

FAST CONFORMER WITH LINEARLY SCALABLE ATTENTION FOR EFFICIENT SPEECH RECOGNITION

Dima Rekesh¹, Nithin Rao Koluguri¹, Samuel Kriman¹, Somshubra Majumdar¹, Vahid Noroozi¹, He Huang¹, Oleksii Hrinchuk¹, Krishna Puvvada¹, Ankur Kumar², Jagadeesh Balam¹, Boris Ginsburg¹

¹NVIDIA, USA,

²Department of Computer Science, University of California, Los Angeles

ABSTRACT

Conformer-based models have become the dominant end-to-end architecture for speech processing tasks. With the objective of enhancing the conformer architecture for efficient training and inference, we carefully redesigned Conformer with a novel downsampling schema. The proposed model, named *Fast Conformer(FC)*, is $2.8\times$ faster than the original Conformer, supports scaling to Billion parameters without any changes to the core architecture and also achieves state-of-the-art accuracy on Automatic Speech Recognition benchmarks. To enable transcription of long-form speech up to 11 hours, we replaced global attention with limited context attention post-training, while also improving accuracy through fine-tuning with the addition of a global token. Fast Conformer, when combined with a Transformer decoder also outperforms the original Conformer in accuracy and in speed for Speech Translation and Spoken Language Understanding.

Index Terms: speech recognition, speech translation, spoken language understanding

1. INTRODUCTION

Conformer is a Transducer (RNNT) model for automatic speech recognition (ASR) proposed by Gulati et al [1]. Conformer models obtain state-of-the-art results on multiple speech benchmarks[2] thanks to their encoder architecture which combines depth-wise convolutional layer for local features and self-attention layer for global context. Conformers have also been rapidly adopted in industry, especially for streaming ASR on-device and in the cloud. ¹ At the same time, Conformer models use more compute and memory than convolution-only ASR models, e.g. Quartznet [3], because self-attention layers have quadratic time and memory complexity vs. input sequence length. The quadratic complexity imposes a severe limitation on the maximum audio length which can be processed by Conformer. Scaling Conformer models require modifying conv kernel sizes[4] in Conformer

¹F. Beaufays, “Google Cloud launches new models for more accurate Speech AI”, 2022. C. Barnes, “Run Google Cloud Speech AI locally, no internet connection required”, 2022

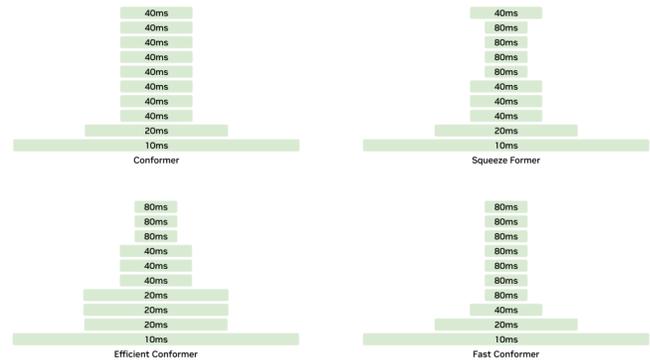


Fig. 1. The downsampling schemas for Conformer, SqueezeFormer, Efficient Conformer and Fast Conformer. Fast Conformer increases sampling rate from 10ms to 80 ms using 3 depth-wise convolutional sub-sampling layers. The additional 2x reduction in the encoder output length versus Conformer yields further compute-memory savings in RNNT decoder.

blocks to stabilize large model training. We redesign Conformer to address these challenges. Namely:

1. We redesign the downsampling schema and sub-sampling block to increase downsampling to 8x (see Fig.1)
2. We optionally replace original self-attention layers with a combination of local attention and global context token post-training, similarly to LongFormer [5].

The proposed encoder has $2.9\times$ fewer multiply-add operations with global attention, and can be made to scale linearly with sequence length post-training. It is $2.8\times$ faster than an equivalent Conformer during inference. It can scale to 1 Billion parameters without any changes to the core architecture. We hereby refer to this model as **Fast Conformer**. At the same time, Fast Conformer maintains highly competitive word error rates (WER) on ASR benchmarks.

We experiment with a Longformer-based attention mechanism, and find that by using limited context attention in combination with a single global attention token we are able

to achieve good results on long-form ASR, while being more than $3\times$ faster during inference. We then test Fast Conformer model on two additional speech processing tasks: speech translation (ST) and speech language understanding (SLU). On ST, with transformer and transducer decoders, we obtain strong scores for En-De translation with significant speedup compared to Conformer. On SLU we obtain state-of-the-art results on the Speech Intent Classification and Slot Filling task, and discuss the relatively modest speedup compared to the Conformer baseline. The models and training recipes have been open sourced through NVIDIA NeMo.²

2. FAST CONFORMER ARCHITECTURE

2.1. Downsampling schema

Conformer encoder is comprised of a stack of alternating multi-head attention [6], depth-wise separable convolutional [7] and fully connected layers with residual connections. The encoder starts with a sub-sampling module, which increases the frame rate from 10 ms to 40 ms. The 4x decrease in the sequence length helps reduce computational and memory costs of the attention layers in all following blocks. This sub-sampling module is relatively expensive, accounting for over 20 % of the computation time for each forward pass of the model for the "Large" Conformer (120 M parameters) [8].

A straightforward way to accelerate Conformer is to increase downsampling rate from 4x to 8x. For example, *EfficientConformer* [9] reduces the sequence length of the speech features by 8x using progressive downsampling: first, 2x down-sampling in the first layer, then another 2x in the middle of encoder, and final 2x in the last encoder layer. One of the drawbacks of progressive subsampling is a computation imbalance between different attention layers. The initial attention layers work on much longer sequences, so they have 16x more computational cost as compared to final attention layers that operate on much shorter sequences (see Fig. 1). Squeezeformer [8] combines progressive downsampling with *Temporal U-Net structure*. Squeezeformer adds extra downsampling at the middle of encoder, and an upsampling layer at the end of encoder to recover 4x time resolution (see Fig. 1). A similar strategy was used in Uconv-Conformer [10].

One of the reasons why previous works limit final encoder downsampling by 4x is related to the usage of the Conformer encoder with a CTC loss [11]. CTC requires that the input to the loss function must be longer than the target sequence length.³ Encoder output may become too short after 8x downsampling if the model uses character tokenization. For example, we found that most of training samples in Librispeech

²<https://github.com/NVIDIA/NeMo>

³Note that this constraint does not apply to the RNNT loss [12], and we are free to use any tokenization scheme as necessary. However, as the RNNT decoder is autoregressive, it is far more efficient at reducing the required number of calls to the Transducer Decoder and Joint by using large sub-word vocabularies.

Table 1. Downsampling schemas and subsampling layer type for Conformer, EfficientConformer, Squeezeformer, and Fast Conformer. **K** is kernel size in convolutional filters.

Model	Subsampling schema	Type	K
Conformer[2]	2/4	2D Conv	31
Squeezeformer[8]	progressive 2/4/8/4	Depth-wise sep	31
Eff. Conformer[9]	progressive 2/4/8	Depth-wise sep	15
Fast Conformer	2/4/8	Depth-wise sep	9

(LS) [13] will not satisfy CTC condition after 8X subsampling if we use character tokenization. To enable 8x downsampling for Conformer-CTC, we switch from character tokenization to Sentencepiece Byte Pair Encoding (BPE) [14] with vocabulary sizes ranging from 128 to 1024 tokens.

To accelerate Conformer, we made the following novel changes in the original design:

1. 8x downsampling at the start of the encoder, thereby reducing the compute cost of subsequent attention layers by 4x
2. Replacement of the original convolution sub-sampling layers with depthwise separable convolutions [15]
3. Reduction of the number of convolutional filters in the downsampling block to 256
4. Reduction of the convolutional kernel size to 9.

A detailed comparison of Fast Conformer downsampling with previous works is presented in Table 1.

To determine the contribution of each change to the model accuracy, we took Conformer-RNNT Large (115 M parameters) as a baseline and gradually applied each design change. First, we added another 2x convolutional subsampling layer. Next, we used depthwise-separable convolutions in the second and third subsampling layers. Then we reduced the number of channels in the subsampling layers from 512 to 256. Finally, we reduced the convolutional kernel size in the conformer blocks from 31 to 9. Encoder inference speed was measured with a batch size of 128 on A100/80G GPU using 20s speech samples. Results are shown in Table 2. The encoder speed increased 2.8x, while maintaining model accuracy.

2.2. Long-form audio transcription

While standard multi-head attention layers have been successful in processing short utterances, their scalability to long sequences has been limited due to the quadratic scaling of the self-attention operation with sequence length. For example, Conformer can process at once maximum 15 minutes audio on a single A100 GPU. To address this challenge, several alternative approaches have been explored. A common approach is buffered transcription, where the audio sequence is

Table 2. Accuracy vs speed for each component of Fast Conformer downsampling schema. Models were tested on LibriSpeech test-other incrementally for each modification starting from the original Conformer. Encoder inference speed (samples/sec) was measured with batch size 128 on 20 sec audio samples. The number of parameters (M) is shown for the encoder only.

Encoder	WER, % test-other	Inference, samples/s	Params, M	GMACS
Baseline Conformer	5.19	169	115	143.2
+8X Stride	5.11	303	115	92.5
+Depthwise conv	5.12	344	111	53.2
+256 channels	5.09	397	109	48.8
+Kernel 9	4.99	467	109	48.7

divided into shorter chunks, which are then transcribed separately before being merged to form a complete transcription. Efficient-Conformer used *grouped attention* to reduce the cost of early attention layers on long sequences by grouping neighboring time elements of the sequence along the feature dimension before applying scaled dot-product attention.

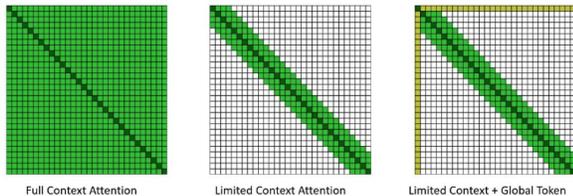


Fig. 2. Fast Conformer combines local attention with a global context token to process long speech samples. The figure is adapted from Longformer [5]

We decided to use an alternative approach inspired by Longformer [5] - local attention augmented with global tokens. We use a single global attention token, which attends to all other tokens and has all other tokens attend to it, using a separate set of query, key and value linear projections, while others attend in a fixed-size window surrounding each token (see Fig. 2). By switching to limited context attention, we extend the maximum duration that the model can process at once on a single A100 GPU by 45x: from 15 minutes for unmodified Conformer to 675 minutes for Fast Conformer with limited context (see Table 3). In order to efficiently compute this attention, we utilize the overlapping chunks approach introduced in Longformer.

3. EXPERIMENTS

3.1. Automatic Speech Recognition

The Fast Conformer model was evaluated on following English ASR benchmarks: LibriSpeech (LS)[13], English part

Table 3. Maximum audio duration (min), which can be transcribed by model with batch size 1, A100 GPU.

Model	Max duration, min
Conformer	15
Fast Conformer	25
Conformer + Limited Context	135
Fast Conf + Limited Context	675

of Multilingual LibriSpeech (MLS) [16], Mozilla Common Voice (MCV) [17], and Wall Street Journal (WSJ) [18]. We used Fast Conformer-RNNT and Fast Conformer-CTC along with baseline Conformer models in Large configuration.

3.1.1. LibriSpeech

First, we trained all models on the LibriSpeech dataset only. We used the SentencePiece unigram tokenizer with 128 tokens for CTC, and 1024 tokens for RNNT models. Training of baseline Conformer-RNNT and -CTC models, was done with AdamW optimizer with the Noam learning rate scheduler and peak learning rates of 0.0025 and 0.001 respectively. We set a linear warmup schedule for 15k steps in all experiments. Fast Conformer models were trained using cosine scheduler with peak learning rates of 0.005 and 0.001 respectively. Both Conformer and Fast Conformer-RNNT models were trained for 80k steps, CTC models, for 380k steps. The models were trained on 32 GPUs with a global batch size of 2048. Last five checkpoints were averaged [19]. We used the Deepspeed profiler⁴ to estimate Multiply Accumulate operations (MACS) on a single 30s audio. The results are shown in Table 4. Fast Conformer has slightly better accuracy than regular Conformer. The proposed encoder is 3x more compute efficient than original Conformer encoder, and is significantly faster than EfficientConformer and SqueezeFormer.

3.1.2. Large-25k hours NeMo ASR Set

To test Fast Conformer capacity with respect to larger dataset, we trained Fast Conformer and Conformer models on 25k hours of speech, composed from LibriSpeech, Mozilla Common Voice, National Singapore Corpus and other public English speech datasets. All RNNT models were trained for 300k steps and CTC models for 1M steps. AdamW was used with cosine scheduler with a 15k linear warmup and learning rates of 0.0025 and 0.001, for RNNT and CTC respectively. The models have been tested on LibriSpeech test-other, MLS, MCV and WSJ-93 test sets. The results are presented in Table 5. Fast Conformer outperformed the original Conformer on most benchmarks.

⁴<https://www.deepspeed.ai/tutorials/flops-profiler/>

Table 4. ASR: Fast Conformer-Large with CTC and RNNT decoders trained on Librispeech. Greedy WER (%) was measured on Librispeech test-other. The number of parameters (M) and compute (Multiply-Acc, GMAC) are shown for encoder only.

Encoder	WER, % test-other	Params, M	Compute, GMACS
<i>RNNT decoder</i>			
Conformer	5.19	115	143.2
Fast Conformer	4.99	109	48.7
<i>CTC decoder</i>			
Conformer	5.74	121	149.2
Eff. Conformer[9]	5.79	125	101.3
SqueezeFormer[8]	6.05	125	91.0
Fast Conformer	5.64	115	51.5

Table 5. ASR: Fast Conformer-Large with CTC and RNNT decoders trained on English ASR set with 25K hours combined from public speech datasets. Greedy WER (%) was measured on Librispeech test-other, MCV 8, MLS and WSJ-92 test sets

Encoder	LS test-other	MCV 8 test	MLS En	WSJ-92 test
<i>RNNT decoder</i>				
Conformer	3.74	7.87	5.77	1.47
Fast Conformer	3.79	8.18	5.76	1.42
<i>CTC decoder</i>				
Conformer	4.50	9.40	6.60	1.70
Fast Conformer	4.19	9.00	6.42	1.59

3.2. Speech Translation

Next, we analyze the efficacy of Fast Conformer on Speech Translation (ST) from English to German. We trained two architectures with the same Conformer-like encoder and different autoregressive decoders: either RNNT, or 6-layer Transformer trained with cross entropy loss.

In all the experiments, we initialized encoder with the corresponding weights from ASR RNNT models trained on 25k hours of speech. The parameters of decoder and joint module were initialized randomly. Our vocabulary consists of 16384 YouTokenToMe⁵ byte-pair-encodings trained on German text. Our models have been trained on all available datasets from IWSLT22 [20] competition which corresponds to 4k hours of speech. Some of the datasets did not contain German translations, so we generated them ourselves with text-to-text machine translation model trained on WMT21 [21] and in-

⁵<https://github.com/VKCOM/YouTokenToMe>

Table 6. Speech Translation, MUST-C V2 tst-COMMON dataset. SacreBLEU, total inference time, and relative inference speed-up were measured with a batch size 32 for two speech translation models with Conformer-based encoder and either RNNT, or Transformer decoder.

Encoder	BLEU	Time, sec	Speed-up
<i>Transformer decoder</i>			
Conformer	31.0	267	1X
Fast Conformer	31.4	161	1.66X
<i>RNNT decoder</i>			
Conformer	23.2	83	1X
Fast Conformer	27.9	45	1.84X

Table 7. Speech intent classification and slot filling on SLURP dataset. ESPNet-SLU and SpeechBrain-SLU models use a HuBERT [23] encoder pre-trained via a self-supervised objective on LibriLight-60k [24]. Inference time and relative speed-up against Conformer are measured with batch size 32.

Model	Intent Acc.	SLURP F1	Inference, sec	Rel. Speed-up
SpeechBrain-SLU	87.70	76.19	-	-
ESPnet-SLU	86.52	76.91	-	-
Conformer/Fast Conformer+Transformer Decoder				
Conformer	90.14	82.27	210	1X
Fast Conformer	90.68	82.04	191	1.1X

domain finetuned on Must-C v2[22].

The results for all models are shown in Table 6. We note that RNNT loss is generally not suitable for speech translation due to its implicit monotonic alignment assumption. Surprisingly, Fast Conformer-RNNT translation model gets BLEU score of 27.89. In addition, the inference of this model is up to 1.84× faster compared to Conformer.

3.3. Spoken Language Understanding

Next, we apply the pre-trained Fast Conformer encoder to spoken language understanding (SLU). We study the *Speech Intent Classification and Slot Filling* (SICSF) task, which should detect user intents and extract the corresponding lexical fillers for detected entity slots [25]. An intent can be a composition of a scenario type and an action type. Slots and fillers are represented by key-value pairs. The ground-truth intents and slots of input are organized as a Python dictionary, represented as a string. The SICSF task is to predict this structured Python dictionary as a string, based on the input audio. Experiments are conducted using the SLURP [25] dataset, where intent accuracy and SLURP-F1 are used as the evaluation metric.

We use as the baseline Conformer encoder, initialized

Table 8. Fast Conformer versus Conformer on long audio. We evaluated four versions of Fast Conformer: (1) FC with full context attention (2) FC with limited context (3) FC with limited context and global token. Models have been evaluated on two long-audio benchmarks: TED-LIUM v3 and Earnings 21. Normalized greedy WER(%).

Model	TED-LIUM v3	Earnings21
Conformer	9.18	18.26
Fast Conformer (buffered)	9.15	17.65
+ Limited Context	8.25	16.08
+ Global Token	7.5	11.85

from pretrained ASR model, with a Transformer decoder. We also compare Fast Conformer against two state-of-the-art models from ESPNet-SLU [26] and SpeechBrain [27]. ESPNet-SLU and SpeechBrain both use a HuBERT [23] encoder pre-trained via a self-supervised objective on the entire LibriLight-60k [24] dataset. ESPNet-SLU further finetunes the encoder on LibriSpeech before training on SLURP.

The results for SLURP experiments are shown in Table 7. The model with pre-trained Fast Conformer encoder significantly surpasses the ESPNet-SLU and SpeechBrain models, which have been pre-trained on nearly 60K hours of speech. Fast Conformer attains very high accuracy quite close to Conformer. Its decoding is also 10(%) faster than Conformer. The speed-up is not as high as for ASR for following reason: the ratio of acoustic signal length (after 8x downsampling) to target token length is roughly **1:2.22** for the SICSF task. The execution cost for encoder is dwarfed by slow autoregressive Transformer decoder, and therefore we used batch 32 to balance the cost of the encoder against the decoder to show speed-up.

3.4. Long-form audio transcription

Fast Conformer with limited context and global token for long audio was trained in the following way. This model has shared query, key and value projection layers that are used for global and local attention. The encoder is initialized with the checkpoint pre-trained on our internal 25k hour set with full context. We then do additional fine-tuning on the same dataset with limited context attention for 10k steps. We used learning rate warmup of 1k steps, with maximum learning rate 1e-6 and cosine rate annealing to 0. The size of the limited context was set to 128 steps on each side of a token, which corresponds to around 10 seconds.

The performance of Fast Conformer was evaluated using two long-form audio datasets: TED-LIUM v3 [28], and Earnings-21 [29]. We compare our limited context model with full context Fast Conformer as well as with base Conformer trained on the same dataset. We used the Whisper normalizer on both transcripts and predictions to evaluate the

Table 9. Table presents parameter modifications required for constructing the FC-L, -XL, and -XXL models. Increased the hidden dimension (d_model), encoder layers and decoder RNNT layers to build XL from L. Keeping other model parameters constant from XL, we increased the number of encoder layers to construct FC-XXL.

Model	Hidden Dimension	Encoder Layers	RNNT Layers	Model Parameters
L	512	17	1	120 M
XL	1024	24	2	600 M
XXL	1024	42	2	1.1 B

models. For Conformer and Fast Conformer with full context attention we used 20 second buffers. For Fast Conformer with limited context we processed the full audio in one forward pass. Fast Conformer with new attention mechanism significantly outperforms Conformer and Fast Conformer with global attention on both long-form ASR benchmark sets (see Table 8).

4. SCALING FAST CONFORMER MODEL

To show scaling capacity of Fast Conformer models, we designed three model sizes: Large (L), Extra Large (XL) and Extra Extra Large (XXL), similar to Conformer scaling in [4] as shown in Table: 9.

We observed that when scaling FC models from XL to XXL, pretraining the encoder with Self Supervised learning helps stabilize training and enable high learning rates. We adopted the pretraining and finetuning method of SSL models based on Wav2Vec 2.0 [30]. Unlike Conformer models[4], we didn't change the conformer blocks and relative attention while scaling up models. From -L to -XXL core architecture of all FC models remains the same.

Table 10. Comparison of XL and XXL models on ASR benchmark datasets. Table illustrates the performance comparison between Conformer-XL versus Fast Conformer-XL RNNT models, as well as the improvement of Fast Conformer-XXL over Fast Conformer-XL. All models are trained on 25k hrs of ASR Set. Greedy WER(%).

Model	LS Test-clean	LS Test-other	MLS Test	GMACS
Conformer-XL	1.49	2.80	5.32	686
FC-XL	1.50	2.88	4.90	253
FC-XXL	1.38	2.52	4.58	441

The XL and XXL models yield superior results within fewer training steps when compared to L models. As shown in Table 10 when evaluated on HF-Leaderboard[31] evalua-

Table 11. The performance of FC-XL and FC-XXL models was compared with an augmented training set. The table illustrates the performance improvement of both models as we augment the training dataset by an additional 40,000 hours (ASR Set ++) on HF-audio leaderboard test sets [31] after whisper text normalization [32] applied. Greedy WER (%).

Model	Decoder	Train Dataset	LS Test-clean	LS Test-other	TED-LIUM V3	Vox Populi	MCV 9 Test	AMI Test	Earnings 22	SPGI Speech	Giga Speech
FC-XL	RNNT	ASR Set	1.50	2.88	4.49	5.74	7.26	18.28	16.37	4.40	11.58
		ASR Set ++	1.63	3.06	3.86	6.05	8.07	17.55	14.78	3.47	10.07
FC-XL	CTC	ASR Set	1.73	3.47	4.71	6.09	7.51	18.41	17.89	5.04	11.84
		ASR Set ++	1.87	3.76	3.78	7.00	10.57	16.3	14.14	4.11	10.35
FC-XXL	RNNT	ASR Set	1.38	2.52	4.74	5.56	6.07	18.81	16.66	4.98	11.95
		ASR Set ++	1.46	2.47	3.92	5.39	5.79	17.1	14.11	3.11	9.96
FC-XXL	CTC	ASR Set	1.69	3.4	4.64	6.45	8.31	17.62	16.44	4.91	11.61
		ASR Set ++	1.83	3.54	3.54	6.53	9.02	15.62	13.69	4.20	10.27

tion test sets ⁶, performance analysis showcases the effectiveness of the FC-XL and FC-XXL models trained on 25k hrs of NeMo ASR Set. The -XL models presented were trained for 70k steps, while -XXL the models were trained for 100k steps with an effective batch size of 2048. For the XL model, we use Adamw optimizer with Noam learning rate scheduler[33] with peak learning rate of 6e-4 and 15k linear warmup steps, whereas FC-XXL models were linearly warmed up for 25k steps. All XL and XXL models were initialized with pre-trained SSL checkpoints. We also observed that finetuning a RNNT FC-XL model with CTC just for 40k steps showed similar performance to training FC-XL CTC model for 200k steps from scratch.

4.1. Scaling Dataset

The FC-XXL RNNT model, which was trained on 25k hours of ASR Set, achieves similar state-of-the-art performance on LS-test other as [4], while also achieving the best performance on other benchmark datasets. However, in order to effectively leverage large models, it is imperative to proportionally augment dataset sizes in accordance with the model sizes. Hence, we trained these large models by incorporating an additional 40,000 hours of internal dataset (ASR Set ++). The integration of these supplementary datasets facilitated enhanced accuracy and noise robustness in both XL and XXL models. Table 11 shows successful utilization of large amount of data by effectively training the 1B parameter model. Furthermore, Figure 3 illustrates the noise robustness of the FC-XXL models across different signal-to-noise ratio (SNR) levels on the clean set of the Librispeech dataset.

5. CONCLUSIONS

In this paper, we propose Fast Conformer, a redesigned Conformer with a novel downsampling schema. Fast Conformer

⁶https://huggingface.co/spaces/hf-audio/open_asr_leaderboard

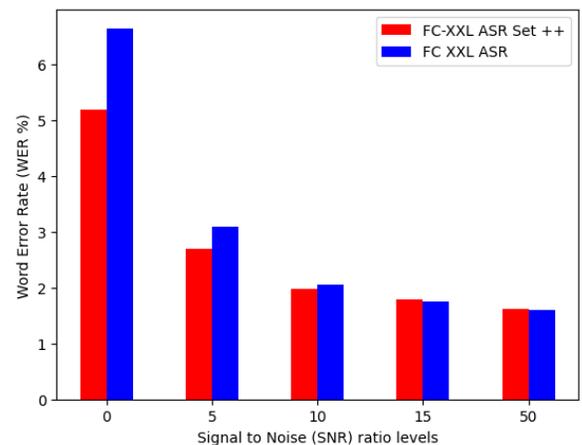


Fig. 3. Figure demonstrates the noise robustness of FC-XXL models on the LS Clean evaluation set, showcasing their performance across various signal-to-noise ratio (SNR) levels.

uses 2.9x less compute while delivering roughly the same WER as the original Conformer. Evaluations on tasks such as speech translation and spoken language understanding show a strong model accuracy while achieving significant speed-ups in the encoder computation. When the attention module is replaced with local attention, we show that the greater efficiency enables long-form transcription of an 11-hour audio segment in a single forward pass. The results on long-form audio are improved further by adding a single global attention token. We finally show that Fast Conformer architecture can be easily scaled to 1B parameters which enables us to further improve accuracy while achieving noise robustness when training on larger datasets.

6. REFERENCES

[1] 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 *Interspeech*, 2020.
- [2] Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, Hirofumi Inaguma, Naoyuki Kamo, Chenda Li, Daniel Garcia-Romero, Jiantong Shi, et al., “Recent developments on Espnet toolkit boosted by Conformer,” in *ICASSP*, 2021.
- [3] Samuel Krizan, Stanislav Beliaev, Boris Ginsburg, Jocelyn Huang, Oleksii Kuchaiev, Vitaly Lavrukhin, Ryan Leary, Jason Li, and Yang Zhang, “QuartzNet: Deep automatic speech recognition with 1D time-channel separable convolutions,” in *ICASSP*, 2020.
- [4] Yu Zhang, James Qin, Daniel S Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V Le, and Yonghui Wu, “Pushing the limits of semi-supervised learning for automatic speech recognition,” *arXiv:2010.10504*, 2020.
- [5] Iz Beltagy, Matthew E. Peters, and Arman Cohan, “Longformer: The long-document Transformer,” *arXiv:2004.05150*, 2020.
- [6] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [7] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *CVPR*, 2017.
- [8] Sehoon Kim, Amir Gholami, Albert Shaw, Nicholas Lee, Karttikeya Mangalam, Jitendra Malik, Michael W Mahoney, and Kurt Keutzer, “Squeezeformer: An efficient Transformer for automatic speech recognition,” in *NeurIPS*, 2022.
- [9] Maxime Burchi and Valentin Vielzeuf, “Efficient Conformer: Progressive downsampling and grouped attention for automatic speech recognition,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2021.
- [10] Andrei Andrusenko, Rauf Nasretidinov, and Aleksei Romanenko, “Uconv-conformer: High reduction of input sequence length for end-to-end speech recognition,” in *ICASSP*, 2022.
- [11] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *ICML*, 2006.
- [12] A. Graves, “Sequence transduction with recurrent neural networks,” in *ICML*, 2012.
- [13] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an ASR corpus based on public domain audio books,” in *ICASSP*, 2015.
- [14] Taku Kudo and John Richardson, “SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” *arXiv:1808.06226*, 2018.
- [15] François Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *CVPR*, 2017.
- [16] Vineel Pratap, Qiantong Xu, Anuroop Sriram, Gabriel Synnaeve, and Ronan Collobert, “MLS: A large-scale multilingual dataset for speech research,” *Interspeech*, 2020.
- [17] “Mozilla: A journey to less than 10% word error rate,” <https://hacks.mozilla.org/2017/11/a-journey-to-10-word-error-rate/>.
- [18] D. B. Paul and J. M. Baker, “The design for the Wall Street Journal based CSR corpus,” in *Proc. of the workshop on Speech and Natural Language*. ACL, 1992, pp. 357–362.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [20] Antonios Anastasopoulos, Loïc Barrault, Luisa Bontivogli, Marcelly Zanon Boito, Ondřej Bojar, Roldano Cattoni, Anna Currey, Georgiana Dinu, Kevin Duh, Maha Elbayad, et al., “Findings of the iwslt 2022 evaluation campaign,” in *IWSLT*, 2022.
- [21] Sandeep Subramanian, Oleksii Hrinchuk, Virginia Adams, and Oleksii Kuchaiev, “NVIDIA NeMo Neural Machine Translation Systems for English-German and English-Russian News and Biomedical Tasks at WMT21,” *arXiv:2111.08634*, 2021.
- [22] Roldano Cattoni, Mattia Antonino Di Gangi, Luisa Bontivogli, Matteo Negri, and Marco Turchi, “MuST-C: A multilingual corpus for end-to-end speech translation,” *Computer Speech & Language*, vol. 66, pp. 101155, 2021.
- [23] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Tsai, Kushal Lakhota, 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*, vol. 29, pp. 3451–3460, 2021.
- [24] Jacob Kahn, Morgane Rivière, Weiyi Zheng, Evgeny Kharitonov, Qiantong Xu, Pierre-Emmanuel Mazaré,

- Julien Karadayi, Vitaliy Liptchinsky, Ronan Collobert, Christian Fuegen, et al., “Libri-Light: A benchmark for ASR with limited or no supervision,” in *ICASSP*, 2020.
- [25] Emanuele Bastianelli, Andrea Vanzo, Pawel Swietojanski, and Verena Rieser, “SLURP: A Spoken Language Understanding Resource Package,” in *EMNLP*, 2020.
- [26] Siddhant Arora, Siddharth Dalmia, Pavel Denisov, Xunkai Chang, Yushi Ueda, Yifan Peng, Yuekai Zhang, Sujay Kumar, Karthik Ganesan, Brian Yan, et al., “ESPnet-SLU: Advancing spoken language understanding through ESPnet,” in *ICASSP*, 2022.
- [27] Yingzhi Wang, Abdelmoumene Boumadane, and Abdelwahab Heba, “A fine-tuned Wav2vec 2.0/HuBERT benchmark for speech emotion recognition, speaker verification and spoken language understanding,” *arXiv:2111.02735*, 2021.
- [28] François Hernandez, Vincent Nguyen, Sahar Ghannay, Natalia Tomashenko, and Yannick Esteve, “Ted-lium 3: Twice as much data and corpus repartition for experiments on speaker adaptation,” in *Speech and Computer: SPECOM 2018*. Springer, 2018.
- [29] Miguel Del Rio, Natalie Delworth, Ryan Westerman, Michelle Huang, Nishchal Bhandari, Joseph Palakapilly, Quinten McNamara, Joshua Dong, Piotr Żelasko, and Miguel Jetté, “Earnings-21: A practical benchmark for ASR in the wild,” in *Interspeech*, 2021.
- [30] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” in *NeurIPS*, 2020.
- [31] Vaibhav Srivastav, Somshubra Majumdar, Nithin Koluguri, Adel Moumen, Sanchit Gandhi, Hugging Face Team, Nvidia NeMo Team, and Speech-Brain Team, “Open automatic speech recognition leaderboard,” <https://huggingface.co/spaces/huggingface.co/spaces/open-asr-leaderboard/leaderboard>, 2023.
- [32] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, “Robust speech recognition via large-scale weak supervision,” 2022.
- [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.