

# Residual Language Model for End-to-end Speech Recognition

Emiru Tsunoo<sup>1</sup> Yosuke Kashiwagi<sup>1</sup> Chaitanya Narisetty<sup>2</sup> Shinji Watanabe<sup>2</sup>

<sup>1</sup> Sony Group Corporation, Japan  
<sup>2</sup> Carnegie Mellon University, USA

emiru.tsunoo@sony.com

## Abstract

End-to-end automatic speech recognition suffers from adaptation to unknown target domain speech despite being trained with a large amount of paired audio–text data. Recent studies estimate a linguistic bias of the model as the internal language model (LM). To effectively adapt to the target domain, the internal LM is subtracted from the posterior during inference and fused with an external target-domain LM. However, this fusion complicates the inference and the estimation of the internal LM may not always be accurate. In this paper, we propose a simple external LM fusion method for domain adaptation, which considers the internal LM estimation in its training. We directly model the residual factor of the external and internal LMs, namely the residual LM. To stably train the residual LM, we propose smoothing the estimated internal LM and optimizing it with a combination of cross-entropy and mean-squared-error losses, which consider the statistical behaviors of the internal LM in the target domain data. We experimentally confirmed that the proposed residual LM performs better than the internal LM estimation in most of the cross-domain and intra-domain scenarios.

**Index Terms:** speech recognition, language model, attention-based encoder–decoder, internal language model estimation

## 1. Introduction

End-to-end (E2E) automatic speech recognition (ASR) has attracted interest as a method of directly integrating acoustic models (AMs) and language models (LMs) because of its simple training and efficient decoding procedures. In recent years, various approaches have been studied, including connectionist temporal classification (CTC) [1, 2, 3], attention-based encoder–decoder models [4, 5], hybrid models [6, 7], and transducers [8, 9, 10].

E2E ASR requires pairs of audio and text data for training. Even with a large amount of paired data, Del Rio et al. demonstrated that training with 960 h of Librispeech read speech does not result in sufficient performance in the mismatched domain of earnings calls [11]. If the target domain has paired data, adaptation techniques can be adopted [12, 13, 14, 15]. However, in most scenarios, orders-of-magnitude more text-only data of the target domain are available, and it is more efficient to shift the linguistic bias of the E2E ASR model towards the domain of interest using such data.

Many researchers have studied fusion methods using an external LM trained with text-only data. Shallow fusion, which linearly interpolates the E2E ASR model with an external LM, is the most popular approach [4, 16, 17]. More structural integration can be observed in deep fusion [18], cold fusion [19], and component fusion [20], which require additional training. Fundamentally, the probability estimation of the LMs relies on softmax computation. Therefore, although there are efficient

log-sum-exp calculation tricks, it incurs a higher computational cost as the vocabulary size increases. The density ratio approach [21] focuses more on domain adaptation by assuming that the source and target domains are acoustically consistent, and it adapts the E2E ASR model with LMs trained in each domain by following Bayes’ rule. Recently, the estimation of an internal LM, a linguistic bias of E2E ASR, has been investigated, and by subtracting from the ASR posterior, it improves performance in both cross-domain and intra-domain scenarios [22, 23, 24, 25]. However, both the density ratio and internal LM estimation complicate the inference computation. In addition, due to the domain mismatch, the estimation of the internal LM may not always be accurate.

In this paper, we propose a simple external LM fusion for domain adaptation, which considers the internal LM. Instead of subtracting the estimated internal LM and fusing with an external target-domain LM, we directly model the residual factor of them. The difference of the probability distributions, namely the residual LM, is trained with a target-domain text-only dataset, considering the estimated internal LM in the specific domain. Thus, the residual LM not only conveys the linguistic characteristic of the dataset, but also aggregates the estimation results of the internal LM in the target-domain corpus into the model, thereby alleviating the domain mismatch problem. In addition, because the distribution is no longer a probability, the residual LM can omit costly softmax computations in the output layer. We propose a training approach that applies smoothing to the internal LM probability, and combines cross-entropy with mean squared errors (MSEs) for the loss function. The trained residual LM can be simply fused in the same manner as shallow fusion. We performed experiments to determine the effectiveness of the proposed residual LM in cross-domain and intra-domain scenarios using various corpora. The results show that the proposed residual LM improves performance in cross-domain scenarios by 4.0% relative word error rate (WER) in the Librispeech–TEDLIUM3 adaptation, with faster inference. Additionally, the residual LM fusion method performs robustly in intra-domain scenarios.

## 2. Formulation and Related Studies

ASR is a problem in determining the most probable token sequence  $\mathbf{y}$  given an input audio  $\mathbf{x}$ . In a scenario in which the training and target domains differ, and the text corpus in the target domain is easily accessible, it is useful to combine the ASR model with an external LM trained using the target corpus. With a Bayesian interpretation, classical hybrid ASR systems determine the highest probability by combining with the external LM, as follows:

$$p(\mathbf{y}|\mathbf{x}) = p(\mathbf{x}|\mathbf{y}; \theta_{AM}) \cdot p(\mathbf{y}; \theta_{LM}) \cdot \frac{1}{p(\mathbf{x})}, \quad (1)$$

where  $\theta_{AM}$  and  $\theta_{LM}$  denote parameters of the AM and LM, respectively. In E2E systems,  $p(\mathbf{x}|\mathbf{y};\theta_{AM})$  is replaced by E2E neural networks, which can be further decomposed into the following terms using the Bayesian theorem:

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S)}{p(\mathbf{y};\theta_{E2E}^S)} \cdot p(\mathbf{x};\theta_{E2E}^S) \cdot p(\mathbf{y};\theta_{LM}^T) \cdot \frac{1}{p(\mathbf{x})}, \quad (2)$$

where  $\theta_{E2E}$  is a parameter set of an E2E ASR model. For domain adaptation, the E2E models are trained in a source domain, and the external LMs are trained in a target domain; thus,  $*^S$  and  $*^T$  represent source and target domains, respectively. By omitting  $p(\mathbf{x})$  and  $p(\mathbf{x};\theta_{E2E}^S)$ , which are not required to search for the highest probability of  $\mathbf{y}$ , the score function for recognition can be expressed in a logarithmic scale as

$$\text{Score}(\mathbf{y}|\mathbf{x}) = \log p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S) - \log p(\mathbf{y};\theta_{E2E}^S) + \log p(\mathbf{y};\theta_{LM}^T). \quad (3)$$

The first term is the output posterior of the E2E ASR neural network, and the second term is the implicit linguistic bias (prior) of the trained E2E model. The last term is an external LM trained in the target-domain text corpus.

### 2.1. Shallow fusion

E2E ASR models are often used with an external LM trained with a text-only corpus of the target domain. This is reasonable because, unlike the E2E ASR requiring a large corpus of paired audio–text data for training, external LM training requires only text data, which can be easily obtained with a large scale. Shallow fusion is a common method for integrating an E2E ASR model with an LM [4, 16, 17]. In practice, the second term in Eq. (3) is omitted because it is intractable, and in shallow fusion, only the external LM term is combined by introducing an LM weight,  $\lambda_{LM}$  as:

$$\text{Score}(\mathbf{y}|\mathbf{x}) = \log p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S) + \lambda_{LM} \log p(\mathbf{y};\theta_{LM}^T). \quad (4)$$

### 2.2. Density ratio approach

The density ratio [21] assumes that the source and target domains are acoustically consistent. Assuming that an LM trained with text-only data in the source domain can represent the linguistic prior of the E2E model, i.e., the second term of Eq. (3), the density ratio uses the ratio of LMs trained in both domains for adaptation. Thus, the score in the density ratio approach can be expressed using the source-domain and target-domain LM weights,  $\lambda_{DR}$  and  $\lambda_{LM}$ , as follows:

$$\text{Score}(\mathbf{y}|\mathbf{x}) = \log p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S) - \lambda_{DR} \log p(\mathbf{y};\theta_{LM}^S) + \lambda_{LM} \log p(\mathbf{y};\theta_{LM}^T). \quad (5)$$

### 2.3. Internal language model estimation

Internal LM estimation (ILME) [22, 23, 24, 25] attempts to estimate the second term in Eq. (3), i.e.,  $\log p(\mathbf{y};\theta_{E2E}^S)$ . As a result, the score for decoding using the internal LM is expressed as follows.

$$\text{Score}(\mathbf{y}|\mathbf{x}) = \log p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S) - \lambda_{ILM} \log p(\mathbf{y};\theta_{E2E}^S) + \lambda_{LM} \log p(\mathbf{y};\theta_{LM}^T), \quad (6)$$

where  $\lambda_{ILM}$  is a weight parameter for the estimated internal LM.

A common method of estimating the internal LM,  $p(\mathbf{y};\theta_{E2E}^S)$ , is to replace the encoder output with zero-filled vectors and infer only with the decoder. This is because, particularly in attention-based encoder–decoder and transducer architectures, decoders function similarly to LMs as they estimate next  $i$ -th token  $y_i$  with a given previous output  $\mathbf{y}_{1:i-1}$ . This estimation approach is statistically reasonable in the source domain because the encoder output is most likely normalized to zero-mean vectors with a normalization layer. However, this estimation is performed in every inference, which complicates the computation. In addition, when it is estimated with the mismatched domain speech, the behavior of the internal LM becomes unpredictable and the estimation may not always be accurate.

## 3. Residual Language Model

### 3.1. Definition of the residual language model

Instead of estimating the internal LM in every inference, we propose to directly model the residual factor of the target-domain external LM and the estimated internal LM. The model predicts the difference between the second and third terms (internal and external LMs, respectively) in Eq. (3), which we defined as the residual LM. By directly modeling the residual term, we can simplify the inference computation to shallow fusion, as follows:

$$\text{Score}(\mathbf{y}|\mathbf{x}) = \log p(\mathbf{y}|\mathbf{x};\theta_{E2E}^S) + \lambda_{LM} f(\mathbf{y};\theta_{res}^T), \quad (7)$$

where  $f(\mathbf{y})$  is the proposed residual LM. The residual LM conveys both the second and third terms in Eq. (3); thus Eq. (7) strictly follows the score calculation (3) derived by Bayes’ interpretation. The main differences from ILME are as follows.

- The residual LM models the residual factor of an external LM and the internal LM, which simplify the inference procedure and reduce computational cost.
- The difference of two LMs is no longer a probability distribution; thus it can further omit the log-softmax operation, which requires costly log-sum-exp calculation.
- The residual LM conveys statistical behavior of the estimated internal LM in the target-domain text corpus.

The residual LM is trained using the target-domain text-only data. To model the residual terms of the external and internal LMs, we can define the training target of the model output,  $r(\mathbf{y})$  as follows:

$$r(\mathbf{y}) = \log q_{\mathbf{y}}^* - \gamma \log p(\mathbf{y};\theta_{E2E}^S), \quad (8)$$

where  $q_{\mathbf{y}}^*$  is a reference label, and  $\gamma$  is a tunable parameter. To avoid log-zero computation, a smoothed label as in [26] is adopted for the reference  $q_{\mathbf{y}}^*$ , as

$$q_{y,k}^* = (1 - \omega)\delta_{y,k} + \frac{\omega}{K}, \quad (9)$$

where  $k$  is the vocabulary index,  $\delta_{y,k}$  is Dirac delta with respect to  $k$ ,  $K$  is the vocabulary size, and  $\omega$  is the smoothing weight. Throughout the training data, the internal LM is estimated using Eq. (8). While ordinary ILME (Sec. 2.3) is performed only with the input speech sample, the residual LM statistically considers its behavior in the entire target-domain training data, which may be effective during inference with the target-domain data.

### 3.2. Smoothing of the internal language model

Because the estimation of the internal LM is not always reliably preformed in the target domain, the probability distribution  $p(\mathbf{y}; \theta_{\text{E2E}}^S)$  may be inaccurate. If the probability of the token of interest is incorrectly estimated to be significantly low, the target distribution  $r(\mathbf{y})$  diverges to infinity. To prevent this, we use a temperature  $T > 1$ , as introduced in [27], to soften its distribution.

$$\tilde{p}(\mathbf{y}; \theta_{\text{E2E}}^S)_i = \frac{\exp(z_i/T)}{\sum_{i=1}^K \exp(z_i/T)} + \epsilon, \quad (10)$$

where  $z_i$  is the output of the decoder of the E2E ASR model. We further introduce a small value,  $\epsilon$  to avoid log zero in Eq. (8). By replacing  $p(\mathbf{y})$  with  $\tilde{p}(\mathbf{y})$ , the softened target is used instead as

$$\tilde{r}(\mathbf{y}) = \log q_{\mathbf{y}}^* - \gamma \log \tilde{p}(\mathbf{y}; \theta_{\text{E2E}}^S). \quad (11)$$

### 3.3. Training of the residual language model

The residual LM is trained to minimize the distance between the target distribution  $\tilde{r}(\mathbf{y})$  and the model output  $f(\mathbf{y})$ . A straightforward approach is to minimize the L1 norm or MSEs between the model output and the target distribution. However, in our preliminary experiments, the trained model did not reasonably perform. We assume that this was because the Euclidean-based optimization attempted to minimize the distances of all vocabulary entries equally, which did not contribute to improving recognition performance.

To stably train the residual LM, we decompose the target function by introducing a normalization term, as follows:

$$\tilde{r}(\mathbf{y}) = \log q(\mathbf{y}) + \log N(\mathbf{y}), \quad (12)$$

where  $q(\mathbf{y})$  and  $N(\mathbf{y})$  are defined as

$$\begin{aligned} q(\mathbf{y}) &= \text{softmax} \left( \log q_{\mathbf{y}}^* - \gamma \log \tilde{p}(\mathbf{y}; \theta_{\text{E2E}}^S) \right) \\ &= \frac{q_{\mathbf{y}}^*}{\tilde{p}(\mathbf{y}; \theta_{\text{E2E}}^S)^\gamma} \cdot \frac{1}{N(\mathbf{y})}, \end{aligned} \quad (13)$$

$$N(\mathbf{y}) = \sum_{k=1}^K \frac{q_{\mathbf{y},k}^*}{\tilde{p}(\mathbf{y}; \theta_{\text{E2E}}^S)_k^\gamma}. \quad (14)$$

Thus, the target function  $\tilde{r}(\mathbf{y})$  is decomposed into the probabilistic term,  $\log q(\mathbf{y})$ , and the bias term,  $\log N(\mathbf{y})$ . We propose to separately minimize the distances pertaining to these terms, with a combination of cross-entropy of probabilities and the Euclidean distance of the biases.

#### 3.3.1. Cross-entropy loss for the probabilistic term

LMs are generally optimized by minimizing the negative log likelihood by computing the cross entropy, which is equivalent to minimizing the perplexity. Inspired by this, we apply cross-entropy to optimize the probabilistic term,  $\log q(\mathbf{y})$ , in Eq. (12). To extract the probabilistic factor from the residual LM, we marginalize the model output by applying softmax function.

$$p(\mathbf{y}; \theta_{\text{res}}^T) = \text{softmax}(f(\mathbf{y}; \theta_{\text{res}}^T)). \quad (15)$$

Then a cross-entropy loss is accumulated in the target-domain dataset  $\mathcal{D}_T$  as

$$\mathcal{L}_{\text{CE}} = - \sum_{\mathbf{y} \in \mathcal{D}_T} \sum_i q(y_i) \log p(y_i; \theta_{\text{res}}^T). \quad (16)$$

#### 3.3.2. Mean-squared-error loss for the bias term

We also aim to minimize the distance between bias terms of the target and model. The bias of the residual LM can be derived as

$$f(\mathbf{y}; \theta_{\text{res}}^T) - \log p(\mathbf{y}; \theta_{\text{res}}^T) = \log \sum_{k=1}^K \exp(f(y_k; \theta_{\text{res}}^T)). \quad (17)$$

The bias term of the target,  $\log N(\mathbf{y})$  in Eq. (12), is also a scalar value, defined in Eq. (14).

We further simplify those terms to define an objective function. In the case that the label  $q_{y_i}^*$  is a hard label, i.e.,  $\omega = 0$ , Eq. (14) becomes

$$N(y_i) = \frac{1}{\tilde{p}(y_i; \theta_{\text{E2E}}^S)^\gamma}. \quad (18)$$

Therefore, we approximate the bias using Eq. (18) and define an MSE loss as

$$\begin{aligned} \mathcal{L}_{\text{MSE}} &= \frac{1}{|\mathcal{D}_T|} \sum_{\mathbf{y} \in \mathcal{D}_T} \sum_i (f(y_i; \theta_{\text{res}}^T) - \log p(y_i; \theta_{\text{res}}^T) \\ &\quad + \gamma \log \tilde{p}(y_i; \theta_{\text{E2E}}^S))^2. \end{aligned} \quad (19)$$

#### 3.3.3. Integrated objective function

Finally, we integrate aforementioned optimization in a hybrid manner. The cross-entropy and MSE losses are combined with a weighted sum using a parameter  $\eta$  as follows.

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \eta \mathcal{L}_{\text{MSE}}. \quad (20)$$

Note that both losses contain statistic terms of the internal LM, in  $q(y_i)$  of Eq. (16) and in Eq. (19). Therefore, the trained residual LM considers the statistical behaviors of the internal LM in the target-domain data.

## 4. Experiments

### 4.1. Cross-domain evaluation

To evaluate the effectiveness of the proposed residual LM on domain adaptation, we evaluated cross-domain scenarios in English and Japanese.

#### 4.1.1. Experimental setup

For the English evaluation, we trained an E2E ASR model with the Librispeech dataset [28], a read speech corpus, and applied it to the TED-LIUM 3 [29] dev/test set, which is a spontaneous lecture style. For Japanese, we trained an E2E ASR model with the lecture style CSJ corpus [30], and then applied it to the LaboroTV dev set [31], a corpus of TV programs. We trained streaming Transformer E2E ASR models following [32]. The input acoustic features were 80-dimensional filter bank features. The transformer architecture consisted of 12 encoder blocks and six decoder blocks, with four-head 256-unit attention layers and 2048-unit feed-forward layers. Contextual block encoding [33] was applied to the encoder with a block size of 40, a shift size of 16, and a look-ahead size of 8. The models were trained using multitask learning with CTC loss, as in [6], with a weight of 0.3. We used the Adam optimizer and Noam learning rate decay, and applied SpecAugment [34].

External LMs for baseline shallow fusion [16] as well as the proposed residual LMs were trained using the text-only

Table 1: ASR results in cross-domain adaptation scenarios.

	LS $\rightarrow$ TEDLIUM3 (WER)			CSJ $\rightarrow$ LaboroTV (CER)	
	Dev	Test	Dec. Speed	Dev	Dec. Speed
Shallow Fusion [16]	13.2	12.6	x1.0	24.6	x1.0
Density Ratio [21]	12.9	12.7	x0.92	<b>21.9</b>	x0.97
ILME [22]	12.9	12.2	x0.58	23.7	x0.58
<i>w/ Smoothing</i>	12.9	12.2	—	24.3	—
Residual LM	<b>12.6</b>	<b>12.1</b>	<b>x1.08</b>	22.7	<b>x1.04</b>
<i>w/o Smoothing</i>	12.9	12.2	—	24.1	—

data of the training set in the target corpora, i.e. TED-LIUM3 and LaboroTV. Both LMs were four-layer unidirectional LSTM with 1024 units for the English task and two-layer unidirectional LSTM with 2048 units for Japanese. We applied the byte-pair encoding subword tokenization with 5000 token classes for English LMs. The tokens for Japanese LMs had 3262 character classes. The training weight in Eq. (8) was set as  $\gamma = 0.3$ . We set the temperature in Eq. (10) as  $T = 2$  and the loss integration weight in Eq. (20) as  $\eta = 0.1$  for all experiments.

In addition to shallow fusion, we compared residual LMs with density ratio [21] and ILME [22]. After the parameter search, the LM weight was set to  $\lambda_{LM} = 0.6$ , and the weight for density ratio in Eq. (5) was  $\lambda_{DR} = \{0.1, 0.3\}$ , for the respective dataset. The internal LM weight in Eq. (6) was set to  $\lambda_{ILM} = 0.3$  for both languages. The beam size for decoding was 10. The internal LM was estimated by replacing the encoder output with a zero-filled tensor as in [22, 24]. We also measured the inference speed using randomly sampled 100 utterances from each dev set of the target domain.

#### 4.1.2. Experimental results

The experimental results are listed in Table. 1. In both the English and Japanese scenarios, the density ratio approach and ILME achieved lower WERs than the baseline shallow fusion. The proposed residual LM performed better than ILME and achieved the best performance in the English adaptation with a WER of 12.1 % on the test set (4.0 % WER relative improvement over the shallow fusion). Although, for the Japanese task, the residual LM performed poorer than the density ratio approach, there was an improvement from the baseline shallow fusion and ILME. We assume that the residual LM can consider the statistical behavior of the internal LM, as discussed in Sec. 3.1.

The results of the decoding speed are shown relative to shallow fusion as a base. ILME required almost twice a decoding duration as shallow fusion, because the decoders of the E2E ASR models were required to compute twice, once for ASR and once for ILME. The density ratio was also slightly slower than the shallow fusion because it was required to compute two LMs. The proposed residual LM was slightly faster than the baseline shallow fusion, because it can omit softmax operation, particularly in the larger vocabulary size in English setup.

We performed further ablation studies on the proposed residual LM. When we replaced the smoothed target  $\tilde{r}(\mathbf{y})$  with the regular target  $r(\mathbf{y})$ , as defined in Eq. (8), we observed a significant performance drop from the smoothed target in the Japanese case. We assume that the internal LM estimation is not always accurate in the target domain. On the other hand, applying smoothing did not aid the ILME to improve performance. We assume that using the smoothed soft labels in training gains positive effect similarly to the knowledge distillation learning [35].

Table 2: Chinese/Japanese CERs in intra-domain scenarios.

	AISHELL-1		CSJ		
	Dev	Test	eval1	eval2	eval3
Shallow Fusion [16]	5.7	6.3	6.0	4.4	5.1
ILME [22]	5.8	6.4	6.3	4.5	5.3
Residual LM	<b>5.6</b>	<b>6.1</b>	<b>5.5</b>	<b>4.2</b>	<b>4.6</b>

Table 3: Librispeech WERs in an intra-domain scenario.

	test-clean	test-other	Dec. speed
Shallow Fusion [16]	<b>2.3</b>	<b>5.2</b>	x1.0
ILME [22]	<b>2.3</b>	<b>5.2</b>	x0.51
Residual LM	<b>2.3</b>	5.3	<b>x1.13</b>

## 4.2. Intra-domain evaluation

### 4.2.1. Chinese/Japanese corpora

We performed an intra-domain evaluation to determine if the proposed residual LM did not have an adverse impact in the matched conditions. The residual LMs were evaluated using AISHELL-1 [36] and in the CSJ evaluation set. The experiments followed the configuration in Sec. 4.1.1, except the LMs for AISHELL-1 consisted of two LSTM layers with 650 units, whose output is 4233 character classes. The LM weights were set to  $\lambda_{LM} = \{0.2, 0.3\}$  in the respective evaluation sets, the weights for the ILME were  $\lambda_{ILM} = \{0.1, 0.3\}$ , and the training parameters were  $\gamma = \{0.1, 0.3\}$  respectively.

The character error rate (CER) results are presented in Table. 2. In our reproduction of ILME, even with the best effort to search for the parameters, we observed degradation in both the AISHELL-1 and CSJ evaluation sets. The literature reported that ILME had a positive effect even in the intra-domain scenarios [22], but they were evaluated only in the English dataset, and were not tested in other languages. In contrast, our proposed residual LM performed robustly throughout the corpora and even had lower CERs particularly in the Japanese scenario.

### 4.2.2. Transformer LM evaluation using Librispeech

Lastly, we evaluated the state-of-the-art architecture using the Librispeech dataset. For the E2E ASR model, we adopted 12 conformer encoder blocks [9] with 512-unit eight-head attention and 2048 unit feed forward layers, followed by six transformer decoder with 512-unit eight attention heads and 2048 unit feed forward layers. The LMs were 16-layer transformers. We set the parameters as  $\lambda_{LM} = 0.6$ ,  $\lambda_{ILM} = 0.3$  and the beam size was 30. Inference speed was also evaluated using the test-clean set.

Table. 3 shows the results. No significant difference was observed between the shallow fusion and ILME in this setup. Although we observed slight degradation in the test-other set, both ILME and the proposed residual LM performed robustly with the large model architecture. We observed similar tendency of inference speed as in Table. 1.

## 5. Conclusion

We propose a simple external LM fusion method for domain adaptation, which considers the internal LM estimation in its training. We directly modeled the ratio of an external target domain LM to an internal LM of the E2E ASR model, which is called residual LM. The residual LM was stably trained using a combination of cross-entropy and MSE losses. The experimental results indicated that the proposed residual LM performed better than the internal LM estimation in most of the cross-domain and intra-domain scenarios.

## 6. References

- [1] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proc. of 23rd International Conference on Machine Learning*, 2006, pp. 369–376.
- [2] Y. Miao, M. Gowayyed, and F. Metze, “EESSEN: End-to-end speech recognition using deep RNN models and WFST-based decoding,” in *Proc. of ASRU Workshop*, 2015, pp. 167–174.
- [3] D. Amodei *et al.*, “Deep Speech 2: End-to-end speech recognition in English and Mandarin,” in *Proc. of 33rd International Conference on Machine Learning*, vol. 48, 2016, pp. 173–182.
- [4] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, “Attention-based models for speech recognition,” in *Proc. of NIPS*, 2015, pp. 577–585.
- [5] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *Proc. of ICASSP*, 2016, pp. 4960–4964.
- [6] S. Watanabe, T. Hori, S. Kim, J. R. Hershey, and T. Hayashi, “Hybrid CTC/attention architecture for end-to-end speech recognition,” *Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1240–1253, 2017.
- [7] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang *et al.*, “A comparative study on transformer vs RNN in speech applications,” in *Proc. of ASRU Workshop*, 2019, pp. 449–456.
- [8] A. Graves, A.-R. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *Proc. of ICASSP*, 2013, pp. 6645–6649.
- [9] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. of Interspeech*, 2020, pp. 5036–5040.
- [10] Q. Zhang, H. Lu, H. Sak, A. Tripathi, E. McDermott, S. Koo, and S. Kumar, “Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss,” in *Proc. of ICASSP*, 2020, pp. 7829–7833.
- [11] M. Del Rio, N. Delworth, R. Westerman, M. Huang, N. Bhandari, J. Palakapilly, Q. McNamara, J. Dong, P. Zelasko, and M. Jetté, “Earnings-21: A practical benchmark for ASR in the wild,” in *Proc. of Interspeech*, 2021, pp. 3465–3469.
- [12] K. Yao, D. Yu, F. Seide, H. Su, L. Deng, and Y. Gong, “Adaptation of context-dependent deep neural networks for automatic speech recognition,” in *Spoken Language Technology Workshop (SLT), 2012 IEEE*, 2012, pp. 366–369.
- [13] M. Delcroix, K. Kinoshita, C. Yu, A. Ogawa, T. Yoshioka, and T. Nakatani, “Context adaptive deep neural networks for fast acoustic model adaptation in noisy conditions,” in *Proc. of ICASSP*, 2016, pp. 5270–5274.
- [14] O. Klejch, J. Fainberg, and P. Bell, “Learning to adapt: A meta-learning approach for speaker adaptation,” in *Proc. of Interspeech*, 2018, pp. 867–871.
- [15] E. Tsunoo, Y. Kashiwagi, S. Asakawa, and T. Kumakura, “End-to-end adaptation with backpropagation through WFST for on-device speech recognition system,” in *Proc. of Interspeech*, 2019, pp. 764–768.
- [16] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Sathesh, S. Sengupta, A. Coates *et al.*, “Deep speech: Scaling up end-to-end speech recognition,” *arXiv preprint arXiv:1412.5567*, 2014.
- [17] A. Kannan, Y. Wu, P. Nguyen, T. N. Sainath, Z. Chen, and R. Prabhavalkar, “An analysis of incorporating an external language model into a sequence-to-sequence model,” in *Proc. of ICASSP*, 2018, pp. 5824–5828.
- [18] C. Gulcehre, O. Firat, K. Xu, K. Cho, L. Barrault, H.-C. Lin, F. Bougares, H. Schwenk, and Y. Bengio, “On using monolingual corpora in neural machine translation,” *arXiv preprint arXiv:1503.03535*, 2015.
- [19] A. Sriram, H. Jun, S. Sathesh, and A. Coates, “Cold fusion: Training seq2seq models together with language models,” in *Proc. of Interspeech*, 2018, pp. 387–391.
- [20] C. Shan, C. Weng, G. Wang, D. Su, M. Luo, D. Yu, and L. Xie, “Component fusion: Learning replaceable language model component for end-to-end speech recognition system,” in *Proc. of ICASSP*, 2019, pp. 5361–5635.
- [21] E. McDermott, H. Sak, and E. Variani, “A density ratio approach to language model fusion in end-to-end automatic speech recognition,” in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 434–441.
- [22] Z. Meng, S. Parthasarathy, E. Sun, Y. Gaur, N. Kanda, L. Lu, X. Chen, R. Zhao, J. Li, and Y. Gong, “Internal language model estimation for domain-adaptive end-to-end speech recognition,” in *2021 IEEE Spoken Language Technology Workshop (SLT)*, 2021, pp. 243–250.
- [23] Z. Meng, N. Kanda, Y. Gaur, S. Parthasarathy, E. Sun, L. Lu, X. Chen, J. Li, and Y. Gong, “Internal language model training for domain-adaptive end-to-end speech recognition,” in *Proc. of ICASSP*, 2021, pp. 7338–7342.
- [24] M. ZeinEdein, A. Glushko, W. Michel, A. Zeyer, R. Schlüter, and H. Ney, “Investigating methods to improve language model integration for attention-based encoder-decoder asr models,” in *Proc. of Interspeech*, 2021, pp. 2856–2860.
- [25] A. Zeyer, A. Merboldt, W. Michel, R. Schlüter, and H. Ney, “Librispeech transducer model with internal language model prior correction,” in *Proc. of Interspeech*, 2021, pp. 2052–2056.
- [26] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proc. of CVPR*, 2016, pp. 2818–2826.
- [27] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531* 2.7, 2015.
- [28] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “LibriSpeech: an ASR corpus based on public domain audio books,” in *Proc. of ICASSP*, 2015, pp. 5206–5210.
- [29] F. Hernandez, V. Nguyen, S. Ghannay, N. Tomashenko, and Y. Esteve, “Ted-lium 3: twice as much data and corpus repartition for experiments on speaker adaptation,” in *International conference on speech and computer*. Springer, 2018, pp. 198–208.
- [30] K. Maekawa, H. Koiso, S. Furui, and H. Isahara, “Spontaneous speech corpus of Japanese,” in *Proc. of the International Conference on Language Resources and Evaluation (LREC)*, 2000, pp. 947–9520.
- [31] S. Ando and H. Fujihara, “Construction of a large-scale japanese asr corpus on tv recordings,” in *Proc. of ICASSP*, 2021, pp. 6948–6952.
- [32] E. Tsunoo, C. Narisetty, M. Hentschel, Y. Kashiwagi, and S. Watanabe, “Run-and-back stitch search: novel block synchronous decoding for streaming encoder-decoder ASR,” *arXiv preprint arXiv:2201.10190*, 2022.
- [33] E. Tsunoo, Y. Kashiwagi, T. Kumakura, and S. Watanabe, “Transformer ASR with contextual block processing,” in *Proc. of ASRU Workshop*, 2019, pp. 427–433.
- [34] 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. of Interspeech*, 2019.
- [35] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.
- [36] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, “AIShell-1: An open-source Mandarin speech corpus and a speech recognition baseline,” in *Oriental COCOSA*, 2017, pp. 1–5.