

Block Refinement Learning for Improving Early Exit in Autoregressive ASR

Naotaka Kawata*, Shota Orihashi*, Satoshi Suzuki*, Tomohiro Tanaka*, Mana Ihori*, Naoki Maikishima*,
Taiga Yamane*, and Ryo Masumura*

* NTT Human Informatics Laboratories, NTT Corporation, Japan
E-mail: naotaka.kawata@ntt.com

Abstract—While the autoregressive transformer models of automatic speech recognition (ASR) are highly accurate, the inference time is long because of their sequential decoding. Early exit is a technique that aims to speed up the inference process by terminating it early on the basis of output from intermediate decoder layers. When an inference result with a high-confidence token is obtained in low intermediate layers, inference can be terminated at that point and the computational complexity is reduced compared to not terminating. In order to terminate inference early, it is necessary to improve the accuracy of the low intermediate layers. However, early exit often focuses on improving the accuracy of high intermediate layers because they determine the upper accuracy limit of ASR. As a result, it is difficult to terminate inference in low intermediate layers because the confidence of those layers is low. To solve this problem, we propose block refinement learning (BRL), which is a re-training method of existing early-exiting models. BRL trains the low intermediate layers while maintaining the overall accuracy of the model by considering gradients of not only low intermediate layers but also high intermediate layers. In this way, low intermediate layers can be trained while maintaining the accuracy of the high intermediate layers. We demonstrated the effectiveness of BRL on Japanese discourse ASR tasks.

I. INTRODUCTION

Reducing inference time is important for automatic speech recognition (ASR). Previous studies have reduced the inference time by reducing the model parameters and computational complexity through the use of knowledge distillation [1]–[3], quantization [4]–[6], and pruning [7]–[9].

On the other hand, early exit [10], [11], which aims to speed up the inference process of ASR, has attracted attention in recent years due to its high affinity with ASR models that have multi-layer decoder. Early exit trains a model that has output layer associated with each of its intermediate decoder layers as well as the final decoder layer. In this paper, this model is referred to as the “early-exiting model”. Also, the pair of a decoder layer and an output layer attached to it is referred to as the “intermediate layer”. The early-exiting model can reduce the inference time by dynamically selecting inference paths for each input on the basis of the outputs of the intermediate layers. For instance, by terminating inference in low intermediate layers when inference result in low intermediate layers is high confidence and continuing inference until high intermediate layers when inference result in low intermediate layers is low confidence.

In training the early-exiting model, it often focuses on the accuracy of the high intermediate layers because it defines

the upper limit of the ASR accuracy. However, the training policy that focuses on the accuracy of the high intermediate layers neglects the accuracy of low intermediate layers. As a result, the effect of reducing the inference time is smaller because it becomes hard for early exit to terminate inference in low intermediate layers. Given this background, a method is needed to improve the accuracy of low intermediate layers while maintaining the accuracy of high intermediate layers.

In this paper, we propose block refinement learning (BRL), which is a training method of the early-exiting model, to improve the ASR accuracy of low intermediate layers. An overview of BRL is shown in Figure 1. The key advantage of BRL is that it improves the ASR accuracy of low intermediate layers while maintaining the ASR accuracy of high intermediate layers in the early-exiting model. To this end, it trains the parameters of low intermediate layers, while keeping the parameters of all other layers fixed. At the same time, BRL trains low intermediate layers by exploiting not only the gradients for the corresponding attached output layers, but also the gradients for all output layers including the high intermediate ones. We demonstrated the effectiveness of BRL on Japanese discourse ASR tasks.

Our contributions are summarized as follows:

- We established a method to improve the ASR accuracy of low intermediate layers while maintaining the ASR accuracy of high intermediate layers in the early-exiting model.
- By our method, we experimentally demonstrated faster inference while maintaining ASR performance.

II. RELATED WORK

This section details previous studies on early exit. The previous studies used two different training methods to train the early-exiting models which has multiple output layers: two stage training [12]–[14], in which only the added output layers are trained, and joint training [15]–[17], in which the entire early-exiting model is trained in a single training. With two-stage training, only the added output layers are trained, while the parameters of the other component are fixed. Therefore, it is possible to maintain the performance of the final output layer. On the other hand, with joint training, all parameters included in the early-exiting model are trained in a single training. It has been pointed out that the accuracy of low intermediate layers is poor in two-stage training because it

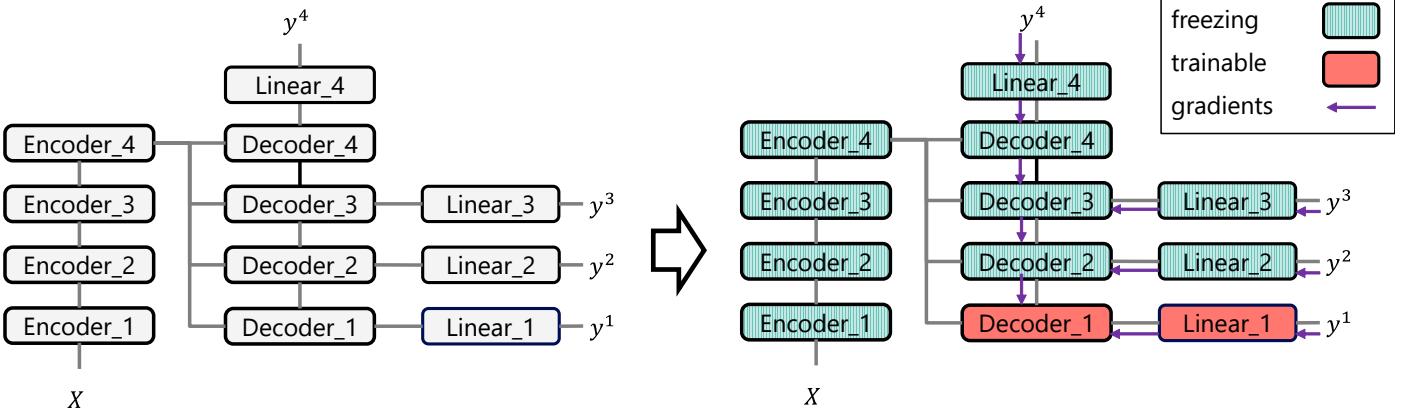


Fig. 1: Overview of the proposed method. The left side shows the early-exiting model, and the right side shows the training of BRL for the early-exiting model. As an example, this figure shows the early-exiting model with four encoder layers and four decoder layers, where the first layer is an intermediate layer trained by BRL.

trains only the added output layers [18]. Also, the accuracy of low intermediate layers is high in joint training because joint training trains all layers, not just the added ones. In this paper, we focus on joint training for its advantage of reducing the inference time by ensuring the high accuracy of the low intermediate layers.

In inference, the early-exiting model achieves a speedup by terminating the inference early based on exit criteria. The possible exit criteria include entropy-based [19], probability-based [16], [20], and patience-based methods [21]. We investigated the effectiveness of BRL by using probability-based and patience-based exit criteria as representative examples.

III. THE TRANSFORMER-BASED EARLY-EXITING MODEL

A. Transformer-based ASR model

This section details the ASR model based on transformer [22]. This model estimates a text $\mathbf{Y} = (y_1, \dots, y_T)$, that is the inference result from acoustic features $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$. y_t is the t -th token of the text, and T is the length of output sequence. \mathbf{x}_n is the n -th frame of the features, and N is the frame length. The ASR model estimates the following probability P of generating \mathbf{Y} for \mathbf{X} :

$$P(\mathbf{Y}|\mathbf{X}; \Theta) = \prod_{t=1}^T P(y_t|y_{1:t-1}, \mathbf{X}; \Theta), \quad (1)$$

where $y_{1:t-1} = (y_1, \dots, y_{t-1})$ and Θ denotes the ASR model parameter. The probability P obtained by the transformer-based ASR model is as follows:

$$\mathbf{h} = \text{TransformerEnc}(\mathbf{X}; \Theta_{\text{enc}}), \quad (2)$$

$$\mathbf{e}_{t-1} = \text{Embedding}(y_{t-1}; \Theta_{\text{emb}}), \quad (3)$$

$$\mathbf{c}_{t-1}^l = \begin{cases} \text{TransformerDec}(\mathbf{h}, \mathbf{e}_{1:t-1}; \Theta_{\text{dec}}^l) & \text{if } l = 1, \\ \text{TransformerDec}(\mathbf{h}, \mathbf{c}_{1:t-1}^{l-1}; \Theta_{\text{dec}}^l) & \text{otherwise,} \end{cases} \quad (4)$$

$$P(y_t|y_{1:t-1}, \mathbf{X}; \Theta) = \text{Softmax}(\mathbf{c}_{t-1}^L; \Theta_{\text{linear}}^L), \quad (5)$$

where $\text{TransformerEnc}(\cdot)$ is a transformer encoder that consists of an embedding layer, positional encoding layer, and multiple transformer encoder layers. The transformer encoder layer consists of a multi-head self-attention layer and a feed-forward layer. Θ_{enc} denotes its parameters. $\text{Embedding}(\cdot)$ is an embedding layer that consists of an embedding layer and a positional encoding layer. Θ_{emb} denotes its parameters. $\text{TransformerDec}(\cdot)$ is a transformer decoder that consists of a masked-multi-head attention layer, encoder-decoder multi-head-attention layer, and feed-forward layer. Θ_{dec}^l denotes the parameter of the l -th transformer decoder, where $l \in [1, L]$. L is the number of decoder layers. $\text{Softmax}(\cdot)$ is a softmax layer with a linear transformation. Θ_{linear}^l denotes the parameter of the l -th softmax layer connected to the l -th transformer decoder. Note that the softmax layer is only connected to the last transformer decoder in the basic ASR model. The parameter $\Theta = \{\Theta_{\text{enc}}, \Theta_{\text{emb}}, \Theta_{\text{dec}}^1, \dots, \Theta_{\text{dec}}^L, \Theta_{\text{linear}}^L\}$ is optimized with a cross-entropy loss function \mathcal{L} , defined as

$$\mathcal{L}(\Theta) = - \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{D}} \sum_{t=1}^T \log P(y_t|y_{1:t-1}, \mathbf{X}; \Theta). \quad (6)$$

\mathcal{D} denotes training data containing multiple \mathbf{X} and \mathbf{Y} pair data.

B. The early-exiting model

This section details the transformer-based early-exiting ASR model. The model can obtain probabilities P^l from l -th intermediate decoder layer. The probabilities P^l are obtained as follows:

$$P^l(y_t|y_{1:t-1}, \mathbf{X}; \Theta_E) = \text{Softmax}(\mathbf{c}_{t-1}^l; \Theta_{\text{linear}}^l). \quad (7)$$

In the early-exiting model, one $\text{Softmax}(\cdot)$ is connected to each $\text{TransformerDec}(\cdot)$. The parameter $\Theta_E = \{\Theta_{\text{enc}}, \Theta_{\text{emb}}, \Theta_{\text{dec}}^1, \dots, \Theta_{\text{dec}}^L, \Theta_{\text{linear}}^1, \dots, \Theta_{\text{linear}}^L\}$ is optimized

Algorithm 1 Inference with probability-based exit criteria
(Input: \mathbf{X} , Output: \mathbf{Y})

```

1:  $\mathbf{h} \leftarrow \text{TransformerEnc}(\mathbf{X}; \boldsymbol{\theta}_{\text{enc}})$ 
2: for  $t = 1$  to  $T$  do
3:    $\mathbf{e}_{t-1} = \text{Embedding}(y_{1:t-1}; \boldsymbol{\theta}_{\text{emb}})$ 
4:   for  $l = 1$  to  $L$  do
5:     if  $l = 1$  then
6:        $\mathbf{c}_{t-1}^l \leftarrow \text{TransformerDec}(\mathbf{h}, \mathbf{e}_{1:t-1}; \boldsymbol{\theta}_{\text{dec}}^l)$ 
7:     else
8:        $\mathbf{c}_{t-1}^l \leftarrow \text{TransformerDec}(\mathbf{h}, \mathbf{c}_{1:t-1}^{l-1}; \boldsymbol{\theta}_{\text{dec}}^l)$ 
9:     end if
10:     $P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E) \leftarrow \text{Softmax}(\mathbf{c}_{t-1}^l; \boldsymbol{\theta}_{\text{linear}}^l)$ 
11:     $s_t^l \leftarrow \text{Score}(P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E))$ 
12:    if  $s_t^l \geq \lambda^l$  then
13:       $y_t \leftarrow \arg \max_{y_t} P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E)$ 
14:      break
15:    end if
16:   end for
17: end for

```

with a cross-entropy loss function \mathcal{L}_E , defined as

$$\mathcal{L}_E(\boldsymbol{\Theta}_E) = - \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{D}} \sum_{t=1}^T \sum_{l=1}^L w^l \log P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E). \quad (8)$$

The weights w^l set in Schuster et al. [23] favor higher layers, as follows:

$$w^l = \frac{l}{\sum_{k=1}^L k}. \quad (9)$$

C. Inference of the early-exiting model

This section details two exit criteria of early exit, probability-based [20] and patience-based (PABEE) [21]. The probability-based exit criteria terminates inference when the early-exiting model outputs high-confidence token from its intermediate layers. Specifically, terminating inference occurs when the confidence score $s_t^l \in [0, 1]$ exceeds the thresholds $\lambda^l \in [0, 1]$. The confidence score $s_t^l = \text{Score}(\cdot)$ is defined by the difference between the top two values of $\text{Softmax}(\cdot)$. The thresholds λ^l are determined by setting arbitrary objective functions and solving an optimization problem. The flow of inference with the probability-based exit criteria is shown in Algorithm 1.

PABEE terminates inference when the early-exiting model outputs same token continuously from its intermediate layers. It uses a patience counter p to store the number of times that outputs do not change and the threshold $\hat{p} \in [1, L-1]$. Terminating inference occurs when p is equal to \hat{p} . The flow of inference with the patience-based exit criteria is shown in Algorithm 2.

Note that due to the structure of the transformer self-attention mechanism, all previous hidden states $\mathbf{c}_{1:t-1}^{l-1}$ are es-

Algorithm 2 Inference with patience-based exit criteria (Input: \mathbf{X} , Output: \mathbf{Y})

```

1:  $\mathbf{h} \leftarrow \text{TransformerEnc}(\mathbf{X}; \boldsymbol{\theta}_{\text{enc}})$ 
2: for  $t = 1$  to  $T$  do
3:    $p \leftarrow 0$ 
4:    $\mathbf{e}_{t-1} = \text{Embedding}(y_{1:t-1}; \boldsymbol{\theta}_{\text{emb}})$ 
5:   for  $l = 1$  to  $L$  do
6:     if  $l = 1$  then
7:        $\mathbf{c}_{t-1}^l \leftarrow \text{TransformerDec}(\mathbf{h}, \mathbf{e}_{1:t-1}; \boldsymbol{\theta}_{\text{dec}}^l)$ 
8:     else
9:        $\mathbf{c}_{t-1}^l \leftarrow \text{TransformerDec}(\mathbf{h}, \mathbf{c}_{1:t-1}^{l-1}; \boldsymbol{\theta}_{\text{dec}}^l)$ 
10:     end if
11:      $P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E) \leftarrow \text{Softmax}(\mathbf{c}_{t-1}^l; \boldsymbol{\theta}_{\text{linear}}^l)$ 
12:     if  $\arg \max_{y_t} P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E)$ 
13:        $= \arg \max_{y_t} P^{l-1}(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E)$  then
14:          $p \leftarrow p + 1$ 
15:       else
16:          $p \leftarrow 0$ 
17:       end if
18:       if  $p = \hat{p}$  then
19:          $y_t \leftarrow \arg \max_{y_t} P^l(y_t|y_{1:t-1}, \mathbf{X}; \boldsymbol{\Theta}_E)$ 
20:         break
21:       end if
22:     end for
22:   end for

```

sential to the calculation of \mathbf{c}_{t-1}^l . Therefore, when terminating inference occurred in l -th layer in the estimation of t -th token, hidden states $\{\mathbf{c}_{t-1}^{l+1}, \dots, \mathbf{c}_{t-1}^L\}$ are complemented by copying \mathbf{c}_{t-1}^l .

IV. PROPOSED METHOD

This section details BRL, an effective method of training the early-exiting model. The key advantage of BRL is that it improves the ASR accuracy of specified low intermediate layers, while maintaining the ASR accuracy of high intermediate layers in the early-exiting model. BRL has two features:

- Fixing the parameters of high intermediate layers and training the parameters of low intermediate layers.
- Calculating gradients from not only low intermediate layers, but also high intermediate layers.

As shown in Figure 1, the initial parameter of BRL is the existing early-exiting model. In BRL, only the weights $\{\boldsymbol{\theta}_{\text{dec}}^1, \dots, \boldsymbol{\theta}_{\text{dec}}^j, \boldsymbol{\theta}_{\text{linear}}^1, \dots, \boldsymbol{\theta}_{\text{linear}}^j\}$ in $\boldsymbol{\Theta}_E$ are updated using Equation (8). $j \in [1, L]$ denotes the number of specified low intermediate layers to be trained.

BRL sets uniform weights for w^l as follows:

$$w^l = \frac{1}{L}. \quad (10)$$

TABLE I: Dataset details.

	Data size (hours)	Number of utterances	Number of characters
Train	512.6	413,240	13,349,780
Valid	1.9	1,292	47,970
Test	1.3	1,385	32,089

BRL sets uniform weights w^l because BRL focuses on improving the ASR accuracy of low intermediate layers, while maintaining the ASR accuracy of high intermediate layers.

V. EXPERIMENT

A. Dataset

We evaluated the effectiveness of BRL on Japanese discourse ASR tasks by using the corpus of spontaneous Japanese (CSJ) [24]. We used CSJ as the training set (Train), validation set (Valid), and test set (Test). Valid and Test are Eval2 and Eval3 provided by CSJ, respectively. Details of these datasets are given in Table I. We used 80 log mel-scale filterbank coefficients as acoustic features and text segmented into character units. To exclude speech data that was too long, we removed from the dataset any data that contained text with more than 300 tokens or more than 1,500 audio frames.

B. Settings

The early-exiting model is a transformer-based encoder-decoder model with intermediate layers. We used two early-exiting models with different numbers of decoder layers: EE-base and EE-deep. EE-base has 8 encoder layers and 6 decoder layers. Every intermediate decoder layer in EE-base has a corresponding output layer. EE-deep has 8 encoder layers and 12 decoder layers, i.e., more decoder layers than in EE-base. EE-deep has an output layer for each even-numbered intermediate decoder layer. EE-base and EE-deep were created under the following conditions: the dimensions of the input and output continuous representations were set to 512, the dimensions of the inner outputs in the position-wise feed-forward networks were set to 2,048, and the number of heads in the multi-head attention was set to 4. To stabilize training, the dropout rate in the transformer encoder and decoder layers was set to 0.1. For the activation function, we used swish activation layer [25].

BRL was applied to EE-base and EE-deep to yield BRL-base and BRL-deep. BRL-base was trained from the first through the third of its intermediate layers ($j = 3$) by utilizing BRL. Also, BRL-deep was trained from the first through the sixth of its intermediate layers ($j = 6$). The number of intermediate layers to be trained, j , was chosen to maximize the sum of the ASR accuracies obtained from all the intermediate layers.

In the training, our models were optimized by using RAdam [26]. We set a mini-batch size of 64 and the learning rate of 0.0001. We used label smoothing [27] for regularization. For data augmentation, we applied SpecAugment with frequency masking and time masking [28], where the number

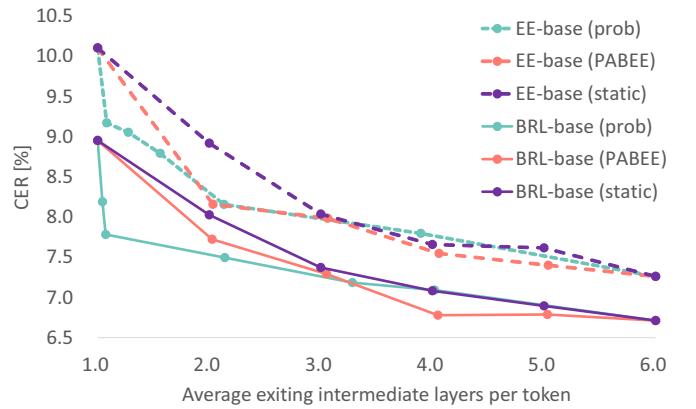


Fig. 2: CER and average exiting layer of models trained by BRL in EE-base.

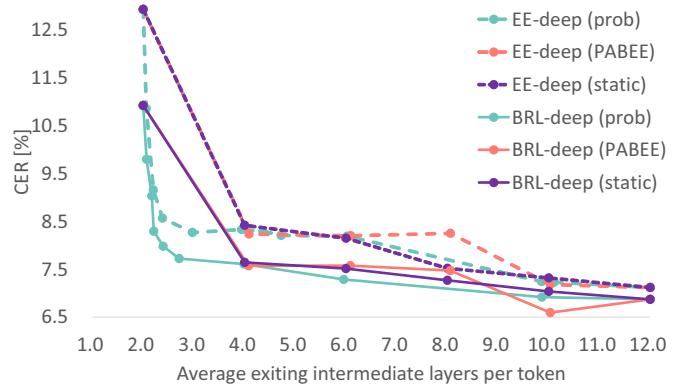


Fig. 3: CER and average exiting layer of models trained by BRL in EE-deep.

of frequency masks and time-step masks were each set to 2, the frequency-masking width was randomly chosen from 0 to 27 frequency bins, and the time-masking width was randomly chosen from 0 to 100 frames. Early stopping was applied if no best model was found in the validation set after 10 epochs.

In the decoding, we used a beam search with a beam size of 4. In this experiment, we tested three exit criteria: probability-based (prob), patience-based (PABEE), and static (static). In prob, we set an optimal threshold λ^l that simultaneously minimized the character error rate (CER) and the average exiting intermediate layer per token in the validation set. To find the optimal threshold λ^l , we used NSGA-II [29], which is a multi-objective genetic algorithm for solving multi-objective optimization problems. We searched for the optimal λ^l in 300 epochs. In PABEE, we set a patience counter \hat{p} from 1 to $L-1$. Static is the results of forced termination of the inference in each intermediate layer, which are presented to validate the training effect of BRL. We used the CPU, which was INTEL XEON Silver 4310, for the decoding and batch size was 1. The Intel oneAPI Math Kernel Library [30] was enabled for the computations, and 24 threads were used to perform parallel computations.

TABLE II: The impact of early exit on actual inference speed.

	\downarrow CER[%]	\uparrow Speedup
EE-base (static by final output)	7.26	$\times 1.00$
EE-base (prob)	7.80	$\times 1.38$
	8.16	$\times 2.27$
	9.17	$\times 3.30$
BRL-base (prob)	7.09	$\times 1.39$
	7.49	$\times 2.56$
	7.78	$\times 3.43$

The accuracy of the early-exiting model was measured by CER. Reduction in inference time was measured by the average exiting intermediate layers per token.

C. Results

Figure 2 shows the trade-off between inference speed and CER in EE-base. From a comparison of EE-base (static) and BRL-base (static), we can find that BRL improved the ASR accuracy of the individual low intermediate layers. From a comparison of EE-base (prob) and BRL-base (prob), we can also find that BRL-base (prob) achieved a comparable CER to that of EE-base (prob) with fewer average exiting intermediate layers. Furthermore, we can find the same trend even when the exit criteria were changed to PABEE. Figure 3 shows the trade-off between inference speed and CER in EE-deep. It confirms that BRL has a similar effect in EE-deep as in EE-base. These results indicate that BRL can improve the ASR accuracy to train a low intermediate layer while maintaining the ASR accuracy of high intermediate layers.

To show the actual effect of BRL, we evaluated it in terms of computation speed and CER. Table II shows the impact of early exit on the actual inference speed. EE-base (static by final output) shows the inference speed and CER when always exiting at the final decoder layer. Speedup is a ratio based on EE-base (static by final output). The effect of speedup can be compared by Speedup under close CER. By comparing EE-base (prob) and BRL-base (prob) in Table II, we can find that close CER to EE-base (static by final output) are achieved by Speedup of $\times 3.43$ in BRL-base (prob), while $\times 1.38$ in EE-base (prob). In other words, BRL-base (prob) achieves 2.49 times faster compared to EE-base (prob) at close CER.

Further, we investigated the effect of changing the number of trained low intermediate layers in BRL. We adjusted j from 1 to 6 in BRL. Figure 4 shows the result. Even when the number of intermediate layers to be trained changed, BRL remained effective under all experimental conditions. $j = 6$ is the result of re-training all decoder layers including the final one. $j = 6$ and EE-base (prob) have different settings for loss weights. The loss weights of $j = 6$ obey Equation (9) and the loss weights of EE-base (prob) obey Equation (10). The result that $j = 6$ exceeds the ASR performance of EE-base (prob) indicate that BRL can also be used for the early-exiting model with insufficient training in the intermediate layers.

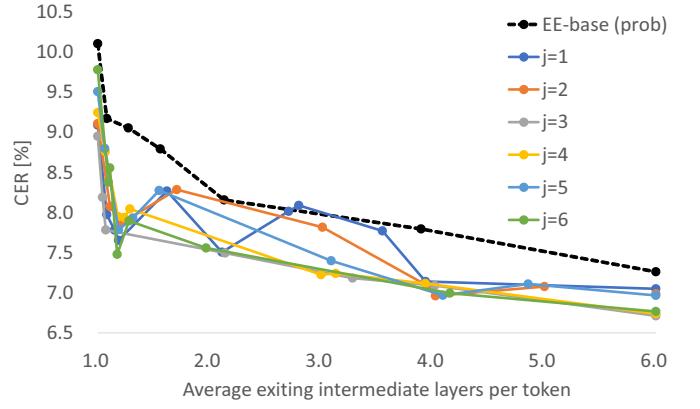


Fig. 4: All model used probability-based exit criteria. Result of varying the number of low intermediate layers. $j = 3$ indicates BRL-base (prob).

VI. CONCLUSION

We presented a block refinement learning (BRL) that achieves faster inference for the early-exiting model in automatic speech recognition (ASR). The key advantage of BRL is that it improves the ASR accuracy of low intermediate layers while maintaining the ASR accuracy of high intermediate layers in the existing early-exiting model. To this end, BRL trains the parameters of the low intermediate layers, while keeping the parameters of all other layers fixed. In addition, it exploits the gradients from not only the corresponding attached output layers, but also from all output layers including those associated with the high intermediate layers. We experimentally showed that BRL can improve the ASR accuracy of low intermediate layers while maintaining the ASR accuracy of high intermediate layers in the early-exiting model.

REFERENCES

- [1] M. A. Haidar, C. Xing, and M. Rezagholizadeh, “Transformer-based asr incorporating time-reduction layer and fine-tuning with self-knowledge distillation,” in *Proc. Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2021, pp. 2102–2106.
- [2] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS) Workshop on Deep Learning*, 2015.
- [3] L. Li, Y. Lin, D. Chen, S. Ren, P. Li, J. Zhou, and X. Sun, “Cascadebert: Accelerating inference of pre-trained language models via calibrated complete models cascade,” in *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2021, pp. 475–486.
- [4] S. Kim, A. Gholami, Z. Yao, N. Lee, P. Wang, A. Nrusimha, B. Zhai, T. Gao, M. W. Mahoney, and K. Keutzer, “Integer-only zero-shot quantization for efficient speech recognition,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 4288–4292.
- [5] S. Shen, Z. Dong, J. Ye, L. Ma, Z. Yao, A. Gholami, M. W. Mahoney, and K. Keutzer, “Q-bert: Hessian based ultra low precision quantization of bert,” in *Proc. AAAI Conference on Artificial Intelligence (AAAI)*, 2020, pp. 8815–8821.
- [6] W. Zhang, L. Hou, Y. Yin, L. Shang, X. Chen, X. Jiang, and Q. Liu, “TernaryBERT: Distillation-aware ultra-low bit BERT,” in *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 509–521.

- [7] C.-I. J. Lai, Y. Zhang, A. H. Liu, S. Chang, Y.-L. Liao, Y.-S. Chuang, K. Qian, S. Khurana, D. Cox, and J. Glass, “Parp: Prune, adjust and re-prune for self-supervised speech recognition,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [8] J. Frankle and M. Carbin, “The lottery ticket hypothesis: Finding sparse, trainable neural networks,” in *Proc. International Conference on Learning Representations (ICLR)*, 2019.
- [9] S. Han, J. Pool, J. Tran, and W. J. Dally, “Learning both weights and connections for efficient neural networks,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2015, pp. 1135–1143.
- [10] R. Tang, K. Kumar, J. Xin, P. Vyas, W. Li, G. Yang, Y. Mao, C. Murray, and J. Lin, “Temporal early exiting for streaming speech commands recognition,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 7567–7571.
- [11] D. Berrebbi, B. Yan, and S. Watanabe, “Avoid overthinking in self-supervised models for speech recognition,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023, pp. 1–5.
- [12] S. Scardapane, M. Scarpiniti, E. Baccarelli, and A. Uncini, “Why should we add early exits to neural networks?” *Cognitive Computation*, vol. 12, no. 5, pp. 954–966, Sep 2020.
- [13] J. Xin, R. Tang, Y. Yu, and J. Lin, “BERxit: Early exiting for BERT with better fine-tuning and extension to regression,” in *Proc. European Chapter of the Association for Computational Linguistics (EACL)*, 2021, pp. 91–104.
- [14] G. A. Wright, U. Cappellazzo, S. Zaiem, D. Raj, L. O. Yang, D. Falavigna, M. N. Ali, and A. Brutti, “Training early-exit architectures for automatic speech recognition: Fine-tuning pre-trained models or training from scratch,” in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 7567–7571.
- [15] S. Bae, J. Ko, H. Song, and S.-Y. Yun, “Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding,” in *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023, pp. 5910–5924.
- [16] S. Teerapittayanon, B. McDanel, and H. Kung, “Branchynet: Fast inference via early exiting from deep neural networks,” in *Proc. International Conference on Pattern Recognition (ICPR)*, 2016, pp. 2464–2469.
- [17] Y. Ji, J. Wang, J. Li, Q. Chen, W. Chen, and M. Zhang, “Early exit with disentangled representation and equiangular tight frame,” in *Proc. Association for Computational Linguistics (ACL)*, 2023.
- [18] J. Xin, R. Tang, J. Lee, Y. Yu, and J. Lin, “DeeBERT: Dynamic early exiting for accelerating BERT inference,” in *Proc. Association for Computational Linguistics (ACL)*, 2020, pp. 2246–2251.
- [19] S. Geng, P. Gao, Z. Fu, and Y. Zhang, “Romebert: Robust training of multi-exit bert,” *arXiv preprint arXiv:2101.09755*, 2021.
- [20] T. Schuster, A. Fisch, J. Gupta, M. Dehghani, D. Bahri, V. Q. Tran, Y. Tay, and D. Metzler, “Confident adaptive language modeling,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2022, pp. 17456–17472.
- [21] W. Zhou, C. Xu, T. Ge, J. McAuley, K. Xu, and F. Wei, “Bert loses patience: Fast and robust inference with early exit,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [22] 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 *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [23] Z. Liu, Z. Xu, H.-J. Wang, T. Darrell, and E. Shelhamer, “Anytime dense prediction with confidence adaptivity,” in *Proc. International Conference on Learning Representations (ICLR)*, 2022.
- [24] K. Maekawa, H. Koiso, S. Furui, and H. Isahara, “Spontaneous speech corpus of Japanese,” in *Proc. International Conference on Language Resources and Evaluation (LREC)*, 2000.
- [25] P. Ramachandran, B. Zoph, and Q. V. Le, “Searching for activation functions,” *arXiv preprint arXiv:1710.05941*, 2017.
- [26] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, “On the variance of the adaptive learning rate and beyond,” in *Proc. International Conference on Learning Representations (ICLR)*, 2020.
- [27] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2818–2826.
- [28] 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. Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2019, pp. 2613–2617.
- [29] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, “A fast and elitist multiobjective genetic algorithm: NSGA-II,” *IEEE Transactions on Evolutionary Computation*, pp. 182–197, 2002.
- [30] M. Krainiuk, M. Goli, and V. R. Pasquazi, “oneapi open-source math library interface,” in *Proc. International Workshop on Performance, Portability and Productivity in HPC (P3HPC)*, 2021, pp. 22–32.