

ADAMER-CTC: CONNECTIONIST TEMPORAL CLASSIFICATION WITH ADAPTIVE MAXIMUM ENTROPY REGULARIZATION FOR AUTOMATIC SPEECH RECOGNITION

SooHwan Eom¹, Eunseop Yoon¹, Hee Suk Yoon¹,
Chanwoo Kim², Mark Hasegawa-Johnson³, Chang D. Yoo^{1†}

¹ Korea Advanced Institute of Science and Technology, Daejeon, Republic of Korea

² Korea University, Seoul, Republic of Korea

³ University of Illinois Urbana-Champaign

ABSTRACT

In Automatic Speech Recognition (ASR) systems, a recurring obstacle is the generation of narrowly focused output distributions. This phenomenon emerges as a side effect of Connectionist Temporal Classification (CTC), a robust sequence learning tool that utilizes dynamic programming for sequence mapping. While earlier efforts have tried to combine the CTC loss with an entropy maximization regularization term to mitigate this issue, they employed a constant weighting term on the regularization during the training, which we find may not be optimal. In this work, we introduce Adaptive Maximum Entropy Regularization (AdaMER), a technique that can modulate the impact of entropy regularization throughout the training process. This approach not only refines ASR model training but ensures that as training proceeds, predictions display the desired model confidence.

Index Terms— Automatic Speech Recognition, Connectionist Temporal Classification, Entropy Maximization

1. INTRODUCTION

Deep learning-based Automatic Speech Recognition (ASR) has significantly advanced due to the use of the Connectionist Temporal Classification (CTC) loss [1]. CTC loss allows training on large datasets without requiring explicit alignment between speech-transcript pairs. Central to this alignment flexibility is the addition of the *blank symbol* to the existing output label set, providing a mechanism to handle varying sequence lengths and consecutive identical characters.

However, training with CTC loss often results in a peaky distribution as shown in Fig. 1. This peaky behavior is a manifestation of the blank symbol design and the entropy minimization property inherent to CTC loss [2, 3, 4]. Randomly initialized ASR model exhibits a tendency to predict only

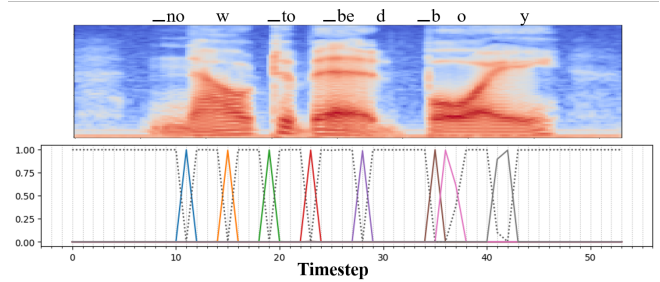


Fig. 1: Illustration of CTC peaky distribution output. The dotted line denotes the output probability of the blank symbol.

blank symbols, as it searches for a single point to “jump” to non-blank symbols in the target sequence, resulting in sharp spikes in its prediction. The entropy-minimizing nature of this training then leads to the predictions during training being stuck or converging into such peaky sub-optimal solutions, hindering the training process.

Prior works [5, 6] have attempted to address this issue by combining the CTC loss with constant entropy maximization as an auxiliary regularization loss. However, our research suggests that such constant weighting on uncertainty might be counterproductive in the latter stages of training, as confidence in predictions becomes more crucial. In this paper, we propose Adaptive Maximum Entropy Regularization (AdaMER), which dynamically adjusts the effect of the entropy regularization throughout training. Through extensive experiments using the LibriSpeech [7] corpus, we show that AdaMER can effectively improve the CTC training of ASR models.

2. RELATED WORK

2.1. Overconfident Problem in CTC

The Connectionist Temporal Classification (CTC) framework [1] has significantly impacted numerous end-to-end sequence learning tasks, revolutionizing areas such as speech recognition [8, 9], text recognition [10], and video segmentation

[†]Corresponding author

This work was partly supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) [NO.2022-0-00184, Development and Study of AI Technologies to Inexpensively Conform to Evolving Policy on Ethics] and SAMSUNG Research, Samsung Electronics Co., Ltd.

[11], due to its ability to model the sequential output without the need for explicit alignment information between input and output sequences. However, despite the substantial breakthroughs made possible by the use of CTC, there exist a few persisting challenges.

One such challenge is the tendency of the CTC-trained model to generate highly peaky distributions [1, 2, 3], a characteristic often interpreted as a sign of overconfidence. This ‘overconfidence’ refers to scenarios in which the model predicts certain outputs with high certainty, which, while positive if the prediction is correct, can lead to pronounced errors if the prediction is off the mark. The ‘peaky behavior’ of CTC is highly related to the ‘blank’ symbol, which is an additional label for handling consecutive identical labels. According to the analysis from [2, 4], while blank symbols are crucial for CTC loss to ensure convergence, it is also the cause of the localized peaky predictions. They specifically point out that models that are initialized uniformly and then trained using gradient descent methods—a common practice in machine learning—are susceptible to converging to suboptimal local optima. These suboptimal solutions are often marked by overconfidence or peaky behavior, thereby suggesting that the issue is intrinsically linked to the model’s training dynamics rather than the model architecture itself.

2.2. Regularization in CTC

Model overconfidence is a common challenge across the machine learning landscape, prompting the creation of a multitude of regularization techniques to handle it. One solution to the overconfidence problem is to directly penalize the model confidence, which is often expressed as the entropy of the model prediction. Label smoothing regularization [12], which uses soft targets instead of hard one-hot targets for cross-entropy labels, can be seen as maximizing relative entropy between the model prediction and uniform distribution [13, 14]. In ASR, prior works including [5] have applied label smoothing on the CTC criterion.

Confidence regularization is also connected to the maximum entropy principle [15], which states that the distribution that best describes the current state is the one that leaves the largest amount of uncertainty (or entropy) consistent with the constraint. Strategies such as maximum entropy regularization are widely used within the reinforcement learning domain [16, 17, 18], which have been deployed to foster exploration and deter early convergence. For CTC, EnCTC [6] have proposed a maximum conditional entropy-based regularization for CTC, applied to the optical character recognition task. This method considers the conditional distribution of potential paths given the input sequence and label sequence. Based on the fact that the error signal of CTC is proportional to the likelihood of the path, EnCTC is designed to prevent the entropy of feasible paths from declining too rapidly, thus mitigating the impact of CTC’s convergence to a single path

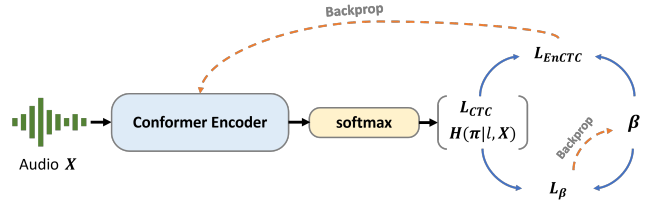


Fig. 2: CTC with Maximum Entropy Regularization. An entropy-based regularization term is added to the original CTC objective to increase the entropy while training. This boosts the probability of the paths nearby, enhancing exploration throughout the training process.

and encouraging exploration during training.

3. CTC WITH MAXIMUM ENTROPY REGULARIZATION

In this section we explain the standard CTC method, the maximum entropy regularization, and our AdaMER method.

3.1. Problem Formulation

Consider a speech dataset \mathcal{D} . This dataset consists of pairs of input sequences and their corresponding target label sequences. A pair is denoted as (X, l) , where $X = [X_1, \dots, X_T]$ is an input sequence of length T and $l = [l_1, \dots, l_L]$ is the corresponding target label sequence of length L , such that $L \leq T$. All possible elements of l is inside the fixed-sized alphabet set \mathcal{S} . In this study, X represents a speech sample, and l is the corresponding transcription. The alignment between X and l is not known but in a sequential manner, which is a common scenario in many speech recognition problems.

3.2. Connectionist Temporal Classification

Connectionist Temporal Classification (CTC) provides a way to handle the alignment problem in sequence prediction tasks [1]. It introduces a new objective function, referred to as the CTC loss. The CTC loss is essentially the negative log-likelihood of the correct label sequence given the input sequence. This is mathematically represented as:

$$\mathcal{L}_{\text{CTC}} = -\log p(l|X_{1:T}) \quad (1)$$

CTC defines the conditional likelihood $p(l|X_{1:T})$ as the sum of probabilities of all possible valid alignments between the input and label sequence. Such alignment is called the path π , and the many-to-one mapping from the feasible paths to the label sequence l is denoted as $\mathcal{B}(\pi)$. The conditional likelihood then can be written as:

$$p(l|X_{1:T}) = \sum_{\pi \in \mathcal{B}^{-1}(l)} p(\pi|X_{1:T}) \quad (2)$$

and

$$p(\pi|X_{1:T}) = \prod_{t=1}^T y_{\pi_t}^t \quad (3)$$

where y_k^t is the probability for symbol k on timestamp t . This probability is usually calculated by the softmax function:

$$y_k^t = \frac{e^{a_k^t}}{\sum_{k'} e^{a_{k'}^t}} \quad (4)$$

where $\{a_k^t | t \in [1, T], k \in \mathcal{S}'\}$ is the output right before the softmax activation layer and $\mathcal{S}' = \mathcal{S} \cup \emptyset$ is the alphabet set including the additional ‘blank’ symbol denoted as \emptyset . This additional ‘blank’ label is introduced to handle varying lengths between the input sequence and the output sequence and also deals with the unknown alignment problem. The total number of output labels for CTC becomes $|\mathcal{S}| + 1$.

In practical implementations, CTC uses a dynamic programming algorithm, akin to the forward-backward algorithm in Hidden Markov Models (HMM), to efficiently compute these probabilities. The final objective is to maximize the total probability of the correct label sequence over all possible alignments.

3.3. Maximum Entropy Regularization

To resolve the overconfident peaky prediction, an entropy-based regularization method, named Entropy-Regularized Connectionist Temporal Classification (EnCTC), has previously been used for the optical character recognition task [6]. Maximum Entropy Regularization is a widely used method in reinforcement learning, such as in Soft Actor-Critic [17, 18], which encourages the agent to explore the broader policy space and prevents the agent from converging to sub-optimal policy. Inspired by this, EnCTC introduces a maximum conditional entropy-based regularization term which can be represented as follows:

$$L_{\text{EnCTC}} = L_{\text{CTC}} - \beta H(\pi|l, X) \quad (5)$$

where β is a factor that regulates the intensity of the maximum conditional entropy regularization.

The entropy of the feasible paths, given the input sequence X and the target sequence l , can be denoted as:

$$\begin{aligned} H(\pi|l, X) &= \sum_{\pi \in \mathcal{B}^{-1}(l)} -p(\pi|l, X) \log p(\pi|l, X) \\ &= -\frac{1}{p(l|X)} \sum_{\pi \in \mathcal{B}^{-1}(l)} p(\pi|X) \log p(\pi|X) \quad (6) \\ &\quad + \log p(l|X) \end{aligned}$$

A rapid convergence towards a single feasible path implies a swift reduction in the conditional entropy. By including the term in the objective function that fosters the maximization of the entropy, it subsequently boosts the probability of the paths nearby, enhancing exploration throughout the

training process. In this paper, we show that the Maximum Entropy Regularization can benefit the training process for speech recognition.

3.4. Adaptive Maximum Entropy Regularization

During the initial phase of CTC training, the randomly initialized network tends to output only blank symbols[3, 4]. The maximum entropy regularization from the previous section can encourage the network to localize non-blank symbols by penalizing low-entropy blank-only predictions. On the other hand, in the later stage of training, the good and bad paths become distinctive to the model. In this case, the model’s prediction should be more deterministic, and thus, fixed entropy regularization can actually act as a hindrance. Therefore, we show the alternative method of entropy regularization, which allows the model to adjust its ambiguity in its prediction.

To begin with, inspired by the improved version of Soft Actor-Critic [18], we can formulate the following constrained optimization problem:

$$\begin{aligned} \max_{\theta} \log p_{\theta}(l|X) &= \max_{\theta} \log \sum_{\pi \in \mathcal{B}^{-1}(l)} p_{\theta}(\pi|X), \quad (7) \\ \text{s.t. } \mathbb{E}_{p_{\theta}(\pi|X, l)}[-\log p_{\theta}(\pi|X, l)] &\geq \mathcal{H}, \forall X \end{aligned}$$

where \mathcal{H} is the desired minimum entropy of the path π .

Although our primal CTC objective is not strictly convex, we can approximate the optimization process using a similar Lagrangian dual technique as in [18]. By introducing the Lagrangian dual variable β , we can formulate the following dual function:

$$\min_{\beta \geq 0} \max_{\theta} (\log \sum_{\pi \in \mathcal{B}^{-1}(l)} p_{\theta}(\pi|X) + \beta[H(\pi|l, X) - \mathcal{H}]) \quad (8)$$

The above dual problem can be solved via gradient descent on the learnable parameters θ and β . Therefore, we will incorporate this new β loss function in addition to the loss function in Eq. 5:

$$L_{\beta} = \beta[H(\pi|l, X) - \mathcal{H}] \quad (9)$$

The above loss is used to update β through the training process. In practical implementation, we use the stop gradient function for both L_{EnCTC} and L_{β} since we do not want the gradient of L_{EnCTC} to flow through β or vice versa. While [18] used dual gradient descent, we jointly trained θ and β for efficient end-to-end training.

4. EXPERIMENTAL SETTINGS

In this section, we present the benchmarks employed for evaluating our approach, such as the datasets, model, and metrics.

| Loss | test-clean | test-other | test-all |
|--------------|-------------|--------------|-------------|
| L_{CTC} | 5.86 | 14.18 | 10.02 |
| L_{EnCTC} | 5.28 | 12.74 | 9.01 |
| L_{AdaMER} | 5.13 | 12.62 | 8.88 |

Table 1: WER (%) (lower is better) on LibriSpeech. test-all denotes the aggregation of test-clean and test-other. The values reported are the average of three runs with different random seeds.

4.1. Datasets

LibriSpeech corpus [7] was used for both training and evaluation. LibriSpeech corpus is the widely-used English speech dataset for ASR. For training, we use train-960 subset, which contains 460 hours of clean speech and 500 hours of noisy speech data. For per-epoch evaluation, we use dev-clean set, and for evaluation, we use test-clean and test-other each to compare the performance boost with and without the noise.

4.2. Model and Metric

In this paper we used Conformer for our baseline model [19]. Conformer is the ASR model which combines a convolutional neural network and transformer to capture both local and global dependencies. We implement conformer-small in NVIDIA NeMo toolkit¹, whose encoder comprises 16 conformer layers. The output of the encoder is trained using CTC objective and decoded by greedy searching. For L_{AdaMER} and L_{EnCTC} , we chose initial β as 0.2 and 1.0 respectively. Target entropy \mathcal{H} was set to $1.1U$ where U is the target length. We used AdamW optimizer and Noam learning rate scheduler with 10,000 warmup steps. The models were trained 100 epochs each. Word error rate (WER) is used as the evaluation metric for all experiments. Experiment results are recorded and rendered from WandB dashboard.²

5. RESULTS

5.1. Librispeech Results

Table 1 demonstrates the Word Error Rate (WER) across different training objectives (i.e., L_{AdaMER} , L_{EnCTC} , or L_{CTC}). We find that using the EnCTC and AdaMER performance improvement is more pronounced compared to standard CTC training for both test-clean and test-other subsets. Specifically, when comparing the WER of CTC with EnCTC and AdaMER, the improvement was 1.01% and 1.14%, respectively. This clearly suggests while entropy regularization can help automatic speech recognition tasks, our adaptive entropy regularization can handle it more effectively.

¹<https://developer.nvidia.com/nemo>

²<https://wandb.ai/site>

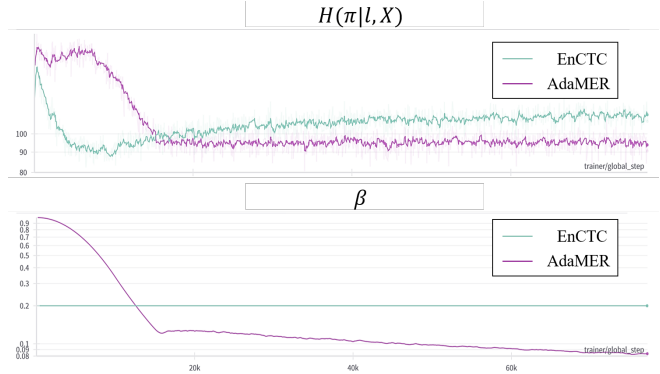


Fig. 3: Change of $H(\pi|l, X)$ (top) and β (down) during training the model with L_{AdaMER} (purple) and L_{EnCTC} (green). The horizontal axis is the training step. The results are recorded and rendered from WandB dashboard.

Fig. 3 shows the change of the value of entropy $H(\pi|l, X)$ and the weight β throughout training. Notably, the entropy associated with the EnCTC model demonstrates a continual increase over the course of training. In contrast, the AdaMER method consistently maintains the entropy at a relatively stable level. Furthermore, the trajectory of β in AdaMER is characterized by a decreasing trend, which persists until the entropy reaches a predefined target value. Beyond this point, the reduction in β becomes more gradual. This behavior evidences the adaptive capability of the AdaMER approach in modulating the intensity of entropy regularization.

6. CONCLUSIONS

In conclusion, our study addresses a pivotal challenge observed in ASR systems trained via Connectionist Temporal Classification (CTC): the issue of overly confident peaky distribution, especially attributed to the blank symbol. Delving into the entropy minimization property inherent to CTC, we recognized the tendency of the model to predominantly predict blank symbols, often leading to sub-optimal outcomes. We introduced the Adaptive Maximum Entropy Regularization (AdaMER) technique, marking a significant departure from traditional methods that attempted to combine CTC loss with a constant entropy maximization term. AdaMER’s uniqueness lies in its dynamic entropy-based scheduler, which modulates the effect of regularization as training progresses. This ensures that while the system is explorative initially, it converges with confident outputs towards the latter stages of training—an essential balance for optimized performance. Our comprehensive experiments demonstrated potential in enhancing the performance of ASR models trained with CTC, especially in tackling the issue of peaky distributions. Our work underscores the importance of adaptive regularization on model confidence in ASR training.

7. REFERENCES

- [1] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd International Conference on Machine Learning*, New York, NY, USA, 2006, ICML ’06, p. 369–376, Association for Computing Machinery.
- [2] Albert Zeyer, Ralf Schlüter, and Hermann Ney, “Why does ctc result in peaky behavior?,” *arXiv preprint arXiv:2105.14849*, 2021.
- [3] Théodore Bluche, *Deep Neural Networks for Large Vocabulary Handwritten Text Recognition. (Réseaux de Neurones Profonds pour la Reconnaissance de Texte Manuscrit à Large Vocabulaire)*, Ph.D. thesis, University of Paris-Sud, Orsay, France, 2015.
- [4] Théodore Bluche, Hermann Ney, Jérôme Louradour, and Christopher Kermorvant, “Framewise and ctc training of neural networks for handwriting recognition,” in *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, 2015, pp. 81–85.
- [5] Suyoun Kim, Michael L. Seltzer, Jinyu Li, and Rui Zhao, “Improved training for online end-to-end speech recognition systems,” 2018.
- [6] Hu Liu, Sheng Jin, and Changshui Zhang, “Connectionist temporal classification with maximum entropy regularization,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. 2018, vol. 31, Curran Associates, Inc.
- [7] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.
- [8] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al., “Deep speech 2: End-to-end speech recognition in english and mandarin,” in *International conference on machine learning*. PMLR, 2016, pp. 173–182.
- [9] Ying Zhang, Mohammad Pezeshki, Philémon Brakel, Saizheng Zhang, Cesar Laurent Yoshua Bengio, and Aaron Courville, “Towards end-to-end speech recognition with deep convolutional neural networks,” *arXiv preprint arXiv:1701.02720*, 2017.
- [10] Alex Graves, Marcus Liwicki, Santiago Fernández, Roman Bertolami, Horst Bunke, and Jürgen Schmidhuber, “A novel connectionist system for unconstrained handwriting recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 5, pp. 855–868, 2009.
- [11] Mengxi Lin, Nakamasa Inoue, and Koichi Shinoda, “Ctc network with statistical language modeling for action sequence recognition in videos,” in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, New York, NY, USA, 2017, Thematic Workshops ’17, p. 393–401, Association for Computing Machinery.
- [12] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [13] Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey Hinton, “Regularizing neural networks by penalizing confident output distributions,” 2017.
- [14] Clara Meister, Elizabeth Salesky, and Ryan Cotterell, “Generalized entropy regularization or: There’s nothing special about label smoothing,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Online, July 2020, pp. 6870–6886, Association for Computational Linguistics.
- [15] E. T. Jaynes, “Information theory and statistical mechanics,” *Phys. Rev.*, vol. 106, pp. 620–630, May 1957.
- [16] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in *International conference on machine learning*. PMLR, 2016, pp. 1928–1937.
- [17] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” *CoRR*, vol. abs/1801.01290, 2018.
- [18] Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine, “Soft actor-critic algorithms and applications,” 2019.
- [19] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.