

# CIF-RNNT: STREAMING ASR VIA ACOUSTIC WORD EMBEDDINGS WITH CONTINUOUS INTEGRATE-AND-FIRE AND RNN-TRANSDUCERS



Wen Shen TEO, Yasuhiro MINAMI (UEC)

## Introduction

### Conventional CIF

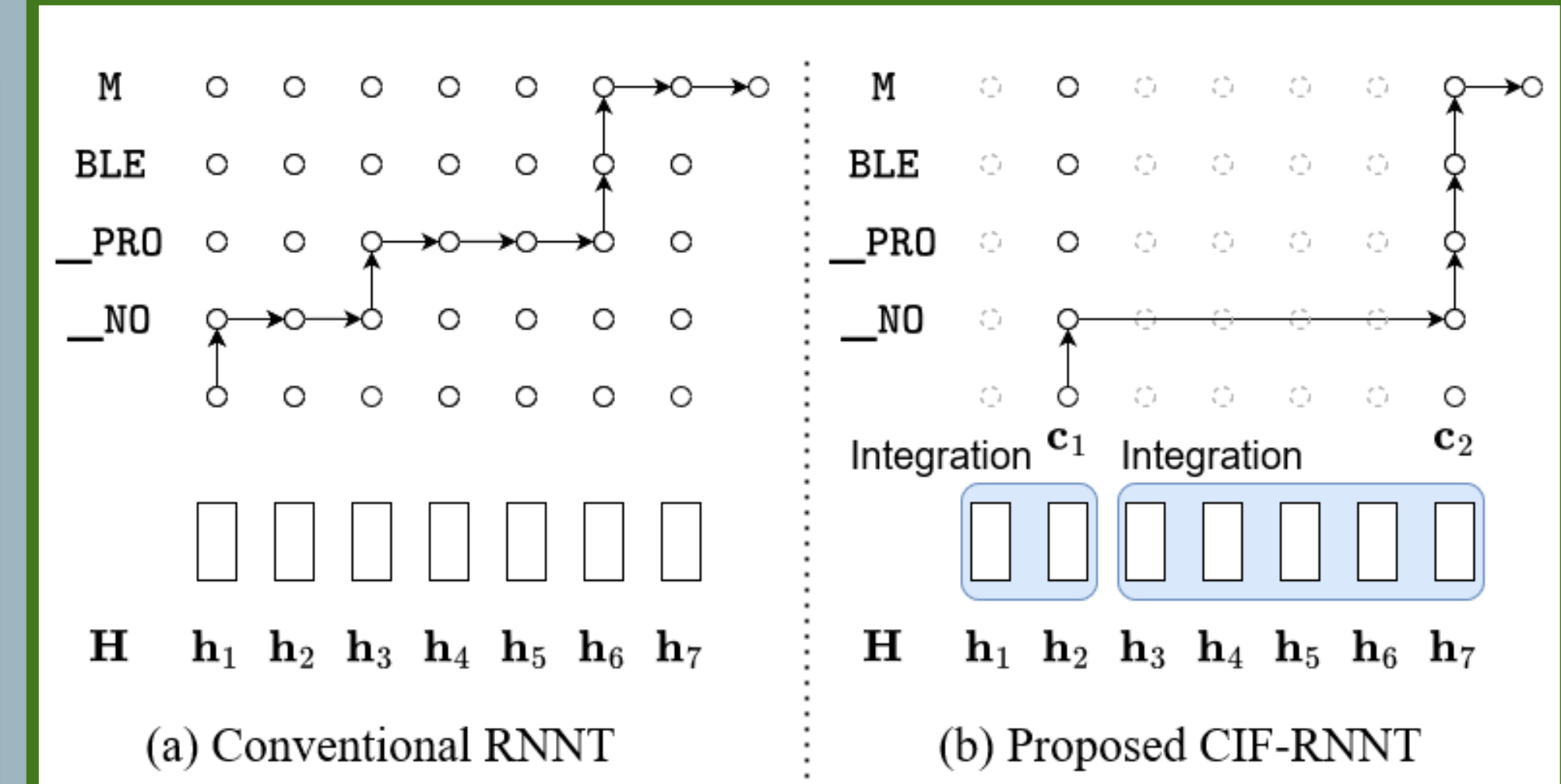
But!!

Bridge between speech and text

Speech and text must share same vocabulary!

### Idea

- CIF : Compress encoder outputs into acoustic word embeddings (AWEs)
- + RNN-T : Spell the word from the AWE using a limited output token set



## Continuous Integrate-and-Fire [1]

Encoder outputs  $H = (h_0, \dots, h_t, \dots, h_{T-1}) \xrightarrow{\text{CIF}} \text{Compressed AWEs } C = (c_0, \dots, c_m, \dots, c_{M-1})$

$$\alpha'_t = \begin{cases} \frac{M^*}{\hat{M}} \alpha_t, & \text{if training} \\ \alpha_t, & \text{otherwise} \end{cases}$$

$M^*$ : target length  
 $\hat{M} = \sum_{t=0}^{T-1} \alpha_t$ : predicted length

$$\mathcal{L}_{qua} = |M^* - \hat{M}|$$

Optimized with Cross Entropy (CE) loss against target output tokens  $Y^*$ .  
 $\rightarrow M^* = |Y^*|$

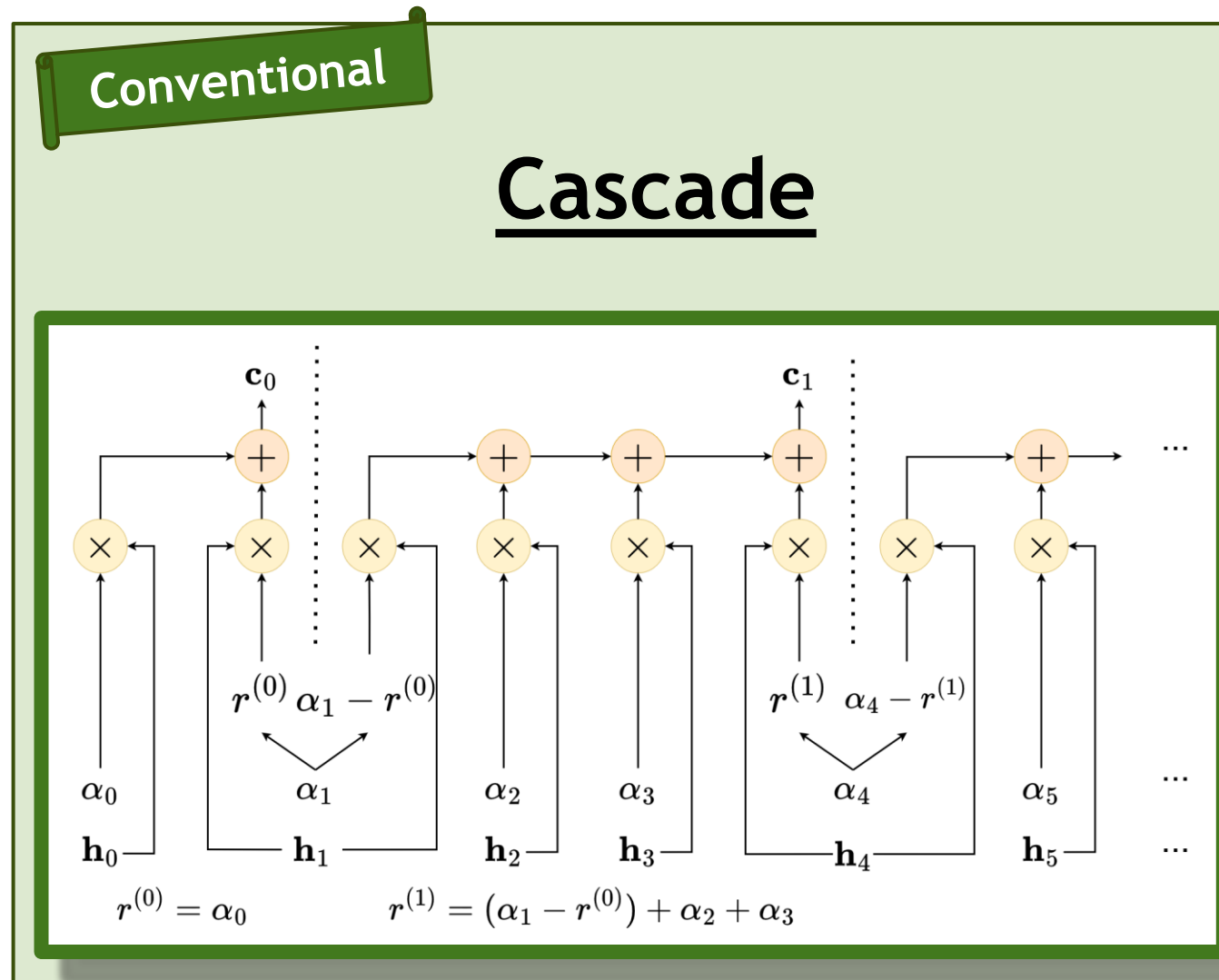
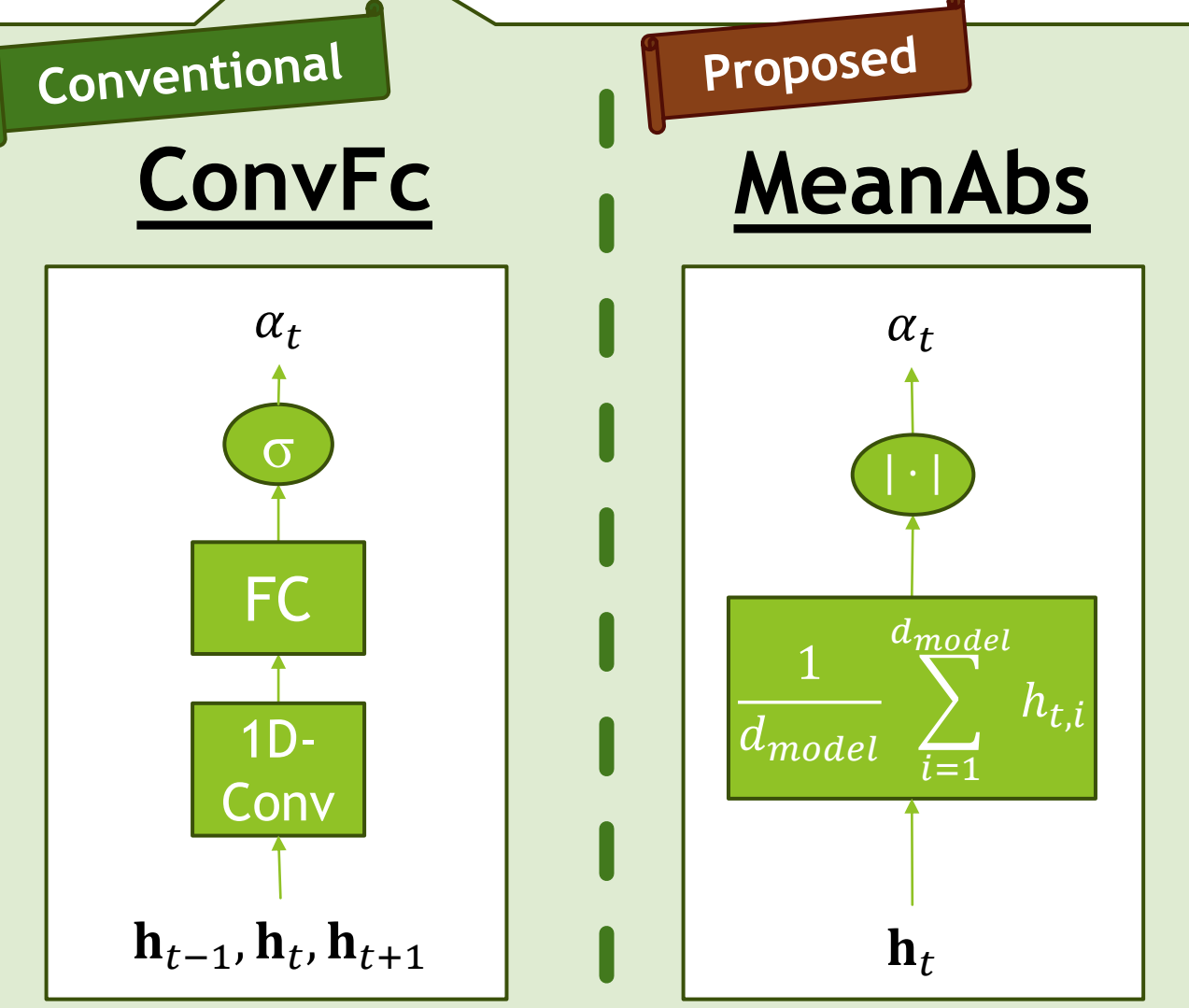
Weight Prediction ( $\omega$ )

Weight Scaling

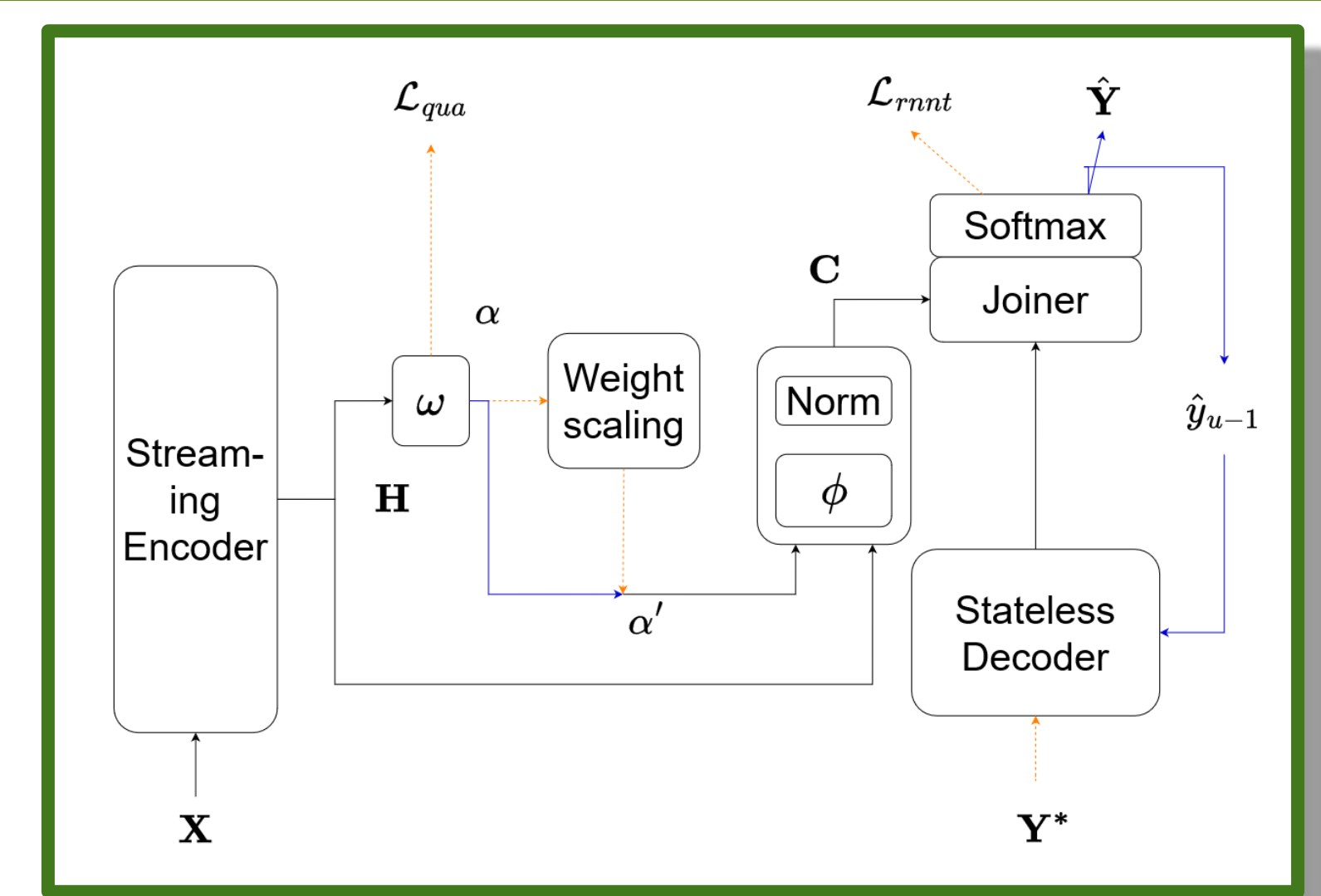
Quantity Loss

Integration ( $\phi$ )

Training Objective



## CIF-RNNT



$$\mathcal{L} = \mathcal{L}_{rnnt} + \mathcal{L}_{qua}$$

Modifications

Detach CE loss from CIF

Insert CIF into RNN-T

Advantages

Detach  $M$  from  $|Y|$

No limit to representable words

Less decoding steps

## Experiment Setup

### Datasets

- LibriSpeech**
- Training data: 'clean-100' and 'clean-360'
  - Output token set: 500 BPE tokens
  - $M^*$ : Number of words
- CSJ**
- Training data: All
  - Output token set: 3266 characters
  - $M^*$ : Number of morphemes
  - Transcript: Fluent

### Setup

Framework: Next Generational Kaldi (NGK)  
Encoder: pruned transducer stateless7 streaming [2] (removed final downsampling layer)  
Decoder: stateless decoder [3] (extended context size to 4)  
Pad: 30 frames at the end of utterances

### Training

Epochs: 40  
RNNT loss: Pruned RNNT loss [4]  
Prune range: (LibriSpeech) 16 (CSJ) 8  
Augmentation: Spec- and noise augmentation

### Decoding

Beam size: 4  
Chunk size: 640ms  
Search method: 'batched beam search'  
Max token per AWE: (LibriSpeech) 9 (CSJ) 5

## Error Rates & RTF

Design		LibriSpeech			CSJ (fluent)			
$\omega$	$\phi$	clean	other	RTF	eval1	eval2	eval3	RTF
CIF-less		4.27	13.06	0.0119	3.99	3.39	3.61	0.0064
ConvFc	Cascade	5.28	14.74	0.0045	4.98	4.22	4.30	0.0026
MeanAbs	Cascade	5.03	14.33	0.0042	4.96	4.25	4.35	0.0025
ConvFc	R.Attn	4.62	13.38	0.0042	4.45	3.75	4.29	0.0032
MeanAbs	R.Attn	4.41	13.09	0.0046	4.50	3.74	4.02	0.0028

Note:  
1. Among combinations of averaged checkpoints from previous epochs up to epoch 40, models with the best (W/C)ERs on the validation set were selected to evaluate the test sets.  
2. "CIF-less" refers to the model trained with the original "pruned transducer stateless7 streaming" recipe.

## $\phi$ -WER

### Definition

- Measures capability to locate word boundaries
- Splits output of greedy path per  $\phi$  for evaluation

(Relaxed)

- To isolate identified but wrongly located  $\phi$ s
- Splits word with "\_" if lone  $\phi$  is encountered

Ref	Hyp
_THERE $\phi$ _IS $\phi$ _NO $\phi$ _OPENING $\phi$ _FOR $\phi$ _YOU $\phi$	_THERE _IS $\phi$ $\phi$ _NO $\phi$ _OPEN $\phi$ ING $\phi$ _FOR _YOU $\phi$
( _THERE $\phi$ $\rightarrow$ _THERE _IS $\phi$ ) ( _IS $\phi$ $\rightarrow$ $\phi$ ) _NO $\phi$ ( _OPENING $\phi$ $\rightarrow$ _OPEN $\phi$ ) (* $\rightarrow$ ING $\phi$ ) ( _FOR $\phi$ $\rightarrow$ _FOR _YOU $\phi$ ) ( _YOU $\phi$ $\rightarrow$ *)	( _THERE $\phi$ _IS $\phi$ _NO $\phi$ ( _OPENING $\phi$ $\rightarrow$ _OPEN $\phi$ ) (* $\rightarrow$ ING $\phi$ ) ( _FOR $\phi$ $\rightarrow$ _FOR _YOU $\phi$ ) ( _YOU $\phi$ $\rightarrow$ *)

### Results

- Only LibriSpeech
- Because Japanese lacks whitespace

Design		clean		other	
$\omega$	$\phi$	Strict	Relax	Strict	Relax
ConvFc	Cascade	24.13	12.55	37.33	24.21
MeanAbs	Cascade	23.26	13.02	35.56	23.32
ConvFc	R.Attn	21.99	8.22	32.22	17.99
MeanAbs	R.Attn	30.20	7.92	38.55	17.19

## Conclusion

### CIF-RNNT: Incorporated CIF into RNN-Ts

- ✓ Streamingly compressed acoustic information into meaningful word units (AWEs)
- ✓ Sped up decoding by reducing decoding operations.
- ✓ Minimized accuracy degradation with novel CIF mechanisms.

## References

- [1] Linhao Dong and Bo Xu, "CIF: Continuous Integrate-and-Fire for End-to-End Speech Recognition", ICASSP 2020.
- [2] Daniel Povey, et. al., [https://github.com/k2-fsa/icefall/egs/librispeech/ASR/pruned\\_transducer\\_stateless7\\_streaming/zipformer.py](https://github.com/k2-fsa/icefall/egs/librispeech/ASR/pruned_transducer_stateless7_streaming/zipformer.py)
- [3] Mohammadreza Ghodsi, et. al., "RNN Transducers with Stateless Prediction Network", ICASSP 2020.
- [4] Fangjun Kuang, et. al., "Pruned RNN-T for Fast Memory-Efficient ASR Training", Interspeech 2022.

teouenshen@gmail.com

teowenshen

Teo Wen Shen

Read my paper here!



Let's connect on LinkedIn!