

INCREMENTAL LEARNING FOR END-TO-END AUTOMATIC SPEECH RECOGNITION

Li Fu, Xiaoxiao Li, Libo Zi, Zhengchen Zhang, Youzheng Wu, Xiaodong He, Bowen Zhou

JD AI Research, Beijing, China

ABSTRACT

In this paper, we propose an incremental learning method for end-to-end Automatic Speech Recognition (ASR) which enables an ASR system to perform well on new tasks while maintaining the performance on its originally learned ones. To mitigate catastrophic forgetting during incremental learning, we design a novel explainability-based knowledge distillation for ASR models, which is combined with a response-based knowledge distillation to maintain the original model’s predictions and the “reason” for the predictions. Our method works without access to the training data of original tasks, which addresses the cases where the previous data is no longer available or joint training is costly. Results on a multi-stage sequential training task show that our method outperforms existing ones in mitigating forgetting. Furthermore, in two practical scenarios, compared to the target-reference joint training method, the performance drop of our method is 0.02% Character Error Rate (CER), which is 97% smaller than the drops of the baseline methods.

Index Terms— incremental learning, automatic speech recognition, knowledge distillation

1. INTRODUCTION

While end-to-end Automatic Speech Recognition (ASR) models have been widely used in many applications [1, 2], they usually suffer from performance degradation when applied to new tasks with new accents, new words, or new acoustic environments, etc. Thus, it is worth adapting a pre-trained ASR model to new tasks while maintaining its performance on original tasks. One real-world case is to enable an ASR model trained on native speakers to also recognize speeches with foreign accents. Another case is to add the most recent internet slang to a legacy model’s recognition ability. Note that for both scenarios, any modifications to the pre-trained ASR models should not degrade their performance on the previously learned tasks.

Naively, one can jointly train the acoustic model on the union of the new task data and the original task data. However, with the continually emergence of new tasks for ASR, joint training will be required repeatedly which is practically infeasible due to the usually high training cost. In addition, the original task data may be unavailable due to security

and privacy issues [3]. An alternative solution is to apply a specifically designed language model for decoding schemes or post corrections; but this approach is limited by some essential problems arising from the acoustic model which are still unsolved [4]. Fine-tuning is usually used to adapt a pre-trained model to new tasks by modifying its parameters solely based on training data from new tasks [5]. However, the performance of a fine-tuned model on original tasks may be degraded due to the absence of mechanisms preventing the model from forgetting.

Incremental learning studies the problem of gradually learning on new tasks while maintaining the existing knowledge [6, 7]. Based on incremental learning, we propose a novel training method for end-to-end ASR models in adaptation to new tasks without obvious forgetting. In particular, our method only uses the pre-trained model and the data from new tasks, but *does not* use any data from original tasks, which reduces training cost and addresses data privacy issues. The advantages of our method over directly using the pre-trained model, joint training, and fine-tuning are summarized in Table 1.

To maintain a model’s previous knowledge without access to any data from the original task, one can align the model’s *behavior* with the model pre-trained on the original task. Thus, we adopt a teacher-student framework [8], where the teacher model is the pre-trained ASR model; while the student model is initialized with the pre-trained ASR model. Ideally, one can force the outputs of the two models to be similar for a sufficiently abundant input set with high modality. However, in our incremental learning scenario, alignment of the two models can only exploit the data from the new task, which has limited similarity with the data from the original task in terms of data distribution. Empirically, such similarity can be exploited by Response-based Knowledge Distillation (RBKD) mechanisms to help the student model align the predictions of the teacher model on original tasks [9]. However, due to the difference in distribution between the new task data and the original task data, we believe that only aligning the predictions of the teacher model and the student model using RBKD is not sufficient for the student model to adequately learn the teacher model’s behaviour [10]. Thus, we propose a novel Explainability-based Knowledge Distillation (EBKD) to help the student model also learn the “reason” for the predictions produced by the teacher model.

Table 1. Advantages of Incremental Learning (IL) without access to Original Task (OT) data over: directly using the pre-trained model, joint training, and fine-tuning, in terms of access to OT data or New Task (NT) data, training cost, performance on OTs or NTs. Generally, IL requires no access to OT data, requires minimum training cost, while showing good performance on both OTs and NTs.

	Pre-training	Joint training	Fine-tuning	IL
Access to OT data	Yes	Yes	No	No
Access to NT data	No	Yes	Yes	Yes
Training cost	High	High	Low	Low
Performance on OTs	Good	Good	Poor	Good
Performance on NTs	Poor	Good	Good	Good

Specifically, to train the student model, we propose a novel training loss function involving the following three terms: 1) a term associated with a conventional Connectionist Temporal Classification (CTC) loss [11], which is adopted to help the ASR model learn on the new task; 2) a term associated with a novel EBKD proposed for the ASR incremental learning; 3) a term associated with a RBKD from the teacher model. In particular, the EBKD term and the RBKD term are designed to maintain the pre-trained model’s predictions and the “reason” for the predictions. Experimentally, for ASR incremental learning tasks, the novel EBKD term significantly improves the performance of our method in comparison with the baseline methods involving only the RBKD term [12]. Our main contributions are shown as follows:

1. We propose a new incremental learning method for end-to-end ASR to improve the model’s performance on new tasks while maintaining its previous knowledge. Compared with the joint training method, ours does not require the training set of the pre-trained model, thus is computationally more efficient.
2. We propose a novel EBKD for ASR incremental learning which largely reduces the CERs compared with the baseline methods [12] on a multi-stage sequential training task.
3. We also evaluate the effectiveness of our method in two practical scenarios: for the ASR model to incrementally learn a new accent and new words, respectively; our method outperforms the baseline methods [12] with obvious CER decreasing.

2. RELATED WORK

Existing incremental learning methods mainly fall into three categories, based on how the original task data is used: a) original task data is involved in generating synthetic data for training; b) original task data is sampled when constructing a new data set for training; c) original task data is *not* used for training [13]. Although methods from the first two categories may mitigate the forgetting by leveraging the knowledge of original task data, sophisticated algorithms for data synthesis

and/or sampling are required. Moreover, these methods fail when the original task data is not available due to data privacy issues.

In this paper, we focus on the category of incremental learning methods *without access to* original task data. Existing methods usually adopt RBKD and/or Elastic Weight Consolidation (EWC) for maintaining the previous knowledge [14]. In comparisons, for ASR domain extension tasks, RBKD usually outperforms EWC [12]. One possible reason is that EWC applies overly tight constraints to the model parameters, which likely causes impaired learning for new tasks [15]. While we focus on end-to-end ASR models, incremental learning methods for hybrid ASR models based on RBKD and/or EWC are also explored and have achieved clear improvements over fine-tuning [16]. Also for hybrid ASR models, [17] proposed to modify the model architecture by adding new parameters during training. However, increment of model complexity usually causes higher time consumption during inference.

3. INCREMENTAL LEARNING FOR ASR

3.1. Problem description

Given an ASR model pre-trained on original task data D_1 and a data set D_2 associated with a new task (different from D_1 in terms of accents, words, or acoustic environments, etc), the incremental learner aims to train a new ASR model that:

- performs well on the new task associated with D_2 ;
- has similar performance with the pre-trained ASR model on the original task associated with D_1 .

Note that the learner has access to the pre-trained ASR model and D_2 , but *has no access to* D_1 used for model pre-training. For convenience, we denote $D_2 = \{\mathbf{x}^i, \mathbf{y}^i | i \in \{1, \dots, N\}\}$, with N the number of samples; $\mathbf{x}^i \in \mathbf{R}^{F \times S_i}$ is the i^{th} sample of D_2 which is a sequence of F -dimensional acoustic features with length S_i ; $\mathbf{y}^i \in \mathbf{L}^{U_i}$ is the associated label sequence with length U_i , with \mathbf{L} the finite label alphabet.

3.2. Incremental learning method

We propose a novel loss function for incremental learning, which contains a CTC loss [11], a RBKD loss, and an EBKD loss. The CTC loss is introduced to help the model learn on the new task; while the RBKD loss and the EBKD loss are designed to maintain the model’s performance on the original task. Note that we adopt CTC loss and associated ASR model architectures here only for an example. Our work can be easily extended to alternative end-to-end ASR model frameworks, such as Recurrent Neural Network Transducer (RNN-T) [18], Listen Attend and Spell (LAS) [19], etc. Similar to [20], we consider ASR models involving Self-Attention Blocks (SABs). For example, our experiments use a

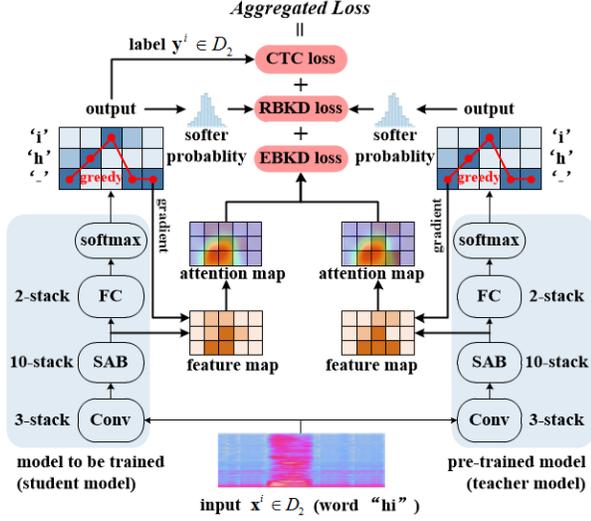


Fig. 1. Outline of our incremental learning method for end-to-end ASR. We adopt a teacher-student framework, where the teacher model (right) is the pre-trained ASR model; while the student model (left) is initialized with the pre-trained ASR model. The training loss of our method is aggregated from three terms: a CTC loss solely based on the model to be trained (student model), which is introduced to help the model learn on the new task; two distilling losses (consisting of a RBKD loss and an EBKD loss) leveraging the pre-trained model (teacher model), which are designed to maintain the model’s performance on the original task.

model architecture with a successive stack of 3 Convolutional (Conv) layers, 10 SABs, 2 Fully Connected (FC) layers, and a feature-wise softmax activation. In this section, we drop the model parameters for simplicity, which are described in detail in Section 4.1. The outline of our method is summarized in Figure 1.

CTC loss: For any labeled sample $\{x^i, y^i\} \in D_2$, denote the softmax output of the model as $[\pi_v^1, \dots, \pi_v^k, \dots, \pi_v^K]$, where K is the sequence length and $\pi_v^k \in \mathbf{R}^M$ is the probability mass with M the alphabet size plus one blank symbol; $v = 1$ and $v = 2$ denote the pre-trained model and the model to be trained, respectively. A valid CTC path with K symbol sequence $c = (c_1, \dots, c_K)$ is a variant of the true label y^i that allows occurrences of blank symbols and repetitions. Then, the CTC loss can be written as [11]:

$$L_{\text{CTC}} = -\log\left(\sum_{c \in \mathcal{C}(x^i, y^i)} \prod_{k=1}^K \pi_{2,c_k}^k\right) \quad (1)$$

where $\mathcal{C}(x^i, y^i)$ is the set of all valid CTC paths; π_{v,c_k}^k is the probability corresponding to the symbol c_k in the vector π_v^k .

RBKD loss and EBKD loss: To maintain a model’s previous knowledge without access to any data from the original task, one can align the model’s *behavior* with the model pre-

trained on the original task. As discussed in Section 1, we employ a RBKD and an EBKD mechanism to help the student model to learn the prediction and the “reason” behind the prediction of the teacher model.

To guide the student model to produce similar predictions as the teacher model, we use the same RBKD loss as in [12]:

$$L_{\text{RBKD}} = -\sum_{k=1}^K \sum_{m=1}^M p_{1,m}^k \log(p_{2,m}^k) \quad (2)$$

where $p_{v,m}^k = (\pi_{v,m}^k)^{1/T} / \sum_{m=1}^M (\pi_{v,m}^k)^{1/T}$, and T is a temperature scalar that produces the softer probability.

To help the student model also learn the “reason” for the predictions produced by the teacher model, we propose a novel EBKD loss for ASR incremental learning to train the student model to have similar attention maps [22] to those of an already trained teacher model. Specifically, for input x^i , we choose model v ’s output of the last SAB layer as the feature map, which is denoted as $\mathbf{A}_v \in \mathbf{R}^{d_h \times K}$ with hidden dimension d_h . Then we calculate the importance of each element in \mathbf{A}_v to the model’s most-probable prediction by taking the gradient of the model’s greedy prediction probability $p_v = \prod_{k=1}^K \max_m (\pi_{v,m}^k)$ over \mathbf{A}_v , yields the importance map $\alpha_v \in \mathbf{R}^{d_h \times K}$ [22]:

$$\alpha_v = \partial \log(p_v) / \partial \mathbf{A}_v \quad (3)$$

Note that one can also consider other choices of the feature map, or even a combination of multiple feature maps. Then, the attention map $\mathbf{Q}_v \in \mathbf{R}^{d_h \times K}$ is defined as

$$\mathbf{Q}_v = f_{\text{ReLU}}(\alpha_v \odot \mathbf{A}_v) \quad (4)$$

where $f_{\text{ReLU}}(\cdot)$ is the element-wise ReLU function, which introduces a positive influence on the result of interest [22]; and \odot denotes element-wise multiplication. In our method, the attention map provides a reasoning associated with the high-level features of interest for the model’s most-probable prediction, which is proved to be useful for the student model learning the behaviour of the teacher model. Finally, the attention map is normalized feature-wise which yields our EBKD loss:

$$L_{\text{EBKD}} = \frac{1}{K} \sum_{k=1}^K \left\| \frac{\mathbf{Q}_2^k}{\|\mathbf{Q}_2^k\|_2} - \frac{\mathbf{Q}_1^k}{\|\mathbf{Q}_1^k\|_2} \right\|_2 \quad (5)$$

where $\mathbf{Q}_v^k \in \mathbf{R}^{d_h}$ is the k^{th} vector of \mathbf{Q}_v .

Aggregated loss: Loss for our incremental learning for end-to-end ASR is:

$$L = L_{\text{CTC}} + \beta L_{\text{RBKD}} + \gamma L_{\text{EBKD}} \quad (6)$$

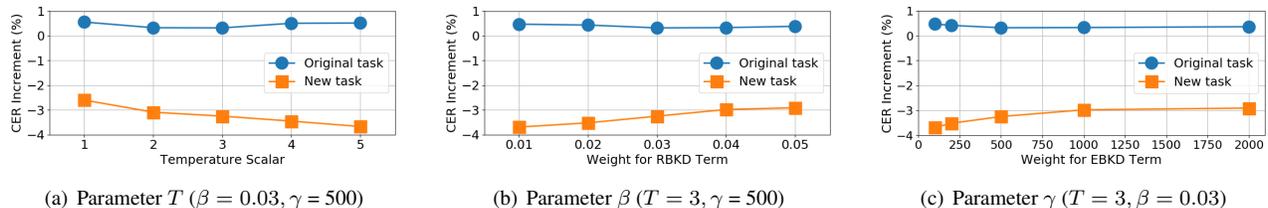


Fig. 2. Compared to the pre-trained model, the CER increments of the model trained by our incremental learning method on both original tasks and new tasks (Performance improvements are indicated by negative CER increments), for a range of choices of (a) the temperature scalar T ; (b) the weight β for the RBKD term; and (c) the weight γ for the EBKD term.

where hyperparameters β and γ are used to balance these terms – their choices are experimentally studied in Section 4.3 (see Figure 2).

Deeper insight of the proposed EBKD loss: For our incremental learning framework, there is no guarantee that the samples \mathbf{x}^i for the new task follow exactly the same distribution as those for the original task. But the attention map generated by the teacher model for any sample \mathbf{x}^i highlights the important internal layer features that are decisive to its output. Thus, our EBKD term, by its formulation, is designed to guide the student model paying more attention to these features “reasoned” by the teacher model. With the EBKD term, the student model is then trained to give the same responses to these important features. Formally, Q_1 should be similar to Q_2 .

Although visualization of the above mentioned effect for speech features is much more difficult than for the image domain [22], we find that there is a significant correlation between knowledge forgetting and EBKD loss in our incremental learning tasks. As shown in Figure 3 (as an example derived from the main results in Table 4), the recognition error on the original task and the value of the proposed EBKD loss show a strong correlation: 0.77 on average for our method and the baseline method involving only the RBKD term [12]. Compared with the baseline method, the knowledge forgetting (CER increment on the original task) of our method is mitigated as our EBKD loss decreases with convergence – **the EBKD loss functions as expected.**

4. RESULTS AND DISCUSSION

4.1. Experimental setup

Data preparation: To evaluate our proposed method more thoroughly, we use a combinatory data set containing 12,301 hours of Mandarin speech from a variety of data subsets with different accents, contents/topics of the text, and total duration [23–27]. Details for these data subsets are shown in Table 2. For the public data, we use their original train-test split. For the other data subsets, we hold out 10% speech samples for testing. Each speech is pre-processed as a sequence of 80-dimensional filter banks, with each element of

Table 2. Details for each composition of the combinatory Mandarin speech data set used in our experiments.

Composition name	Duration (h)	Description
JD in-house data	10,000	Conversational telephony speech
Public data*	2000	Standard Mandarin speech
New accents data	300	Accent telephony speech
New words data	1	Various new words of internet bank speech

*The public data consists of AISHELL-1 (170h) [23], AISHELL-2-10S (1000h) [24], THCHS-30 (40h) [25], Primewords (100h) [26], ST-CMDS (100h) [27], etc.

the sequence obtained from a 20ms window (with 10ms shift) of the speech [28].

ASR Model: The model considered here contains 3 Convs, 10 SABs, and 2 FCs in sequential order (see Figure 1). For the three Convs, the filter sizes are (41,11), (21,11), and (21,11), respectively; the channel counts are 32, 32, and 96, respectively; and the strides are (2,2), (2,1), and (2,1), respectively [11]. All 10 SABs have 8 multi-heads and are 256-dimensional. The output dimensions of the penultimate FC and the last FC are 1024 and 7229, respectively, where 7229 is the number of Chinese characters ever occurring in the entire data set plus a symbol blank. Following the convention of CTC-based ASR, a 4-gram language model, which is trained using the KenLM toolkit [29], is used for inference. The pre-training of the ASR model uses the CTC loss as in (1) [11]. Both the pre-training and our incremental training are trained with a mini batch of 128 and using the same learning rate scheduler as in [20].

4.2. Main experiment: multi-stage sequential training

A 5-stage sequential training scenario: Consider the following 5 data subsets: B_1 : AISHELL-1, B_2 : Primewords, B_3 : ST-CMDS, B_4 : new accents data, B_5 : new words data. In the 1st stage, an ASR model is trained on the training set of B_1 based on the descriptions in Section 4.1. In the 2nd stage, the model pre-trained for B_1 should be modified, given the training set of B_2 , to perform well on both B_1 (the original task) and B_2 (the new task). The subsequent three stages are

Table 3. CERs of our method compared with fine-tuning, RBKD, RBKD with EWC, and joint training, on the Original Tasks (OTs) ever visited and the current New Task (NT), for each stage of the 5-stage sequential training which visits data sets B_1 , B_2 , B_3 , B_4 , and B_5 sequentially. Results show that our method consistently outperforms the two baseline methods with lower average CERs on the OTs and the NTs in all stages, especially in stages 4 and 5 where the distribution and amount of NT data are more different from the OT data.

	Stage 1		Stage 2		Stage 3			Stage 4				Stage 5				
	Initial	OT	NT	OT		NT	OT			NT	OT				NT	
	B_1	B_1	B_2	B_1	B_2	B_3	B_1	B_2	B_3	B_4	B_1	B_2	B_3	B_4	B_5	
Joint training	5.10	4.10	2.09	4.27	2.32	1.78	4.06	2.48	1.99	6.30	4.07	2.48	1.98	6.32	4.06	
Fine-tuning	5.10	8.22	2.14	10.11	13.24	1.66	28.10	34.00	28.77	6.56	27.36	36.56	28.87	10.27	3.89	
RBKD [12]	5.10	4.53	2.17	5.76	5.21	1.67	12.44	20.15	13.19	6.73	13.83	20.32	14.15	10.06	3.98	
RBKD+EWC [12]	5.10	4.12	2.58	4.75	2.79	2.20	6.22	11.80	7.62	8.94	6.47	11.85	7.90	9.64	4.68	
Our method	5.10	4.14	2.18	4.55	2.72	1.77	5.31	8.78	5.41	6.90	5.10	9.11	6.28	8.32	3.97	

similar to the 2nd stage, where in each stage, the model from the previous stage should be modified to adapt to a new task (i.e. B_3 , B_4 , and B_5 , sequentially) given its associated training set.

Performance evaluation protocol: We use CER as the metric for performance evaluation. For stage $s \in \{2, 3, 4, 5\}$, the learned model’s performance on original tasks is jointly reflected by its CER on test set for B_1, \dots, B_{s-1} ; while its performance on the new task is reflected by the CER on the current test set for B_s .

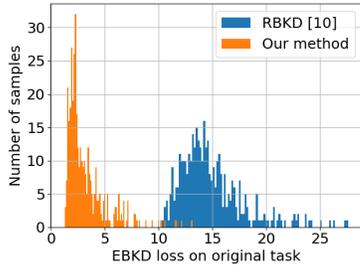
Methods and implementation: We compare our method with the two incremental learning methods in [12]: a vanilla RBKD method and a RBKD method combined with EWC – to our knowledge, these are the only two existing incremental learning methods for ASR without access to original task data. We also compare our method with the fine-tuning method mentioned in Section 1. For all these methods, in each sequential training stage, *only* the training set for the new task is exploited. We also consider the joint training method, which, in each stage, exploits training data accumulated from all previous and current stages; thus, it is supposed to demonstrate a target-reference performance for both new tasks and original tasks [30]. For our method, we set the hyperparameters as $T = 3$, $\beta = 0.03$ and $\gamma = 500$, while their choices will not significantly influence the performance of our method, as will be shown in Section 4.3.

Results and discussion: In Table 3, for each of the 5 sequential training stages, we show CERs of the learned model on the test set for tasks ever visited, for the 5 methods. For better visual comparison between these methods, in Figure 4, for stage 2, 3, 4, and 5, we show the average CER for each method over all tasks ever visited, and also the CER for the task associated with the current stage. Clearly, fine-tuning suffers from catastrophic forgetting, as its average CERs on original tasks are uniformly large for all stages. RBKD alleviates the forgetting to some extent compared with fine-tuning, but its performance on original tasks is still not satisfactory. RBKD combined with EWC shows even lower average CERs than the vanilla RBKD (though still not comparable to our method) on original tasks; but its CERs on the

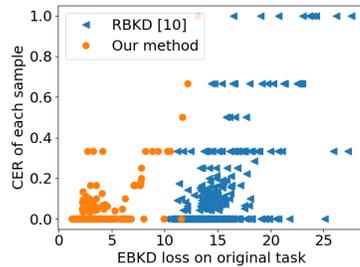
new tasks are larger than other methods for all stages. This is consistent with our discussion in Section 2 that EWC may over-constrain the model parameters, which possibly leads to impaired learning for new tasks. In comparison, on the original tasks, our method clearly outperforms fine-tuning, RBKD, and RBKD with EWC; on the new tasks, our method shows similar performance as RBKD and fine-tuning, while outperforming RBKD with EWC clearly. In particular, the effectiveness of our novel EBKD mechanism in maintaining the model’s performance on original tasks is clearly shown as our method achieves lower average CERs than the RBKD-based methods.

4.3. Choice of hyperparameters

In principle, when T is large, predictions made by the pre-trained model with high confidence will be influential to our model training [31]. Since the pre-trained model usually performs poorly (and also with low confidence) on new tasks, a large T will likely block the “wrong guidance” from the pre-trained model regarding new tasks. For β and γ , since they are weights for the terms devoted to maintaining the trained model’s performance on original tasks, setting them overly large may prevent the model from learning for new tasks. Here, we study the choice of these three hyperparameters. For simplicity, we consider a single-stage incremental learning on the new accents data, with an ASR model pre-trained on the combined training set from the JD in-house data and the public data. Let $T = 3$, $\beta = 0.03$ and $\gamma = 500$ (used in Section 4.2) be the baseline choice of the hyperparameters. Each time, we vary one hyperparameter with the other two fixed to the baseline choice for performance evaluation. Specifically, we consider $T \in \{1, 2, 3, 4, 5\}$, $\beta \in \{0.01, 0.02, 0.03, 0.04, 0.05\}$, and $\gamma \in \{100, 200, 500, 1000, 2000\}$, respectively. As shown in Figure 2, with the above-mentioned hyperparameter choices, the model trained by our method achieves a significant CER decreasing (2.7%–3.8%) on the new tasks, and only a small CER increment (less than 0.5%) on the original tasks when compared with the pre-trained model. In general, within a



(a) Distributions of EBKD loss on original task



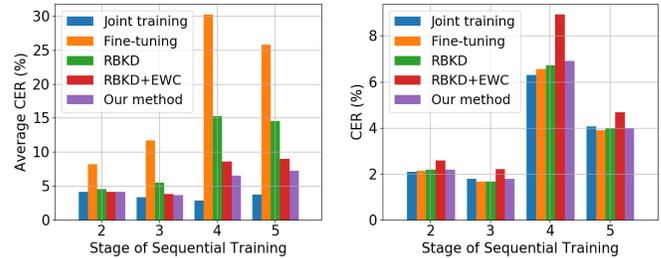
(b) Relations between CER of each sample and EBKD loss on original task

Fig. 3. EBKD losses on original task (400 samples randomly selected) of our method and RBKD [10]: (a) shows the EBKD loss of our method decreases with convergence; (b) shows the correlations between CER of each sample and EBKD loss on the original task are 0.77 for these samples with recognition errors (CER>0) and 0.63 for all selected samples.

proper range, the choice of the hyperparameters does not significantly affect the performance of our method. In practice, one can choose these hyperparameters using a small, held-out validation set.

4.4. Results of two practical scenarios

In practice, an ASR model pre-trained on a large data set for some tasks may need to be adapted to a new task. Here, we compare the three incremental learning methods (including ours), fine-tuning, and joint training (the target-reference method), on adapting a pre-trained ASR model to a new task on the new accents data (scenario I) and a new task on the new words data (scenario J), respectively. Model pre-training is performed on a combined training set from the JD in-house data and the public data. Table 4 shows the CER performance of the above-mentioned methods (and also the pre-trained model itself for reference) on the original tasks and the new tasks for both scenarios I and J . Similar to the experiment in Section 4.2, for both (close-to-practice) scenarios, fine-tuning, RBKD, and RBKD with EWC more or less suffer from forgetting the original tasks. In comparison, the performance of our method on the original tasks is strong for both



(a) Performance on the original tasks (b) Performance on the new tasks

Fig. 4. Average CERs on the original tasks and CERs on the new tasks for stage 2-5 of the 5-stage sequential training, for the 5 methods.

Table 4. CERs of the three incremental learning methods (including ours), fine-tuning, and joint training on the Original Tasks (OTs) and New Tasks (NTs), for adapting a pre-trained ASR model to a new accent (scenario I) and new words (scenario J), respectively. The performance of the pre-trained model itself is included for reference.

	CER on OTs		CER on NTs		Avg. CER
	Sc. I	Sc. J	Sc. I	Sc. J	
Pre-training	4.46	4.46	10.07	18.47	9.37
Joint training	4.79	4.79	6.70	1.63	4.48
Fine-tuning	6.02	7.69	7.19	1.49	5.60
RBKD [12]	5.62	5.76	6.97	2.37	5.18
RBKD+EWC [12]	4.77	4.90	9.24	6.99	6.48
Our method	4.78	4.70	6.82	1.52	4.46

scenarios – even comparable with the performance of joint training which directly uses the training set for both new tasks and original tasks. On the new tasks which the pre-trained model should be adapted to, our method clearly outperforms the two baseline methods, and is comparable with fine-tuning and joint training for the two scenarios in general. Numerically, compared to the target-reference joint training method, the performance drop (i.e. the average CER increment for all test sets in both scenarios) of our method is 0.02% CER, which is 97% smaller than the drops of the baseline methods.

5. CONCLUSIONS

We propose an incremental learning method for end-to-end ASR models' adaptation to new speech recognition tasks without obviously degrading its performance on the original task it is trained for. To overcome forgetting, our method proposes a training loss containing a novel EBKD loss and a RBKD loss based on the pre-trained model. Our method does not require access to the original task data of the pre-trained model; and achieves the best CER performance on both original tasks and new tasks when compared with the baseline methods.

6. REFERENCES

- [1] R. Prabhavalkar, K. Rao, T. Sainath, et al., "A comparison of sequence-to-sequence models for speech recognition," in *18th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2017, pp. 939-943.
- [2] E. Battenberg, J. Chen, R. Child, et al., "Exploring neural transducers for end-to-end speech recognition," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2017, pp. 206-213.
- [3] T. Ha, T. Dang, H. Le, et al., "Security and privacy issues in deep learning: A brief review," *SN Computer Science*, 2020, 1(5):1-15.
- [4] X. Liu, Y. Wang, X. Chen, et al., "Efficient lattice rescoring using recurrent neural network language models," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014, pp. 4908-4912.
- [5] G. Kurata and K. Audhkhasi, "Guiding CTC posterior spike timings for improved posterior fusion and knowledge distillation," in *20th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2019, pp. 1616-1620.
- [6] F. Castro, M. Marin-Jimenez, N. Guil, et al., "End-to-end incremental learning," in *IEEE Conference on European Conference on Computer Vision (ECCV)*, 2018, pp. 233-248.
- [7] Z. Li and D. Hoiem, "Learning without forgetting," in *IEEE Conference on European Conference on Computer Vision (ECCV)*, 2016, pp. 614-629.
- [8] S. Li, L. Li, Q. Hong, et al., "Improving transformer-based speech recognition with unsupervised pre-training and multi-task semantic knowledge learning," in *21th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2020, pp. 5006-5010.
- [9] J. Gou, B. Yu, S. Maybank, et al., "Knowledge distillation: A survey," *International Journal of Computer Vision*, 2021.
- [10] P. Dhar, R. Singh, K. Peng, et al., "Learning without memorizing," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5133-5141.
- [11] S. Amodei, R. Ananthanarayanan, J. Anubhai, et al., "Deep speech 2: End-to-end speech recognition in English and Mandarin," in *International Conference on Machine Learning (ICML)*, 2016, pp. 173-182.
- [12] S. Ghorbani, S. Khorram, and J. Hansen, "Domain expansion in DNN-based acoustic models for robust speech recognition," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 252-259.
- [13] Y. Wu, Y. Chen, L. Wang, et al., "Large scale incremental learning," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 374-382.
- [14] K. Sim, F. Beaufays, A. Benard, et al., "Personalization of end-to-end speech recognition on mobile devices for named entities," in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2019, pp. 23-30.
- [15] J. Schwarz, J. Luketina, W. Czarnecki, et al., "Progress & Compress: A scalable framework for continual learning," in *International Conference on Machine Learning (ICML)*, 2018, pp. 4528-4537.
- [16] B. Houston and K. Kirchhoff, "Continual learning for multi-dialect acoustic models," in *21th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2020, pp. 576-580.
- [17] S. Sadhu and H. Hermansky, "Continual learning in automatic speech recognition," in *21th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2020, pp. 1246-1250.
- [18] A. Graves, "Sequence transduction with recurrent neural networks," in *International Conference on Machine Learning (ICML)*, 2012.
- [19] W. Chan, N. Jaitly, Q. Le, et al., "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 4960-4964.
- [20] J. Salazar, K. Kirchhoff, and Z. Huang, "Self-attention networks for connectionist temporal classification in speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 7115-7119.
- [21] S. Zagoruyko and N. Komodakis, "Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer," in *International Conference on Learning Representations (ICLR)*, 2017.
- [22] R. Selvaraju, M. Cogswell, A. Das, et al., "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 618-626.

- [23] H. Bu, J. Du, X. Na, et al., "AISHELL-1: An open-source Mandarin speech corpus and a speech recognition baseline," in *20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment (O-COCOSDA)*, 2017, pp. 1-5.
- [24] J. Du, X. Na, X. Liu, et al., "AISHELL-2: Transforming Mandarin ASR research into industrial scale," *arXiv preprint arXiv: 1808.10583*, 2018.
- [25] D. Wang and X. Zhang, "THCHS-30: A free Chinese speech corpus," *arXiv preprint arXiv: 1512.01882*, 2015.
- [26] Primewords Information Technology Co. Ltd., Primewords Chinese Corpus Set 1, 2018, <https://www.primewords.cn>.
- [27] ST-CMDS-20170001.1, Free ST Chinese Mandarin Corpus.
- [28] L. Narayana and S. Kopparapu, "Choice of Mel filter bank in computing MFCC of a resampled speech," in *International Conference on Information Science, Signal Processing and their Applications (ISSPA)*, 2010, pp. 121-124.
- [29] K. Heafield, "KenLM: Faster and smaller language model queries," in *Proceedings of the Sixth Workshop on Statistical Machine Translation*, 2011 pp. 187–197.
- [30] M. Delange, R. Aljundi, M. Masana, et al. "A continual learning survey: Defying forgetting in classification tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [31] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *NIPS Deep Learning and Representation Learning Workshop*, 2015.