

DEFORMER: COUPLING DEFORMED LOCALIZED PATTERNS WITH GLOBAL CONTEXT FOR ROBUST END-TO-END SPEECH RECOGNITION



INTERSPEECH 2022

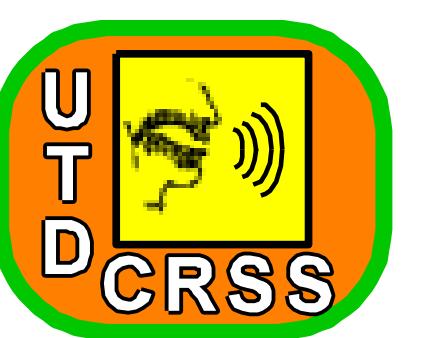
September 18 - 22 • Incheon Korea



International Speech Communication Association

Jiamin Xie and John H. L. Hansen

Jiamin.Xie@utdallas.edu, John.Hansen@utdallas.edu



Center for Robust Speech Systems (CRSS)

Erik Jonsson School of Engineering & Computer Science
The University of Texas at Dallas
Richardson, Texas 75080-3021, U.S.A.



INTERSPEECH 2022 Sept. 18-22, 2022 Incheon, Korea

Introduction

Background

- Convolutional neural networks (CNN) have improved speech recognition greatly by exploiting localized T-F patterns
- Patterns are assumed to exist in a **rigid** and **symmetric** kernel
- What about **asymmetric** kernels? How will the localized pattern learned in this way interact with each other in a global context?

Proposition

- Use the **Deformable CNN (DCNN)** to replace regular CNNs
- Analyze localized patterns obtained and its global interaction by modifying the popular Conformer [1] architecture
- Experiment different initialization methods for the DCNN

[1] A.Gulati, J.Qin, C.-C.Chiu, N.Parmar, Y.Zhang, J.Yu, W.Han, S.Wang, Z.Zhang, Y.Wu et al., "Conformer: Convolution-augmented transformer for speech recognition," *Interspeech 2020*.

INTERSPECH 2022 September 18 - 22 • Incheon Korea

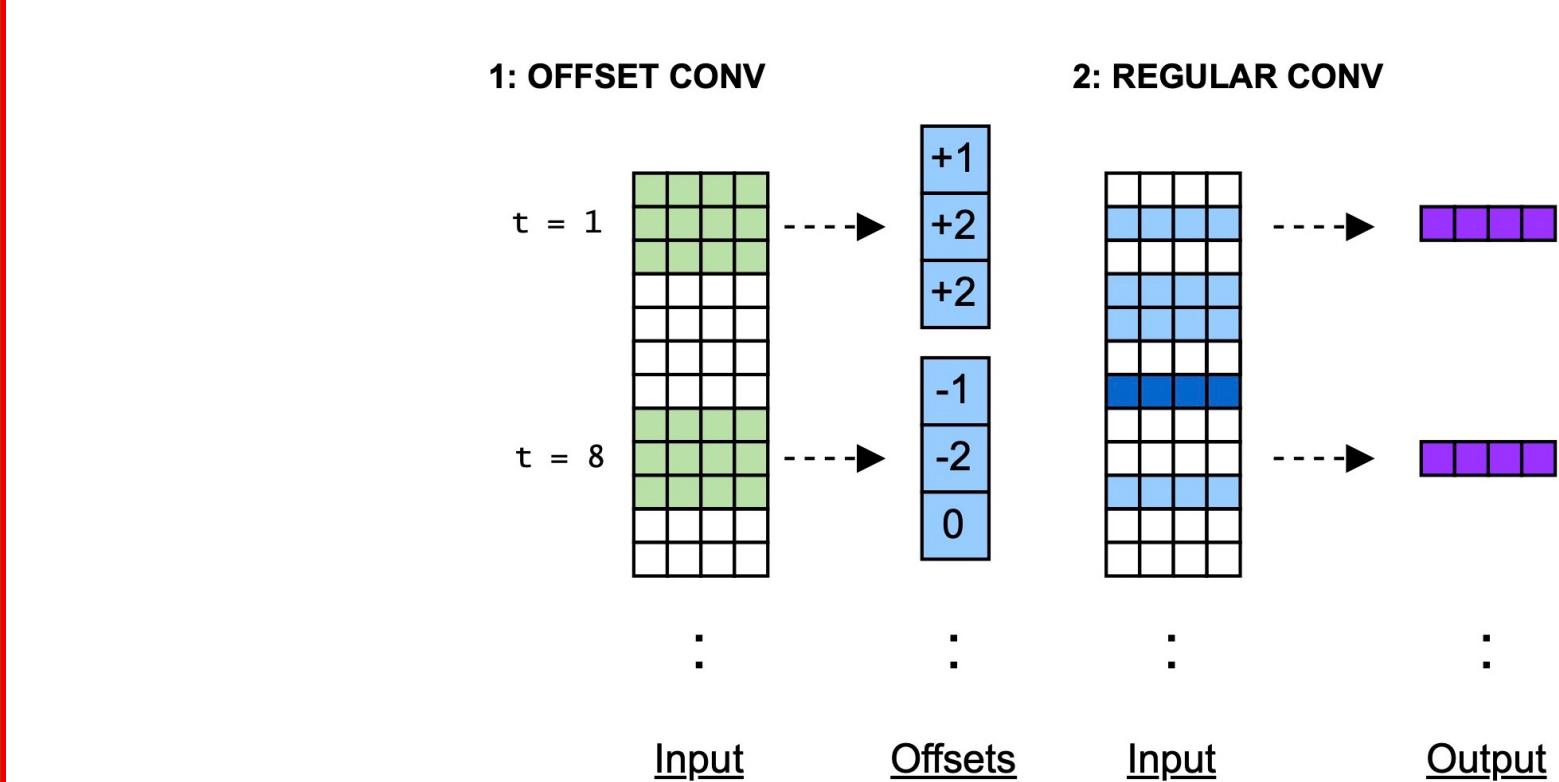
Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 1 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Motivation

Learnable kernel by the DCNN

- The DCNN consists of two regular convolution steps to compute
- One to compute the kernel offsets from the input, another for the output
- The learned positions can overlap and focus on certain input region



INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 2 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Proposed method

Encoder using 1-D DCNN

- We replace the 1D depth-wise convolution in the convolution module of the Conformer by its deformable variant

Formulation

$$\Delta p = \text{Conv}1D_{\text{offset}}(X, p_0) \quad X(p') = X([p']) * ([p'] - p' + 1) \\ p' = p_0 + \Delta p \quad X([p']) = X([p']) + X([p'] + 1) * (p' - [p']) \\ Y = \text{Conv}1D_{\text{output}}(X, p')$$

[2] K. An, Y. Zhang, and Z. Ou, "Deformable tdnn with adaptive receptive fields for speech recognition," *Interspeech 2021*.

INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 3 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Proposed method

Formulation cont.

- 1-D deformable depth-wise convolution

$$X \in \mathbb{R}^{T \times F} \quad Y \in \mathbb{R}^{T \times N} \\ \downarrow \text{Slice} \quad \uparrow \text{Concatenate} \\ X_g \in \mathbb{R}^{T \times F/g} \quad \text{DCNN} \quad Y_g \in \mathbb{R}^{T \times N/g}$$

- Which convolution step to be made depth-wise or both?

- Intuitive to make Conv1D_{output} depth-wise*
- Debatable to make Conv1D_{offset} depth-wise*

Proposed method

Deformer architecture

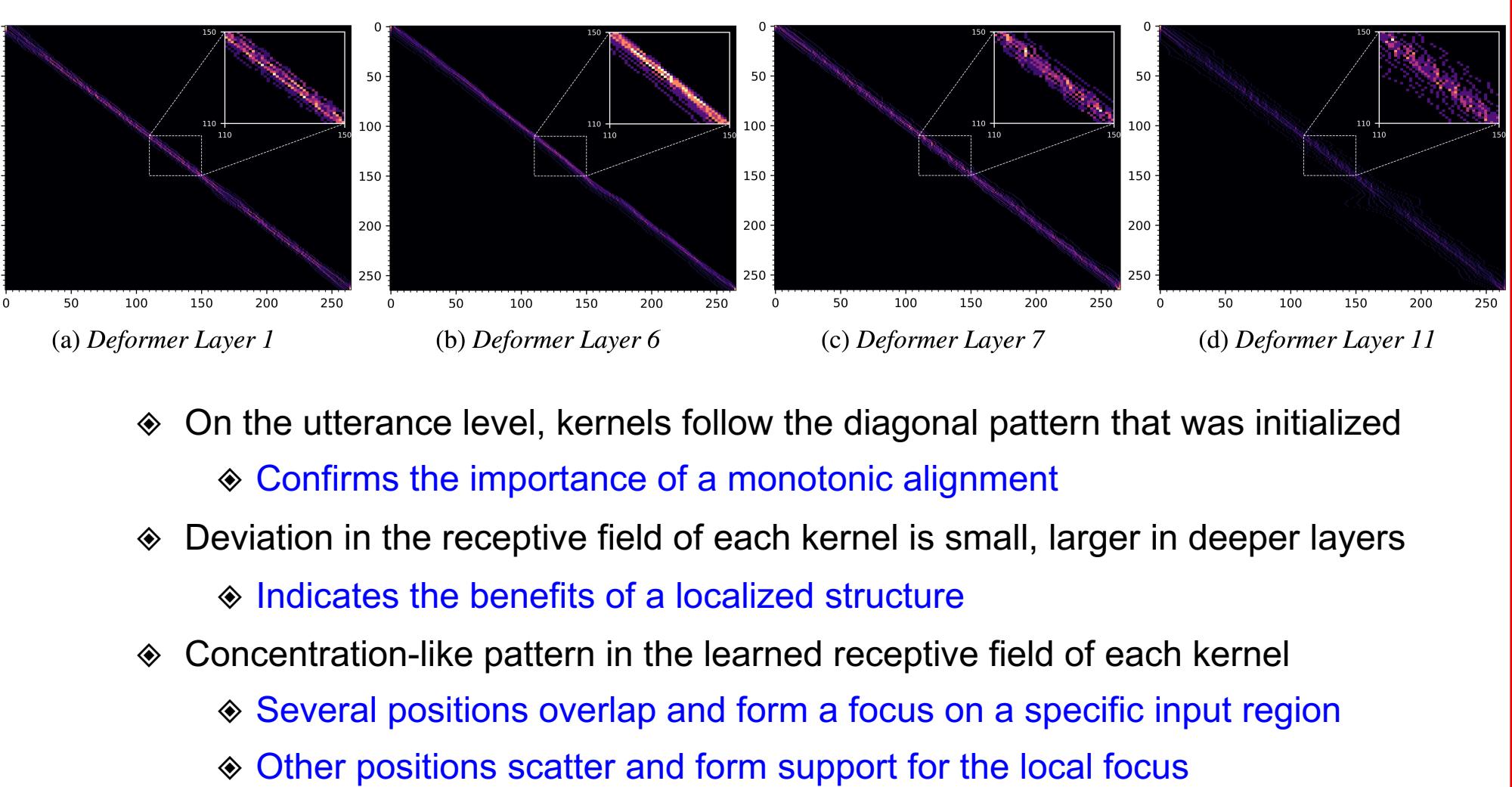
- 12-layer Deformer encoder and 6-layer Transformer decoder
- A mix of deformable and non-deformable layers works the best
- Deformation in the early layer can help as well

Configurations	Deformable Layers	Non-deformable Layers
Layer Index	{1, 6, 7, 10, 11}	{0, 2-5, 8, 9}
Layer Dimensions	256	256
Attention Heads	4	4
Kernel Size	15	15
Dilation	1	1
Stride	1	1
Convolution Groups	256	256
Deformable Groups	1	1

Table 1: Deformer Encoder Configuration

Pattern Analysis

Learned localized patterns



- On the utterance level, kernels follow the diagonal pattern that was initialized
- Confirms the importance of a monotonic alignment*
- Deviation in the receptive field of each kernel is small, larger in deeper layers
- Indicates the benefits of a localized structure*
- Concentration-like pattern in the learned receptive field of each kernel
- Several positions overlap and form a focus on a specific input region
- Other positions scatter and form support for the local focus

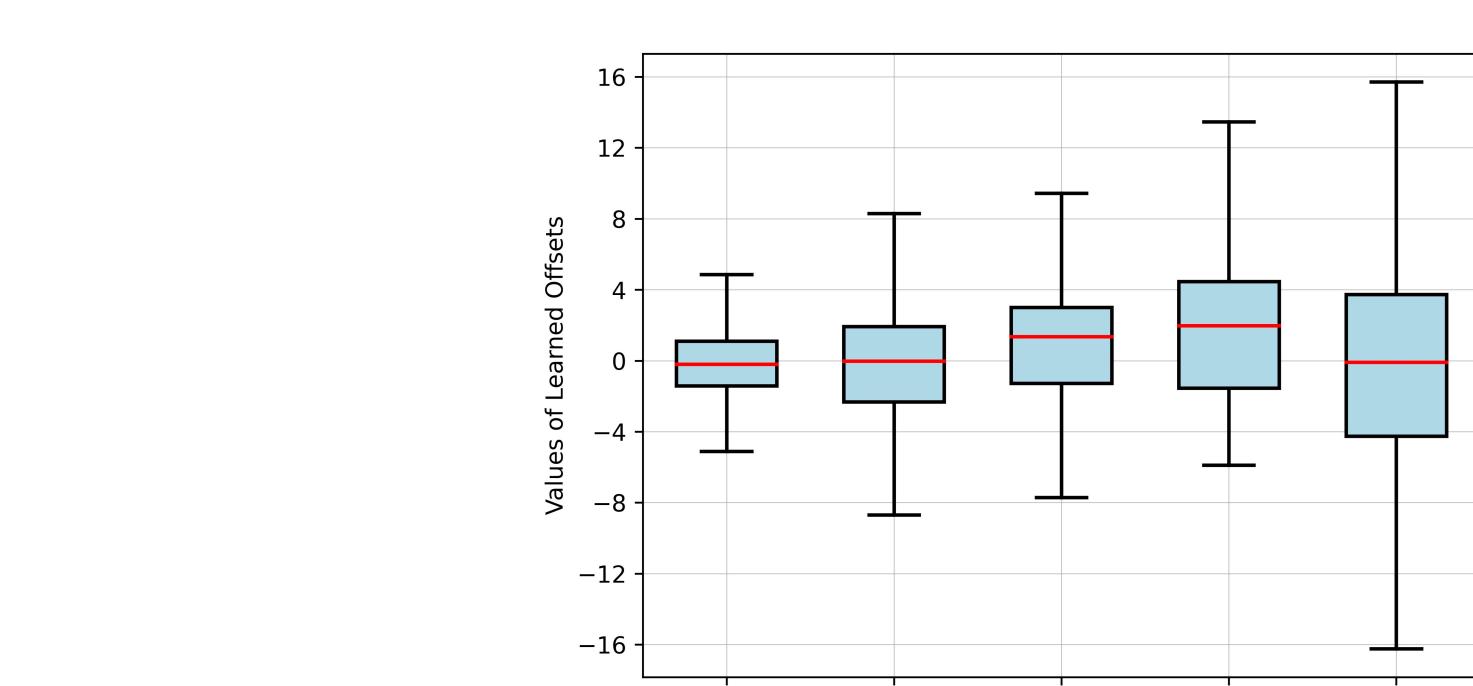
INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 6 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Pattern Analysis

Localized offset statistics



- Increasing spread of distribution over layers
- Deeper layers learn both larger offsets and receptive fields for the kernel
- The distribution has a long tail and is symmetric around zero median
- Verifies a persisting concentration-like pattern even in different utterances*

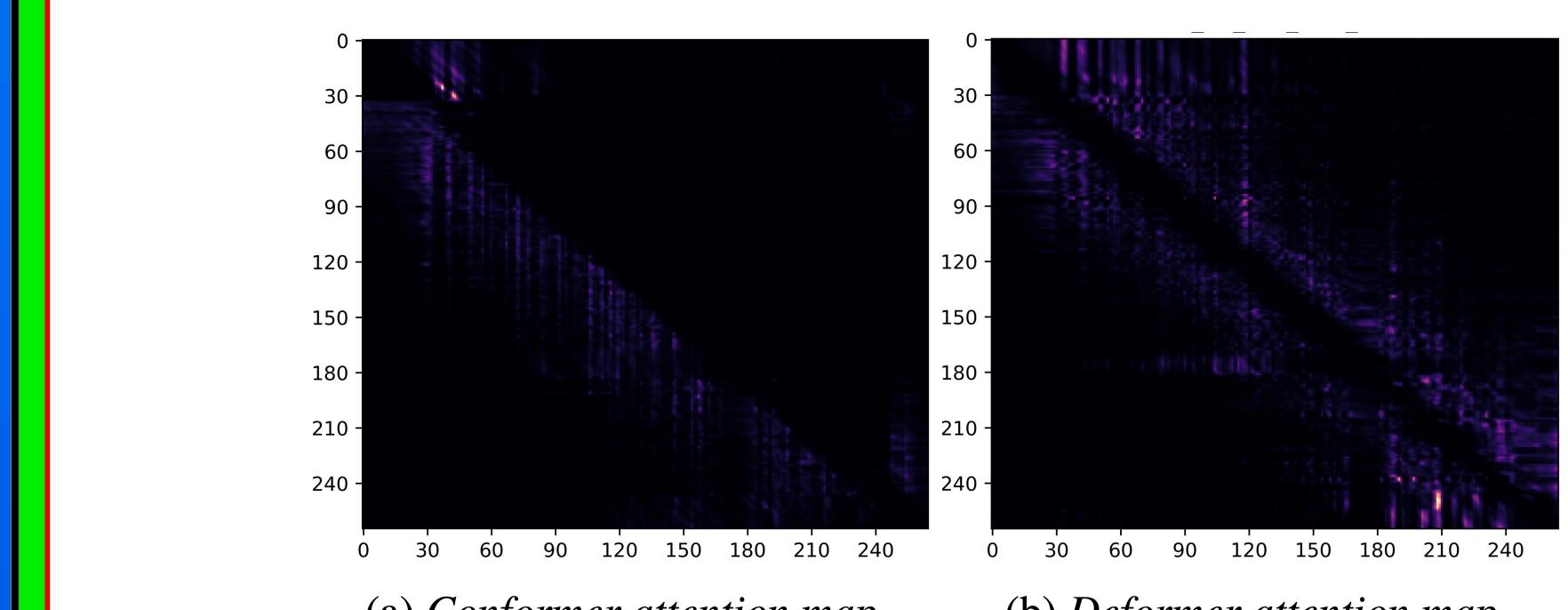
INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 7 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Pattern Analysis

Global pattern (interaction within the utterance context)



- Compared both attention maps of a single head in the last layer
- The Deformer attends additionally to the future above the diagonal
- Indicates deformation brings more relevance among features, but is structured and kept in a small context

[3] S.-w. Yang, A. T. Liu, and H.-y. Lee, "Understanding self-attention of self-supervised audio transformers," *Interspeech 2020*.

INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 8 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Pattern Analysis

Quantitative evaluation of pattern interaction

Model	Globalness	Verticity	Diagonality
Conformer	(0, 7.48]	[-7.48, 0]	[-0.75, 0]
Deformer	4.57	-6.95	-0.18

- Evaluate attention heads of the Deformer using three metrics from [3]
- The globalness increased by +3.3% with a slight decrease of 0.3% in verticity and 1.3% in diagonality, within each respective range

- Verifies a boost in global relevance among features
- Indicates the boost retains the original attention structure, for example, a vertical or a diagonal structure

[3] S.-w. Yang, A. T. Liu, and H.-y. Lee, "Understanding self-attention of self-supervised audio transformers," *Interspeech 2020*.

Experimental Setup

Front-end

- Fbank+Pitch features

Data split

- WSJ (train: si284, dev: dev93, eval: eval92)

Systems

- Baseline: 12-layer Conformer encoder + 6-layer Transformer decoder
- Model: 12-layer Deformer encoder + 6-layer Transformer decoder
- The training setup follows the exact recipe in the Esptnet [5]

Experiments

- Initializing DCNN with Xavier random weights or weights of zeros
- Varying 0.5x or 1x learning rates for the DCNN compared to the rest of the network
- Changing the number of deformable groups to verify the depthwise computation

[4] D. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E.D. Cubuk, and Q.V. Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," *arXiv:1904.08779*, 2019.

[5] S.Watanabe, T.Hori, S.Karita, T.Hayashi, J.Nishitoba, Y.Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen et al., "Esptnet: End-to-end speech processing toolkit," *arXiv:1804.00015*, 2018.

[6] INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 10 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Results

Parameter initialization

Model	#Params	Initialization	WER (%) dev eval	CER (%) dev eval
Conformer Base	43.05M	Xavier	11.2 8.9	3.9 3.0
Deformer (mult=0.5)	43.34M	Xavier	10.7 8.4	3.7 2.8
Deformer (mult=1.0)	43.34M	Xavier	11.1 9.0	3.9 3.0
Deformer (mult=0.5)	43.34M	Zero	10.8 8.8	3.7 3.0
Deformer (mult=1.0)	43.34M	Zero	10.5 8.4	3.7 2.9

- Effectively initializing offsets from zero performs the best overall and gives a more calibrated system
- Improves the Conformer baseline by a +5.6% relative
- Different learning rate multiplier yield results differently, but the fluctuation is smaller with zero intialization

INTERSPECH 2022 September 18 - 22 • Incheon Korea

Email: (Jiamin.Xie, John.Hansen)@utdallas.edu

Slide 11 ISCA Interspeech 2022, Incheon, Korea, Sept. 18-22, 2022

Results

Deformable groups

Model	Deformable Groups	WER (%) dev eval	CER (%) dev eval
Deformer (mult=0.5, zero init.)	256	11.2 9.1	3.9 3.0
Deformer (mult=1.0, zero init.)	256	11.1 8.9	3.9 3.0
Deformer (mult=0.5, zero init.)	2		