

WeNet: Production First and Production Ready End-to-End Speech Recognition Toolkit

Binbin Zhang¹, Di Wu¹, Chao Yang¹, Xiaoyu Chen¹, Zhendong Peng¹, Xiangming Wang¹, Zhuoyuan Yao², Xiong Wang², Fan Yu², Lei Xie², Xin Lei¹

¹Mobvoi Inc., Beijing, China

²Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xi'an, China

binbinzhang@mobvoi.com

Abstract

In this paper, we present a new open source, production first and production ready end-to-end (E2E) speech recognition toolkit named WeNet. The main motivation of WeNet is to close the gap between the research and the production of E2E speech recognition models. WeNet provides an efficient way to ship ASR applications in several real-world scenarios, which is the main difference and advantage to other open source E2E speech recognition toolkits. This paper introduces WeNet from three aspects, including model architecture, framework design and performance metrics. Our experiments on AISHELL-1 using WeNet, not only give a promising character error rate (CER) on a unified streaming and non-streaming two pass (U2) E2E model but also show reasonable RTF and latency, both of these aspects are favored for production adoption. The toolkit is publicly available at <https://github.com/mobvoi/wenet>.

Index Terms: WeNet, Production Ready, U2

1. Introduction

The E2E models, including CTC[1, 2], recurrent neural network transducer (RNN-T)[3, 4] and attention based encoder-decoder (AED)[5, 6, 7], have gained more and more attention to speech recognition over the last few years. Compared with the hybrid ASR framework, the most attractive merit of E2E models is the extremely simplified training procedure. Recent work [8, 9, 10] also shows that E2E systems have surpassed conventional hybrid ASR systems in the standard of word error rate (WER). Considering the foregoing mentioned advantages of E2E models, deploying the emerging ASR framework into real-world productions becomes necessary. However, deploying an E2E system is trivial and there are a lot of problems to be solved.

First, **the streaming problem**. Streaming inference is essential for many scenarios that require the ASR system has to respond quickly with low latency. However, it is difficult for some E2E models to be streaming, such as AED, either great effort is required or big accuracy loss is introduced to make such model work in a streaming fashion[11, 12, 13].

Second, **unifying streaming and non-streaming modes**. Unified streaming and non-streaming model can reduce the development effort, training cost, and deployment cost, especially for small companies, which is also preferred for production adoption[14, 15].

Third, **the production problem**, which is also the most important problem we care about during WeNet design. Great efforts are required to promote the E2E model into a real production application even though we already have a unified model which has reasonable performance on both streaming and non-

streaming applications. Firstly, we have to convert the research model to a production model, which is painful for a dynamic graph based deep learning toolkit, such as PyTorch[16]. Secondly, we have to carefully design the inference workflow in terms of the model architecture, applications and runtime platforms. For the model architecture, most E2E models first do an encoder forward computation then an autoregressive beam search. The workflow is more complicated than a simple neural network forward and problems become even more complicated if streaming is required. For both cloud and on-device applications, computation and memory cost should be seriously considered. Especially for on-device models, inference optimization and model quantization play very important roles. As for runtime platforms, although there are various platforms could be used to do neural network inference, such as ONNX (Open Neural Network Exchange), LibTorch in PyTorch, TensorRT[17], OpenVINO, MNN[18], and NCNN, it still requires both speech specific and advanced deep learning optimization knowledge to select the best one for your application.

In this work, we present WeNet to address the above problems. “We” in WeNet is inspired by “WeChat”, which means connection and share, “Net” is from Espnet[19] since we have referred to a lot of the excellent design in Espnet. The main motivation of WeNet is to close the gap between research and production of E2E speech recognition models, to reduce the effort of producing E2E models explore better E2E models for production. Therefore, WeNet is designed for production in nature, which makes our WeNet distinctly different from other toolkits.

On the production first and production ready principles, WeNet adopts the following implementations. First, WeNet adopts the Unified two Pass (U2)[20] framework to solve the streaming and unified problems. Second, from model training to deployment, WeNet only depends on PyTorch and its ecosystem. Furthermore, WeNet also provides a off-the-shell pipeline for both cloud server and on-device (Android) deployment. The key advantages of WeNet are:

1. **Production first and production ready:** The Python code of WeNet meets the requirements of TorchScript, so the model trained by WeNet can be directly exported by Torch JIT and use LibTorch for inference. There is no gap between the research model and the production model. Neither model conversion nor additional code is required for model inference.
2. **Unified solution for streaming and non-streaming ASR:** WeNet adopts the U2 framework to achieve an accurate, fast and unified E2E model, which is favorable for industry adoption.

3. **Portable runtime:** Several runtimes will be provided to show how to host WeNet trained models on different platforms, including server (x86) and embedded (ARM in Android platforms).
4. **Light weight:** WeNet is designed specifically for E2E speech recognition with clean and simple code. It is all based on PyTorch and its corresponding ecosystem. It has no dependencies on Kaldi, which simplifies installation and usage.

As our experiments show, WeNet is a simple, accurate speech recognition toolkit with an end-to-end solution from research to production.

2. Related Works

EspNet is the most popular open source platform for end-to-end speech research, it mainly focuses on end-to-end automatic speech recognition (ASR), and adopts widely-used dynamic neural network toolkits Chainer and PyTorch as the main deep learning engine. It provides E2E implementations including CTC, AED, RNN-T, and RNN language model rescoring. While EspNet is useful and widely-used for research, it is hard to directly use the model trained by EspNet in production, and there is no production consideration or support in EspNet design.

Currently, there is no such toolkit that focuses on production-level E2E speech recognition.

3. WeNet

3.1. Model Architecture

As we aim to address the streaming problem, the unified problem, the production problem, and the solution should be simple, easy to implement and train, with good performance as well as easy to be productized at runtime.

U2, a unified two-pass gives a great solution to the problems, as shown in Figure 1, which is a joint CTC/AED model. It consists of three parts, a *Shared Encoder*, a *CTC Decoder*, and an *Attention Decoder*. The *Shared Encoder* consists of multiple Transformer[21] or Conformer[22] encoder layers. The *CTC Decoder* consists of a linear layer, which transforms the *Shared Encoder* output to the CTC activation. The *Attention Decoder* consists of multiple Transformer decoder layers. The *Shared Encoder* only sees limited right contexts, and the *CTC Decoder* runs in a streaming mode in the first pass, and the *Attention Decoder* is used in the second pass to give a more accurate result.

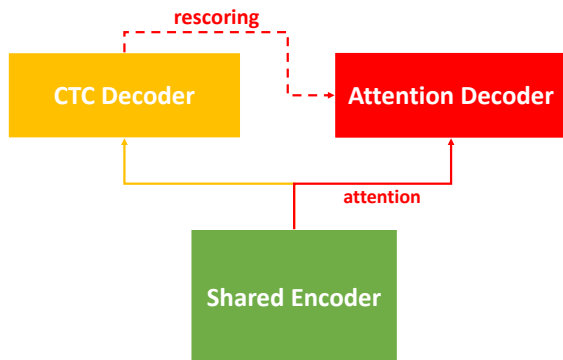


Figure 1: Two pass CTC and AED joint architecture

3.1.1. Training

A combined loss with CTC loss and AED loss is adopted in training, as shown in Equation 1, where \mathbf{x} is the acoustic feature, \mathbf{y} is the corresponding annotation, $\mathbf{L}_{\text{CTC}}(\mathbf{x}, \mathbf{y})$, $\mathbf{L}_{\text{AED}}(\mathbf{x}, \mathbf{y})$ are the CTC and AED loss respectively, λ is a hyperparameter which balances the importance of CTC and AED loss.

$$\mathbf{L}_{\text{combined}}(\mathbf{x}, \mathbf{y}) = \lambda \mathbf{L}_{\text{CTC}}(\mathbf{x}, \mathbf{y}) + (1 - \lambda) \mathbf{L}_{\text{AED}}(\mathbf{x}, \mathbf{y}) \quad (1)$$

A dynamic chunk training technique is applied in the training to unify the none streaming and streaming model and enable latency control. First, The input is split into several chunks by a fixed chunk size C with inputs $[t+1, t+2, \dots, t+C]$, every chunk attends on itself and all the previous chunks, and the whole latency for the *CTC Decoder* in the first pass only depends on the chunk size. If the chunk size is limited, it works in a streaming way, otherwise, it works in a non-streaming way. Second, the chunk size is varied dynamically from 1 to the max length of the current training utterance in the training, so the trained model learns to predict with arbitrary chunk size. Empirically, a larger chunk size gives better results with higher latency, so we can easily balance the accuracy and latency by tuning the chunk size at runtime.

3.1.2. Decoding

For Python decoding in the research stage, to compare and evaluate different parts of the joint CTC/AED model, WeNet supports four decoding modes as follows:

1. attention: apply standard autoregressive beam search on the AED part of the model.
2. ctc_greedy_search: apply CTC greedy search on the CTC part of the model, CTC greedy search is super faster than other modes.
3. ctc_prefix_beam_search: apply CTC prefix beam search on the CTC part of the model, which can give the n-best candidates.
4. attention_rescoring: first apply CTC prefix beam search on the CTC part of the model to generate n-best candidates, and then rescore the n-best candidates on the AED decoder part with corresponding encoder output.

For decoding in the runtime stage, the only attention_rescoring is supported since it's our ultimate solution for production.

3.2. System Design

The overall design stack of WeNet is as Figure 2. The whole framework is fully based on PyTorch and it's an ecosystem as the bottom stack, as what you will see in the following section, TorchScript is used for developing model, TorchAudio is used for on-the-fly feature extraction, DistributedDataParallel is used for distributed training, torch JIT (Just In Time) is used for model exportation, PyTorch quantization is used for the quantized model, and LibTorch is used for production runtime.

The second stack consists of two parts. Python (TorchScript) Research is for developing a research model, TorchScript is used to ensure the model could be correctly exported as a production model. LibTorch Production is for hosting production model, which is designed to support various hardware and platforms like CPU, GPU (CUDA) Linux, Android, and IOS.

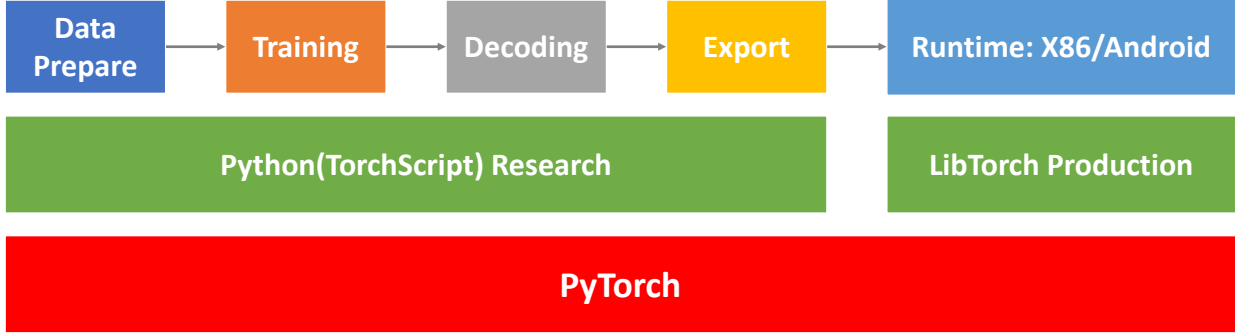


Figure 2: WeNet system design

The third stack shows typical research to production pipeline in WeNet, the following sections will go through the detailed design of these modules.

3.2.1. Data Prepare

Data prepare in WeNet is pretty simple, a Kaldi style label file and wave list file, the model unit dictionary file which maps model unit to corresponding integer id, are all you need. There is no need for any feature extraction in the data preparation stage since we use on-the-fly feature extraction in training.

3.2.2. Training

The training stage in WeNet has the following key features:

On-the-fly feature extraction: this is based on TorchAudio, which can generate the same FBANK feature as Kaldi. Since the feature is extracted on-the-fly from the raw PCM data, we can do data augmentation on raw PCM at a time level, frequency level, and final feature level at the same time, which enrich the diversity of data and data augmentation.

Joint CTC/AED training: joint training speeds up the convergence of the training, improves the stability of the training, as well as gives better recognition results.

Distributed training: WeNet supports multiple GPUs training with DistributedDataParallel in PyTorch.

3.2.3. Decoding

A set of Python tools are provided to recognize the wave files and compute the accuracy. These tools help users validate and debug the model before deploying it in production. All the decoding algorithms in Section 3.1.2 are supported.

3.2.4. Export

For WeNet model is implemented in TorchScript, it could be exported by torch JIT to the production model directly and safely. Then this exported model can be hosted by using the LibTorch library in runtime. Both float model and quantized int8 model are supported. Using a quantized model could double the inference speed or even more when hosted on Android devices.

3.2.5. Runtime

Currently, we support hosting WeNet production model on two mainstream platforms, namely x86 as server runtime and Android as on-device runtime. A C++ API Library and runnable demos for both platforms are provided. The User could also implement their customized system by using the C++ library. We

carefully evaluated the three key metrics of the ASR system, namely accuracy, the real-time factor (RTF), and latency. The performance is suitable for many ASR applications such as service API and on-device voice assistants. The result is reported in section 4.2.

4. Experiments

We carry out our experiments on the open-source Chinese Mandarin speech corpus AISHELL-1[23], which contains a 150-hour training set, a 10-hour development set, and a 5-hour test set. The test set contains 7176 utterances in total. We use 80 dimensional log Mel-filter bank (FBANK) computed on-the-fly by TorchAudio with a 25ms window with a 10ms shift as the feature. SpecAugment[24] is applied with 2 frequency masks with maximum frequency mask ($F = 10$), and 2 time masks with maximum time mask ($T = 50$). Two convolution sub-sampling layers with kernel size 3*3 and stride 2 are used in the front of the encoder, namely 4 times sub-sampling in all. We use 12 transformer layers for the encoder and 6 transformer layers for the decoder. Adam optimizer and a learning rate schedule with 25000 warm-up steps are used in training. Moreover, we get our final model by averaging the top-K best models which have a lower loss on the dev set during training.

4.1. Unified model evaluation

We first trained a non-streaming model (M1) as our baseline, the model is trained and inferenced by full attention. Then we trained a unified mode (M2) with a dynamic chunk strategy. M2 is inferenced with different chunk sizes full/16/8/4 at decoding, full is for full attention non-streaming case, 16/8/4 is for the streaming case.

Table 1: Unified model evaluation

decoding method	M1	M2			
		full	16	8	4
attention	5.76	6.13	6.43	6.59	6.8
ctc_greedy_search	6.21	6.75	7.85	8.41	9.44
ctc_prefix_beam_search	6.21	6.74	7.85	8.41	9.43
attention_rescoring	5.47	5.79	6.50	6.89	7.49

First, as shown in Table 1, the unified model not only shows comparable results to the non-streaming model on the full attention case but also gives promising results on streaming case with limited chunk size 16/8/4, which shows the effectiveness

of dynamic chunk training strategy.

Second, by comparing the four decoding modes, we can see the attention_rescoring mode always improves CTC results in both the non-streaming model and the unified model. The ctc_greedy_search and ctc_prefix_beam_search have almost the same performance, and they degrade significantly as the chunk size decreases. While the attention mode degrades slightly, and the attention_rescoring mode alleviates the degradation of the ctc_prefix_beam_search results. And as the U2 paper shows, the attention_rescoring mode is faster and it has a better RTF than the attention mode since attention mode is an autoregressive procedure while the attention_rescoring mode is not. Overall, the attention_rescoring not only shows promising results but also has a lower RTF.

So the dynamic chunk based unified model with attention_rescoring decoding is our choice for production, resulting in only the attention_rescoring mode is supported in our runtime.

4.2. Runtime benchmark

This section shows the quantization, RTF, and latency benchmark on the unified model M2 described above. We do our benchmark on a server x86 platform, an on-device ARM Android platform respectively.

For the cloud x86, the CPU is 4 cores Intel(R) Xeon(R) CPU E5-2620 v4 @ 2.10GHz, the memory is 16G in total. Only one thread is used for CPU threading and TorchScript inference¹ for each utterance since the cloud service requires parallel processing, and a single thread avoids performance degradation in parallel processing.

For the on-device Android, the CPU is 4 cores Qualcomm Snapdragon 865, the memory is 8.00GB. A single thread is used for the on-device inference.

4.2.1. Quantization

Here we just compare the CER difference before and after quantization. The RTF of quantization is shown in the following section. As in Table 2, the CER is comparable when quantization is applied. The CER of the float model is slightly different from what we listed in Table 1, since Table 1 is tested by Python research tools while the results here are tested by runtime tools.

Table 2: CER before and after quantization

quantization/decoding_chunk	full	16	8	4
NO (float32)	5.87	6.49	6.88	7.46
YES (int8)	5.89	6.54	6.89	7.51

4.2.2. RTF

As in Table 3, we can see the RTF increases as the chunk size decreases since the smaller chunk requires more iterations for the forward computation. Further, quantization yields about 2 times speedup on on-device (Android) and a slight improvement on the server (x86).

4.2.3. Latency

For latency benchmark, we create a WebSocket server/client to simulate a real streaming application. This benchmark is only

Table 3: RTF benchmark

model/decoding_chunk	full	16	8	4
server (x86) float32	0.079	0.095	0.128	0.186
server (x86) int8	0.072	0.081	0.098	0.134
on-device (Android) float32	0.164	0.251	0.350	0.505
on-device (Android) int8	0.082	0.114	0.130	0.201

carried out on the server x86 platform. The average latency we evaluated is described here:

1. **model latency (L1):** the wait time introduced by the model structure. for our chunk based decoding, the average wait time is half of the chunk. And the total model latency of our model is $(chunk/2 * 4 + 6) * 10$ (ms), where 4 is the subsampling rate, 6 is the lookahead introduced by the first two CNN layers in the encoder, 10 is the frame shift.
2. **rescoring cost (L2):** the time cost on the second pass attention rescoring.
3. **final latency (L3):** the user (client) perceived latency, which is the time difference between the user stopping speaking and getting the recognition result. When our ASR server receives the speech stopping signal, it first forwards the left speech for CTC searching, then does the second pass attention rescoring, so rescoring cost is part of the final latency. The network latency should be also taken into account for a real production, however since we tested the server/client in the same machine, so the network latency is negligible.

Table 4: Latency benchmark

decoding_chunk	L1 (ms)	L2 (ms)	L3 (ms)
16	380	115	142
8	220	115	135
4	140	114	130

As we see in Table 4, First, the rescoring costs are almost the same for different chunk sizes, this is reasonable since rescoring computation is invariant to chunk size. Second, the final latency is dominated by the rescoring cost, which means we can further reduce the final latency by reducing the rescoring cost. Third, the final latency increases slightly as the decoding chunk varies from 4 to 8 to 16.

5. Conclusions

We present a new open source E2E speech recognition toolkit, which is production first and production ready, provides a unified solution for streaming and non-streaming application, benchmarks the accuracy, RTF and latency. The whole toolkit is well designed, lightweight, and it shows great performance.

¹https://pytorch.org/docs/stable/notes/cpu_threading_torchscript_inference.html

6. References

- [1] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 369–376.
- [2] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen et al., “Deep speech 2: End-to-end speech recognition in english and mandarin,” in *International conference on machine learning*, 2016, pp. 173–182.
- [3] A. Graves, “Sequence transduction with recurrent neural networks,” *arXiv preprint arXiv:1211.3711*, 2012.
- [4] A. Graves, A.-r. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.
- [5] J. Chorowski, D. Bahdanau, K. Cho, and Y. Bengio, “End-to-end continuous speech recognition using attention-based recurrent nn: First results,” *arXiv preprint arXiv:1412.1602*, 2014.
- [6] W. Chan, N. Jaitly, Q. V. Le, and O. Vinyals, “Listen, attend and spell,” *arXiv preprint arXiv:1508.01211*, 2015.
- [7] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, “Attention-based models for speech recognition,” in *Advances in neural information processing systems*, 2015, pp. 577–585.
- [8] R. Prabhavalkar, K. Rao, T. N. Sainath, B. Li, L. Johnson, and N. Jaitly, “A comparison of sequence-to-sequence models for speech recognition,” in *Interspeech*, 2017, pp. 939–943.
- [9] T. N. Sainath, R. Pang, D. Rybach, Y. He, R. Prabhavalkar, W. Li, M. Visontai, Q. Liang, T. Strohman, Y. Wu et al., “Two-pass end-to-end speech recognition,” *arXiv preprint arXiv:1908.10992*, 2019.
- [10] S. Kim, T. Hori, and S. Watanabe, “Joint ctc-attention based end-to-end speech recognition using multi-task learning,” in *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2017, pp. 4835–4839.
- [11] C. Raffel, M.-T. Luong, P. J. Liu, R. J. Weiss, and D. Eck, “Online and linear-time attention by enforcing monotonic alignments,” *arXiv preprint arXiv:1704.00784*, 2017.
- [12] C.-C. Chiu and C. Raffel, “Monotonic chunkwise attention,” *arXiv preprint arXiv:1712.05382*, 2017.
- [13] H. Inaguma, M. Mimura, and T. Kawahara, “Enhancing monotonic multihead attention for streaming asr,” *arXiv preprint arXiv:2005.09394*, 2020.
- [14] J. Yu, W. Han, A. Gulati, C.-C. Chiu, B. Li, T. N. Sainath, Y. Wu, and R. Pang, “Universal asr: Unify and improve streaming asr with full-context modeling,” *arXiv preprint arXiv:2010.06030*, 2020.
- [15] A. Tripathi, J. Kim, Q. Zhang, H. Lu, and H. Sak, “Transformer transducer: One model unifying streaming and non-streaming speech recognition,” *arXiv preprint arXiv:2010.03192*, 2020.
- [16] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in neural information processing systems*, 2019, pp. 8026–8037.
- [17] H. Vanholder, “Efficient inference with tensorsrt,” 2016.
- [18] X. Jiang, H. Wang, Y. Chen, Z. Wu, L. Wang, B. Zou, Y. Yang, Z. Cui, Y. Cai, T. Yu, C. Lv, and Z. Wu, “Mnn: A universal and efficient inference engine,” in *MLSys*, 2020.
- [19] S. Watanabe, T. Hori, S. Karita, T. Hayashi, J. Nishitoba, Y. Unno, N. E. Y. Soplin, J. Heymann, M. Wiesner, N. Chen et al., “Espnet: End-to-end speech processing toolkit,” *arXiv preprint arXiv:1804.00015*, 2018.
- [20] B. Zhang, D. Wu, Z. Yao, X. Wang, F. Yu, C. Yang, L. Guo, Y. Hu, L. Xie, and X. Lei, “Unified streaming and non-streaming two-pass end-to-end model for speech recognition,” *arXiv preprint arXiv:2012.05481*, 2020.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [22] 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,” *arXiv preprint arXiv:2005.08100*, 2020.
- [23] H. Bu, J. Du, X. Na, B. Wu, and H. Zheng, “Aishell-1: An open-source mandarin speech corpus and a speech recognition baseline,” in *2017 20th Conference of the Oriental Chapter of the International Coordinating Committee on Speech Databases and Speech I/O Systems and Assessment (O-COCOSDA)*. IEEE, 2017, pp. 1–5.
- [24] D. S. 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 preprint arXiv:1904.08779*, 2019.