

# BOUNDARY AND CONTEXT AWARE TRAINING FOR CIF-BASED NON-AUTOREGRESSIVE END-TO-END ASR

Fan Yu<sup>1</sup>, Haoneng Luo<sup>1</sup>, Pengcheng Guo<sup>1</sup>, Yuhao Liang<sup>1</sup>, Zhuoyuan Yao<sup>1</sup>,  
Lei Xie<sup>1\*</sup>, Yingying Gao<sup>2</sup>, Leijing Hou<sup>2</sup>, Shilei Zhang<sup>2</sup>

<sup>1</sup>Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science,  
Northwestern Polytechnical University, Xi'an, China

<sup>2</sup>China Mobile Research Institute

## ABSTRACT

Continuous integrate-and-fire (CIF) based models, which use a soft and monotonic alignment mechanism, have been well applied in non-autoregressive (NAR) speech recognition with competitive performance compared with other NAR methods. However, such an alignment learning strategy may suffer from an erroneous acoustic boundary estimation, severely hindering the convergence speed as well as the system performance. In this paper, we propose a boundary and context aware training approach for CIF based NAR models. Firstly, the connectionist temporal classification (CTC) spike information is utilized to guide the learning of acoustic boundaries in the CIF. Besides, an additional contextual decoder is introduced behind the CIF decoder, aiming to capture the linguistic dependencies within a sentence. Finally, we adopt a recently proposed Conformer architecture to improve the capacity of acoustic modeling. Experiments on the open-source Mandarin AISHELL-1 corpus show that the proposed method achieves a comparable character error rates (CERs) of 4.9% with only 1/24 latency compared with a state-of-the-art autoregressive (AR) Conformer model. Furthermore, when evaluating on an internal 7500 hours Mandarin corpus, our model still outperforms other NAR methods and even reaches the AR Conformer model on a challenging real-world noisy test set.

**Index Terms**— Non-autoregressive, end-to-end speech recognition, continuous integrate-and-fire

## 1. INTRODUCTION

End-to-end (E2E) models have achieved great success in automatic speech recognition (ASR) due to their effectiveness in sequence-to-sequence modeling [1–6]. Compared with the traditional hybrid systems [7, 8], E2E models can not only simplify the model training process but also achieve competitive or even better recognition accuracy. However, most state-of-the-art E2E models follow an autoregressive (AR) fash-

ion and recursively generate an output token based on previously generated tokens and the input sequence. Thus, it will take at least  $L$  iterations to generate an  $L$ -length output sequence, resulting in a complex computation and a large latency with the increment of sequence length. In contrast, the non-autoregressive (NAR) models make a conditional independence assumption among the output tokens and no longer rely on the temporal relationship from left to right, which can generate an  $L$ -length sequence in parallel with a constant  $K$  ( $\ll L$ ) iterations.

NAR models are originally proposed in neural machine translation (NMT) tasks and have been well applied in ASR recently. The major difficulties of NAR models consist of the following two aspects: the accurate length prediction of a target sequence and the parallel inference of decoder. Introducing a length prediction network behind the encoder [9–11] or estimating an empirical target length according to the source sequence [12] are two typical length prediction approaches. However, both methods will bring redundant computations and can not guarantee the accuracy of the predicted lengths, especially in intricate real-world scenarios. Since connectionist temporal classification (CTC) [13] is good at learning frame-wise latent alignments between the input speech and output tokens, more and more studies are trying to get rid of the cumbersome length estimation and focus on incorporating the CTC into NAR models. In [14, 15], conditional chain based methods were proposed for NAR multi-speaker ASR, which inferred the output of each speaker one-by-one using a NAR CTC model. With this design, the total inference steps will be restricted to the number of mixed speakers. In [16], Tian *et al.* proposed a spike-triggered based NAR method, which used the encoded states corresponding to the CTC spikes as the decoder input. Although the decoder can easily perform parallel computation with this design, the length mismatch between the CTC spikes and target sequence may lead to a computation problem of cross-entropy (CE) loss. Inspired by the masked language modeling, Mask-CTC [17, 18] initialized the input target sequence with masked ground-truth during the training and masked token-level CTC outputs

\*Lei Xie is the corresponding author, lxie@nwpu.edu.cn

during the inference, respectively. The idea of Mask-CTC was to infer the masked tokens by the decoder and iteratively refine them. In addition to Mask-CTC, Imputer [19] effectively modeled context dependencies and CASS-NAT [20] generated token-level acoustic embeddings through CTC alignments, which were also confronting the redundant computation problem caused by a longer sequence.

Recently, a novel soft and monotonic alignment mechanism was proposed for sequence transduction, named continuous integrate-and-fire (CIF) [21]. By accumulating the weight of the vector representation in each encoder step, CIF effectively locates the acoustic boundaries and extracts the acoustic information corresponding to each target token. Although it is more accurate in estimating the length of the target sequence, there still exists a deviation in the prediction of the acoustics boundaries. In this paper, in order to improve the accuracy of acoustic boundary prediction (e.g. character boundaries in Mandarin) for the CIF based models, we propose a novel auxiliary objective function to constrain the acoustic boundary of each label by using the spike information and alignment states of CTC. Moreover, a recently proposed Conformer architecture [22] is adopted to enhance the speech representation learning, which has better local and global modeling capabilities than Transformer. Finally, we integrate a new contextual decoder to model the linguistic and contextual dependencies within the target sequences. Evaluating on the open-source Mandarin corpora AISHELL-1, our proposed method achieves a comparable character error rates (CERs) of 4.9% with only 1/24 latency compared with a state-of-the-art AR Conformer model. We also conduct experiments on a large scale 7500 hours internal Mandarin corpus and our model shows similar trend on in-domain test set.

The rest of this paper is organized as follows. Section 2 introduces non-autoregressive continuous integrate-and-fire (CIF) model. Section 3 describes our proposed method. Section 4 presents our experimental setup and results. The conclusions and future work will be given in Section 5.

## 2. NON-AUTOREGRESSIVE CONTINUOUS INTEGRATE-AND-FIRE (CIF)

Given an input speech sequence  $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_T]$ , where  $T$  means the number of frames, conventional AR models iteratively produce output tokens  $\mathbf{y} = [y_1, \dots, y_L]$  as:

$$P_{\text{AR}}(\mathbf{y}|\mathbf{x}) = \prod_{l=1}^L P(y_l|y_{<l}, \mathbf{x}), \quad (1)$$

where  $L$  refers to the length of target sequence. With the previously generated tokens  $y_{<l}$ , AR models estimate the output token  $y_l$  one-by-one, which makes it hard to compute in parallel and results in a large latency during the inference. By contrast, the NAR models aim to get rid of such temporal dependency and perform parallel computation. During the infer-

ence, NAR models will generate the probability distribution of  $\mathbf{y}$  within a constant number of iterations that is not constrained to the sequence length. Mathematically, a NAR model can be defined as:

$$P_{\text{NAR}}(\mathbf{y}|\mathbf{x}) = \prod_{l=1}^L P(y_l|\mathbf{x}). \quad (2)$$

As a NAR model, CIF [21] utilizes a soft and monotonic alignment mechanism between the encoder and decoder, and the encoder outputs, replacing traditional token embeddings, are directly used as the input of decoder to achieve parallel NAR computation. The left part of Fig. 1 shows a detailed structure of the CIF model. At each encoder step, CIF receives the current encoder output  $h_u$  and a corresponding weight  $\alpha_u$ , where  $u \in U$  means the index of encoder step. The weight  $\alpha_u$  scales the amount of acoustic information and is accumulated until reaching a threshold  $\beta = 1$ , which means an acoustic boundary of a specific target token. At this boundary point, the  $h_u$  will be divided into two parts: one for completing the integration of current token, and the other is used for the accumulation of acoustic information for the next token. Then, the weighted summation of encoder outputs is fired as the input of decoder. Therefore, CIF is able to provide a soft alignment between acoustic frames and target labels.

In order to improve the accuracy of sequence length prediction by CIF, a quantity loss is also presented, forcing the model to predict the quantity of integrated embeddings closer to the quantity of targeted label sequence. Quantity loss  $\mathcal{L}_{\text{Qua}}$  can be defined as:

$$\mathcal{L}_{\text{Qua}} = \left| \sum_{u=1}^U \alpha_u - \tilde{S} \right|, \quad (3)$$

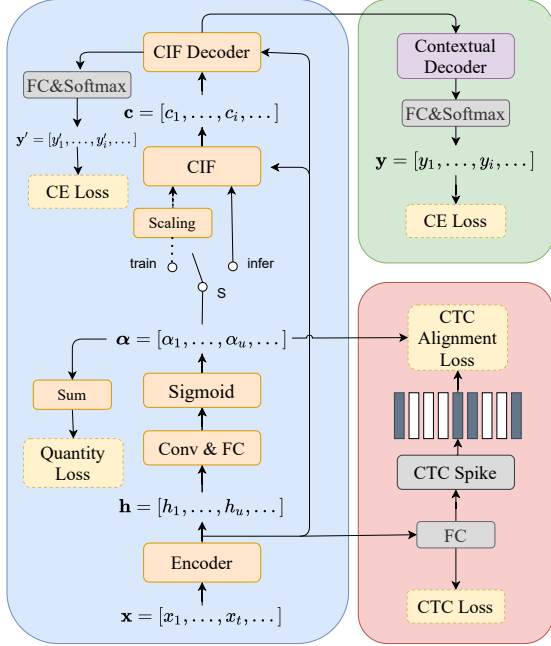
where  $\tilde{S}$  is the ground-truth length of the target sequence  $\mathbf{y}$ .

## 3. PROPOSED METHOD

Our proposed approach is based on the CIF model with substantial improvements to obtain more accurate acoustic boundaries and explicitly capture the linguistic contexts. Specifically, we use a CTC module to guide the learning of acoustic boundaries and introduce a novel context decoder behind the CIF decoder to model the token relationship, as shown in Fig. 1. Besides, since Conformer encoder is proposed to learn both local and global acoustic correlations synchronously, we also adopt it to extract better acoustic hidden representations. Moreover, unlike the original CIF which only uses acoustic embedding of CIF module as in Eq. (4), we fire both the acoustic embedding  $\mathbf{c}$  and encoder output  $\mathbf{h}$  to the decoder to predict the probability of output tokens  $\mathbf{y} = [y_1, \dots, y_i, \dots]$  in parallel as in Eq. (5). We notice that this can lead to a further improvement of performance.

$$P(\mathbf{y}|\mathbf{c}) = \text{Decoder}(\mathbf{c}), \quad (4)$$

$$P(\mathbf{y}|\mathbf{c}, \mathbf{h}) = \text{Decoder}(\mathbf{c}, \mathbf{h}). \quad (5)$$



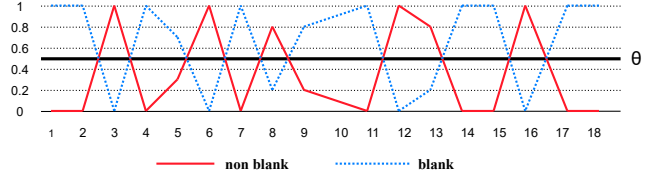
**Fig. 1.** The architecture of our proposed CIF-based non-autoregressive model for the ASR task. The left part illustrates the structure of CIF module, and the right part shows the contextual decoder and CTC alignment loss.

### 3.1. Conformer Encoder

Our encoder block follows the same Conformer architecture as in [22, 23], which includes a multihead self-attention (MHSA) module, a convolution (Conv) module, and a pair of positionwise feed-forward (FFN) module in the Macaron-Net style. Meanwhile the MHSA module can learn long-range global context, the Conv module is good at extracting fine-grained local features. By adopting Conformer, we expect to improve the prediction performance of CTC and CIF modules as well. In addition, incorporating the relatively positional embeddings to MHSA modules further improves the generalization of the model on variable lengths of input sequences.

### 3.2. CTC alignment loss

Due to the speaker’s speaking rate, accent, silence and noise, which often happen in real-world applications, the alignment learning strategy of CIF may result in an inaccurate acoustic boundary estimation, which severely influences the convergence speed. To figure out these problems and improve the performance, we present an additional CTC alignment loss function (as shown in Fig 1) to guide the CIF-based model to predict the acoustic boundary closer to the actual boundary by making full use of the CTC spike information. Since the CTC spike is usually within the acoustic boundary of the label, we can roughly determine the boundary of a label by two consec-



**Fig. 2.** The spike-like trigger probability curve. The red curve represents the non-blank label probability, and the solid black line represents the trigger threshold  $\theta$ . The probability of a non-blank token greater than  $\theta$  is a spike.

utive CTC spikes. The CTC spike module is shown in Fig. 2. Given a CTC spike sequence  $P_s = [0, 0, 1, 0, 0, 1, 0, 1, 0, \dots]$ , where 1 means the probability of non-blank token is greater than a specific threshold  $\theta$ , we constrain the CIF weights  $\alpha$  by the label boundary sequence  $P_b = [0, 2, 5, 7, \dots]$ <sup>1</sup>, where each number is the index of non-blank token. The additional alignment loss  $\mathcal{L}_{Ali}$  can be formulated as:

$$\mathcal{L}_{Ali} = \sum_{i=0}^L \left| \sum_{j=P_b(i)}^{P_b(i+1)} \alpha(j) - 1 \right|. \quad (6)$$

where  $L$  is the ground-truth length of the targets sequence  $\mathbf{y}$ . By providing the alignment label boundary constraints, we combine the boundary information from both CTC spikes and CIF to promote the model to locate boundaries in a parameter-efficient way.

### 3.3. Contextual decoder

To further help the CIF decoder to capture the token relationship, we propose a contextual decoder after CIF decoder to capture the contextual relationship within a sequence. Although the outputs of the CIF decoder can be regarded as an integrated representation of acoustic and linguistic information corresponding to the tokens, the context correlation between tokens is weak. The contextual decoder leverages a stack of self-attention blocks, and its input of the query, key, and value is the high-level representations from the CIF decoder. Meanwhile, contextual decoder does not require the acoustic output representation of encoder since it already contains enough acoustic information of each token. Moreover, considering the deep layers of the model, we calculate the cross-entropy (CE) loss of between the outputs of CIF decoder and contextual-decoder to assist the backward update of the gradient and speed up the convergence. The cross-entropy loss can be formulated as  $\mathcal{L}_{CE} = \mathcal{L}_{CE}(\mathbf{y}') + \mathcal{L}_{CE}(\mathbf{y}'')$ , where  $\mathbf{y}'$  and  $\mathbf{y}''$  indicate CIF decoder output and contextual decoder output respectively. Note that our model with contextual decoder needs to train more epochs or use a pre-trained encoder to speed up the model convergence.

<sup>1</sup>The first element of  $P_b$  is set to 0 for identifying the beginning of the spike sequence.

### 3.4. Training strategy

As described in Section 2, considering the accuracy of sequence length prediction for CIF module and the lack of left-to-right constraints for the attention model, our model adopts the quantity loss and CTC loss. In addition, we introduce a CTC alignment loss  $\mathcal{L}_{\text{ali}}$  to further promote the learning of acoustic boundary as described in Section 3.2. Finally, we also compute the CE loss  $\mathcal{L}_{\text{CE}}$  for the CIF decoder and the contextual decoder to boost the model learning. Therefore, our model is trained with a combination of four different losses:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda_1 \mathcal{L}_{\text{Ali}} + \lambda_2 \mathcal{L}_{\text{CTC}} + \lambda_3 \mathcal{L}_{\text{Qua}}, \quad (7)$$

where  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  are interpolation factors. In this study, these hyper-parameters are set to 1 and the effects of these auxiliary functions will be investigated in the following experiments.

## 4. EXPERIMENTS

### 4.1. Dataset

In this paper, we evaluate our Conformer-CIF based model on two Mandarin speech recognition corpora: an open-source AISHELL-1 corpus [24] and a large-scale internal 7500 hours corpus [25]. The AISHELL-1 corpus contains 150 hours of speech recorded by 340 speakers as training set, 18 hours of speech recorded by 40 speakers as development set, and about 10 hours speech as test set. The speakers of different sets are not overlapped and each set is recorded in a relatively quiet environment. The 7500 hours internal corpus contains reading speech data in various domains, such as entertainment, journalism, literature, technology, free conversation, etc. For experiments on the 7500 hours corpus, we report character error rate results (CERs) on four test sets, namely AISHELL-1 test set (TA1) [24], AISHELL-2 test set (TA2) [26], an internal voice input test set (VI), and a voice assistant (VA) test set. The VI test set consists of about 3.4 hours data with 3063 sentences covering lots of proper nouns and named entities, which is used to verify the language generalization ability of the model. The VA test set is collected under various noisy scenarios, containing about 3.9 hours data with 5000 voice assistant commands.

### 4.2. Experimental Setup

In our work, the 80-dimensional log Mel-filter bank feature (Fbank) plus 3-dimensional pitch feature are used as the input feature. The window size is 25 ms with a shift of 10 ms. We use 4231 characters extracted from the training transcriptions as the modeling units. Our CIF-based NAR model is comprised of 12 encoder layers, 6 CIF decoder layers, and 6 contextual decoder layers. Particularly, when performing experiments on the 7500 hours corpus, the number of encoder

**Table 1.** The character error rates (CERs) of the systems on AISHELL-1. Real-time factor (RTF) is computed as the ratio of the total inference time to the total duration of the test set.

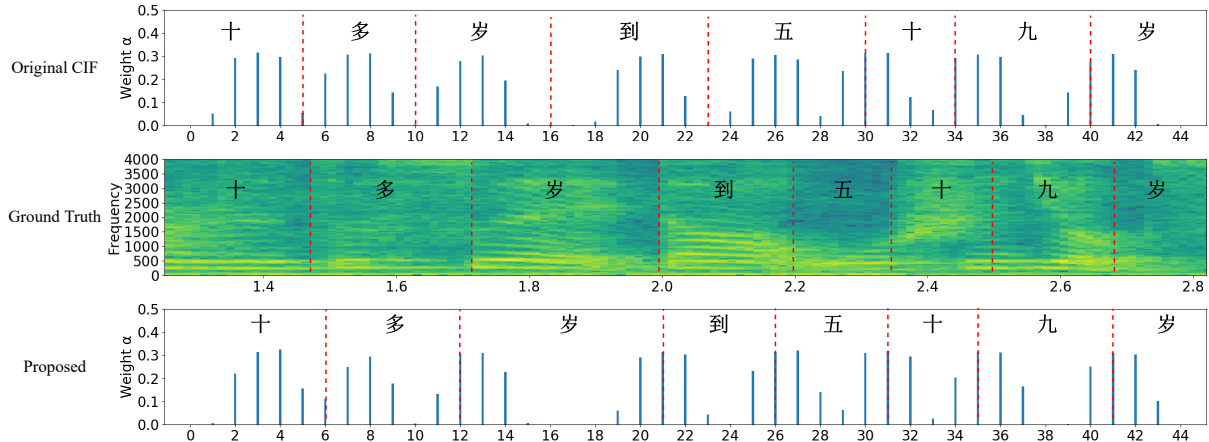
Model	CER (%)		RTF
	Dev	Test	
<i>Autoregressive</i>			
Kaldi chain [24]	-	7.6	-
Conformer	4.8	5.1	0.2768
+ Contextual decoder	4.4	4.8	0.4306
<i>Non-autoregressive</i>			
Transformer-A-FMLM [2]	6.2	6.7	-
Transformer-Insertion [31]	6.1	6.7	-
Transformer-LASO [12]	5.8	6.4	-
Transformer-CASS-NAT [20]	5.3	5.8	-
Conformer-CTC	5.1	5.7	0.0125
Conformer-Spike-Triggered [16] <sup>†</sup>	5.0	5.6	0.0152
<i>Non-autoregressive (proposed)</i>			
Conformer-CIF	4.8	5.3	0.0166
+ Contextual decoder	4.6	5.0	0.0181
+ CTC alignment loss	<b>4.5</b>	<b>4.9</b>	0.0182

<sup>†</sup>: This models is re-implemented by ourselves with the same parameter structure as our model.

layers and contextual decoder layers is set to 20 and 4, respectively. The common parameters of the encoder and decoder layers are:  $d^{\text{head}} = 4$ ,  $d^{\text{attn}} = 256$ ,  $d^{\text{ff}} = 2048$  for the number of attention heads, dimension of attention module, and dimension of feedforward layer, respectively. We use the Adam optimizer [27] and the warmup scheduler [5] to train the model for 80 epochs. The warmup step is set to 25k iterations. SpecAugment [28] and speed perturbation [29] with the factor of 0.0, 1.0, and 1.1 are also applied to avoid overfitting. For a more fair comparison, we re-implement the CTC and Spike-Triggered [16] models with the same parameters and Conformer structure as used in our model. All the experiments are conducted using the open-source ESPnet [30] toolkit on a single NVIDIA RTX 2080Ti GPU.

### 4.3. Results on AISHELL-1

To explore the efficiency of our proposed Conformer-CIF based model, we first evaluate it on the open-source AISHELL-1 corpus. As show in Table 1, we can find that our model outperforms the hybrid Kaldi chain system and even reaches the state-of-the-art AR Conformer model, achieving CERs of 4.8% and 5.3% on dev and test sets, respectively. Comparing with other NAR models, our Conformer-CIF based model is still better than other methods with a slight increment of RTF. When incorporating with the contextual decoder and CTC alignment loss, we can obtain further improvement, achieving CERs of 4.5% and 4.9% with about  $24\times$  faster in inference speech.



**Fig. 3.** Visualization on the prediction of weight  $\alpha$  for each encoded state for the original CIF model (upper) and our CIF model (lower) with contextual decoder and CTC alignment loss. These figures are drawn for the utterance index BAC009S0764W0356 in the AISHELL-1 test set. We find that the original CIF model (upper), as mentioned in [21], is more prone to be located ahead of time, but our CIF model (lower) predicts the boundary of each token more accurately. Red dotted vertical line: acoustic boundary of the token.

#### 4.4. Results on 7500 hours

We further evaluate the proposed model on a large-scale internal 7500 hours corpus, as shown in Table 2. In addition to the clean test sets TA1 and TA2, we also involve two more challenging test sets VI and VA. From the table, we can see a quite similar trend as experiments on the AISHELL-1 corpus. The proposed Conformer-CIF based model outperforms the cascade RNN-T model on 3 test sets and even shows good potential on the most challenging sets VI and VA, where VA is corrupted by real-world noise or reverberation. Comparing with the state-of-the-art AR Conformer model, our model achieves competitive results on TA1 and TA2 tests with only 1/18 latency. Since VI and VA are collected under more complicated acoustic conditions from a mismatched domain, comparing with AR model, the performance of the NAR model is reduced due to lack of context information within the sequences. Moreover, by comparing the results with Conformer-CTC NAR models, our model yields up to 9.5%, 8.0%, 5.8%, and 1.5% relative CERs reductions on four test sets, respectively.

#### 4.5. Result Analysis

The results on Section 4.3 and Section 4.4 have demonstrated that the proposed method significantly improves the performance with a negligible RTF increment. In this Section, we will investigate the effectiveness of the CTC alignment loss and the contextual decoder.

Fig. 3 shows a visualization of the encoded state weight  $\alpha$  for the original CIF model and our proposed Conformer-CIF model with contextual decoder and CTC alignment loss. For the upper and lower parts, the red dash lines are acoustic

**Table 2.** The character error rates (CERs) of the systems on the the 7500-hour corpus. Real-time factor (RTF) is computed as the ratio of the total inference time to the total duration of the test set.

Model	CER(%)				Param.(M)/RTF
	TA1	TA2	VI	VA	
<i>Autoregressive</i>					
Cascade RNN-T [25]	4.6	9.2	8.7	28.1	95.5M / -
Conformer	5.1	8.4	6.7	25.7	71.2M / 0.4487
<i>Non-autoregressive</i>					
Conformer-CTC	6.1	9.7	8.6	27.6	55.8M / 0.0201
Conformer-CIF	<b>5.5</b>	<b>8.9</b>	<b>8.1</b>	<b>27.2</b>	73.6M / 0.0253

boundaries estimated by the models, which means the accumulated weight reaches 1.0 at this point. The middle part is the spectrogram of the audio and the acoustic boundaries are marked manually. We can see that the original CIF model is prone to be located ahead of time, as mentioned in [21], while our model gives a more accurate estimation of the acoustic boundaries. By using the CTC spike information to guide the learning of acoustic boundaries, our model can not only converge faster, but also predict the length of the target sequence correctly, which will dramatically reduce the insertion and deletion errors.

Moreover, we also notice that most of the inference errors are substitution errors with similar pronunciation, as shown in Fig 4. The upper part shows the ground truth sequence, while the middle and lower parts are the output hypotheses obtained by different models. Considering the lack of context dependencies within a sentence, we adopt an additional contextual decoder to improve the linguistic and contextual relationship and aim to recover such substitution errors. From figure 4, we

### Ground truth

重点突破棉花油菜甘蔗收获机械化瓶颈

### CIF-based inference

重点突破棉花油菜干着收获机械化瓶颈

### CIF-based with contextual decoder

重点突破棉花油菜甘蔗收获机械化瓶颈

**Fig. 4.** Decoding example for BAC009S0903W0239 in the AISHELL-1 test set. Red indicates characters with errors and blue indicates ones recovered with a contextual decoder.

can see that the proposed contextual decoder can successfully eliminate substitution errors.

## 5. CONCLUSIONS

In this paper, we propose a boundary and context aware training approach for CIF based non-autoregressive models to improve the accuracy of the acoustic boundary prediction. Our Conformer-CIF based model utilizes the CTC spike information to guide the learning of acoustic boundaries in the CIF, and integrates a new contextual decoder to model the linguistic and contextual dependencies within the target sequence. Meanwhile, we also adopt a recently proposed Conformer architecture to improve the capacity of acoustic modeling. Evaluating on the open-source Mandarin corpora AISHELL-1 and an internal 7500 hours Mandarin corpus, our proposed approach achieves comparable character error rates with only 1/24 latency compared with a state-of-the-art autoregressive Conformer model. In the future, we will integrate the external language model into our proposed model to improve the language generalization ability.

## 6. ACKNOWLEDGEMENT

This work was supported by MoE-CMCC “Artificial Intelligence” Project (MCM20190701).

## 7. REFERENCES

- [1] Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio, “Attention-based models for speech recognition,” in *Proc. NIPS*, 2015, pp. 577–585.
- [2] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *Proc. ICASSP*. IEEE, 2016, pp. 4960–4964.
- [3] Suyoun Kim, Takaaki Hori, and Shinji Watanabe, “Joint CTC-attention based end-to-end speech recognition using multi-task learning,” in *Proc. ICASSP*. IEEE, 2017, pp. 4835–4839.
- [4] Chung-Cheng Chiu, Tara N Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjali Kannan, Ron J Weiss, Kanishka Rao, Ekaterina Gonnina, et al., “State-of-the-art speech recognition with sequence-to-sequence models,” in *Proc. ICASSP*. IEEE, 2018, pp. 4774–4778.
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, et al., “Attention is all you need,” in *Proc. NeurIPS*, 2017, pp. 5998–6008.
- [6] Linhao Dong, Shuang Xu, and Bo Xu, “Speech-Transformer: A no-recurrence sequence-to-sequence model for speech recognition,” in *Proc. ICASSP*. IEEE, 2018, pp. 5884–5888.
- [7] Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur, “A time delay neural network architecture for efficient modeling of long temporal contexts,” in *Proc. INTERSPEECH*. ISCA, 2015, pp. 3214–3218.
- [8] Daniel Povey, Vijayaditya Peddinti, Daniel Galvez, Pegah Ghahremani, Vimal Manohar, Xingyu Na, Yiming Wang, and Sanjeev Khudanpur, “Purely sequence-trained neural networks for asr based on lattice-free mmi,” in *Proc. INTERSPEECH*. ISCA, 2016, pp. 2751–2755.
- [9] Jason Lee, Elman Mansimov, and Kyunghyun Cho, “Deterministic non-autoregressive neural sequence modeling by iterative refinement,” in *Proc. EMNLP*. ACL, 2018, pp. 1173–1182.
- [10] Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK Li, and Richard Socher, “Non-autoregressive neural machine translation,” in *Proc. ICLR*, 2018.
- [11] Xuezhe Ma, Chunting Zhou, Xian Li, Graham Neubig, and Eduard Hovy, “FlowSeq: Non-autoregressive conditional sequence generation with generative flow,” in *Proc. EMNLP-IJCNLP*, 2019, pp. 4273–4283.
- [12] Ye Bai, Jiangyan Yi, Jianhua Tao, Zhengkun Tian, Zhengqi Wen, and Shuai Zhang, “Listen attentively, and spell once: Whole sentence generation via a non-autoregressive architecture for low-latency speech recognition,” in *Proc. INTERSPEECH*. ISCA, 2020, pp. 3381–3385.
- [13] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks,” in *Proc. ICML*. PMLR, 2006, pp. 369–376.
- [14] Jing Shi, Xuankai Chang, Pengcheng Guo, Shinji Watanabe, Yusuke Fujita, Jiaming Xu, Bo Xu, and Lei

- Xie, “Sequence to multi-sequence learning via conditional chain mapping for mixture signals,” in *Proc. NeurIPS*, 2020, pp. 3735–3747.
- [15] Pengcheng Guo, Xuankai Chang, Shinji Watanabe, and Lei Xie, “Multi-speaker ASR combining non-autoregressive Conformer CTC and conditional speaker chain,” in *Proc. INTERSPEECH*. ISCA, 2021.
- [16] Zhengkun Tian, Jiangyan Yi, Jianhua Tao, Ye Bai, Shuai Zhang, and Zhengqi Wen, “Spike-triggered non-autoregressive Transformer for end-to-end speech recognition,” in *Proc. INTERSPEECH*. ISCA, 2020, pp. 5026–5030.
- [17] Yosuke Higuchi, Shinji Watanabe, Nanxin Chen, Tetsuji Ogawa, and Tetsunori Kobayashi, “Mask CTC: Non-autoregressive end-to-end ASR with CTC and mask predict,” in *Proc. Interspeech*. ISCA, 2020, pp. 3655–3659.
- [18] Yosuke Higuchi, Hirofumi Inaguma, Shinji Watanabe, Tetsuji Ogawa, and Tetsunori Kobayashi, “Improved Mask-CTC for non-autoregressive end-to-end ASR,” in *Proc. ICASSP*. IEEE, 2021, pp. 8363–8367.
- [19] William Chan, Chitwan Saharia, Geoffrey Hinton, Mohammad Norouzi, and Navdeep Jaitly, “Imputer: Sequence modelling via imputation and dynamic programming,” in *Proc. ICML*. PMLR, 2020, pp. 1403–1413.
- [20] Ruchao Fan, Wei Chu, Peng Chang, and Jing Xiao, “CASS-NAT: CTC alignment-based single step non-autoregressive Transformer for speech recognition,” in *Proc. ICASSP*. ISCA, 2021, pp. 5874–5878.
- [21] Linhao Dong and Bo Xu, “CIF: Continuous integrate-and-fire for end-to-end speech recognition,” in *Proc. ICASSP*. IEEE, 2020, pp. 6079–6083.
- [22] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, et al., “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. Interspeech*. ISCA, 2020, pp. 5036–5040.
- [23] Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jia-tong Shi, et al., “Recent developments on ESPnet toolkit boosted by Conformer,” in *Proc. ICASSP*. IEEE, 2021, pp. 5874–5878.
- [24] 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 *Proc. O-COCOSDA*. IEEE, 2017, pp. 1–5.
- [25] Xiong Wang, Zhuoyuan Yao, Xian Shi, and Lei Xie, “Cascade RNN-Transducer: Syllable based streaming on-device Mandarin speech recognition with a syllable-to-character converter,” in *Proc. SLT*. IEEE, 2021, pp. 15–21.
- [26] Jiayu Du, Xingyu Na, Xuechen Liu, and Hui Bu, “Aishell-2: Transforming Mandarin ASR research into industrial scale,” *arXiv preprint arXiv:1808.10583*, 2018.
- [27] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *Proc. ICLR*, 2015.
- [28] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Proc. Interspeech*, 2019, pp. 2613–2617.
- [29] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur, “Audio augmentation for speech recognition,” in *Proc. Interspeech*, 2015, pp. 3586–3589.
- [30] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, Yuya Unno, Nelson-Enrique Yalta Soplín, Jahn Heymann, Matthew Wiesner, Nanxin Chen, et al., “ESPnet: End-to-End speech processing toolkit,” in *Proc. Interspeech*, 2018, pp. 2207–2211.
- [31] Yuya Fujita, Shinji Watanabe, Motoi Omachi, and Xuankai Chan, “Insertion-based modeling for end-to-end automatic speech recognition,” in *Proc. Interspeech*, 2020, pp. 3660–3664.