

INTERMPL: MOMENTUM PSEUDO-LABELING WITH INTERMEDIATE CTC LOSS

Yosuke Higuchi^{1,2}, Tetsuji Ogawa², Tetsunori Kobayashi², Shinji Watanabe¹

¹Carnegie Mellon University, USA ²Waseda University, Japan

ABSTRACT

This paper presents InterMPL, a semi-supervised learning method of end-to-end automatic speech recognition (ASR) that performs pseudo-labeling (PL) with intermediate supervision. Momentum PL (MPL) trains a connectionist temporal classification (CTC)-based model on unlabeled data by continuously generating pseudo-labels on the fly and improving their quality. In contrast to autoregressive formulations, such as the attention-based encoder-decoder and transducer, CTC is well suited for MPL, or PL-based semi-supervised ASR in general, owing to its simple/fast inference algorithm and robustness against generating collapsed labels. However, CTC generally yields inferior performance than the autoregressive models due to the conditional independence assumption, thereby limiting the performance of MPL. We propose to enhance MPL by introducing intermediate loss, inspired by the recent advances in CTC-based modeling. Specifically, we focus on self-conditional and hierarchical conditional CTC, that apply auxiliary CTC losses to intermediate layers such that the conditional independence assumption is explicitly relaxed. We also explore how pseudo-labels should be generated and used as supervision for intermediate losses. Experimental results in different semi-supervised settings demonstrate that the proposed approach outperforms MPL and improves an ASR model by up to a 12.1% absolute performance gain. In addition, our detailed analysis validates the importance of the intermediate loss.

Index Terms— pseudo-labeling, intermediate loss, semi-supervised learning, end-to-end speech recognition, deep learning

1. INTRODUCTION

End-to-end (E2E) automatic speech recognition (ASR) [1–3] has achieved remarkable improvements in performance thanks to innovative sequence-to-sequence modeling techniques [4–7] with sophisticated neural network architectures [8–10]. While E2E ASR has shown promising results on a variety of benchmarks [11–13], training E2E ASR is generally data-hungry: its performance often relies on the availability of abundant labeled (transcribed) speech data [14], which is not always achievable due to high annotation costs.

In order to mitigate the heavy requirement on labeled data, semi-supervised learning has been actively studied in E2E ASR, utilizing a large quantity of unlabeled speech-only data to enhance the model performance. Among diverse semi-supervised learning approaches, pseudo-labeling (PL) [15] has been gathering attention due to its simple yet effective training algorithm [16–24]. In typical PL, a teacher (seed) model is first trained on labeled data and used to generate pseudo-labels by transcribing unlabeled data. A student model is then trained using the labeled and pseudo-labeled data, with the aim of performing better than the teacher. Shallow fusion is often performed during the labeling process, which utilizes an external language model (LM) to generate higher-quality pseudo-labels [19, 25]. Data augmentation also plays an important role in providing the student with informative training signals [17, 18, 21].

In addition to the techniques above, iterating the PL process has shown to be essential for improving ASR performance, periodically updating pseudo-labels as the model training proceeds [20–23]. Momentum PL (MPL) [26–29] is one of the recent iterative methods, inspired by the mean teacher framework [30]. MPL trains a pair of offline (\approx teacher) and online (\approx student) models that interact and learn from each other. In each training step, pseudo-labels are generated on the fly by the offline model via greedy decoding and used as targets to train the online model. The offline model maintains an exponential moving average of the online model weights to stabilize the label generation. Through the interaction between the two models, MPL enables continuous updates on pseudo-labels and, concurrently, improves the label quality.

Connectionist temporal classification (CTC) [4] is a promising approach for conducting iterative PL, particularly MPL, compared to other autoregressive models equipped with a recurrent decoder (e.g., attention-based encoder-decoder [2, 3] and transducer [5]). The non-autoregressive formulation in CTC allows a model to transcribe unlabeled speech data efficiently with its simple, fast and parallelized inference algorithm, a crucial property for the on-the-fly label generation in MPL. Furthermore, CTC is robust against generating collapsed pseudo-labels, which is frequently caused by autoregressive decoding (e.g., word skipping or repeating [25, 31]). This does not necessarily require a model to apply heuristic filtering techniques for excluding erroneous labels that hinder semi-supervised training [22]. However, the ASR performance of CTC often lags behind those of the autoregressive models [11]. This is attributed to the fact that CTC assumes output tokens are conditionally independent of each other, making a model less capable of capturing contextual information.

Hence, our work aims to further enhance the MPL performance while maintaining the advantages of CTC-based modeling. To this end, we propose to introduce intermediate CTC loss [32, 33] to MPL, given the recent advances in non-autoregressive E2E ASR [34]. We consider adopting self-conditional CTC (SC-CTC) [35] and hierarchical conditional CTC (HC-CTC) [36] for MPL. SC-CTC applies auxiliary CTC losses to intermediate model layers and utilizes each intermediate prediction as a condition for subsequent layers. This induces the contextualization of representations, which is beneficial for relaxing the conditional independence assumption. HC-CTC extends SC-CTC by gradually increasing the granularity of each output sequence in a hierarchical manner [37, 38], where the hierarchical structure allows a model to learn the progressive generation of a target sequence. Through SC-CTC and HC-CTC, a model is capable of generating multiple pseudo-labels from its intermediate layers. We thereby explore how pseudo-labels should be generated and used as supervision for intermediate losses during MPL training.

The key contributions of our work are summarized as follows: 1) We propose InterMPL, which enhances MPL-based semi-supervised ASR by intermediate CTC loss; 2) Experimental results show that InterMPL significantly outperforms MPL in various semi-supervised scenarios. We also present a detailed

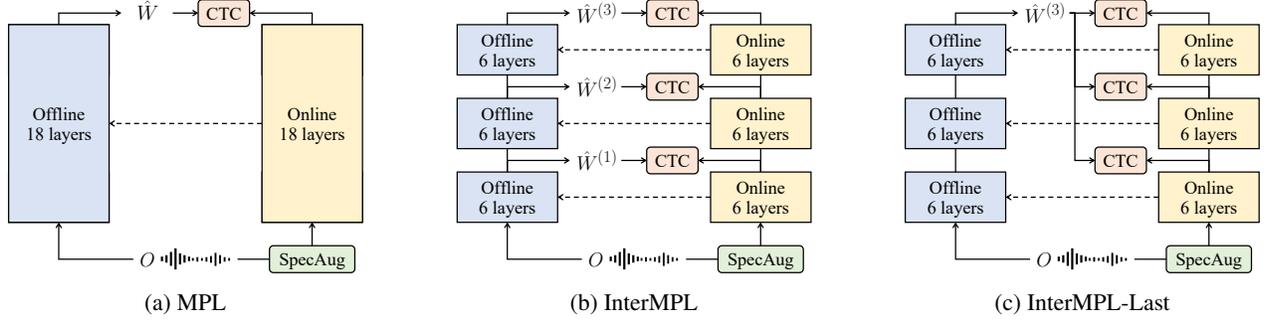


Fig. 1. Comparisons between conventional MPL and proposed InterMPL for semi-supervised ASR. A dashed line (\leftarrow) indicates the momentum update of the offline model using the online model parameters. The number of total losses is set to three for InterMPL (i.e., $|\mathcal{K}| = 3$).

analysis to validate the significance of the intermediate loss; and 3) The codes and recipes are made publicly available at <https://github.com/YosukeHiguchi/espnet/tree/intermpl>.

2. BACKGROUND

2.1. CTC-based End-to-End ASR with Intermediate Loss

E2E ASR is defined as a sequence-mapping problem between a T -length input sequence $O = (\mathbf{o}_t \in \mathbb{R}^F | t = 1, \dots, T)$ and U -length output sequence $W = (w_u \in \mathcal{V} | u = 1, \dots, U)$, where \mathbf{o}_t is an F -dimensional acoustic feature at frame t , w_u is an output token at position u , and \mathcal{V} is a vocabulary. For embedding O into a latent representation space, we construct a Conformer-based model [10] consisting of K encoder layers. The k -th layer takes as input a previous sequence $H^{(k-1)}$ and outputs a sequence $H^{(k)}$ as

$$H^{(k)} = \text{Encoder}^{(k)}(H^{(k-1)}), \quad (1)$$

where $H^{(0)} = O$, and $H^{(k)} = (\mathbf{h}_t^{(k)} \in \mathbb{R}^D | t = 1, \dots, T)$ is a sequence of D -dimensional hidden vectors with the same length as that of O . For simplicity, we define H as the last sequence $H^{(K)}$.

CTC CTC [1] formulates E2E ASR by considering all possible alignments between O and W . To align the sequences at the input frame level, CTC augments an output sequence by allowing consecutive identical tokens and inserting a blank symbol ϵ . Let $A = (a_t \in \mathcal{V} \cup \{\epsilon\} | t = 1, \dots, T)$ be an augmented output sequence, which we refer to an alignment path between O and W . CTC trains a model to predict the paths by minimizing the following loss:

$$\mathcal{L}_{\text{ctc}}(W|H) = -\log \sum_{A \in \mathcal{B}^{-1}(W)} \prod_t p(a_t|H), \quad (2)$$

where $\mathcal{B}(\cdot)$ is the collapsing function [4] that maps A to W by suppressing repeated tokens and removing blank symbols, and $\mathcal{B}^{-1}(W)$ is a set of all possible paths compatible with W .

Intermediate CTC with conditioning Intermediate CTC [33] applies auxiliary CTC losses to the intermediate hidden-state of the encoder. In addition to the original CTC loss applied to the last layer (Eq. (2)), the intermediate losses are calculated as

$$\mathcal{L}_{\text{ic}}(W|O) = \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{ctc}}(W|H^{(k)}), \quad (3)$$

where \mathcal{K} is a set of layer indices where the CTC losses are computed, and we equally distribute the weight across the losses. Note that Eq. (3) always includes the last loss (i.e., $K \in \mathcal{K}$) and is equal to Eq. (2) when $\mathcal{K} = \{K\}$. Self-conditional CTC (SC-CTC) [35] extends Intermediate CTC by conditioning the encoder using a sequence predicted from each intermediate layer. Specifically, after calculating $H^{(k)}$ from Eq. (1) at an intermediate layer, SC-CTC adds a sequence of posterior probability distributions

$$P^{(k)} = (p(a_t|H^{(k)}) \in [0, 1]^{|V|+1} | t = 1, \dots, T) \text{ as}$$

$$H^{(k)} \leftarrow H^{(k)} + \text{Linear}_{|V|+1 \rightarrow D}(P^{(k)}), \quad (4)$$

where $k \in \mathcal{K}$. This has been shown to further improve Intermediate CTC by relaxing the conditional independence assumption in Eq. (2) ($a_t \perp a_{\neq t} | H$). Hierarchical-conditional CTC (HC-CTC) [36] extends SC-CTC by hierarchically increasing the output unit size of each intermediate prediction, using different subword vocabularies.

2.2. Semi-Supervised ASR with Momentum Pseudo-Labeling

In semi-supervised ASR, a seed model is first trained on labeled data $\mathcal{D}_{\text{lab}} = \{(O_n, W_n) | n = 1, \dots, N\}$ using the CTC loss from Eq. (2). Momentum pseudo-labeling (MPL) [26] is then applied to the seed model to improve the performance using unlabeled speech-only data $\mathcal{D}_{\text{unlab}} = \{O_m | m = N+1, \dots, N+M\}$. Figure 1(a) illustrates the training process of MPL based on a pair of *online* and *offline* models. Let ξ and ϕ denote the parameters of the online and offline models, which are initialized with the pre-trained seed model parameters.

Online model training Given the m -th unlabeled sample $O_m \in \mathcal{D}_{\text{unlab}}$ and its encoded sequence H_m , the online model is trained on pseudo-labels \hat{W}_m generated on the fly by the offline model with ϕ :

$$\hat{W}_m = \mathcal{B}(\arg\max_{a_t} p(a_t|H_m, \phi) | t = 0, \dots, T). \quad (5)$$

With the pseudo-labeled sample (O_m, \hat{W}_m) , the online model with ξ is trained via a gradient descent optimization based on the CTC loss $\mathcal{L}_{\text{ctc}}(\hat{W}_m|H_m, \xi)$ from Eq. (2). Here, the input speech is augmented by SpecAugment [39] (as shown in Fig. 1(a)) to facilitate the model training on pseudo-labels. Note that MPL also uses the n -th labeled sample $(O_n, W_n) \in \mathcal{D}_{\text{lab}}$ and trains the online model with supervised loss $\mathcal{L}_{\text{ctc}}(W_n|H_n, \xi)$, which helps the online model stabilize and promote learning from unlabeled data.

Offline model training After every update of the online model, the offline model accumulates the parameters of the online model as $\phi \leftarrow \alpha\phi + (1-\alpha)\xi$, an exponential moving average with a momentum coefficient $\alpha \in (0, 1)$. This momentum update makes the offline model serve as an ensemble of the online models at different training steps, stabilizing and reinforcing the label generation in Eq. (5).

Through the above interaction between the two models, MPL realizes stable and continuous ASR training on unlabeled data, concurrently improving the quality of pseudo-labels.

3. INTERMPL

We propose a semi-supervised ASR method that introduces the intermediate CTC loss to MPL. The conventional MPL is founded on the CTC-based modeling, whose performance can be limited due to the conditional independence assumption (cf. Eq. (2)). To fur-

ther enhance MPL, we adopt SC-CTC or HC-CTC for constructing a seed model, which is expected to facilitate better CTC training/decoding and thus promote the succeeding semi-supervised process with higher-quality pseudo-labels. Given labeled data \mathcal{D}_{lab} , a seed model is trained by a supervised loss $\mathcal{L}_{\text{ic}}(W_n|O_n)$ from Eq. (3) along with the conditioning mechanism in Eq. (4).

Initialized with the seed model trained by SC-CTC or HC-CTC, the online model can accept intermediate supervision using pseudo-labels, and the offline model can generate multiple pseudo-labels from its intermediate layers. This motivated us to consider two approaches, namely **InterMPL** (Fig. 1(b)) and **InterMPL-Last** (Fig. 1(c)), for fully utilizing the intermediate mechanism in MPL.

InterMPL In InterMPL, the offline model generates pseudo-labels from each prediction layer, as shown in Fig. 1(b) with three different outputs. These pseudo-labels are used to calculate a loss for the corresponding layer of the online model. Given the m -th unlabeled sample $O_m \in \mathcal{D}_{\text{unlab}}$, the k -th offline encoder layer emits hidden vectors $H_m^{(k)}$ and generates a k -th prediction $\hat{W}_m^{(k)}$ as in Eq. (5) as

$$\hat{W}_m^{(k)} = \mathcal{B}(\operatorname{argmax}_{a_t} p(a_t|H_m^{(k)}, \phi)|t = 0, \dots, T), \quad (6)$$

where $k \in \mathcal{K}$. With the unlabeled input and multiple pseudo-labels $\langle O_m, \{\hat{W}_m^{(k)}\}_{k \in \mathcal{K}} \rangle$, the objective function of the online model is defined based on Eq. (3) as

$$\mathcal{L}_{\text{ic}}(\{\hat{W}_m^{(k)}\}_{k \in \mathcal{K}}|O_m, \xi) = \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{ctc}}(\hat{W}_m^{(k)}|H^{(k)}, \xi). \quad (7)$$

This training strategy is compatible with both SC-CTC and HC-CTC-based InterMPL, which we assume particularly effective for HC-CTC with varying output units. HC-CTC trains an ASR model to learn a progressive generation of a target sequence, using the intermediate loss with increasing subword vocabulary size. We expect pseudo-labels generated at different granularities to facilitate semi-supervised learning by providing ancillary training signals.

InterMPL-Last For SC-CTC, InterMPL may not be an optimal choice, as SC-CTC calculates intermediate losses using the same sequence targeted in the last layer. Hence, we design another variant called InterMPL-Last. Different from InterMPL (Fig. 1(b) vs. Fig. 1(c)), InterMPL-Last utilizes only the final hypothesis of the offline model as pseudo-labels for calculating all the losses in the online model. Given the m -th unlabeled sample $O_m \in \mathcal{D}_{\text{unlab}}$ and the last pseudo-labels generated by the offline model $W_m^{(K)}$, the objective function of the online model is defined based on Eq. (3) as

$$\mathcal{L}_{\text{ic}}(\hat{W}_m^{(K)}|O_m, \xi) = \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{ctc}}(\hat{W}_m^{(K)}|H^{(k)}, \xi). \quad (8)$$

InterMPL-Last enables the online model to be trained on the most accurate pseudo-labels predicted by the offline model, which permits more effective use of SC-CTC for semi-supervised training.

4. EXPERIMENTS

4.1. Experimental Setting

We used the ESPnet toolkit [40] for conducting the experiments, and all the codes and recipes are made publicly available (see Sec. 1).

Data We used LibriSpeech (LS) [41] and TED-LIUM3 (TED3) [42]. LS is a corpus of read English speech, containing 960 hours of training data (split into *train-clean-100*, *train-clean-360*, and *train-other-500*). TED3 is a corpus of English Ted Talks consisting of 450 hours of training data (*train-ted3*). We used the standard development and test sets of each dataset for tuning hyper-parameters and evaluating performance, respectively. As input speech features, we extracted 80 mel-scale filterbank coefficients with three-dimensional pitch

Table 1. WERs [%] for models trained on fully labeled data. A*, B*, and C* indicate the oracle results for each semi-supervised setting. The LibriSpeech results are divided into “test-{clean / other}” sets.

Setting	Model	LibriSpeech	TED-LIUM3
		Test WER (\downarrow)	Test WER (\downarrow)
LS-100	S1 CTC	8.4 / 23.1	26.7
	S2 SC-CTC	7.5 / 21.3	24.2
	S3 HC-CTC	7.4 / 20.4	23.8
LS-100 / LS-360	A1 CTC	4.6 / 13.5	–
	A2 SC-CTC	3.9 / 12.0	–
	A3 HC-CTC	4.0 / 11.6	–
LS-100 / LS-860	B1 CTC	3.5 / 8.8	–
	B2 SC-CTC	3.1 / 7.8	–
	B3 HC-CTC	3.2 / 7.7	–
LS-100 / TED3	C1 CTC	–	7.5
	C2 SC-CTC	–	6.8
	C3 HC-CTC	–	7.1

features using Kaldi [43]. We used SentencePiece [44] to construct subword vocabularies from the *train-clean-100* transcriptions.

Semi-supervised settings We regarded *train-clean-100* (LS-100) as the labeled data \mathcal{D}_{lab} . Based on a seed model trained on LS-100, we simulated three semi-supervised settings using different unlabeled data $\mathcal{D}_{\text{unlab}}$: LS-100/LS-360, an in-domain setting using unlabeled *train-clean-360* (LS-360); LS-100/LS-860, an in-domain setting using unlabeled *train-{clean-360, other-500}* (LS-860); and LS-100/TED3, an out-of-domain setting using unlabeled *train-ted3*.

Model architecture We used the Conformer architecture [10, 45] consisting of two convolutional neural network layers followed by a stack of 18 encoder blocks (i.e., $K = 18$). The number of heads, dimension of a self-attention layer, dimension of a feed-forward network, and kernel size were set to 4, 256, 1024, and 7, respectively. Following [28], we replaced batch normalization in the convolution module with group normalization with the group size of 4.

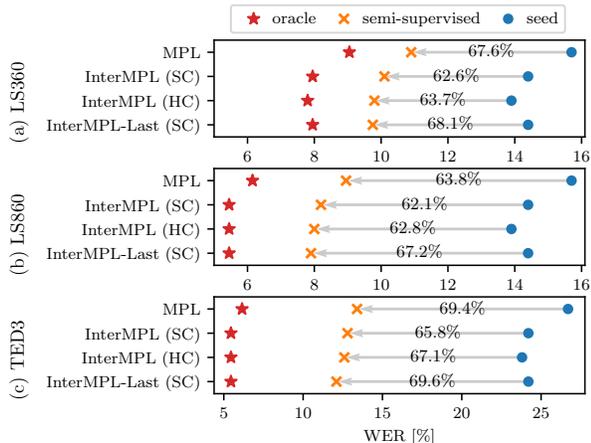
Training and decoding configurations We trained the seed model for 150 epochs using the Adam optimizer [46] with $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-9}$, and Noam learning rate scheduling [47]. We used 25k warmup steps and a learning rate factor of 5.0. The MPL training was iterated up to 200 epochs, and the online model was trained using the Adam optimizer with an initial learning rate of 10^{-3} , $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$. The momentum coefficient α , introduced in Sec. 2.2, was decided following [26]. The subword vocabulary size of CTC was set to 1024. SC-CTC and HC-CTC applied the intermediate loss to the 6th and 12th encoder layers (i.e., $\mathcal{K} = \{6, 12, 18\}$ in Eq. (3)). The output vocabulary size for each loss was set to (1024, 1024, 1024) for SC-CTC and (256, 1024, 4096) for HC-CTC. A final model was obtained for evaluation by averaging model parameters over ten checkpoints that gave the best validation performance. For the MPL-based methods, we followed [26] and used the online model for evaluation. During decoding, we carried out the best path decoding of CTC [4].

4.2. Supervised Baseline and Oracle Results

Table 1 shows the word error rate (WER) of seed models trained on LS-100 (S*) and oracle models trained on fully labeled data in each semi-supervised setting (A*, B*, and C*). Overall, SC-CTC and HC-CTC outperformed CTC by relaxing the conditional independence assumption [35, 36]. The quality of pseudo-labels is crucial for effective semi-supervised training, and we can expect MPL to benefit from the seed models trained with the intermediate loss.

Table 2. WER [%] on in-domain LibriSpeech (LS) settings.

Setting	Method	Init.	Test WER (\downarrow)	
			clean	other
LS-100 /LS-360	X1 MPL	S1 (CTC)	6.3	15.4
	X2 InterMPL	S2 (SC-CTC)	5.7	14.5
	X3 InterMPL	S3 (HC-CTC)	5.5	14.1
	X4 InterMPL-Last	S2 (SC-CTC)	5.4	14.1
LS-100 /LS-860	Y1 MPL	S1 (CTC)	6.0	11.9
	Y2 InterMPL	S2 (SC-CTC)	5.4	11.0
	Y3 InterMPL	S3 (HC-CTC)	5.3	10.7
	Y4 InterMPL-Last	S2 (SC-CTC)	5.1	10.7

**Fig. 2.** Visualization of WRR [%] in each semi-supervised setting.

4.3. Main Results

In-domain setting Table 2 shows WER on the LS settings, comparing the conventional MPL [28] against the proposed InterMPL and InterMPL-Last. In Fig. 2, we also compare the performance of each semi-supervised training method in the WER recovery rate (WRR) [48], which shows how much performance gain is obtained relative to the improvement from the seed to oracle WERs. WRRs for LS are averaged on the clean and other sets. Note that the seed models (S^*) from Table 1 were used for the initialization in each method. Looking at the results on the LS-360 setting (X^*) in Table 2, both InterMPL and InterMPL-Last led to distinct improvements over MPL ($X1$ vs. $X2, X3, X4$), indicating the effectiveness of using the well-trained seed models and applying intermediate CTC loss during semi-supervised training. Comparing SC-CTC and HC-CTC-based InterMPL, HC-CTC resulted in better performance by benefiting from using the pseudo-labels at different granularity ($X2$ vs. $X3$). InterMPL-Last was better suited for SC-CTC-based training than InterMPL ($X2$ vs. $X4$), as it was hypothesized that higher-quality labels are more appropriate for intermediate supervision. Overall, HC-CTC-based InterMPL and InterMPL-Last similarly achieved the best performance, while InterMPL-Last gave higher WRRs in Fig. 2(a).

In the LS-860 setting with more unlabeled data (Y^*) in Table 2, the general trend was consistent with what was observed in the LS-360 setting. In terms of WRR in Fig. 2, InterMPL-Last had the most significant gain, which was even higher than those of MPL. Both InterMPL and InterMPL-Last were scalable to larger amounts of unlabeled data.

Out-of-domain setting Table 3 lists results on the out-of-domain

Table 3. WER [%] on out-of-domain LS-100/TED3.

Method	Init.	Test WER (\downarrow)
Z1 MPL	S1 (CTC)	13.4
Z2 InterMPL	S2 (SC-CTC)	12.8
Z3 InterMPL	S3 (HC-CTC)	12.6
Z4 InterMPL-Last	S2 (SC-CTC)	12.1

Table 4. Ablation study on LS-100/LS-360.

Method	Test WER (\downarrow)		Test WRR (\uparrow)	
	clean	other	clean	other
InterMPL ($X2$)	5.7	14.5	51.9	73.3
w/o inter. loss	6.4	15.5	30.9	62.6
InterMPL ($X3$)	5.5	14.1	55.3	72.0
w/o inter. loss	5.8	14.8	45.8	63.8
InterMPL-Last ($X4$)	5.4	14.1	59.2	77.0
w/ init. from S1	5.9	14.4	46.4	74.1

TED3 setting. Both InterMPL and InterMPL-Last outperformed MPL ($Z1$ vs. $Z2, Z3, Z4$), demonstrating stable training on unlabeled data under the domain-mismatched condition. In contrast to the in-domain results, SC-CTC and HC-CTC-based InterMPL resulted in a similar performance ($Z2$ vs. $Z3$), and InterMPL-Last achieved lower WERs than InterMPL ($Z4$ vs. $Z2, Z3$). HC-CTC was less significant in the out-of-domain semi-supervised scenario, including the oracle results in Table 1, which we attribute to inferior generalization capability. Subword vocabularies are constructed from the small LS-100 text set, and the large vocabulary size used in HC-CTC (i.e., 4096) was not generalized well to the TED3 domain.

4.4. Ablation Study on Intermediate Loss

Table 4 shows an ablation study validating the effectiveness of InterMPL. We initialized a model using parameters of SC-CTC ($S2$) or HC-CTC ($S3$) and performed standard MPL without intermediate loss. Compared to the InterMPL results ($X2, X3$), we observed that removing intermediate loss led to worsen WERs with degraded WRRs. Interestingly, MPL based on SC-CTC initialization resulted in a similar performance as that of vanilla MPL ($X1$). This indicates the importance of applying intermediate loss during semi-supervised training. We also performed InterMPL-Last initialized from CTC ($S1$), which gave better results than those of MPL ($X1$). The results suggest the importance of applying intermediate loss to both the seed model and semi-supervised training.

5. CONCLUSION

We proposed InterMPL, a semi-supervised ASR method enhancing MPL using intermediate CTC loss. We adopted SC-CTC or HC-CTC for training a seed model and explored how pseudo-labels can be generated and used during semi-supervised training. The experimental results and analysis revealed that InterMPL substantially outperforms MPL by fully using the intermediate loss mechanism. Future work should explore using an external language model (LM) in InterMPL, such as combining with LM-based PL [25, 28] and applying shallow fusion to intermediate predictions [49].

Acknowledgement This work was supported in part by JST ACT-X (JPMJAX210J) and JSPS KAKENHI (JP21J23495). This work was also based on results obtained from a project, JPNP20006, commissioned by the New Energy and Industrial Technology Development Organization (NEDO).

6. REFERENCES

- [1] A. Graves and N. Jaitly, "Towards end-to-end speech recognition with recurrent neural networks," in *Proc. ICML*, 2014, pp. 1764–1772.
- [2] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho *et al.*, "Attention-based models for speech recognition," in *Proc. NeurIPS*, 2015, pp. 577–585.
- [3] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *Proc. ICASSP*, 2016, pp. 4960–4964.
- [4] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006, pp. 369–376.
- [5] A. Graves, "Sequence transduction with recurrent neural networks," *arXiv preprint arXiv:1211.3711*, 2012.
- [6] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proc. NeurIPS*, 2014, pp. 3104–3112.
- [7] D. Bahdanau *et al.*, "Neural machine translation by jointly learning to align and translate," in *Proc. ICLR*, 2014.
- [8] L. Dong, S. Xu, and B. Xu, "Speech-Transformer: A no-recurrence sequence-to-sequence model for speech recognition," in *Proc. ICASSP*, 2018, pp. 5884–5888.
- [9] S. Kriman, S. Beliaev, B. Ginsburg, J. Huang *et al.*, "QuartzNet: Deep automatic speech recognition with 1D time-channel separable convolutions," in *Proc. ICASSP*, 2020, pp. 6124–6128.
- [10] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar *et al.*, "Conformer: Convolution-augmented Transformer for speech recognition," in *Proc. Interspeech*, 2020, pp. 5036–5040.
- [11] C.-C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar *et al.*, "State-of-the-art speech recognition with sequence-to-sequence models," in *Proc. ICASSP*, 2018, pp. 4774–4778.
- [12] C. Lüscher, E. Beck, K. Irie, M. Kitza *et al.*, "RWTH ASR systems for LibriSpeech: Hybrid vs attention," in *Proc. Interspeech*, 2019, pp. 231–235.
- [13] S. Karita, N. Chen, T. Hayashi, T. Hori *et al.*, "A comparative study on Transformer vs RNN in speech applications," in *Proc. ASRU*, 2019, pp. 449–456.
- [14] J. Kahn, M. Rivière, W. Zheng, E. Kharitonov *et al.*, "Libri-Light: A benchmark for ASR with limited or no supervision," in *Proc. ICASSP*, 2020, pp. 7669–7673.
- [15] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Proc. ICML*, 2013.
- [16] B. Li, T. N. Sainath, R. Pang, and Z. Wu, "Semi-supervised training for end-to-end models via weak distillation," in *Proc. ICASSP*, 2019, pp. 2837–2841.
- [17] R. Masumura, M. Ithori, A. Takashima, T. Moriya *et al.*, "Sequence-level consistency training for semi-supervised end-to-end automatic speech recognition," in *Proc. ICASSP*, 2020, pp. 7054–7058.
- [18] F. Weninger, F. Mana, R. Gemello, J. Andrés-Ferrer *et al.*, "Semi-supervised learning with data augmentation for end-to-end ASR," in *Proc. Interspeech*, 2020.
- [19] W.-N. Hsu, A. Lee, G. Synnaeve, and A. Hannun, "Semi-supervised speech recognition via local prior matching," *arXiv preprint arXiv:2002.10336*, 2020.
- [20] Q. Xu, T. Likhomanenko, J. Kahn, A. Hannun *et al.*, "Iterative pseudo-labeling for speech recognition," in *Proc. Interspeech*, 2020.
- [21] Y. Chen, W. Wang, and C. Wang, "Semi-supervised ASR by end-to-end self-training," in *Proc. Interspeech*, 2020, pp. 2787–2791.
- [22] D. S. Park, Y. Zhang, Y. Jia, W. Han *et al.*, "Improved noisy student training for automatic speech recognition," in *Proc. Interspeech*, 2020, pp. 2817–2821.
- [23] T. Likhomanenko, Q. Xu, J. Kahn, G. Synnaeve *et al.*, "slim-IPL: Language-model-free iterative pseudo-labeling," in *Proc. Interspeech*, 2021, pp. 741–745.
- [24] N. Moritz, T. Hori, and J. Le Roux, "Semi-supervised speech recognition via graph-based temporal classification," in *Proc. ICASSP*, 2021, pp. 6548–6552.
- [25] J. Kahn, A. Lee, and A. Hannun, "Self-training for end-to-end speech recognition," in *Proc. ICASSP*, 2020.
- [26] Y. Higuchi, N. Moritz, J. Le Roux, and T. Hori, "Momentum pseudo-labeling for semi-supervised speech recognition," in *Proc. Interspeech*, 2021, pp. 726–730.
- [27] V. Manohar, T. Likhomanenko, Q. Xu, W.-N. Hsu *et al.*, "Kaizen: Continuously improving teacher using exponential moving average for semi-supervised speech recognition," in *Proc. ASRU*, 2021, pp. 518–525.
- [28] Y. Higuchi, N. Moritz, J. Le Roux, and T. Hori, "Advancing momentum pseudo-labeling with conformer and initialization strategy," in *Proc. ICASSP*, 2022, pp. 7672–7676.
- [29] —, "Momentum pseudo-labeling: Semi-supervised ASR with continuously improving pseudo-labels," *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1424–1438, 2022.
- [30] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Proc. NeurIPS*, 2017, pp. 1195–1204.
- [31] J. Chorowski and N. Jaitly, "Towards better decoding and language model integration in sequence to sequence models," *arXiv preprint arXiv:1612.02695*, 2016.
- [32] A. Tjandra, C. Liu, F. Zhang, X. Zhang *et al.*, "Deja-vu: Double feature presentation and iterated loss in deep Transformer networks," in *Proc. ICASSP*, 2020, pp. 6899–6903.
- [33] J. Lee and S. Watanabe, "Intermediate loss regularization for CTC-based speech recognition," in *Proc. ICASSP*, 2021, pp. 6224–6228.
- [34] Y. Higuchi, N. Chen, Y. Fujita, H. Inaguma *et al.*, "A comparative study on non-autoregressive modelings for speech-to-text generation," in *Proc. ASRU*, 2021, pp. 47–54.
- [35] J. Nozaki and T. Komatsu, "Relaxing the conditional independence assumption of CTC-based ASR by conditioning on intermediate predictions," in *Proc. Interspeech*, 2021, pp. 3735–3739.
- [36] Y. Higuchi, K. Karube, T. Ogawa, and T. Kobayashi, "Hierarchical conditional end-to-end ASR with CTC and multi-granular subword units," in *Proc. ICASSP*, 2022, pp. 7797–7801.
- [37] R. Sanabria and F. Metzger, "Hierarchical multitask learning with CTC," in *Proc. SLT*, 2018, pp. 485–490.
- [38] K. Krishna *et al.*, "Krishna, kalpesh and toshniwal, shubham and livescu, karen," *arXiv preprint arXiv:1807.06234*, 2018.
- [39] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu *et al.*, "SpecAugment: A simple data augmentation method for automatic speech recognition," in *Proc. Interspeech*, 2019, pp. 2613–2617.
- [40] S. Watanabe, T. Hori, S. Karita, T. Hayashi *et al.*, "ESPnet: End-to-end speech processing toolkit," in *Proc. Interspeech*, 2018.
- [41] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP*, 2015.
- [42] F. Hernandez, V. Nguyen, S. Ghannay, N. Tomashenko *et al.*, "TED-LIUM 3: Twice as much data and corpus repartition for experiments on speaker adaptation," in *Proc. SPECOM*, 2018.
- [43] D. Povey, A. Ghoshal, G. Boulianne, L. Burget *et al.*, "The Kaldi speech recognition toolkit," in *Proc. ASRU*, 2011.
- [44] T. Kudo, "Subword regularization: Improving neural network translation models with multiple subword candidates," in *Proc. ACL*, 2018.
- [45] P. Guo, F. Boyer, X. Chang, T. Hayashi *et al.*, "Recent developments on ESPnet toolkit boosted by Conformer," in *Proc. ICASSP*, 2021, pp. 5874–5878.
- [46] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [47] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit *et al.*, "Attention is all you need," in *Proc. NeurIPS*, 2017.
- [48] J. Ma and R. Schwartz, "Unsupervised versus supervised training of acoustic models," in *Proc. Interspeech*, 2008.
- [49] T. Komatsu, Y. Fujita, J. Lee, L. Lee *et al.*, "Better intermediates improve CTC inference," in *Proc. Interspeech*, 2022, pp. 4965–4969.