

MACRO-BLOCK DROPOUT FOR IMPROVED REGULARIZATION IN TRAINING END-TO-END SPEECH RECOGNITION MODELS

Chanwoo Kim¹, Sathish Indurti¹, Jinhwan Park¹, and Wonyong Sung²

Samsung Research¹, Seoul, South Korea
Seoul National University², Seoul, South Korea

{chanw.com, s.indurti, jh0354.park}@samsung.com, wysung@snu.ac.kr

ABSTRACT

This paper proposes a new regularization algorithm referred to as *macro-block dropout*. The overfitting issue has been a difficult problem in training large neural network models. The dropout technique has proven to be simple yet very effective for regularization by preventing complex co-adaptations during training. In our work, we define a *macro-block* that contains a large number of units from the input to a Recurrent Neural Network (RNN). Rather than applying dropout to each unit, we apply random dropout to each *macro-block*. This algorithm has the effect of applying different drop out rates for each layer even if we keep a constant average dropout rate, which has better regularization effects. In our experiments using Recurrent Neural Network-Transducer (RNN-T), this algorithm shows relatively 4.30 % and 6.13 % Word Error Rates (WERs) improvement over the conventional dropout on LibriSpeech *test-clean* and *test-other*. With an Attention-based Encoder-Decoder (AED) model, this algorithm shows relatively 4.36 % and 5.85 % WERs improvement over the conventional dropout on the same test sets.

Index Terms: neural-network, regularization, macro-block, dropout, end-to-end speech recognition

1. INTRODUCTION

Deep learning technologies have significantly improved speech recognition accuracy recently [1, 2]. There have been series of remarkable changes in speech recognition algorithms during the past decade. These improvements have been obtained by the shift from Gaussian Mixture Model (GMM) to the Feed-Forward Deep Neural Networks (FF-DNNs), FF-DNNs to Recurrent Neural Network (RNN) such as the Long Short-Term Memory (LSTM) networks [3]. Thanks to these advances, voice assistant devices such as Google Home [4], Amazon Alexa and Samsung Bixby are

widely used at home environments. Tremendous amount of research has been conducted in the process of switching from a conventional speech recognition system consisting of an Acoustic Model (AM), a Language Model (LM), and a decoder based on a Weighted Finite State Transducer (WFST) to a complete end-to-end all-neural speech recognition system [5].

Such notable shifts happened not only in research on model architectures, but also in research on model robustness as well. In conventional approaches, researchers have focused on cleaning or transforming speech signals [6, 7, 8, 9, 10, 11] and features [12, 13]. However, it has been recently observed that data augmentation [14, 15] is especially powerful in enhancing the model robustness. Data augmentation during the training may be helpful in reducing the environmental mismatch between the training and testing time. The Small Power Boosting (SPB) algorithm [16] may be considered as an extreme example of reducing environmental mismatch by intentionally transforming portions of inputs which are more susceptible to noise or reverberation.

A large number of end-to-end speech recognition systems are based on the Attention-based Encoder-Decoder (AED) [5] and the Recurrent Neural Network-Transducer (RNN-T) [17] algorithms. These complete end-to-end systems have started outperforming conventional WFST-based decoders for large vocabulary speech recognition tasks [2]. Further improvements in these end-to-end speech recognition systems have been possible thanks to a better choice of target units such as Byte Pair Encoded (BPE) and *unigram language model* [18] subword units, and an improved training methodologies such as Minimum Word Error Rate (MWER) training [19]. In training such all neural network structures, over-fitting has been a major issue. For improved regularization in training, various approaches have been proposed [20]. Data-augmentation has been also proved to be useful in improving model training [4, 21, 22, 23]. The dropout approach [24] has been applied to overcome this issue in which both the input and the hidden units are randomly dropped out for regularization. This dropout approach has inspired a number of related approaches [25, 26, 27]. In *DropBlock* [28], the

Thanks to Samsung Electronics for funding this research. The authors are thankful to president Sebastian Seung, Executive Vice President (EVP) Daniel Lee, EVP Seunghwan Cho, and the Language and Voice Team (LVT) members at Samsung Research.

authors claimed that dropping out at random is not effective in removing semantic information when training Convolutional Neural Networks (CNNs) because nearby activations contain closely related information. Note that in CNN, filters are applied to nearby elements, thus activation units are spatially correlated. Motivated by this, they apply a square mask centered around each zero position. However, since this kind of spatial correlation does not hold for fully connected feed-forward layers or RNNs, the application of *DropBlock* is limited to CNNs.

In this paper, we present a new regularization algorithm referred to as *macro-block dropout*. We define a *macro-block* that contains multiple input units to a neural network layer. Rather than applying dropout to each unit, we apply random dropout to each *macro-block*. In our experiments using an RNN-T [17] and an Attention-based Encoder Decoder (AED) in Sec. 5, this simple algorithm has shown a quite significant improvement over the conventional dropout approach. Our contribution in this paper may be summarized as follows. First, by using large macro-blocks, the ratio of dropped input units in each layer has large variation even if we keep a constant dropout rate. This variation in the ratio of dropped units leads to better regularization. Thus, unlike *DropBlock* that relies on spatial correlation in CNNs, we may apply *macro-block dropout* to any kinds of neural networks such as Feed-Forward (FF) networks and RNNs. To the best of our knowledge, our work is the first in using big chunks consisting of large number of input units to RNNs for masking. Second, we propose a new way of input scaling after applying a mask in *macro-block dropout*. This new scaling approach is related to the fact that the portion of dropped units significantly varies in *macro-block dropout*. Third, we proposed a very low-cost regularization algorithm. As will be seen in Sec 5, the best performance is achieved when the number of macro blocks for each layer is only four. This means that we only need to generate four random values for each layer. Thus, this *macro-block dropout* is very simple with very small computational requirement.

2. RELATED WORKS

Dropout is a simple regularization technique to alleviate the overfitting problem by preventing co-adaptations [24]. When the dimension of the input to a neural network layer is d_x , we create a random mask tensor \mathbf{m} with the same dimension d_x . Each element $m \in \mathbf{m}$ follows the Bernoulli distribution $m \sim \text{Bernoulli}(1 - q)$, where q is the dropout rate. Given an input \mathbf{x} , the dropout output \mathbf{x}_{out} is obtained by the following equation:

$$\mathbf{x}_{\text{out}} = \frac{\mathbf{x} \odot \mathbf{m}}{1 - q}, \quad (1)$$

where \odot is a Hadamard product. The scaling by $\frac{1}{1-q}$ is applied to keep the sum of input units the same through this

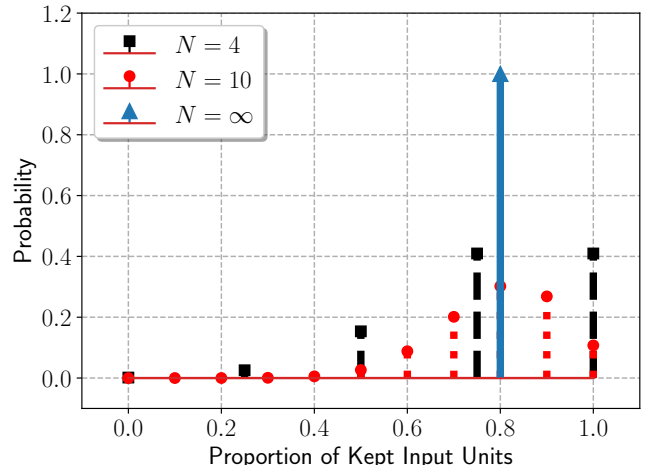
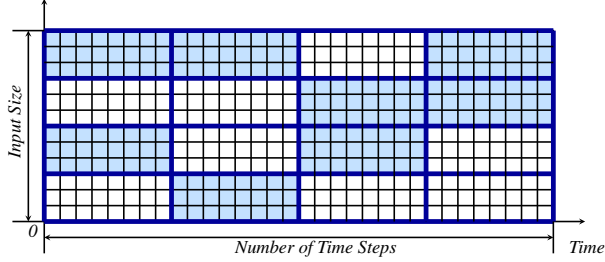


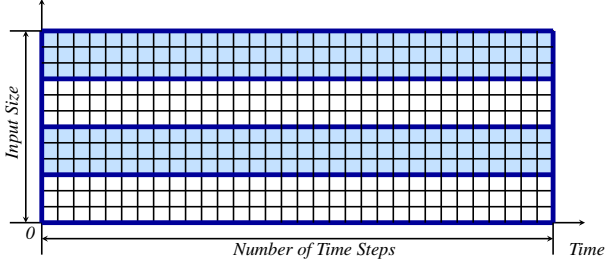
Fig. 1: The Probability Mass Function (PMF) representing the ratio of kept input units when the numbers of macro blocks are $N = 4$ and $N = 10$, respectively. In this case, the dropout rate q is assumed to be 0.2. When the conventional dropout is applied, the total number of input units to an RNN layer is usually from tens of thousands up to several millions. In this case, the Probability Density Function (PDF) can be approximated by a delta function, which is represented by the plot with the legend label of $N = \infty$.

masking process. In (1), we adopt the “*inverted dropout*” approach rather than the original form of *dropout* in [24] where the scaling is performed during the inference time. Before introducing our *macro-block dropout*, let us consider the variance of the dropped ratio of an input layer. When the total number of input units to a single RNN is N , the expected ratio of kept units is given by $\frac{(1-q)N}{N} = 1 - q$, and the standard deviation is given by $\sqrt{\frac{q(1-q)}{N}}$. From the central limit theorem, we know that this distribution can be approximated by a Gaussian distribution. If we use typical values of $q = 0.2$ and $N = 10^5$, the standard deviation becomes 0.00126, which is very small. As shown in Fig. 1, when there are a large number of units, the distribution is very similar to a delta function centered at $1 - q$. Since the ratio of kept units is always very close to $1 - q$, we conclude that $\frac{\sum \mathbf{x} \odot \mathbf{m}}{\sum \mathbf{x}} \approx 1 - q$. Thus we can safely use the fixed scaling of $\frac{1}{1-q}$ to keep the sum of the input units the same after dropout in (1). Dropout has been turned out to be especially useful in improving the training of dense network models for image classification [29], speech recognition [30], and so on. This dropout approach inspired many other related approaches such as *DropConnect* [25], *drop-path* [26], *shake-shake* [27], *ShakeDrop* [31], and *DropBlock* [28] regularizations.

The authors of *DropBlock* [28] claim that dropping out



(a) Application of two-dimensional *macro-block dropout* to the input of an RNN.



(b) Application of one-dimensional *macro-block dropout* to the input of an RNN.

Fig. 2: Application of *macro-block dropout* to the inputs to a Recurrent Neural Network (RNN) layer: (2a) Two-dimensional and (2b) One-dimensional *macro-block dropout* cases. Each tiny rectangle defined by the grid corresponds to each input unit. Larger rectangular chunks are *macro-blocks*. Region in the light blue color represent *macro-blocks* to be dropped out.

inputs to CNNs at random is not effective in removing spatially correlated information. In *DropBlock*, the *zero position*, which is the center of a square mask, is randomly sampled in the same way as the *dropout* in [24]. It has been reported that *DropBlock* significantly outperforms the baseline *dropout* when applied to CNNs. However, unlike our *macro-block dropout*, Bernoulli random numbers are generated for each unit of inputs to neural networks. Thus, the ratio of dropped units has very small variation similar to that of the conventional dropout case, which is one of the biggest differences from the proposed *Macro-block dropout*.

3. MACRO-BLOCK DROPOUT

3.1. Definition of a macro-block

Let us consider a two-dimensional data \mathbf{x} with the dimension of $\mathbf{d}_x = (N_x, N_y)$ that is the input to a neural-network layer. *Macro-blocks* are constructed by equally partitioning this input \mathbf{d}_x along each axis into a new dimension of $\mathbf{d}_{(\text{par})} = (P_x, P_y)$ with the constraint of $N_x = P_x M_x$ and $N_y = P_y M_y$ where M_x and M_y are macro block sizes along the x and y axes.

Fig. 2 illustrates examples of *macro-blocks* when this

algorithm is applied to the input of an RNN. As shown in this figure, the input of an RNN layer is a two-dimensional variable with the dimension of (number_of_time_steps, input_size). Fig. 2a shows the case when this two-dimensional region is partitioned by $\mathbf{d}_{(\text{par})} = (4, 4)$. We refer to this approach as the two-dimensional *macro-block* approach. As another example, we may also consider the case when $\mathbf{d}_{(\text{par})} = (1, 4)$, which is shown in Fig. 2b. In this case, the mask pattern does not change along the time axis. We refer to this approach as the one-dimensional *macro-block* approach.

Unlike the conventional *dropout* or the *DropBlock* approach described in Sec. 2, since there is only a very small number of macro-blocks in the input, the ratio of kept input units in a single layer has very large variation as shown by the stem plot with the legend label of $N = 4$ in Fig. 1. For example, when $\mathbf{d}_{(\text{par})} = (1, 4)$ and $q = 0.2$, with the probability of $0.8^4 = 0.4096$, dropout will not happen at all for a single layer. With the probability of $\binom{4}{2} 0.8^2 0.2^2 = 0.1536$, half of the input units will be dropped out. The standard deviation of kept ratio is given by $\sqrt{\frac{q(1-q)}{4}} = 0.2$. We believe that such large variation is the reason why *macro-block dropout* works for RNNs while *DropBlock* is limited to CNNs.

The Word Error Rates (WERs) using an RNN-T model with the one-dimensional and the two-dimensional *macro-block dropout* approaches are summarized in Table 1. The structure of this RNN-T model and the experimental configuration are described in detail in Sec. 4 and Sec. 5, respectively. From this result, we conclude that the one-dimensional *macro-block dropout* approach is more effective than two-dimensional approach for RNNs. This observation is similar to what is mentioned in [32] for the baseline dropout. The masking using the one-dimensional approach corresponds to “the same weight realization for a single input”, which is consistent with the Bayesian interpretation of *dropout* for RNNs. This improvement has been commonly found in other literatures as well [33].

In obtaining this result, we choose the partition dimensions $\mathbf{d}_{(\text{par})} = (1, 4)$ and $\mathbf{d}_{(\text{par})} = (4, 4)$ for one- and two-dimensional cases, respectively. We choose these dimensions since the best WERs for one- and two-dimensional *macro-block dropout* cases are obtained with these dimensions in our experiments on the *LibriSpeech* corpus.

3.2. Application of dropout to macro-blocks

Having defined the required terms in Sec. 3.1, we proceed to explain the algorithm in detail in this section. The entire algorithm is summarized in Algorithm 1. During the inference time, *macro-block dropout* is not applied as the original dropout. During the training time, we create a random tensor \mathbf{r} whose dimension is $\mathbf{d}_{(\text{par})}$. This tensor is created from the *Bernoulli* distribution with the probability of one given by $1 - q$, where q is the dropout probability. This \mathbf{r} is then resized to match the dimension of the input \mathbf{x} . For simplic-

Algorithm 1 Macro-block Dropout

```
1: Input: Inputs to a layer:  $\mathbf{x}$ , the dimension for partitioning:  $d_{(\text{par})}$ , dropout_rate  $q$ , mode  
2: if mode == Inference then  
3:   return  $\mathbf{x}$   
4: end if  
5: Creates a random tensor  $\mathbf{r}$  with a dimension of  $d_{(\text{par})}$ :  
6:   For each element  $r$  of  $\mathbf{r}$ ,  $r \sim \text{Bernoulli}(1 - q)$ .  
7: Creates a masking tensor  $\mathbf{m}$  by resizing  $\mathbf{r}$  using the nearest-neighbour method to match the dimension of  $\mathbf{x}$ .  
8: Applies the mask:  
9:    $\mathbf{x}_m = \mathbf{x} \odot \mathbf{m}$ .  
10: Obtains the output  $\mathbf{x}_{\text{out}}$  by scaling  $\mathbf{x}_m$  :  
11:    $\mathbf{x}_{\text{out}} = \left| \frac{\sum_{\text{all elements}} \mathbf{x}}{\sum_{\text{all elements}} \mathbf{x} \odot \mathbf{m}} \right| \mathbf{x}_m$ .
```

Table 1: Word Error Rates (WERs) with the RNN-T model described in Sec. 4 using the one-dimensional *macro-block dropout* of $d_{(\text{par})} = (1, 4)$, and the two-dimensional *macro-block dropout* of $d_{(\text{par})} = (4, 4)$. In these experiments, the dropout rate of 0.2 is used since the best WER in each case is obtained at this rate.

Test Set	Baseline Dropout	Macro-Block Dropout	
		1-D : (1, 4)	2-D : (4, 4)
test-clean	3.95 %	3.78 %	3.92 %
test-other	12.23 %	11.48 %	11.50 %

ity, this resize operation is performed using the well-known *nearest-neighborhood* interpolation approach.

The scaling factor r is given by the following equation:

$$r = \left| \frac{\sum_{\text{all elements}} \mathbf{x}}{\sum_{\text{all elements}} \mathbf{x} \odot \mathbf{m}} \right|. \quad (2)$$

We apply the absolute value operation in (2), because the sign of the numerator and the denominator of (2) may be different when \mathbf{x} is the output of an RNN such as an LSTM or a GRU. More specifically, the hidden output of an LSTM is given by the following equation [3, 1]:

$$\mathbf{h}_{[m]} = \mathbf{o}_{[m]} \odot \sigma_h(\mathbf{c}_{[m]}), \quad (3)$$

where m is a time index, \odot is the Hadamard product, $\sigma_h(\cdot)$ is the hyperbolic tangent function, $\mathbf{o}_{[m]}$ is the output-gate value, and $\mathbf{c}_{[m]}$ is the cell value, respectively. From (3), it is obvious that $\mathbf{h}_{[m]}$ may have both positive and negative values, since the range of σ_h is between -1 and 1. In our speech recognition experiments, it is observed that performance is slightly worse if this absolute value operation is not applied in (2). In performing division in (2), we employ “a safe division” implemented by the `tf.math.divide_no_nan` Tensorflow [34] API to prevent division by zero.

We observe that the scaling in (2) is more effective than the simple scaling of $\frac{1}{1-q}$ used in the baseline dropout in (1) since $\frac{\sum \mathbf{x} \odot \mathbf{m}}{\sum \mathbf{x}} \approx 1 - q$ does not hold with *macro-block dropout* because of the large variance in the ratio of kept input units. Table 2 summarizes WERs obtained with the conventional scaling of $\frac{1}{1-q}$ and the scaling in (2) on the *LibriSpeech test-clean* and *test-other* sets. We use the RNN-T model that will be described in Sec. 5. For *macro-block dropout*, we employ the one-dimensional approach with the partition dimension of $d_{(\text{par})} = (1, 4)$. The experimental configuration in obtaining these results will be described in Sec. 5.

Table 2: Word Error Rates (WERs) with the RNN-T model described in Sec. 4 using the scaling suggested by (2) and $\frac{1}{1-q}$. The dropout rate is 0.2 and the partition dimension for the 1-dimensional *macro-block dropout* is $d_{(\text{par})} = (1, 4)$.

Test Set	Baseline Dropout	1-D Macro-Block Dropout: (1, 4)	
		Scaling using (2)	Scaling using $\frac{1}{1-q}$
test-clean	3.95 %	3.78 %	4.04 %
test-other	12.23 %	11.48 %	11.50 %

4. SPEECH RECOGNITION MODEL

For speech recognition experiments, we employed an RNN-T speech recognizer and an Attention-based Encoder Decoder (AED). Our speech recognition system is built *in-house* using Keras models [35] implemented for Tensorflow 2.3 [34]. The RNN-T structures have three major components: an encoder (also known as a transcription network), a prediction network, and a joint network. In our implementation, the

Table 3: Word Error Rates (WERs) with the RNN-T model and the Attention-based Encoder-Decoder (AED) model described in Sec. 4 using the baseline dropout and the one-dimensional *macro-block dropout* approaches. In these experiments, the dropout rate of 0.2 is used since the best WERs for both approaches are obtained at this rate.

Model	Test Set	Baseline Dropout	One-Dimensional Macro-Block Dropout			
			with different $\mathfrak{d}_{(\text{par})}$			
			(1, 3)	(1, 4)	(1, 5)	(1, 10)
RNN-T	test-clean	3.95 %	4.11 %	3.78 %	3.88 %	3.94 %
	test-other	12.23 %	11.57 %	11.48 %	11.52 %	11.50 %
Attention-based	test-clean	3.67 %	3.66 %	3.51 %	3.54 %	3.61 %
Encoder Decoder	test-other	11.62 %	11.20 %	10.94 %	10.98 %	11.07 %

encoder consists of six layers of bi-directional LSTMs interleaved with 2:1 max-pooling layers [36] in the bottom three layers. Thus, the overall temporal sub-sampling factor is 8:1 because of these three 2:1 max-pooling layers. The prediction network consists of two layers of uni-directional LSTMs. The unit size of all these LSTM layers is 1024. *Macro-block dropout* is applied to all the inputs of each LSTM layer of the encoder and the prediction network except the first layer of the encoder. From the transducer output, a linear embedding vector with a dimension of 621 is obtained and fed back into the prediction network. The AED model has three components: an encoder, a decoder, and an attention block. In our implementation, the encoder of the AED model is identical to the encoder structure of the RNN-T model explained above. As a decoder, we use a single layer of uni-directional LSTM whose unit size is 1024.

The loss employed for training the RNN-T model is a combination of the Connectionist Temporal Classification (CTC) loss [37] applied to the *encoder* output and the RNN-T loss [17] applied to the full network, which is represented by the following:

$$\mathbb{L}_{\text{CTC-RNN-T}} = \mathbb{L}_{\text{CTC}} + \mathbb{L}_{\text{RNN-T}}. \quad (4)$$

We refer to this loss in (4) as the *joint CTC-RNN-T loss*. For the AED model, we employ the joint *CTC-CE loss*, which is given by:

$$\mathbb{L}_{\text{CTC-CE}} = \mathbb{L}_{\text{CTC}} + \mathbb{L}_{\text{CE}}, \quad (5)$$

as in [1]. For better stability during the training, we use the gradient clipping by global norm [38], which is implemented in Tensorflow as the `tf.clip_by_global_norm` API.

To obtain further improvement in speech recognition accuracy, we incorporate an improved shallow-fusion technique. In this approach, we subtract log prior probabilities of each label obtained from the transcript of the speech recognition training database. This idea is initially described in [39].

Our formulation is based on our earlier work in [40]:

$$\log p_{\text{sf}}(y_l | \mathbf{x}_{[0:m]}) = \log p(y_l | \mathbf{x}_{[0:m]}, \hat{y}_{0:l}) - \lambda_p \log p_{\text{prior}}(y_l) + \lambda_{\text{lm}} \log p_{\text{lm}}(y_l | \hat{y}_{0:l}), \quad (6)$$

where y_l is the output at the output label index l , $\mathbf{x}_{[0:m]}$ is the input feature vector sequence from the zero-th frame up to the $(m-1)$ -st frame, and $\hat{y}_{0:l}$ is the estimated output label sequence from the output index zero up to $l-1$. λ_p and λ_{lm} are constant weighting coefficients for the log prior probability denoted by $\log p_{\text{prior}}(y_l)$ and the log probability from the Language Model (LM) denoted by $\log p_{\text{lm}}(y_l | \hat{y}_{0:l})$, respectively. $\log p(y_l | \mathbf{x}_{[0:m]}, \hat{y}_{0:l})$ is the log probability from the speech recognition model. In our experiments, we use λ_p of 0.002 and λ_{lm} of 0.48 respectively as in [40].

5. EXPERIMENTAL RESULTS

Table 4: Word Error Rates (WERs) with the Attention-based Encoder Decoder (AED) model described in Sec. 4 with the improved shallow fusion in (6) with a Transformer LM [41].

Test Set	Baseline Dropout	Macro-Block Dropout
test-clean	2.44 %	2.37 %
test-other	7.87 %	7.42 %

In this section, we explain experimental results using the *macro-block dropout* approach with the RNN-T and the AED model described in Sec. 4. For training, we used the entire 960 hours LibriSpeech [42] training set consisting of 281,241 utterances. For evaluation, we used the official 5.4 hours test-clean and 5.1 hours test-other sets consisting of 2,620 and 2,939 utterances respectively. The pre-training stage has some similarities to our previous work in [43]. In this pre-training stage, the number of LSTM layers in the

encoder increased at every 10,000-steps starting from two LSTM layers up to six layers. We use an Adam optimizer [44] with the initial learning rate of 0.0003, which is maintained for the entire pre-training stage and one full epoch after finishing the pre-training stage. After this step, this learning rate decreases exponentially with a decay rate of 0.5 for each epoch. $x[m]$ and y_l are the input *power-mel filterbank* feature of size 40 and the output label, respectively. m is the input frame index and l is the decoder output step index. We use the *power-mel filterbank* feature instead of the more commonly used *log-mel filterbank* feature based on our previous experimental results [40, 45, 43]. In our *power-mel filterbank* feature, we employ the power-law nonlinearity with the power coefficient of $\frac{1}{15}$, which is suggested in [13, 46]. For better regularization in training, we apply the *SpecAugment* as a data-augmentation technique in all the experiments [22]. In our experiments with different dropout rates ranging from 0.0 to 0.5, for both the conventional *dropout* and the *macro-block dropout* approaches, the best WERs are obtained when the dropout rate q is 0.2.

Table 3 summarizes the experimental results with conventional dropout and *macro-block dropout* approaches using the RNN-T and AED models. In the case of the *macro-block dropout* approach, we conducted experiments with four different partition sizes ($d_{(\text{par})} = \{(1, 3), (1, 4), (1, 5), (1, 10)\}$) with the one-dimensional masking pattern that is shown in Fig. 2b. From this table, we observe that the best WER is obtained when the number of blocks is four ($d_{(\text{par})} = (1, 4)$). As shown in this table, when the RNN-T model is employed, the *macro-block dropout* algorithm shows relatively 4.30 % and 6.13 % WER improvements over the conventional dropout approach on the *LibriSpeech test-clean* and *test-other* sets, respectively. In the same table, we observe that the performance improvement using the AED model is also similar. We obtain 4.36 % and 5.85 % relative WER improvements on the same sets, respectively. Finally, we apply the improved shallow fusion in 6 to further improve the performance. Table 4 shows WERs obtained with the AED model using the improved shallow fusion in (6) with a Transformer LM [41]. As shown in this Table, *macro-block dropout* shows 2.86 % and 5.72 % relative WER improvement on the *LibriSpeech test-clean* and *test-other* sets.

6. CONCLUSIONS

In this paper, we describe a new regularization algorithm referred to as *macro-block dropout*. In this approach, rather than applying dropout to each input unit to RNN layers, we apply a random mask to a bigger chunk referred to as a *macro-block*. We propose an improved way of performing scaling for better performance with *macro-block dropout* in Sec. 3.2. In experiments using RNN-T and AED models, we obtain significantly better results with *macro-block dropout* compared

to the conventional dropout approach. We compare the variance of the ratio of input units that are not dropped with the conventional dropout approach and with *macro-block dropout* approach. We observe that this variance is significantly larger with *macro-block dropout*. We believe this larger variance helps regularization during training.

7. REFERENCES

- [1] C. Kim, D. Gowda, D. Lee, J. Kim, A. Kumar, S. Kim, A. Garg, and C. Han, "A review of on-device fully neural end-to-end automatic speech recognition algorithms," in *2020 54th Asilomar Conference on Signals, Systems, and Computers*, Nov. 2020, pp. 277–283.
- [2] C.-C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar, P. Nguyen, Z. Chen, A. Kannan, R. J. Weiss, K. Rao, E. Gonina, N. Jaitly, B. Li, J. Chorowski, and M. Bacchiani, "State-of-the-art speech recognition with sequence-to-sequence models," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 4774–4778.
- [3] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, no. 9, pp. 1735–1780, Nov. 1997.
- [4] C. Kim, A. Misra, K. Chin, T. Hughes, A. Narayanan, T. N. Sainath, and M. Bacchiani, "Generation of large-scale simulated utterances in virtual rooms to train deep-neural networks for far-field speech recognition in google home," in *Proc. Interspeech 2017*, 2017, pp. 379–383. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2017-1510>
- [5] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 577–585. [Online]. Available: <http://papers.nips.cc/paper/5847-attention-based-models-for-speech-recognition.pdf>
- [6] E. A. P. Habets, J. Benesty, I. Cohen, S. Gannot, and J. Dmochowski, "New insights into the mvdr beamformer in room acoustics," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 1, pp. 158–170, Jan 2010.
- [7] Y. Ephraim and H. Van Trees, "A signal subspace approach for speech enhancement," *IEEE Transactions on Speech and Audio Processing*, vol. 3, no. 4, pp. 251–266, 1995.
- [8] C. Kim and R. M. Stern, "Power function-based power distribution normalization algorithm for robust speech recognition," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2009, pp. 188–193.
- [9] C. Kim, K. Chin, M. Bacchiani, and R. M. Stern, "Robust speech recognition using temporal masking and thresholding algorithm," in *INTERSPEECH-2014*, Sept. 2014, pp. 2734–2738.
- [10] S. Mun and S. Shon, "Domain mismatch robust acoustic scene classification using channel information conversion," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 845–849.

- [11] J. S. Chung, J. Huh, and S. Mun, "Delving into VoxCeleb: Environment Invariant Speaker Recognition," in *Proc. Odyssey 2020 The Speaker and Language Recognition Workshop*, 2020, pp. 349–356. [Online]. Available: <http://dx.doi.org/10.21437/Odyssey.2020-49>
- [12] U. H. Yapanel and J. H. L. Hansen, "A new perceptually motivated MVDR-based acoustic front-end (PMVDR) for robust automatic speech recognition," *Speech Communication*, vol. 50, no. 2, pp. 142–152, Feb. 2008.
- [13] C. Kim and R. M. Stern, "Power-Normalized Cepstral Coefficients (PNCC) for Robust Speech Recognition," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, pp. 1315–1329, July 2016.
- [14] C. Kim, T. Sainath, A. Narayanan, A. Misra, R. Nongpiur, and M. Bacchiani, "Spectral distortion model for training phase-sensitive deep-neural networks for far-field speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 5729–5733.
- [15] C. Kim, E. Variiani, A. Narayanan, and M. Bacchiani, "Efficient implementation of the room simulator for training deep neural network acoustic models," in *INTERSPEECH-2018*, Sept 2018, pp. 3028–3032. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2018-2566>
- [16] C. Kim, K. Kumar and R. M. Stern, "Robust speech recognition using small power boosting algorithm," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2009, pp. 243–248.
- [17] A. Graves, A. rahman Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, May 2013, pp. 6645–6649.
- [18] T. Kudo, "Subword regularization: Improving neural network translation models with multiple subword candidates," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics, Jul. 2018, pp. 66–75. [Online]. Available: <https://www.aclweb.org/anthology/P18-1007>
- [19] R. Prabhavalkar, T. N. Sainath, Y. Wu, P. Nguyen, Z. Chen, C. Chiu, and A. Kannan, "Minimum word error rate training for attention-based sequence-to-sequence models," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, April 2018, pp. 4839–4843.
- [20] C. Kim, K. Kim, and S. Indurthi, "Small energy masking for improved neural network training for end-to-end speech recognition," in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2020, pp. 7684–7688.
- [21] C. Kim, A. Garg, D. Gowda, S. Mun, and C. Han, "Streaming end-to-end speech recognition with jointly trained neural feature enhancement," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 6773–6777.
- [22] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition," in *Proc. Interspeech 2019*, 2019, pp. 2613–2617. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2019-2680>
- [23] S. Mun, S. Park, D. K. Han, and H. Ko, "Generative adversarial network based acoustic scene training set augmentation and selection using svm hyper-plane." in *DCASE*, 2017, pp. 93–102.
- [24] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014. [Online]. Available: <http://jmlr.org/papers/v15/srivastava14a.html>
- [25] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, and R. Fergus, "Regularization of neural networks using dropconnect," ser. Proceedings of Machine Learning Research, S. Dasgupta and D. McAllester, Eds., vol. 28, no. 3. Atlanta, Georgia, USA: PMLR, 17–19 Jun 2013, pp. 1058–1066. [Online]. Available: <http://proceedings.mlr.press/v28/wan13.html>
- [26] G. Larsson, M. Maire, and G. Shakhnarovich, "Fractalnet: Ultra-deep neural networks without residuals," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=S1VaB4cex>
- [27] X. Gastaldi, "Shake-shake regularization of 3-branch residual networks," in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. OpenReview.net, 2017. [Online]. Available: <https://openreview.net/forum?id=HkO-PCmYl>
- [28] G. Ghiasi, T.-Y. Lin, and Q. V. Le, "Dropblock: A regularization method for convolutional networks," in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Inc., 2018, pp. 10 727–10 737. [Online]. Available: <http://papers.nips.cc/paper/8271-dropblock-a-regularization-method-for-convolutional-networks.pdf>
- [29] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105. [Online]. Available: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [30] G. E. Dahl, T. N. Sainath, and G. E. Hinton, "Improving deep neural networks for lvcsr using rectified linear units and dropout," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 8609–8613.
- [31] Y. Yamada, M. Iwamura, T. Akiba, and K. Kise, "Shakedrop regularization for deep residual learning," *IEEE Access*, vol. 7, pp. 186 126–186 136, 2019.
- [32] Y. Gal and Z. Ghahramani, "A theoretically grounded application of dropout in recurrent neural networks," in *Advances in Neural Information Processing Systems 29*, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds. Curran Associates, Inc., 2016, pp. 1019–1027. [Online]. Available: <http://papers.nips.cc/paper/6241-a-theoretically-grounded-application-of-dropout-in-recurrent-neural-networks.pdf>
- [33] T. Moon, H. Choi, H. Lee, and I. Song, "Rnndrop: A novel dropout for rnns in asr," in *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, 2015, pp. 65–70.

- [34] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, “Tensorflow: A system for large-scale machine learning,” in *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. Savannah, GA: USENIX Association, 2016, pp. 265–283. [Online]. Available: <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi>
- [35] F. Chollet *et al.*, “Keras,” <https://keras.io>, 2015.
- [36] M. Ranzato, F. J. Huang, Y. Boureau, and Y. LeCun, “Unsupervised learning of invariant feature hierarchies with applications to object recognition,” in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [37] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd International Conference on Machine Learning*, ser. ICML ’06. New York, NY, USA: ACM, 2006, pp. 369–376. [Online]. Available: <http://doi.acm.org/10.1145/1143844.1143891>
- [38] R. Pascanu, T. Mikolov, and Y. Bengio, “On the difficulty of training recurrent neural networks,” in *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28*, ser. ICML’13. JMLR.org, 2013, pp. III–1310–III–1318. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3042817.3043083>
- [39] N. Kanda, X. Lu, and H. Kawai, “Maximum-a-posteriori-based decoding for end-to-end acoustic models,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 5, pp. 1023–1034, 2017.
- [40] C. Kim, M. Shin, A. Garg, and D. Gowda, “Improved vocal tract length perturbation for a state-of-the-art end-to-end speech recognition system,” in *INTERSPEECH-2019*, Graz, Austria, Sept. 2019, pp. 739–743. [Online]. Available: <http://dx.doi.org/10.21437/Interspeech.2019-3227>
- [41] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 5998–6008. [Online]. Available: <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>
- [42] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [43] C. Kim, S. Kim, K. Kim, M. Kumar, J. Kim, K. Lee, C. Han, A. Garg, E. Kim, M. Shin, S. Singh, L. Heck, and D. Gowda, “End-to-end training of a large vocabulary end-to-end speech recognition system,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2019, pp. 562–569.
- [44] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [45] C. Kim, M. Kumar, K. Kim, and D. Gowda, “Power-law non-linearity with maximally uniform distribution criterion for improved neural network training in automatic speech recognition,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, Dec. 2019, pp. 988–995.
- [46] C. Kim, “Signal processing for robust speech recognition motivated by auditory processing,” Ph.D. dissertation, Carnegie Mellon University, Pittsburgh, PA USA, Dec. 2010.