustrated in figure 1. The task of the encoder at the transmitter is to encode information bits  $\mathbf{x}$  of length  $K$  into a binary codeword  $\mathbf{u}$  of length  $N$ . We assume that through binary phase shift keying (BPSK) modulation, the codeword  $\mathbf{u}$  can be mapped to a symbol vector  $\mathbf{s}$ . After transmitting the symbol vector  $\mathbf{s}$  over a channel with the additive white Gaussian noise (AWGN), a noisy version of the codeword  $\mathbf{y}$  is received at the receiver. The received vector  $\mathbf{y}$  can be written as  $\mathbf{y} = \mathbf{s} + \mathbf{n}$ , where  $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$  represents the  $N \times 1$  symbol vector. The task of the decoder is to recover the vector  $\mathbf{y}$  to the corresponding information bits  $\mathbf{x}$ . Here, we denote the recovered bits  $\hat{\mathbf{x}}$ . We hope that after decoding by the RNN decoders, the BER can be as close as possible to the maximum a posteriori (MAP) decoding. For details we refer to [5].

### 2.2. Training

The neural network model for supervised learning consists of two phases, training phase and testing phase. Training samples for training phase are generated by the following rules: We first denote the messages bits of length  $N$  including  $K$  information bits  $\mathbf{x}$  and  $N - K$  frozen bits as  $\mathbf{x}'$ , and then the transmitted codeword  $\mathbf{u}$  of length  $N$  can be obtained by following the encoding rules of polar codes  $\mathbf{u} = \mathbf{x}' \mathbf{G}$ , where  $\mathbf{G} = \mathbf{F}^{\otimes m}$ .  $\mathbf{F}^{\otimes m}$  is the  $m$ -th Kronecker power where  $\mathbf{F} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$  and  $m = \log_2 N$ . Finally, the received symbol vector  $\mathbf{y}$  can be obtained when  $\mathbf{u}$  passes through BPSK modulation and channel with AWGN.

Let  $\mathbf{x} = [x_0, \dots, x_{K-1}]$ ,  $\hat{\mathbf{x}} = [\hat{x}_0, \dots, \hat{x}_{K-1}]$  and  $\mathbf{y} = [y_0, \dots, y_{N-1}]$ . We define all possible  $\mathbf{x}$  and  $\hat{\mathbf{x}}$  by the set  $A$  and that of all possible  $\mathbf{y}$  by the set  $B$ . Randomly picked a training sample  $\mathbf{y}$  from  $B$ , there will be a corresponding label  $\mathbf{x}$  in  $A$ . During training phase, we consider four factors that affect the performance of the RNN decoders.

- *Signal to noise ratio of training samples:* In the actual channel, the signal to noise ratio (SNR) is uncertain and time-varying, we simulate the noise in channel decoding phase to set different SNR of training samples and adopt the normalized validation error (NVE) in [4] to measure the performance of RNN decoders. The NVE defines as

$$\text{NVE}(\rho_t) = \frac{1}{S} \sum_{s=1}^S \frac{\text{BER}_R(\rho_t, \rho_{v,s})}{\text{BER}_M(\rho_{v,s})} \quad (1)$$

where  $\rho_t$  and  $\rho_v$  denote the SNR of training data sets and validation data sets, respectively. Let  $\rho_s$  be the  $s$ -th SNR in different validation data sets of the number  $S$ .  $\text{BER}_R(\rho_t, \rho_{v,s})$  denotes BER achieved by RNN decoder, which trained at  $\rho_t$  on the data with  $\rho_{v,s}$ . Similarly,  $\text{BER}_M(\rho_{v,s})$  is the BER achieved by MAP decoding at  $\rho_v$ . As such, our purpose of setting up different  $\rho_t$  is to find the optimal  $\rho_t$  which results in the least NVE in training phase, and then the optimal  $\rho_t$  will be used for the testing phase.

- *Ratio of codebook set A*: To evaluate the generalization ability of RNN decoders, we choose the information bits  $\mathbf{x}$  which covers only part of A, that is, the ratio  $p$  of codebook set A. In the sequel, we set the ratio  $p = 40\%, 60\%, 80\%, 100\%$ , respectively.
- *Loss function*: Here, We choose the mean squared error (MSE) as the loss function, defined by

$$L_{\text{MSE}} = \frac{1}{K} \sum_{i=0}^{K-1} (x_i - \hat{x}_i)^2 \quad (2)$$

where  $x_i$  is the label and  $\hat{x}_i$  donates the  $i$ -th estimated information bit of RNN decoder.

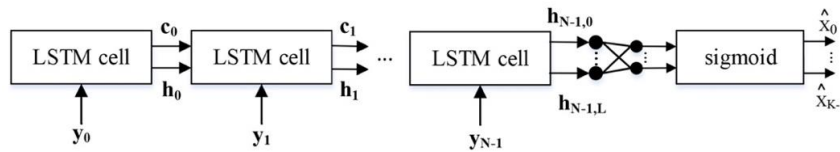
- *Dropout*: There exists a saturation length for RNN decoder. When the saturation length exceeds codeword length  $N$  and  $p < 100\%$ , the RNN decoder will be overfitting. When the saturation length is under the codeword length  $N$ , the RNN decoder will suffer from the problem of underfitting no matter what the  $p$  is. Based on this, we compare the BER performance of the decoders with dropout and without dropout, instead of adding only dropout like [5].

### 3. Proposed RNN decoders

In this section, we describe three types of proposed RNN decoders, which are built upon LSTM, GRU and Bi-GRU, respectively.

#### 3.1. LSTM decoder

In LSTM, gating mechanism is used to control the information flow such that the gradient vanishing problem in RNN is better handled, and the long range dependency is better captured. The flow diagram of our proposed LSTM decoder is shown in figure 2, which mainly consists of two parts, a LSTM cell and a fully connected layer with sigmoid activation function. Let  $M$  denote the Mini-batch size. The input of the LSTM decoder is  $\mathbf{y} = (\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{N-1})$ , where  $\mathbf{y}_t$  ( $t = 0, 1, \dots, N-1$ ) of length  $M$  is the input of LSTM cell at time  $t$ . The memory cell vector  $\mathbf{c}_{t-1}$  and cell output  $\mathbf{h}_{t-1}$  at time  $t-1$  will affect the memory cell vector  $\mathbf{c}_t$  and the cell output  $\mathbf{h}_t$  at time  $t$ .



**Figure 2.** The flow diagram of LSTM decoder.

The LSTM cell in figure 3 includes input gate, forget gate and output gate. The input gate (output gate) uses inputs from other memory cell to decide whether to store (access) certain information in its memory cell. The forget gate determines the importance of memory cell vector  $\mathbf{c}_{t-1}$  to the memory cell vector  $\mathbf{c}_t$ . The memory cell vector  $\mathbf{c}_t$  is updated by adding a new memory content vector  $\hat{\mathbf{c}}_t$  and partially forgetting the existing memory, where  $\hat{\mathbf{c}}_t$  describes the state of current LSTM cell input. The value of the  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ ,  $\mathbf{o}_t$ ,  $\hat{\mathbf{c}}_t$ ,  $\mathbf{c}_t$  and  $\mathbf{h}_t$  are computed by

$$\mathbf{i}_t = \sigma(\mathbf{W}_{yi}\mathbf{y}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i) \quad (3)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{yf}\mathbf{y}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1} + \mathbf{b}_f) \quad (4)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{yo}\mathbf{y}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{b}_o) \quad (5)$$

$$\hat{\mathbf{c}}_t = \tanh(\mathbf{W}_{yc}\mathbf{y}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (6)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t \quad (7)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (8)$$

where  $\sigma(x) = \frac{1}{1+e^{-x}}$ . The initial value  $\mathbf{c}_0 = \mathbf{0}$ ,  $\mathbf{h}_0 = \mathbf{0}$ , where the  $\mathbf{0}$  represents zero vector of length  $L$ .  $\odot$  stands for the element-wise product of two matrices or vectors. The hidden state in LSTM cell is defined as the concatenation of  $(\mathbf{h}_t, \mathbf{c}_t)$ , where the ‘slow’ state  $\mathbf{c}_t$  fights the vanishing gradient problem and the ‘fast’ state  $\mathbf{h}_t$  makes hard decision over short periods of time. Note that

final output of LSTM cell is  $\mathbf{h}_N$ , which is determined by the output gate  $\mathbf{o}_N$  and the memory cell vector  $\mathbf{c}_N$ . Finally, we add a layer of sigmoid activation function to get the decoded value  $\hat{\mathbf{x}}$  of length K by

$$\hat{\mathbf{x}} = \sigma(\mathbf{W}\mathbf{h}_N + \mathbf{b}) \quad (9)$$

Our goal is to train a set of parameter  $\{\mathbf{W}_{yi}, \mathbf{W}_{hi}, \mathbf{W}_{yf}, \mathbf{W}_{hf}, \mathbf{W}_{yo}, \mathbf{W}_{ho}, \mathbf{W}_{yc}, \mathbf{W}_{hc}, \mathbf{W}, \mathbf{b}_i, \mathbf{b}_f, \mathbf{b}_o, \mathbf{b}_c, \mathbf{b}\}$  to achieve the output  $\hat{\mathbf{x}}$  which is as close as possible to the label  $\mathbf{x}$ .

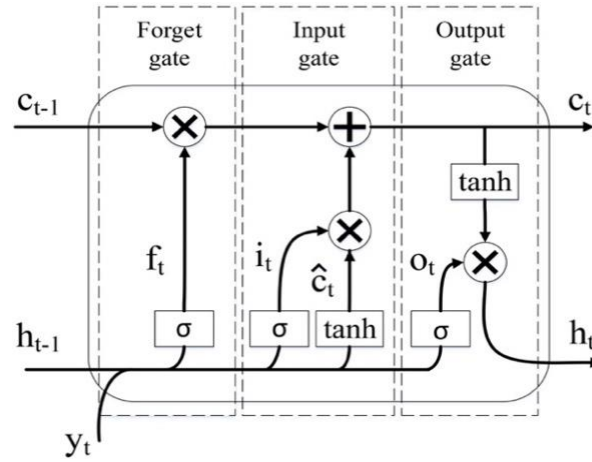


Figure 3. The LSTM cell

### 3.2. GRU decoder

Compared to LSTM decoder, the structure of GRU decoder is much simpler. GRU decoder and LSTM decoder both have gating mechanism to make the flow of information inside the unit. There is no separate memory cells in the GRU, which is diverse from the LSTM. The flow diagram of the proposed GRU decoder is shown in figure 4, which mainly consists of two parts, a GRU cell and a fully connected layer with sigmoid activation function. The input of the GRU decoder is also  $\mathbf{y} = (\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_{N-1})$ . The hidden state  $\mathbf{h}_{t-1}$  at time  $t - 1$  will affect the state of GRU cell at time  $t$ .

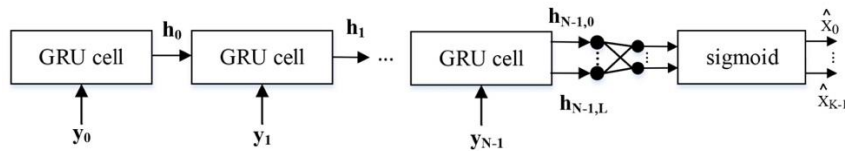


Figure 4. The flow diagram of GRU decoder.

The GRU cell in figure 5 includes two types of gates: update gate and reset gate. The hidden state  $\mathbf{h}_t$  is define as a linear interpolation between previous hidden state  $\mathbf{h}_{t-1}$  and the candidate hidden state  $\hat{\mathbf{h}}_t$  at time  $t$ . The formula is as follows:

$$\mathbf{h}_t = (\mathbf{I} - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \hat{\mathbf{h}}_t \quad (10)$$

where  $\mathbf{I} = (1, 1, \dots, 1)$ , whose dimension is L. An update gate is used to control how much the unit updates its hidden state. We compute  $\mathbf{z}_t$  and the candidate hidden state  $\hat{\mathbf{h}}_t$  by

$$\mathbf{z}_t = \sigma(\mathbf{W}_{hz}\mathbf{h}_{t-1} + \mathbf{W}_{yz}\mathbf{y}_t + \mathbf{b}_z) \quad (11)$$

$$\hat{\mathbf{h}}_t = \tanh(\mathbf{W}_y\mathbf{y}_t + \mathbf{W}_h(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h) \quad (12)$$

where  $\mathbf{r}_t$  decides how much the unit ignores its candidate hidden state  $\hat{\mathbf{h}}_t$ . The smaller the value of the  $\mathbf{r}_t$ , the more ignored.  $\mathbf{r}_t$  is computed similarly to the  $\mathbf{z}_t$ :

$$\mathbf{r}_t = \sigma(\mathbf{W}_{hr}\mathbf{h}_{t-1} + \mathbf{W}_{yr}\mathbf{y}_t + \mathbf{b}_r) \quad (13)$$

We also add a layer with full connection in the end, same with (9).

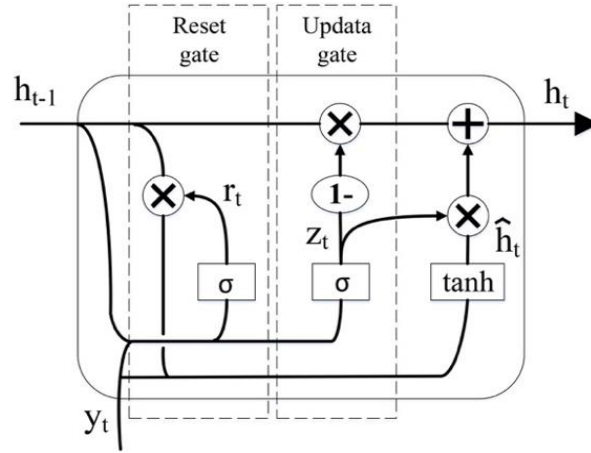


Figure 5. The GRU cell.

### 3.3. Bi-GRU decoder

Bi-GRU decoder is proposed to be able to exceed the LSTM decoder in BER performance. Since the GRU decoder only models the information flow in one direction, we use a Bi-GRU decoder to get information by summarizing from both directions. More specifically, one GRU goes from 1 to N and the other from N to 1 if given a sequence of length N.

As shown in figure 6, the structure of Bi-GRU is the combination of two different directional GRUs. At each time  $t$ , we provide both directional GRU the same input, which is same the input of GRU, and the output is determined by both two GRUs. Denote the hidden states of the forward and backward GRUs as  $\vec{h}_t$  and  $\overleftarrow{h}_t$  at time  $t$ , respectively. The cell output is computed by concatenating the two hidden states at each time.  $\vec{h}_t$  and  $\overleftarrow{h}_t$  are computed by (10)-(12), now we simply replace the above complex calculation with the function  $f$ :

$$\vec{h}_t = f(W_1 y_t + W_2 \vec{h}_{t-1} + \vec{b}) \quad (14)$$

$$\overleftarrow{h}_t = f(W_3 y_t + W_4 \overleftarrow{h}_{t-1} + \vec{b}) \quad (15)$$

$$\mathbf{o}_t = (\vec{h}_t, \overleftarrow{h}_t) \quad (16)$$

where the representation in (16) means to concatenate the forward and backward outputs. As before, we add a layer of with full connection as the last layer to obtain the decoded value of Bi-GRU decoder. The decoded value  $\hat{\mathbf{x}}$  is computed by

$$\hat{\mathbf{x}} = \sigma(W' \mathbf{o}_N + \mathbf{b}') \quad (17)$$

where N is also the total time and the dimension of  $\hat{\mathbf{x}}$  is also K.

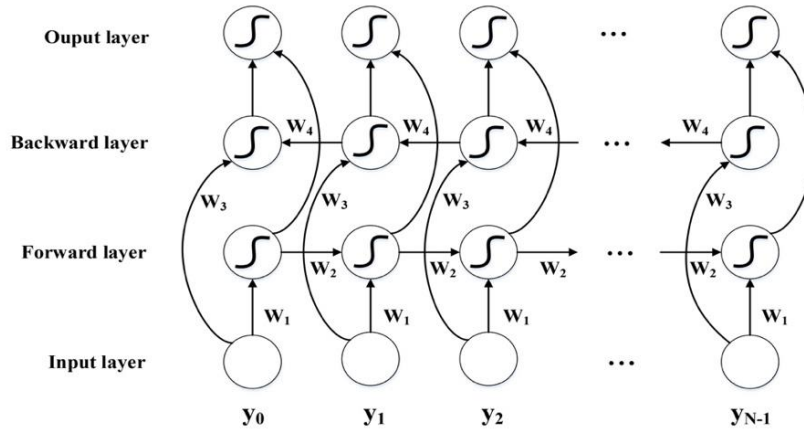


Figure 6. The architecture of Bi-GRU.

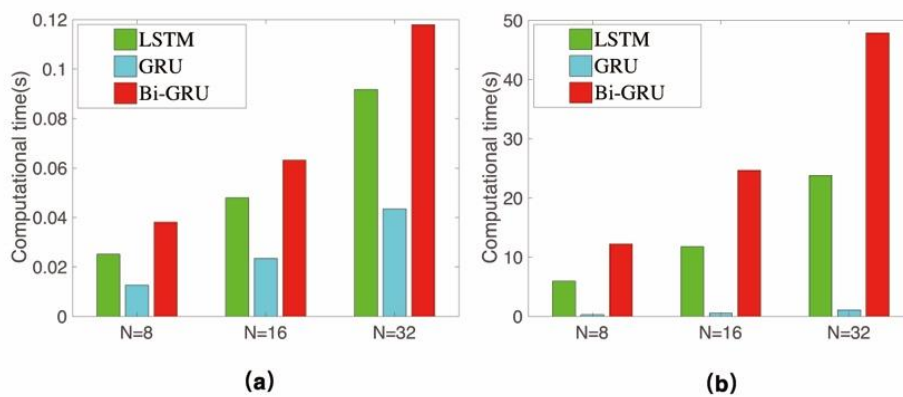
#### 4. Performance evaluation

We compare the performance of proposed three RNN decoders with LSTM, GRU and Bi-GRU of code length  $N=8, 32$ , respectively. We focus on the polar code of rate 1/2 for performance evaluation. In addition, in order to reduce the decoding complexity, we only use a single layer RNN as decoder. For more parameter settings, we select hyperparameters set, which is shown in table 1, for RNN decoders after many trials.

We first observe the computational time of these three decoders in figure 7. The backward propagation time for one training sample of LSTM decoder is almost twice that of the GRU decoder, and Bi-GRU decoder is about three times that of GRU decoder under the same  $N$ . For the forward propagation time for one testing sample under the same  $N$ , the LSTM decoder is almost 21 times longer than GRU decoder, and Bi-GRU decoder is twice as long as LSTM decoder. Therefore, the GRU decoder is the most time saving in computing.

**Table 1.** Hyperparameters setting

Item	Hyperparameters setting
Number of training samples	$10^6$
Number of testing samples	$10^5$
Total training epoch	$10^5$
Mini-batch size	128
Optimization method	Adam optimization
SNRs for training ( $\rho_t$ )	{2,0,2,4,6,8,10,12,14,16,18,20}dB
SNRs for testing ( $\rho_v$ )	{0,0.5,1,1.5,2,2.5,3,3.5,4,4.5,5,5.5,6}dB
Training ratio of codebook $p$	40%,60%,80%,100%
Initialization method	Xavier initialization
L	256

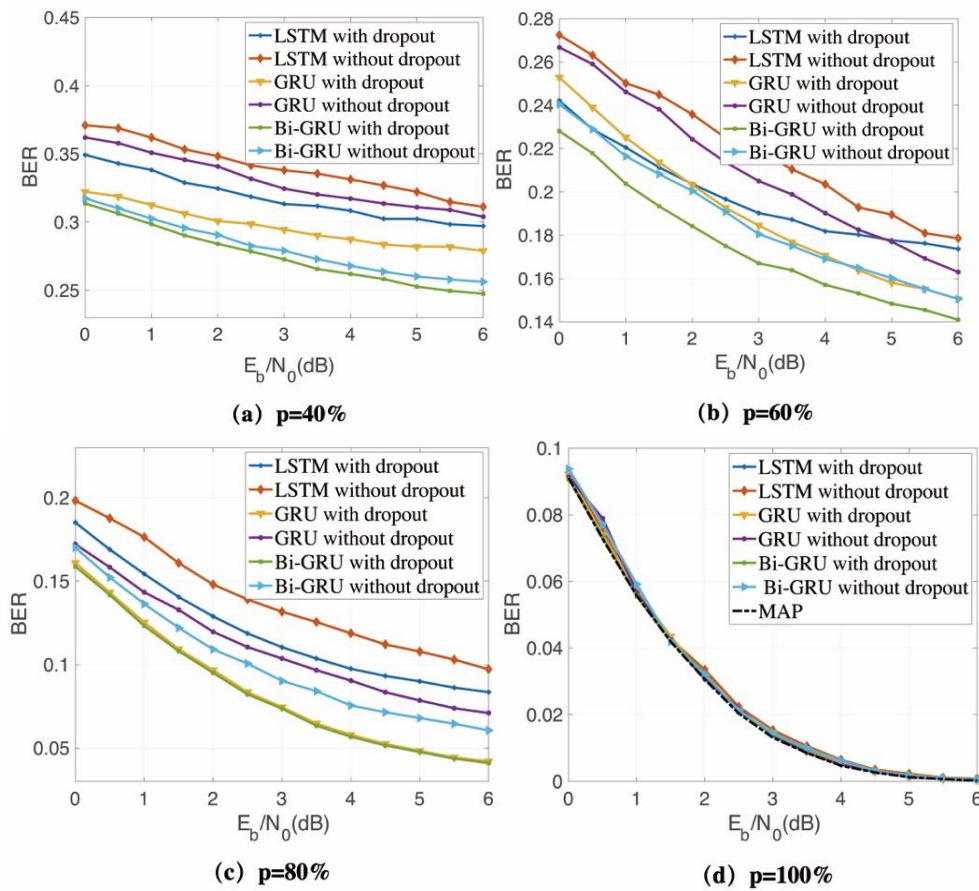


**Figure 7.** The computational time of LSTM, GRU and Bi-GRU under different length of codeword  $N$  with noise. (a) The time for one training sample through backward propagation. (b) The time for one testing sample through forward propagation.

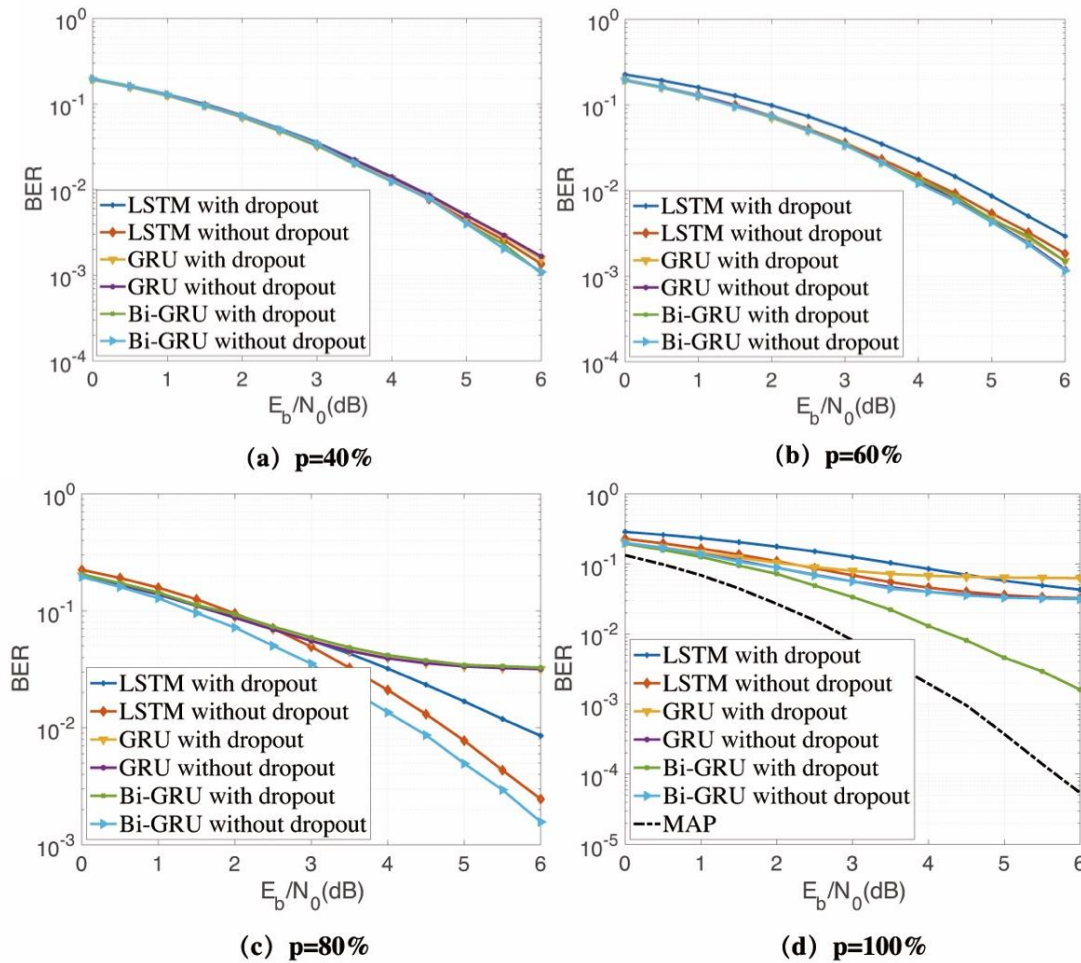
Figure 8. shows the case when  $N = 8$  with noise. We can see that all three RNN decoders can achieve MAP performance no matter dropout or not when  $p = 100\%$ . When  $p < 100\%$ , the BER performance of Bi-GRU decoder is significantly better than LSTM decoder and GRU decoder although they are all overfitting, and three RNN decoders with dropout are better than decoders without dropout.



Figure 9. shows the case when  $N = 32$  with noise. We can see that as  $p$  increases, the overall trend of BER for all RNN decoders increases, which can be interpreted as they are all underfitting. The Bi-GRU decoder always has a slight advantage in BER performance than LSTM decoder and GRU decoder under different  $p$ . The three RNN decoders without dropout are always better than the one with dropout when  $p = 40\%, 60\%, 80\%$ . In figure (d), the LSTM decoder without dropout has better BER performance than LSTM decoder with dropout and the GRU decoder without dropout has better BER performance than GRU decoder with dropout. But the Bi-GRU decoder with dropout makes a breakthrough in BER performance, which is even better than Bi-GRU without dropout, this also shows that Bi-GRU is prone to overfitting.



**Figure 8.** The BER performance achieved by LSTM, GRU and Bi-GRU when  $N = 8$  under different training ratio  $p = 40\%, 60\%, 80\%$  and  $100\%$ .



**Figure 9.** The BER performance achieved by LSTM, GRU and Bi-GRU when  $N = 32$  under different training ratio  $p = 40\%$ ,  $60\%$ ,  $80\%$  and  $100\%$ .

## 5. Conclusion

In this paper, we proposed three types of RNN decoders, named LSTM decoder, GRU decoder and Bi-GRU decoder, respectively. After extensive experiments, we found that the BER performance of GRU decoder, which had a simpler structure and less computational time, was not worse than the LSTM decoder. Although the Bi-GRU decoder always had the best BER performance, it was at the expense of more computational time and was prone to overfitting. We confirmed that the performance of RNN decoders with dropout was better than the decoders without dropout when decoding models were overfitting, while the RNN decoders without dropout was better than decoders with dropout when decoding models were underfitting.

## Acknowledgment

This work was supported by the National Natural Science Foundation of China (No. 61671143).

## References

- [1] 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 (CVPR)*, Las Vegas, Nevada, pp. 770-778, 2016.



- [2] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention based neural machine translation," *Conference on Empirical Methods in Natural Language Process*, Lisbon, Portugal, pp. 1412-1421, 2015.
- [3] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," *IEEE International Conference Acoustics, speech and signal Process (ICASSP)*, Vancouver, Canada, pp. 6645- 6649, 2013.
- [4] T. Gruber, S. Cammerer, J. Hoydis, and S. ten Brink, "On deep learning- based channel decoding," in *The 51st Annual Conference on Information Sciences and Systems (CISS)*. *IEEE*, Nicollet Mall Minneapolis, pp. 1- 6, 2017.
- [5] W. Lyu, Z. Zhang, C. Jiao, et al. "Performance Evaluation of Channel Decoding With Deep Neural Networks," in *IEEE International Conference on Communications (ICC)*, Kansas City, MO, USA, pp. 1-6, 2018.