# CR-CTC: CONSISTENCY REGULARIZATION ON CTC FOR IMPROVED SPEECH RECOGNITION

Anonymous authors

Paper under double-blind review

### ABSTRACT

Connectionist Temporal Classification (CTC) is a widely used method for automatic speech recognition (ASR), renowned for its simplicity and computational efficiency. However, it often falls short in recognition performance compared to transducer or systems combining CTC and attention-based encoder-decoder (CTC/AED). In this work, we propose the Consistency-Regularized CTC (CR-CTC), which enforces consistency between two CTC distributions obtained from different augmented views of the input speech mel-spectrogram. We provide indepth insights into its essential behaviors from three perspectives: 1) it conducts self-distillation between random pairs of sub-models that process different augmented views; 2) it learns contextual representation through masked prediction for positions within time-masked regions, especially when we increase the amount of time masking; 3) it suppresses the extremely peaky CTC distributions, thereby reducing overfitting and improving the generalization ability. Extensive experiments on LibriSpeech, Aishell-1, and GigaSpeech datasets demonstrate the effectiveness of our CR-CTC, which achieves performance comparable to, or even slightly better than, that of transducer and CTC/AED.

029

004

010 011

012

013

014

015

016

017

018

019

021

## 1 INTRODUCTION

End-to-end approaches (Graves et al., 2006; Graves, 2012; Chan et al., 2015), which eliminate the 031 need of pre-aligned speech-text data, have replaced traditional hybrid systems (Bourlard & Morgan, 2012; Hinton et al., 2012) and become dominant methods in automatic speech recognition 033 (ASR). Prominent examples include Connectionist Temporal Classification (CTC) (Graves et al., 034 2006), Transducer (Graves, 2012) (also known as RNN-T), and the method that combines CTC and attention-based encoder-decoder (AED) (Chan et al., 2015), referred to as CTC/AED (Watanabe et al., 2017). To handle the alignment between speech and token sequences, CTC (Graves et al., 2006) introduces a blank token and makes independent predictions at each frame, training the model 037 to maximize the total probability over all valid alignments. Transducer (Graves, 2012) extends CTC by introducing a prediction network and a joint network, explicitly modeling the interdependencies on output labels. CTC/AED (Watanabe et al., 2017) integrates CTC into AED (Chan et al., 2015) for 040 jointly training, while the CTC and AED scores are fused during the decoding process. Among these 041 three methods, CTC is the simplest and most computationally efficient due to its frame-independent 042 assumption, making it a strong candidate for real-world deployment. However, it significantly lags 043 behind transducer and CTC/AED in terms of recognition performance, which limits its applicability. 044

To improve the CTC performance, in this work we propose the Consistency-Regularized CTC (CR-*CTC*), which takes two different augmented views of the same speech mel-spectrogram as input, 046 and enforces consistency between the resulting CTC distributions. We analyze its internal behaviors 047 from three following perspectives. First, it performs self-distillation between sub-models randomly 048 sampled by drop-based techniques (Srivastava et al., 2014; Huang et al., 2016). Second, for positions within time-masked regions, the model is required to predict the target token distributions, forcing it to learn contextual representation based on unmasked context, akin to self-supervised learning 051 methods (Devlin et al., 2019; Baevski et al., 2020; Hsu et al., 2021). Therefore, we especially increase the amount of time masking in CR-CTC to enhance this masked prediction behavior. Third, 052 the consistency regularization suppresses extremely peaky CTC distributions, which mitigates overfitting and improves the model's generalization ability. Inspired by this, we additionally propose an simple method specifically designed to learn smoother CTC distributions (Appendix Section A.1),
 which is experimentally validated to be effective.

We conduct experiments on LibriSpeech, Aishell-1, and GigaSpeech datasets using Zipformer (Yao et al., 2024) as speech encoder. The results demonstrate the superiority of CR-CTC, which significantly outperforms vanilla CTC and achieves results comparable to, or even slightly better than, those of transducer and CTC/AED. In addition, CR-CTC can further improve the performance of transducer and CTC/AED when employed for jointly training. We perform detailed ablation studies on LibriSpeech dataset to investigate the effect of each functional component in CR-CTC and to validate our explanations.

2 RELATED WORK

064

065

**Self-distillation.** Unlike traditional knowledge distillation (Buciluă et al., 2006; Hinton et al., 2015), 066 which transfers knowledge from a larger and high-capacity teacher model to a smaller student model, 067 self-distillation (Furlanello et al., 2018; Zhu et al., 2018; Mobahi et al., 2020; Allen-Zhu & Li, 068 2020) involves learning from a same-architecture model that processes the same training data. This 069 approach enables the model to extract more refined representations and achieve improved performance. For example, BANs (Furlanello et al., 2018) introduces a re-training procedure in which a 071 newly initialized student model is trained to match a pre-trained teacher model, subsequently serv-072 ing as the teacher in the next iteration. Some works also explore constructing the teacher and student 073 models from a shared network, distilling knowledge from deeper layers to shallower layers (Zhang 074 et al., 2019; Kim et al., 2024), or between pairs of sub-models randomly initialized through drop-075 based techniques (Srivastava et al., 2014; Huang et al., 2016), such as R-Drop (Wu et al., 2021) and cosub (Touvron et al., 2023). Our CR-CTC fundamentally conducts self-distillation between ran-076 dom sub-models, sharing similar idea to R-Drop and cosub, while our approach further use different 077 augmented input views, which enriches the diversity of predictions from these sub-models.

079 Masked prediction. Masked prediction has proven highly effective in self-supervised learning (Devlin et al., 2019; Baevski et al., 2019; Joshi et al., 2020; Baevski et al., 2020; Hsu et al., 2021; He 081 et al., 2022; Baevski et al., 2023). In this approach, the model is tasked with predicting masked positions based on the surrounding unmasked context, which encourages the learning of robust 083 contextual representations. Notable methods for speech representation learning include wav2vec 2.0 (Baevski et al., 2020), HuBERT (Hsu et al., 2021), and data2vec 2.0 (Baevski et al., 2023), which 084 primarily differ in their prediction targets. Specifically, wav2vec 2.0 (Baevski et al., 2020) jointly 085 trains a representation quantizer and learns to distinguish the true quantized target from distractors (Oord et al., 2018). HuBERT generates target labels through offline clustering, while data2vec 087 2.0 uses contextualized representations from a teacher model as its targets. Our *CR-CTC* essentially 088 performs masked prediction for positions within time-masked regions, where the target labels are 089 frame-level token distributions generated based on another augmented view of input. 090

Peaky CTC distributions. CTC models are known for predicting extremely peaky distribu-091 tions (Graves et al., 2006; Sak et al., 2015), which can be harmful in certain scenarios, such as 092 forced alignment (Huang et al., 2024) and knowledge distillation (Ding et al., 2020). These peaky distributions lead to inaccurate alignments as the model assigns excessive blanks to non-silence 094 frames. To address this, label priors are employed to suppress the peaky distributions, thereby im-095 proving the accuracy of forced alignment (Huang et al., 2024). As position mismatches of CTC 096 spikes can hinder knowledge distillation performance, some approaches propose to encourage con-097 sistent alignments between the teacher and student (Ding et al., 2020) or to utilize sequence-level 098 distillation (Takashima et al., 2019). Unlike previous works, we demonstrate that peak suppression in CR-CTC can improve the generalization ability of the CTC models. 099

100 Consistency regularization. The technique of consistency regularization has demonstrated effec-101 tiveness in learning generalizable image representations across various learning paradigms, includ-102 ing self-supervised (Chen et al., 2020; Grill et al., 2020; He et al., 2020; Chen & He, 2021), semi-103 supervised (Sajjadi et al., 2016; Laine & Aila, 2016; Sohn et al., 2020), and supervised (Wu et al., 104 2021; Touvron et al., 2023; Heo et al., 2023) learning tasks. Self-supervised methods, such as Sim-105 CLR (Chen et al., 2020), BYOL (Grill et al., 2020), MoCo (He et al., 2020) and SimSiam (Chen & He, 2021), aim to align hidden representations of unlabeled image data from different model 106 branches or different augmented views. They address the training issue of feature collapsing into a 107 constant vector (Chen & He, 2021) through contrastive learning (Chen et al., 2020; He et al., 2020),

momentum encoder (Grill et al., 2020; He et al., 2020), and stop-gradient operation (Chen & He, 2021). In semi-supervised learning, a prominent example leveraging consistency regularization is
FixMatch (Sohn et al., 2020). It generates pseudo-labels based on high-confidence predictions from weakly augmented images, then trains the model to predict these pseudo-labels using the strongly augmented versions of the same images. Additionally, in supervised learning, methods such as R-Drop (Wu et al., 2021) and cosub (Touvron et al., 2023) encourage consistency between predictions of randomly sampled sub-models using drop-based techniques.

115 When employing consistency regularization as unsupervised objective to train transformer encoders 116 on unlabeled speech data, a new training issue arises in the form of the shortcut learning prob-117 lem (Geirhos et al., 2020), which is tackled using reconstruction loss in Speech SimCLR (Jiang et al., 2020) and temporal augmentation in C-Siam (Khorram et al., 2022). Some studies explore leverag-118 ing consistency regularization to enhance model robustness during predicting the pseudo-labels of 119 untranscribed data, which are generated based on different augmentations (Masumura et al., 2020; 120 Weninger et al., 2020; Chen et al., 2021b; Higuchi et al., 2021; Sapru, 2022) or through speech chain 121 reconstruction (Qi et al., 2022). In contrast to these self/semi-supervised ASR works, our work fo-122 cuses on a fully supervised setting, where we introduce consistency loss as a regularization term 123 to improve performance of CTC model trained on labeled data. As the consistency regularization 124 is enforced on CTC distributions, which are stably supervised by the main CTC loss, it inherently 125 avoids the training issues associated with the unsupervised objectives as observed in Speech Sim-126 CLR (Jiang et al., 2020) and C-Siam (Khorram et al., 2022). 127

The idea of R-Drop (Wu et al., 2021) has also been extended to supervised ASR (Gao et al., 2022; 128 Yoon et al., 2024). For example, to improve the CTC/AED system, (Gao et al., 2022) specially 129 designs the spatial-temporal dropout to construct the sub-models, with consistency regularization 130 enforced exclusively on the CTC spike frames. Cons-KD (Yoon et al., 2024) integrates consistency 131 regularization into a knowledge distillation system, enabling the student model to be more robust 132 to inconsistency induced by dropout. In this work, we focus on improving the performance of pure 133 CTC systems and are the first to enable CTC models to match the performance of transducer and 134 CTC/AED systems by a simple yet effective approach. Moreover, we introduce peak suppression as 135 a novel explanatory perspective, demonstrating for the first time that it can mitigate overfitting and 136 enhance the generalization ability of CTC models.

## 3 method

We first introduce the standard CTC algorithm in Section 3.1. Then we present the detailed implementation of our proposed Consistency-Regularized CTC (*CR-CTC*) in Section 3.2, followed by in-depth explanations from different perspectives in Section 3.3.

143 144

> 151 152

158

159 160

161

137

138

## 3.1 PRELIMINARY: CONNECTIONIST TEMPORAL CLASSIFICATION

The ASR task is to convert a sequence of speech frames  $\mathbf{x} = \{x_t\}_1^T$  of length T to a sequence of transcript tokens  $\mathbf{y} = \{y_u \in \mathcal{V}\}_1^U$  of length U, where  $\mathcal{V}$  is the vocabulary and typically  $T \ge U$ . CTC (Graves et al., 2006) extends the vocabulary  $\mathcal{V}$  to  $\mathcal{V}' = \mathcal{V} \cup \{\epsilon\}$  with a blank token  $\epsilon$ , and aims to maximize the total posterior probability of all valid alignments  $\boldsymbol{\pi} = \{\pi_t \in \mathcal{V}'\}_1^T$  between  $\mathbf{x}$  and  $\mathbf{y}$ . Let  $\mathcal{B}(\boldsymbol{\pi})$  denote the many-to-one map that merges repeating tokens and removes all blanks in  $\boldsymbol{\pi}$ , and  $p(\boldsymbol{\pi} | \mathbf{x})$  denote the posterior probability of alignment  $\boldsymbol{\pi}$ , the CTC loss function is formulated as:

$$\mathcal{L}_{\text{CTC}}(\mathbf{x}, \mathbf{y}) = -\log \sum_{\boldsymbol{\pi} \in \mathcal{B}^{-1}(\mathbf{y})} p(\boldsymbol{\pi} | \mathbf{x}).$$

153 Specifically, given the input **x**, it employs an encoder f to estimate the  $|\mathcal{V}'|$ -dimensional probability 154 distributions  $\mathbf{z} = \{z_t\}_1^T$ :  $\mathbf{z} = f(\mathbf{x})^{-1}$ , where f is modeled by a speech encoder network such as 155 Zipformer (Yao et al., 2024) followed by a linear projection layer and a *softmax* function. Note 156 that we now start to use  $\mathcal{L}_{CTC}(\mathbf{z}, \mathbf{y})$  instead of  $\mathcal{L}_{CTC}(\mathbf{x}, \mathbf{y})$  for ease of description in the following 157 sections. Under the frame-independent assumption (Graves et al., 2006),  $p(\pi|\mathbf{x})$  is computed as:

$$p(\boldsymbol{\pi}|\mathbf{x}) = \prod_{t=1}^{T} z_{t,\pi_t},$$
(2)

1)

 $<sup>{}^{1}</sup>T$  is typically downsampled in the encoder f by a factor of 4 for efficiency. This is omitted for the sake of simplicity in expression.



Figure 1: Overall architecture of CR-CTC.

where  $z_{t,\pi_t}$  is the probability of emitting token  $\pi_t$  at frame t.

### 3.2 OUR APPROACH: CONSISTENCY-REGULARIZED CTC

Figure 1 illustrates the overall architecture of our proposed *CR-CTC*. It takes as input two different augmented views,  $\mathbf{x}^{(a)}$  and  $\mathbf{x}^{(b)}$ , both derived from the input speech mel-spectrogram  $\mathbf{x}$ . The two input views are then passed through a shared speech encoder f, which estimates the per-frame distributions:  $\mathbf{z}^{(a)} = f(\mathbf{x}^{(a)}), \mathbf{z}^{(b)} = f(\mathbf{x}^{(b)})$ . In addition to computing the CTC losses on both branches:  $\mathcal{L}_{\text{CTC}}(\mathbf{z}^{(a)}, \mathbf{y})$  and  $\mathcal{L}_{\text{CTC}}(\mathbf{z}^{(b)}, \mathbf{y})$ , we introduce an auxiliary loss (defined in Equation 4) to enforce consistency between  $\mathbf{z}^{(a)}$  and  $\mathbf{z}^{(b)}$ :  $\mathcal{L}_{\text{CR}}(\mathbf{z}^{(a)}, \mathbf{z}^{(b)})$ . The overall loss of the whole model is formulated as:

195 196 197

183

184 185 186

187

$$\mathcal{L} = \frac{1}{2} (\mathcal{L}_{\text{CTC}}(\mathbf{z}^{(a)}, \mathbf{y}) + \mathcal{L}_{\text{CTC}}(\mathbf{z}^{(b)}, \mathbf{y})) + \alpha \mathcal{L}_{\text{CR}}(\mathbf{z}^{(a)}, \mathbf{z}^{(b)}),$$
(3)

where  $\alpha$  is a hyper-parameter that controls the consistency regularization.

J

**Different augmented views.** The two different augmented views,  $\mathbf{x}^{(a)}$  and  $\mathbf{x}^{(b)}$ , are generated by in-200 dependently applying SpecAugment (Park et al., 2019) to two copies of the input mel-spectrogram x. 201 SpecAugment involves warping along time axis, masking blocks of frequency channels, and mask-202 ing blocks of time steps. Since time warping alters feature timing and thus shifts output timestamps, 203 we apply it first before creating the copies to prevent significant timestamp mismatches between 204 the outputs of two branches. Subsequently, random frequency masking and time masking are both 205 applied to the two copies, resulting in  $\mathbf{x}^{(a)}$  and  $\mathbf{x}^{(b)}$ . Note that we also increase the amount of 206 time masking by a factor of 2.5 compared to regular systems. The reason behind this adjustment is 207 explained in Section 3.3, with implementation details provided in Section 4.1. 208

**Consistency regularization loss.** The consistency regularization is applied on each frame t, by minimizing the bidirectional Kullback-Leibler divergence (denoted as  $D_{\text{KL}}$ ) between each pair of distributions  $z_t^{(a)}$  and  $z_t^{(b)}$ :  $D_{\text{KL}}(sg(z_t^{(b)}) || z_t^{(a)})$  and  $D_{\text{KL}}(sg(z_t^{(a)}) || z_t^{(b)})$ , where sg denotes the operation stopping gradient on the target distributions. The regularization loss  $\mathcal{L}_{\text{CR}}(\mathbf{z}^{(a)}, \mathbf{z}^{(b)})$  is formulated as:

214  
215 
$$\mathcal{L}_{CR}(\mathbf{z}^{(a)}, \mathbf{z}^{(b)}) = \frac{1}{2} \sum_{t=1}^{T} D_{KL}(sg(z_t^{(b)}) \| z_t^{(a)}) + D_{KL}(sg(z_t^{(a)}) \| z_t^{(b)}).$$
(4)

# 216 3.3 EXPLANATION 217

We now explain the essential behaviors of our proposed *CR-CTC* from three different perspectives: 1) it performs self-distillation between pairs of sub-models with different input views; 2) it conducts contextual representation learning by predicting the token distributions at masked positions based on unmasked context; 3) it suppresses extremely peaky CTC distributions, mitigating overfitting and enhancing generalization ability. We conduct an empirical investigation through ablation studies in Section 4.3, and the experimental results validate our explanations.

224 Self-distillation. When using model regularization techniques such as dropout (Srivastava et al., 225 2014) and stochastic depth (Huang et al., 2016), which randomly drop parts of the model (neurons or layers), it can be viewed as implicitly training randomly sampled sub-models that are ultimately 226 combined into an ensemble during inference. Similar to R-Drop (Wu et al., 2021) and cosub (Tou-227 vron et al., 2023), in CR-CTC, enforcing consistency regularization between the two branches en-228 ables to perform self-distillation between pairs of randomly sampled sub-models derived from the 229 shared model f, with each sub-model receiving supervision signals in the form of per-frame pre-230 dictions from the other. In addition, feeding different augmented views (with larger amount of time 231 masking) exposes these sub-models to varied aspects of the input data, enhancing their prediction 232 diversity and facilitating richer knowledge transfer as well as complementary representation learn-233 ing. 234

Masked prediction. In CR-CTC, consistency regularization requires frames within the time-masked 235 regions in each branch to predict the corresponding token distributions, which are generated by 236 the other branch on the fly. Similar to masked-based self-supervised models (Devlin et al., 2019; 237 Baevski et al., 2020; Hsu et al., 2021), this behavior encourages the model to capture acoustic in-238 formation on the unmasked context and exploit its implicit language modeling capability. Inde-239 pendently applying random time masking to the two branches reduces the occurrence of positions 240 masked by both branches, thereby improve the quality of the provided target distributions for these 241 masked positions. Furthermore, increasing the amount of time masking in CR-CTC enhances con-242 textual representation learning through the masked prediction behavior.

Peak suppression. In line with previous works (Graves et al., 2006; Sak et al., 2015), we also observe that CTC tends to learn extremely peaky distributions. As shown in Figure 2 (left), almost all non-blank tokens occupy only one frame, while the remaining frames are dominated by the blank token, with both types of emissions occurring with extremely high probabilities. This phenomenon suggests potential overfitting to training data, which limits generalization abality to unseen data.

Enforcing prediction consistency between the two branches in *CR-CTC* guides the model to learn the average of their predictions, ultimately resulting in smoother distributions. The peak suppression behavior reduces overconfidence on training data, thereby improving the model's generalization ability. As presented in Figure 2 (right), *CR-CTC* exhibits reduced token emitting probabilities and an increased occurrence of repeating non-blank tokens. A comparison of concrete statistics on the distribution peakedness between CTC and *CR-CTC* is provided in Table 6.

Inspired by this, we also propose a simple method, called Smooth-Regularized CTC (*SR-CTC*), which incorporates an auxiliary loss into regular CTC, specifically encouraging the model to learn smoother CTC distributions. Appendix Section A.1 presents the details of *SR-CTC*.

<sup>258</sup> 4 EXPERIMENTS

260

261

269

4.1 EXPERIMENTAL SETUP

Datasets. To evaluate the effectiveness of our proposed *CR-CTC*, we conduct experiments on three publicly available ASR datasets: 1) LibriSpeech (Panayotov et al., 2015), which contains 1000 hours of English speech; 2) Aishell-1 (Bu et al., 2017), which consists of 170 hours of Mandarin speech;
3) GigaSpeech (Chen et al., 2021a), comprising 10000 hours of English speech.

Implementation details. Our experiments are performed using the icefall framework<sup>2</sup>, with Lhotse toolkit (Żelasko et al., 2021) for data preparation. For regular ASR recipes in icefall, default parameter settings of SpecAugment (Park et al., 2019) include a time warping factor of 80, 2 frequency

<sup>&</sup>lt;sup>2</sup>https://github.com/k2-fsa/icefall

302

303

304



Figure 2: Visualization of token emitting probabilities for vanilla CTC (left) and our *CR-CTC* (right) on four randomly selected samples from LibriSpeech test set. The gray dashed lines indicate the blank token. Compared to vanilla CTC, the token distributions in *CR-CTC* are smoother with lower emitting probabilities and more repeating non-blank tokens.

305 masking regions with a maximum width of 27, and 10 time masking regions with a maximum width 306 of 100, along with a maximum masking fraction of 15% specifically for time masking. In our CR-307 CTC systems, we utilize larger amount of time masking through increasing both the number of time 308 masking regions and the maximum masking fraction by a factor of 2.5. Speed perturbation (Ko et al., 309 2015) with factors 0.9, 1.0 and 1.1 is applied to LibriSpeech (Panayotov et al., 2015) and Aishell-1 (Bu et al., 2017) datasets. The input features are 80-dimensional mel-spectrograms extracted 310 using 25-ms windows with a 10-ms shift. For LibriSpeech and GigaSpeech datasets, we employ 311 500-class Byte Pair Encoding (BPE) (Sennrich et al., 2016) word pieces as modeling units, while 312 for Aishell-1 dataset, we use 4336-class characters. By default, we set  $\alpha$  in Equation 3 to 0.2. Zip-313 former (Yao et al., 2024), which uses dropout (Srivastava et al., 2014) and stochastic depth (Huang 314 et al., 2016), is used as our speech encoder due to its speed and high performance. Following (Yao 315 et al., 2024), pruned transducer (Kuang et al., 2022), a highly optimized and memory-efficient ver-316 sion of transducer, is employed for comparison. Word-error-rate (WER) and character-error-rate 317 (CER) are employed as ASR metrics for English and Mandarin datasets, respectively. As CR-CTC 318 requires two forward pass during training, we train CR-CTC models with half the batch size and 319 half the number of epochs compared to CTC models, ensuring a fair comparison in terms of train-320 ing cost. Training configuration in terms of number of GPUs and training epochs are provided in 321 Appendix Section A.2. For CTC and CR-CTC systems, we use prefix search decoding (Graves et al., 2006) with a beam size of 4 for comparisons against other state-of-the-art models, and em-322 ploy greedy search decoding for ablation studies. Results comparison between these two decoding 323 methods are provided in Appendix Section A.3. For pruned transducer models, we use beam search

decoding with beam size of 4 (Kang et al., 2023). For CTC/AED systems, we use joint decoding that combines CTC scores and AED scores (Watanabe et al., 2017).

## 4.2 COMPARISON WITH STATE-OF-THE-ART MODELS

In this section, we compare our *CR-CTC* with other state-of-the-art models. For LibriSpeech and GigaSpeech datasets, we also use CR-CTC as an auxiliary loss in CTC/AED and pruned transducer systems for joint training (denoted as CR-CTC/AED and pruned transducer w/ CR-CTC), to further validate the representation learning capability of CR-CTC. Note that for the models that combine *CR-CTC* and pruned transducer, we only utilize the transducer head for decoding, without incor-porating the CTC scores. For the larger GigaSpeech dataset, we additionally use a even larger scale of Zipformer (Zipformer-XL). Model configuration of different scales of Zipformer are pro-vided in Table 14. For Aishell-1 dataset, which is considerably smaller, we conduct experiments on Zipformer-S and Zipformer-M to ensure comparable parameter counts with other models reported in the literature. 

**LibriSpeech dataset.** Table 1 presents the results on LibriSpeech dataset for *CR-CTC* and other state-of-the-art models. Our CR-CTC significantly outperforms the CTC baselines on all three scales of Zipformer encoder. When comparing to CTC/AED models, our CR-CTC achieves lower WER on Zipformer-M/L, while yielding comparable result on Zipformer-S. Similarly, our CR-CTC surpasses pruned transducer on Zipformer-M, and performs comparably on Zipformer-L. It also demonstrates that *CR-CTC* can further enhance the performance of CTC/AED and pruned transducer models when used for jointly training. A notable result is that pruned transducer combined with CR-CTC using Zipformer-L achieves a new state-of-the-art result of 1.88%/3.95% on test-clean/test-other, outperforming both the transducer models with Conformer-L (Gulati et al., 2020) and Stateformer 25L (Fathullah et al., 2023).

Table 1: WER(%) performance of our method on LibriSpeech dataset compared to the best results reported in the literature without using an external language model. 

351			WER	R (%)
353	Model	Params (M)	test-clean	test-other
354	CTC/AED, E-Branchformer-B (Kim et al., 2023)	41.1	2.49	5.61
355	CTC/AED, Branchformer (Peng et al., 2022)	116.2	2.4	5.5
356	CTC/AED, E-Branchformer-L (Kim et al., 2023)	148.9	2.14	4.55
357	Transducer, ContextNet-S (Han et al., 2020)	10.8	2.9	7.0
358	Transducer, ContextNet-M (Han et al., 2020)	31.4	2.4	5.4
359	Transducer, ContextNet-L (Han et al., 2020)	112.7	2.1	4.6
360 361	Transducer, Conformer-S (Gulati et al., 2020) Transducer, Conformer-M (Gulati et al., 2020) Transducer, Conformer-L (Gulati et al., 2020)	10.3 30.7 118.8	2.7 2.3 2.1	6.3 5.0 4.3
363	Transducer, MH-SSM 32L (Fathullah et al., 2023)	140.3	2.01	4.61
364	Transducer, Stateformer 25L (Fathullah et al., 2023)	139.8	1.91	4.36
365	CTC/AED, Zipformer-S (Yao et al., 2024)	46.3	2.46	6.04
366	CTC/AED, Zipformer-M (Yao et al., 2024)	90.0	2.22	4.97
267	CTC/AED, Zipformer-L (Yao et al., 2024)	174.3	2.09	4.59
368 369	Pruned transducer, Zipformer-S (Yao et al., 2024) Pruned transducer, Zipformer-M (Yao et al., 2024) Pruned transducer, Zipformer-L (Yao et al., 2024)	23.3 65.6 148.4	2.42 2.21 2.00	5.73 4.79 4.38
370	CTC, Zipformer-S	22.1	2.85	6.89
371	CTC, Zipformer-M	64.3	2.52	6.02
372	CTC, Zipformer-L	147.0	2.5	5.72
373	<i>CR-CTC</i> , Zipformer-S (ours)	22.1	2.52	5.85
374	<i>CR-CTC</i> , Zipformer-M (ours)	64.3	2.1	4.61
375	<i>CR-CTC</i> , Zipformer-L (ours)	147.0	2.02	4.35
376	<i>CR-CTC</i> /AED, Zipformer-L (ours)	174.3	1.96	4.08
377	Pruned transducer w/ <i>CR-CTC</i> , Zipformer-L (ours)	148.8	<b>1.88</b>	<b>3.95</b>

Aishell-1 dataset. Table 2 presents the results on Aishell-1 dataset. Our *CR-CTC* models not only significantly outperform vanilla CTC by a substantial margin but also achieve better results than all other CTC/AED and pruned transducer models. For example, *CR-CTC* with Zipformer-S surpasses CTC/AED with Zipformer-M while using much fewer parameters.

Table 2: WER(%) performance of our method on Aishell-1 dataset compared to the best results reported in the literature without using an external language model.

Madal	Darama (M)	WER (%)	
Widdel		dev	test
CTC/AED, Conformer in ESPnet (Watanabe et al., 2018)	46.2	4.5	4.9
CTC/AED, Conformer in WeNet (Yao et al., 2021)	46.3	—	4.61
CTC/AED, E-Branchformer in ESPnet (Watanabe et al., 2018)	37.9	4.2	4.5
CTC/AED, Branchformer (Peng et al., 2022)	45.4	4.19	4.43
Pruned transducer, Zipformer-S (Yao et al., 2024)	30.2	4.4	4.67
Pruned transducer, Zipformer-M (Yao et al., 2024)	73.4	4.13	4.4
CTC, Zipformer-S	23.1	4.89	5.26
CTC, Zipformer-M	66.2	4.47	4.8
CTC/AED, Zipformer-S	39.3	4.47	4.8
CTC/AED, Zipformer-M	83.2	4.0	4.32
CR-CTC, Zipformer-S (ours)	23.1	3.9	4.12
CR-CTC, Zipformer-M (ours)	66.2	3.72	4.02

399 400

406

407

408

GigaSpeech dataset. Table 3 shows the results on GigaSpeech dataset. Our *CR-CTC* consistently
 achieves a significantly lower WER than vanilla CTC across all scales of Zipformer. In comparisons
 with CTC/AED or pruned transducer models, our *CR-CTC* demonstrates comparable performance
 on Zipformer L/XL. Additionally, the results indicate that employing *CR-CTC* for joint training can
 further improve the performance of both CTC/AED and pruned transducer models.

Table 3: WER(%) performance of our method on GigaSpeech dataset compared to the best results reported in the literature without using an external language model.

409	Madal	Darama (M)	WEF	R (%)
410	Model	Parallis (M)	dev	test
411	CTC/AED, Transformer (Chen et al., 2021a)	87	12.30	12.30
412	CTC/AED, Conformer in Wenet (Zhang et al., 2022)	113.2	10.7	10.6
413	CTC/AED, Conformer in ESPnet (Chen et al., 2021a)	113.2	10.9	10.8
414	CTC/AED, E-Branchformer in ESPnet (Watanabe et al., 2018)	148.9	10.6	10.5
415	CTC, Zipformer-S	22.1	12.08	11.95
416	CTC, Zipformer-M	64.3	11.23	11.27
417	CTC, Zipformer-L	147.0	11.16	11.16
418	CTC, Zipformer-XL	286.6	10.8	10.87
419	CTC/AED, Zipformer-S	46.3	11.4	11.39
420	CTC/AED, Zipformer-M	90.0	10.57	10.61
404	CTC/AED, Zipformer-L	174.3	10.26	10.38
421	CTC/AED, Zipformer-XL	315.5	10.22	10.33
422	Pruned transducer, Zipformer-S	23.3	10.98	10.94
423	Pruned transducer, Zipformer-M	65.6	10.37	10.42
424	Pruned transducer, Zipformer-L	148.4	10.23	10.28
425	Pruned transducer, Zipformer-XL	288.2	10.09	10.2
426	CR-CTC, Zipformer-S (ours)	22.1	11.68	11.58
427	CR-CTC, Zipformer-M (ours)	64.3	10.62	10.72
428	CR-CTC, Zipformer-L (ours)	147.0	10.31	10.41
429	CR-CTC, Zipformer-XL (ours)	286.6	10.15	10.28
430	CR-CTC/AED, Zipformer-XL (ours)	315.5	9.92	10.07
431	Pruned transducer w/ CR-CTC, Zipformer-XL (ours)	286.6	9.95	10.03

## 432 4.3 ABLATION STUDIES

433 434

466

434 We now perform ablation studies on LibriSpeech dataset using Zipformer-M encoder to investigate 435 the effect of each component in *CR-CTC* (Section 3.2), and to validate our explanations of its be-436 haviors (Section 3.3). Results of tuning  $\alpha$  in Equation 3 and the ratio used to increase the amount of 437 time masking are presented in Table 15.

438 Self-distillation. One self-distillation method in self-supervised learning is to construct a teacher 439 model by tracking the model weights using exponential moving average (EMA) (Grill et al., 2020; 440 He et al., 2020; Baevski et al., 2023). For comparison, we include this approach, referred to as 441 EMA-distilled CTC, which incorporates an auxiliary loss to learn from the CTC distribution of the 442 EMA teacher model. Its details are provided in Appendix Section A.6. As presented in Table 4, CR-*CTC* significantly outperforms EMA-distilled CTC, demonstrating its superiority in self-distillation. 443 For *CR-CTC*, both the lack of increased time masking and the absence of different augmented views 444 lead to WER degradation, indicating the effectiveness of enhancing the input diversity between sub-445 models during self-distillation. Replacing  $D_{\rm KL}$  with hard label-based cross-entropy (CE) function 446 in  $\mathcal{L}_{CR}$  (Equation 4) results in a WER degradation of 0.02%/0.22% on test-clean/test-other. This 447 suggests the advantage of using  $D_{\rm KL}$  which enables a finer-grained self-distillation as it distills over 448 the full CTC lattice, whereas the hard label CE-based method only distills the best alignment. When 449 removing the sg operation in  $\mathcal{L}_{CR}$ , the WER increase by 0.12%/0.35%, which implies that the model 450 might have a tendency towards a degenerated solution (Chen & He, 2021) that is insensitive to the 451 pattern of input masking and model dropout. 452

456	Mathad		WER (%)	
457	Method	test-clean	test-other	
458	CTC baseline	2.51	6.02	
460	EMA-distilled CTC	2.31	5.25	
461	CR-CTC (final)	2.12	4.62	
462	No larger time masking	2.19	4.98	
463	No larger time masking, no different augmented views	2.27	5.11	
160	Use hard-label CE-based $\mathcal{L}_{\mathrm{CR}}$	2.14	4.84	
465	Remove $sg$ in $\mathcal{L}_{CR}$	2.24	4.97	

Table 4: Ablation studies for self-distillation in *CR-CTC* on LibriSpeech dataset using Zipformer-M
encoder and greedy search decoding.

Masked prediction. As reported in Table 5, without increasing the amount of time masking, the 467 WER of *CR-CTC* increases by 0.07%/0.36% on *test-clean/test-other*, suggesting the effectiveness 468 of enhancing the masked prediction behavior for contextual representation learning. Additionally, 469 without using different augmented views, the WER increases further by 0.12%/0.13%. This indi-470 cates the advantage of independently applying random time masking, which improves the quality 471 of the provided target distributions for the masked positions. However, using larger amount of 472 frequency masking leads to a WER degradation of 0.07% on *test-clean*, implying that the perfor-473 mance gain from increasing the amount of time masking is primarily due to the masked prediction 474 behavior, rather than merely increasing the input diversity for the two branches. Furthermore, ap-475 plying larger amount of time masking does not benefit the CTC baseline, as it increases the WER 476 by 0.17%/0.26%. In the final CR-CTC system, excluding frames with time-masked regions in the 477 current branch (self-masked) from  $\mathcal{L}_{CR}$  (Equation 4) leads to a larger WER degradation compared to excluding the remaining unmasked frames (self-unmasked). This highlights the importance of the 478 masked prediction behavior in the overall performance of CR-CTC. 479

Peak suppression. To measure the peakedness of the learned CTC distributions, we compute the averaged duration over all non-blank tokens, as well as the averaged emitting probabilities for the blank token and all non-blank tokens, based on the best alignment obtained through greedy search decoding on the test sets. We also include the method *SR-CTC* (described in Appendix Section A.1) for comparison. As presented in Table 6, compared to the CTC baseline, *CR-CTC* learns smoother distributions and significantly improves the recognition performance. Note that *SR-CTC* also surpasses the CTC baseline by 0.19%/0.8% on *test-clean/test-other*, while exhibiting a notably larger

Mathad	WER (%)		
Method	test-clean	test-other	
CTC baseline	2.51	6.02	
Use larger time masking	2.68	6.28	
CR-CTC (final)	2.12	4.62	
No larger time masking	2.19	4.98	
No larger time masking, no different augmented views	2.27	5.11	
No larger time masking, use larger frequency masking	2.26	4.98	
Exclude self-masked frames in $\mathcal{L}_{CR}$	2.32	5.26	
Exclude self-unmasked frames in $\mathcal{L}_{\mathrm{CR}}$	2.32	5.02	

Table 5: Ablation studies for masked prediction in *CR-CTC* on LibriSpeech dataset using Zipformer M encoder and greedy search decoding.

> average duration of non-blank tokens. This manifests the effectiveness of peak suppression in reducing overfitting and improving generalization performance.

Table 6: Ablation studies for peak suppression in *CR-CTC* on LibriSpeech dataset using ZipformerM encoder and greedy search decoding. We include the averaged duration of all non-blank tokens, as well as the averaged emitting probabilities of the blank token and all non-blank tokens on the best alignments.

Mathad	Non-blank duration   Emit probability (%)		WER (%)		
Method	(frames)	blank	non-blank	test-clean	test-other
CTC baseline	1.04	99.64	98.50	2.51	6.02
SR-CTC	4.25	95.44	90.04	2.32	5.22
CR-CTC	1.28	94.19	89.42	2.12	4.62

Compared to using auxiliary head for jointly training. The straightforward approach to improve the CTC performance is using an auxiliary head of AED (Chan et al., 2015; Watanabe et al., 2017) or pruned transducer (Kuang et al., 2022) for jointly training. As reported in Table 7, *CR-CTC* significantly outperforms these two methods with less model parameters, suggesting the advantage of our method.

Table 7: Comparison between *CR-CTC* and methods using an auxiliary head for jointly training on LibriSpeech dataset using Zipformer-M encoder and greedy search decoding.

Method	Parame (M)	WER (%)		
Wethou		test-clean	test-other	
CTC baseline	64.3	2.51	6.02	
CTC w/ AED head	90.0	2.46	5.57	
CTC w/ pruned transducer head	65.8	2.42	5.4	
CR-CTC	64.3	2.12	4.62	

## 

## 5 CONCLUSION

In this work, we introduce the *CR-CTC* to enhance CTC performance. Specifically, it takes as
input two different augmented views of the same speech mel-spectrogram, and enforce consistency between the two obtained CTC distributions. We explain our method from three different
perspectives: 1) self-distillation between randomly sampled sub-models; 2) masked prediction for
positions within time-masked regions, facilitating the learning of contextual representation; 3) peak
suppression, which reduces overfitting and improves the model's generalization ability. Extensive
experiments on LibriSpeech, Aishell-1, and GigaSpeech datasets demonstrate the effectiveness of *CR-CTC*. Additionally, detailed ablation studies validate our explanations.

540	REFERENCES
541	

556

567

569

570

576

580

581

582

583

- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and 542 self-distillation in deep learning. arXiv preprint arXiv:2012.09816, 2020. 543
- 544 Alexei Baevski, Steffen Schneider, and Michael Auli. vq-wav2vec: Self-supervised learning of discrete speech representations. arXiv preprint arXiv:1910.05453, 2019. 546
- 547 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. Advances in neural information 548 processing systems, 33:12449–12460, 2020. 549
- Alexei Baevski, Arun Babu, Wei-Ning Hsu, and Michael Auli. Efficient self-supervised learning 551 with contextualized target representations for vision, speech and language. In International Con-552 ference on Machine Learning, pp. 1416–1429. PMLR, 2023. 553
- 554 Herve A Bourlard and Nelson Morgan. Connectionist speech recognition: a hybrid approach, vol-555 ume 247. Springer Science & Business Media, 2012.
- Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng. Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline. In 20th conference of the oriental chapter 558 of the international coordinating committee on speech databases and speech I/O systems and 559 assessment (O-COCOSDA), pp. 1–5, 2017.
- 561 Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data 563 mining, pp. 535-541, 2006.
- 564 William Chan, Navdeep Jaitly, Quoc V Le, and Oriol Vinyals. Listen, attend and spell. arXiv 565 preprint arXiv:1508.01211, 2015. 566
- Guoguo Chen, Shuzhou Chai, Guanbo Wang, Jiayu Du, Wei-Qiang Zhang, Chao Weng, Dan Su, 568 Daniel Povey, Jan Trmal, Junbo Zhang, et al. Gigaspeech: An evolving, multi-domain asr corpus with 10,000 hours of transcribed audio. arXiv preprint arXiv:2106.06909, 2021a.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for 571 contrastive learning of visual representations. In International conference on machine learning, 572 pp. 1597–1607. PMLR, 2020. 573
- 574 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In Proceedings of 575 the IEEE/CVF conference on computer vision and pattern recognition, pp. 15750–15758, 2021.
- Zhehuai Chen, Andrew Rosenberg, Yu Zhang, Heiga Zen, Mohammadreza Ghodsi, Yinghui Huang, 577 Jesse Emond, Gary Wang, Bhuvana Ramabhadran, and Pedro J Moreno. Semi-supervision in asr: 578 Sequential mixmatch and factorized tts-based augmentation. In Interspeech, pp. 736–740, 2021b. 579
  - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186, 2019.
- Haisong Ding, Kai Chen, and Qiang Huo. Improving knowledge distillation of ctc-trained acoustic 585 models with alignment-consistent ensemble and target delay. IEEE/ACM transactions on audio, 586 speech, and language processing, 28:2561–2571, 2020.
- 588 Yassir Fathullah, Chunyang Wu, Yuan Shangguan, Junteng Jia, Wenhan Xiong, Jay Mahadeokar, 589 Chunxi Liu, Yangyang Shi, Ozlem Kalinli, Mike Seltzer, et al. Multi-head state space model for 590 speech recognition. arXiv preprint arXiv:2305.12498, 2023.
- Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. 592 Born again neural networks. In International conference on machine learning, pp. 1607–1616. PMLR, 2018.

594	Yingying Gao, Junlan Feng, Tianrui Wang, Chao Deng, and Shilei Zhang. A ctc triggered siamese
595	network with spatial-temporal dropout for speech recognition. arXiv preprint arXiv:2206.08031.
596	2022.
597	

- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
   Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- Alex Graves. Sequence transduction with recurrent neural networks. arXiv preprint
   arXiv:1211.3711, 2012.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist tem poral classification: labelling unsegmented sequence data with recurrent neural networks. In
   *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376, 2006.
- Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena
   Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
   et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.
- Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. Conformer: Convolution-augmented Transformer for Speech Recognition. In *Proc. Interspeech 2020*, pp. 5036–5040, 2020.
- Wei Han, Zhengdong Zhang, Yu Zhang, Jiahui Yu, Chung-Cheng Chiu, James Qin, Anmol Gulati,
  Ruoming Pang, and Yonghui Wu. Contextnet: Improving convolutional neural networks for
  automatic speech recognition with global context. In *Interspeech 2020*, pp. 3610–3614, 2020.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- Byeongho Heo, Taekyung Kim, Sangdoo Yun, and Dongyoon Han. Masking augmentation for supervised learning. *arXiv preprint arXiv:2306.11339*, 2023.
- Yosuke Higuchi, Niko Moritz, Jonathan Le Roux, and Takaaki Hori. Momentum pseudo-labeling
   for semi-supervised speech recognition. *arXiv preprint arXiv:2106.08922*, 2021.
- Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly,
   Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. Deep neural networks
   for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal processing magazine*, 29(6):82–97, 2012.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network.
   *ArXiv*, abs/1503.02531, 2015.
- Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov,
   and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked
   prediction of hidden units. *IEEE/ACM transactions on audio, speech, and language processing*,
   29:3451–3460, 2021.
- Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pp. 646–661, 2016.
- Ruizhe Huang, Xiaohui Zhang, Zhaoheng Ni, Li Sun, Moto Hira, Jeff Hwang, Vimal Manohar,
  Vineel Pratap, Matthew Wiesner, Shinji Watanabe, et al. Less peaky and more accurate ctc forced
  alignment by label priors. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 11831–11835. IEEE, 2024.

677

684

685

686

687

688 689

690

691

648	Dongwei Jiang, Wubo Li, Miao Cao, Wei Zou, and Xiangang Li. Speech simclr: Combining con-
649	trastive and reconstruction objective for self-supervised speech representation learning. arXiv
650	preprint arXiv:2010.13991, 2020.
651	

- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. Span-652 bert: Improving pre-training by representing and predicting spans. Transactions of the association 653 for computational linguistics, 8:64–77, 2020. 654
- 655 Wei Kang, Liyong Guo, Fangjun Kuang, Long Lin, Mingshuang Luo, Zengwei Yao, Xiaoyu Yang, 656 Piotr Żelasko, and Daniel Povey. Fast and parallel decoding for transducer. In ICASSP 2023-2023 657 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 1–5. 658 IEEE, 2023.
- 659 Soheil Khorram, Jaeyoung Kim, Anshuman Tripathi, Han Lu, Qian Zhang, and Hasim Sak. Con-660 trastive siamese network for semi-supervised speech recognition. In ICASSP 2022-2022 IEEE 661 International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 7207–7211. 662 IEEE, 2022. 663
- 664 Eungbeom Kim, Hantae Kim, and Kyogu Lee. Guiding frame-level ctc alignments using self-665 knowledge distillation. arXiv preprint arXiv:2406.07909, 2024.
- 666 Kwangyoun Kim, Felix Wu, Yifan Peng, Jing Pan, Prashant Sridhar, Kyu J Han, and Shinji Watan-667 abe. E-branchformer: Branchformer with enhanced merging for speech recognition. In 2022 668 IEEE Spoken Language Technology Workshop (SLT), pp. 84–91. IEEE, 2023. 669
- 670 Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur. Audio augmentation for 671 speech recognition. In Sixteenth annual conference of the international speech communication 672 association, 2015.
- Fangjun Kuang, Liyong Guo, Wei Kang, Long Lin, Mingshuang Luo, Zengwei Yao, and Daniel 674 Povey. Pruned rnn-t for fast, memory-efficient asr training. arXiv preprint arXiv:2206.13236, 675 2022. 676
- Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. arXiv preprint 678 arXiv:1610.02242, 2016. 679
- Ryo Masumura, Mana Ihori, Akihiko Takashima, Takafumi Moriya, Atsushi Ando, and Yusuke 680 Shinohara. Sequence-level consistency training for semi-supervised end-to-end automatic speech 681 recognition. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and 682 Signal Processing (ICASSP), pp. 7054–7058. IEEE, 2020. 683
  - Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. Advances in Neural Information Processing Systems, 33:3351–3361, 2020.
  - Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748, 2018.
  - Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: an asr corpus based on public domain audio books. In IEEE international conference on acoustics, speech and signal processing (ICASSP), pp. 5206–5210, 2015.
- Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and 693 Quoc V Le. Specaugment: A simple data augmentation method for automatic speech recognition. 694 arXiv preprint arXiv:1904.08779, 2019. 695
- 696 Yifan Peng, Siddharth Dalmia, Ian Lane, and Shinji Watanabe. Branchformer: Parallel mlp-attention 697 architectures to capture local and global context for speech recognition and understanding. In 698 International Conference on Machine Learning, pp. 17627–17643. PMLR, 2022. 699
- Heli Qi, Sashi Novitasari, Sakriani Sakti, and Satoshi Nakamura. Improved consistency training for 700 semi-supervised sequence-to-sequence asr via speech chain reconstruction and self-transcribing. 701 arXiv preprint arXiv:2205.06963, 2022.

- Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. *Advances in neural information processing systems*, 29, 2016.
- Haşim Sak, Andrew Senior, Kanishka Rao, Ozan Irsoy, Alex Graves, Françoise Beaufays, and Johan Schalkwyk. Learning acoustic frame labeling for speech recognition with recurrent neural networks. In 2015 IEEE international conference on acoustics, speech and signal processing (ICASSP), pp. 4280–4284, 2015.
- Ashtosh Sapru. Using data augmentation and consistency regularization to improve semi-supervised
   speech recognition. 2022.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pp. 1715–1725, 2016.
- Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel,
  Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised
  learning with consistency and confidence. *Advances in neural information processing systems*,
  33:596–608, 2020.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
   Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Ryoichi Takashima, Li Sheng, and Hisashi Kawai. Investigation of sequence-level knowledge dis tillation methods for ctc acoustic models. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6156–6160. IEEE, 2019.
- Hugo Touvron, Matthieu Cord, Maxime Oquab, Piotr Bojanowski, Jakob Verbeek, and Hervé Jégou.
   Co-training 2l submodels for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11701–11710, 2023.
- Shinji Watanabe, Takaaki Hori, Suyoun Kim, John R Hershey, and Tomoki Hayashi. Hybrid ctc/attention architecture for end-to-end speech recognition. *IEEE Journal of Selected Topics in Signal Processing*, 11(8):1240–1253, 2017.
- Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson
  Enrique Yalta Soplin, Jahn Heymann, Matthew Wiesner, Nanxin Chen, Adithya Renduchintala,
  and Tsubasa Ochiai. ESPnet: End-to-end speech processing toolkit. In *Proceedings of Inter- speech*, pp. 2207–2211, 2018.
- Felix Weninger, Franco Mana, Roberto Gemello, Jesús Andrés-Ferrer, and Puming Zhan. Semi-supervised learning with data augmentation for end-to-end asr. *arXiv preprint arXiv:2007.13876*, 2020.
- Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu, et al.
  R-drop: Regularized dropout for neural networks. *Advances in Neural Information Processing Systems*, 34:10890–10905, 2021.
- Zengwei Yao, Liyong Guo, Xiaoyu Yang, Wei Kang, Fangjun Kuang, Yifan Yang, Zengrui Jin, Long Lin, and Daniel Povey. Zipformer: A faster and better encoder for automatic speech recognition. In *The Twelfth International Conference on Learning Representations*, 2024.
- Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang, Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu
   Chen, Lei Xie, and Xin Lei. WeNet: Production Oriented Streaming and Non-Streaming End-to End Speech Recognition Toolkit. In *Proc. Interspeech*, pp. 4054–4058, 2021.
- Ji Won Yoon, Hyeonseung Lee, Ju Yeon Kang, and Nam Soo Kim. Cons-kd: Dropout-robust knowledge distillation for ctc-based automatic speech recognition. *IEEE Access*, 2024.
- 755 Piotr Żelasko, Daniel Povey, Jan Trmal, Sanjeev Khudanpur, et al. Lhotse: a speech data representation library for the modern deep learning ecosystem. *arXiv preprint arXiv:2110.12561*, 2021.

756 757 758	Binbin Zhang, Di Wu, Zhendong Peng, Xingchen Song, Zhuoyuan Yao, Hang Lv, Lei Xie, Chao Yang, Fuping Pan, and Jianwei Niu. Wenet 2.0: More productive end-to-end speech recognition toolkit. <i>arXiv preprint arXiv:2203.15455</i> , 2022.
759 760	Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your
761 762	own teacher: Improve the performance of convolutional neural networks via self distillation. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pp. 3713–3722, 2019.
763	Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. Advances
764	in neural information processing systems, 31, 2018.
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
707	
707	
789	
790	
791	
792	
793	
794	
795	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805 800	
000 807	
808	
809	

# 810 A APPENDIX

# 812 A.1 SMOOTH-REGULARIZED CTC

814 Smooth-regularized CTC (*SR-CTC*) discourages peaky distributions by adding an smooth regular-815 ization loss (denoted as  $\mathcal{L}_{SR}$ ) to regular CTC model. Specifically, we first apply a smooth kernel *K* 816 of size 3 to the model prediction z, smoothing it along the time dimension:  $z^{(s)} = smooth(z, K)$ . 817 The smoothing operation is done by using a 1-D depth-wise convolution layer. Then we minimize 818 the  $D_{KL}$  between z and  $z^{(s)}$ , similar to the consistency loss in *CR-CTC* (Equation 4):

$$\mathcal{L}_{\mathrm{SR}}(\mathbf{z}, \mathbf{z}^{(s)}) = \sum_{t=1}^{T} D_{\mathrm{KL}}(sg(z_t^{(s)}) || z_t).$$
(5)

823 The overall loss of *SR*-*CTC* is formulated as:

$$\mathcal{L}' = \mathcal{L}_{\rm CTC}(\mathbf{z}, \mathbf{y}) + \beta \mathcal{L}_{\rm SR}(\mathbf{z}, \mathbf{z}^{(s)}), \tag{6}$$

where  $\beta$  is a hyper-parameter. In this work, we use K = (0.25, 0.5, 0.25) and  $\beta = 0.2$ .

We validate its effectiveness in Section 4.3. Table 6 presents the experimental result.

#### A.2 TRAINING CONFIGURATION

Training configuration, including the number of GPUs and training epochs, on LibriSpeech, Aishell-1 and GigaSpeech datasets are presented in Table 8, Table 9, and Table 10, respectively.

Model	GPUs (80G NVIDIA Tesla A100)	Epochs
CTC, Zipformer-S	1	100
CTC, Zipformer-M	2	100
CTC, Zipformer-L	2	100
<i>CR-CTC</i> , Zipformer-S	1	50
<i>CR-CTC</i> , Zipformer-M	2	50
<i>CR-CTC</i> , Zipformer-L	2	50
<i>CR-CTC</i> /AED, Zipformer-L Pruned transducer w/ <i>CR-CTC</i> , Zipformer-L	22	50 50

Table 8: Training configuration on LibriSpeech dataset.

Table 9: Training configuration on Aishell-1 dataset.

Model	GPUs (80G NVIDIA Tesla A100)	Epochs
CTC, Zipformer-S	1	120
CTC, Zipformer-M	1	120
CTC/AED, Zipformer-S	1	60
CTC/AED, Zipformer-M	1	60
CR-CTC, Zipformer-S	1	60
CR-CTC, Zipformer-M	1	60

#### A.3 RESULTS OF DIFFERENT DECODING METHODS

Results comparison between greedy search decoding and prefix search decoding for CTC and *CR-CTC* on LibriSpeech, Aishell-1 and GigaSpeech datasets are presented in Table 11, Table 12, and Table 13, respectively.

866	Model		(80C NI	GPUs	A 100)	Epochs
000			(800 N		A100)	
000	CTC, Zipformer-S	CTC, Zipformer-S			2	
009	CTC Zipformer-L	CTC, Zipformer-M CTC, Zipformer-L			2	
070	CTC, Zipformer-XL	CTC, Zipformer-XL			4	
872	CTC/AED, Zipformer-S	CTC/AED. Zipformer-S			2	
873	CTC/AED, Zipformer-M	CTC/AED, Zipformer-M			2	
874	CTC/AED, Zipformer-L	CTC/AED, Zipformer-L			2	
875	CTC/AED, Ziptormer-XL	CTC/AED, Zipformer-XL			4	
876	Pruned transducer, Zipformer	Pruned transducer, Zipformer-S			2	
77	Pruned transducer, Zipformer	-M I		$\frac{2}{2}$		30
, 8	Pruned transducer, Zipformer	Pruned transducer, Zipformer-L Pruned transducer, Zipformer, VI			4	
9	CD CTC Zinformon S	<u>AL</u>		2		20
30	<i>CR-CTC</i> , Zipformer-S			$\frac{2}{2}$		30
81	<i>CR-CTC</i> , Zipformer-L			$\frac{2}{2}$		30
2	CR-CTC, Zipformer-XL			4		30
-	CR-CTC/AED, Zipformer-XI	,		4		30
4	Pruned transducer w/ CR-CTC	C, Zipformer-Z	XL	4		30
5						
36	Table 11: WER (%) results o	of different de	ecoding meth	nods on Libr	iSpeec	h dataset.
27						
20	Model	Greedy sear	ch decoding	Prefix searc	h decod	ling
20		test-clean	test-other	test-clean	test-ot	her
0	CTC, Zipformer-S	2.85	6.91	2.85	6.89	)
,	CTC, Zipformer-M	2.51	6.02	2.52	6.02	2
	CIC, Ziptormer-L	2.49	5.7	2.5	5.72	2
	CR-CTC, Zipformer-S	2.57	5.95	2.52	5.85	5
	CR-CTC, Zipformer-M	2.12	4.62	2.1	4.6	1
	CR-CTC, Ziplormer-L	2.03	4.37	2.02	4.33	<u> </u>
	Table 12: WFR (%) results	of different	decoding me	thods on Ai	chall 1	dataset
	Tuble 12. WER (70) results	or unrerent	account int		511011-1	aataset.
		Greedy sear	ch decoding	Prefix sear	ch deco	ding
	Model	dev	test	dev	test	
	CTC Zinformer S	1 88	5.26	1 80	5 76	
	CTC, Zipformer-M	4.46	4.8	4.47	4.8	
		2.0	4.12	2.0	4.10	
	CR- $CTC$ , Zipformer-S	3.9	4.12	3.9	4.12 4.02	
	en-ere, Zipioriner-W	5.15	<b>H.U</b> 2	5.12	T.02	
	Table 13: WFR (%) results of	of different d	ecoding meth	nods on Gig	aSpeec	h dataset
	1000 10.  WER (70) 1000000 0		county mou	ious on Oige	uspece	n uutaset.
		Greedy sea	rch decoding	Prefix searc	h decor	ling
	Model	dev	test	dev	test	ığ
	CTC Zinformer S	12.15	12.03	12.08	11.05	
	CTC, Zipformer-M	11.3	12.05	11.23	11.93	
	CTC, Zipformer-L	11.21	11.19	11.16	11.16	
	CTC, Zipformer-XL	10.85	10.91	10.8	10.87	
	CR-CTC, Zipformer-S	11.85	11.8	11.68	11.58	
	CR-CTC, Zipformer-M	10.78	10.88	10.62	10.72	
	CR-CTC, Zipformer-L	10.42	10.56	10.31	10.41	
	CR-CTC, Zipformer-XL	10.28	10.41	10.15	10.28	

## Table 10: Training configuration on GigaSpeech dataset.

Scale	layer-numbers	embedding-dimensions	feed-forward-dimensions
S	{2,2,2,2,2,2}	{192,256,256,256,256,256}	{512,768,768,768,768,768}
Μ	{2,2,3,4,3,2}	{192,256,384,512,384,256}	{512,768,1024,1536,1024,768}
L	{2,2,4,5,4,2}	{192,256,512,768,512,256}	{512,768,1536,2048,1536,768}
XL	{2,2,4,5,4,2}	{192,384,768,1024,768,384}	{512,1024,2048,3072,2048,1024}

Table 14: Model configuration of Zipformer at four different scales.

A.4 MODEL CONFIGURATION OF DIFFERENT SCALES OF ZIPFORMER

Table 14 presents model configuration of different scales of Zipformer.

A.5 ABLATION STUDIES ON HYPER-PARAMETER TUNING

Table 15 presents results of tuning hyper-parameters, including  $\alpha$  in Equation 3 and the ratio used to increase the amount of time masking for *CR-CTC*.

Table 15: Results of tuning  $\alpha$  that controls  $\mathcal{L}_{CR}$  (Equation 3) and the ratio used to increase the amount of time-masking for *CR-CTC* on LibriSpeech dataset using Zipformer-M encoder and greedy search decoding.

Hyper-parameter	WER (%) test-clean test-other	
$\begin{aligned} \alpha &= 0.1\\ \alpha &= 0.2 \text{ (final)}\\ \alpha &= 0.3 \end{aligned}$	2.19 <b>2.12</b> 2.23	4.8 <b>4.62</b> 4.84
$\begin{array}{c} 1.0\times \text{ time masking} \\ 1.5\times \text{ time masking} \\ 2.0\times \text{ time masking} \\ 2.5\times \text{ time masking (final)} \\ 3.0\times \text{ time masking} \end{array}$	2.19 2.19 2.17 <b>2.12</b> 2.17	4.98 4.73 4.71 <b>4.62</b> 4.81

### A.6 EMA-DISTILLED CTC

In EMA-distilled CTC, the teacher model  $f^{(e)}$  is dynamically constructed for self-distillation. Its weights  $\theta^{(e)}$  are updated using the exponential moving average of the current model's weights  $\theta$ :  $\theta^{(e)} \leftarrow \tau \theta^{(e)} + (1 - \tau)\theta$ , where  $\tau = \min(0.9999, 1 - 10/\max(20, \text{step}))$ . The teacher model  $f^{(e)}$  processes the unmasked input  $\mathbf{x}^{(e)}$ , and produces the CTC distribution  $\mathbf{z}^{(e)} = f^{(e)}(\mathbf{x}^{(e)})$  which serves as distillation target for the current model f. Similar to  $\mathcal{L}_{CR}$  in *CR-CTC* (Equation 4), the distillation loss  $\mathcal{L}_{EMA}$  is defined as:

$$\mathcal{L}_{\text{EMA}}(\mathbf{z}, \mathbf{z}^{(e)}) = \sum_{t=1}^{T} D_{\text{KL}}(sg(z_t^{(e)}) || z_t).$$
(7)

962 The overall loss of EMA-distilled CTC is formulated as:

 $\mathcal{L}'' = \mathcal{L}_{\text{CTC}}(\mathbf{z}, \mathbf{y}) + \gamma \mathcal{L}_{\text{EMA}}(\mathbf{z}, \mathbf{z}^{(e)}), \tag{8}$ 

where  $\gamma$  is a hyper-parameter. In this work, we use  $\gamma = 0.2$ . Table 4 presents the experimental result.