

TORCHAUDIO 2.1: ADVANCING SPEECH RECOGNITION, SELF-SUPERVISED LEARNING, AND AUDIO PROCESSING COMPONENTS FOR PYTORCH

Jeff Hwang^{*1}, Moto Hira^{*1}, Caroline Chen^{*†}, Xiaohui Zhang^{*1}, Zhaoheng Ni^{*1}, Guangzhi Sun², Pingchuan Ma¹, Ruizhe Huang^{†3}, Vineel Pratap¹, Yuekai Zhang⁴, Anurag Kumar¹, Chin-Yun Yu⁵, Chuang Zhu⁴, Chunxi Liu[†], Jacob Kahn¹, Mirco Ravanelli⁶, Peng Sun⁴, Shinji Watanabe⁸, Yangyang Shi¹, Yumeng Tao[†], Robin Scheibler⁷, Samuele Cornell⁸, Sean Kim[†], Stavros Petridis¹

¹Meta, ²University of Cambridge, ³Johns Hopkins University,
⁴NVIDIA, ⁵Queen Mary University of London, ⁶Concordia University,
⁷LY Corporation, ⁸Carnegie Mellon University

ABSTRACT

TorchAudio is an open-source audio and speech processing library built for PyTorch. It aims to accelerate the research and development of audio and speech technologies by providing well-designed, easy-to-use, and performant PyTorch components. Its contributors routinely engage with users to understand their needs and fulfill them by developing impactful features. Here, we survey TorchAudio’s development principles and contents and highlight key features we include in its latest version (2.1): self-supervised learning pre-trained pipelines and training recipes, high-performance CTC decoders, speech recognition models and training recipes, advanced media I/O capabilities, and tools for performing forced alignment, multi-channel speech enhancement, and reference-less speech assessment. For a selection of these features, through empirical studies, we demonstrate their efficacy and show that they achieve competitive or state-of-the-art performance.

Index Terms— Open-Source Toolkit, Speech Recognition, Audio Processing, Self-Supervised Learning

1. INTRODUCTION

With the rapid advancement and increasing pervasiveness of machine learning technologies, usage of open-source toolkits such as Tensorflow [1] and PyTorch [2] for developing machine learning applications has grown significantly. Many modern machine learning applications interface with modalities such as vision, text, and audio. Building such applications, however, requires modality-specific functionality not covered by said general-purpose toolkits.

To address the need for audio and speech facilities in particular, the TorchAudio library has been developed [3]. TorchAudio supplements PyTorch with easy-to-use and performant components for developing audio and speech machine learning models. As a natural extension of PyTorch to the audio domain, TorchAudio embodies many of the same design principles that PyTorch does. Its components support automatic differentiation to facilitate building neural networks and training them end to end. It supports GPU acceleration, which can greatly improve training and inference throughput. It emphasizes composability, simple interfaces shared with PyTorch, and

Version	Contrib.	Commits	Stars	Forks	Dep. repos
0.10 (Sep 2021)	144	1,013	1,428	351	5,420
2.1 (Jul 2023)	204	2,154	2,149	585	31,173

Table 1. TorchAudio’s Github activity statistics, covering unique contributors, commits, stars, forks, and repositories depending on TorchAudio. Statistics for Version 0.10 gleaned from web.archive.org and [3].

minimal dependencies to allow for easily integrating its components into any application that uses PyTorch.

TorchAudio has been widely adopted and actively developed by the PyTorch community, with its Github development statistics having grown substantially since Version 0.10 was presented in [3] (Table 1). The dramatic increase in the number of repositories that depend on TorchAudio in particular strongly reaffirms TorchAudio’s usefulness to the community and success.

This paper begins by summarizing TorchAudio’s design principles and contents. It then expounds significant new features that have been introduced since Version 0.10 [3], covering self-supervised learning (Wav2Vec 2.0 [4], HuBERT [5], XLS-R [6], WavLM [7]), automatic speech recognition (CTC decoder [8], Conformer [9], Emformer [10], audio-visual speech recognition [AV-ASR]), advanced media I/O, CTC-based forced alignment, multi-channel speech enhancement components, and reference-less speech assessment, of which several are technically novel, e.g. real-time AV-ASR, Emformer, CUDA-based CTC decoder, and CUDA-based forced alignment API. It concludes by presenting experimental results for the new features, which demonstrate that they are effective and achieve or exceed parity in run-time efficiency or output quality with public implementations.

2. RELATED WORK

Several popular open-source toolkits implement lower-level audio operations such as I/O, spectrogram generation, and data augmentations. Just as librosa [11] is one such library for Numpy [12] and tfio.audio for Tensorflow, TorchAudio is one such library for PyTorch. The broad applicability of TorchAudio’s data componentry has made it effective in serving more specialized audio data representation li-

* Equal contribution.

† Work done while at Meta.

libraries such as Lhotse [13], which provides abstractions and utilities that streamline data preparation for downstream audio tasks.

Many higher-level audio and speech machine learning toolkits exist in the PyTorch ecosystem, e.g. ESPnet [14], SpeechBrain [15], fairseq [16], and NeMo [17]. These toolkits provide ready-to-use models, training recipes, and components covering audio and speech tasks such as text to speech, speech recognition, speech translation, and speech enhancement. As the aforementioned audio operations are fundamental to such tasks, all of these toolkits rely on TorchAudio.

In addition to lower-level audio components, TorchAudio also provides some of the features offered by these higher-level toolkits. For instance, TorchAudio includes task-specific components such as decoders for speech recognition, multichannel functions, and model architectures, as well as ready-to-use models and training recipes. That being said, as far as such features are concerned, TorchAudio is distinguished from many of these other toolkits in its focus on *stable and established technologies* over the cutting edge. For example, rather than maintaining an extensive model repository and continually updating it with the latest state-of-the-art models, we aim to curate a smaller selection of key models and training recipes to demonstrate the use of TorchAudio’s components and serve as reliable references. Ultimately, we intend for TorchAudio to be first and foremost a library of established components, which allows it to *complement rather than compete with other toolkits in the PyTorch ecosystem*.

3. LIBRARY PRINCIPLES

TorchAudio firmly adheres to several design principles, which we distill from [3] and clarify.

Extend PyTorch to audio. TorchAudio aims to be PyTorch for the audio domain. Its components compose PyTorch operations, share the same abstractions and Tensor-based interfaces with PyTorch, and support foundational PyTorch features such as GPU acceleration and automatic differentiation. Moreover, its only required dependency is PyTorch. As a consequence, it behaves as a natural extension of PyTorch, and its components integrate seamlessly with PyTorch applications.

Be easy to use. TorchAudio is intuitively designed. Each component is implemented closely following C++, Python, and PyTorch best practices.

It is easy to install. TorchAudio’s binaries are distributed through standard Python package managers PyPI and Conda and support major platforms Linux, macOS, and Windows. Optional dependencies are similarly installable via standard package managers. For users who want to use their own custom logic, building TorchAudio from source is straightforward¹.

It is extensively documented. TorchAudio’s official website² comprehensively covers installation directions and the library’s public APIs. Moreover, a wide array of tutorials covering basic and advanced library usages are available on the website and Google Colab. Such resources educate users of all levels of familiarity with audio and speech technologies on how to best use TorchAudio to address their needs.

Favor stability. TorchAudio tends towards mature techniques that are broadly useful. It offers implementations of models and operations

that are or will soon become standards in the field. New features are released following a prototype-beta-stable progression to allow users to preview them without disrupting the official releases. Backwards compatibility breaking changes are released after a minimum of two releases to give users ample time to adapt their use cases. 12,000+ test cases and continuous integration workflows run through Github Actions ensure that the APIs work as expected.

Promote accessibility. TorchAudio is an open source library. Its entire source code is available on Github³, where contributions and feedback are encouraged from all. To enable usage in as many contexts as possible, TorchAudio is released under the permissive BSD-2 license.

4. NEW FEATURES

Relative to Version 0.10 [3], TorchAudio 2.1 includes many significant new features. We elaborate on several of these below. Note that some of these features are technically novel and the first of their kind, e.g. the first AV-ASR model to be capable of real-time inference on CPU, the first public implementation of streaming-capable transformer-based acoustic model Emformer, the first CUDA-based CTC beam search decoder, and the first CUDA-based forced alignment API.

Self-supervised learning. Self-supervised learning (SSL) approaches have consistently improved performance for downstream speech processing tasks. While S3PRL [18] focuses on supporting downstream tasks and benchmarking, TorchAudio focuses on upstream models by providing reliable and production-ready pre-trained models and training recipes. TorchAudio now provides models and pre-trained pipelines for Wav2Vec 2.0 [4], HuBERT [5], XLS-R [6], and WavLM [7]. Each pre-trained pipeline relies on the weights that the corresponding original model uses and thus produces identical outputs. Moreover, each pipeline is easy to use, simply expecting users to call a single method to retrieve a pre-trained model. To facilitate production usage, TorchAudio’s model implementations support TorchScript and PyTorch-native quantization and leverage PyTorch 2.0’s Accelerated Transformers⁴ to speed up training and inference.

TorchAudio also provides end-to-end training recipes that allow for pre-training and fine-tuning HuBERT models from scratch. The training recipes have minimal dependencies beyond PyTorch and TorchAudio and are modularly implemented entirely in imperative code, which makes them conducive to customization and integration with other training flows, as their adoption by other frameworks such as ESPnet [19] demonstrates.

CTC decoder. Beam search is an efficient algorithm that has been used extensively for decoding speech recognition (ASR) model outputs and remains a fast and lightweight alternative to model-based decoding approaches. We have added a CTC beam search decoder that wraps Flashlight Text’s [20] high performance beam search decoder in an intuitive and flexible Python API. The decoder is general purpose, working for both lexicon and lexicon-free decoding as well as various language model types, including KenLM [21] and custom neural networks, and is easily adaptable to different model outputs.

We have also introduced a CUDA-based CTC beam search decoder. By parallelizing computation along the batch, hypothesis, and

¹<https://github.com/pytorch/audio/blob/main/CONTRIBUTING.md>

²<https://pytorch.org/audio/>

³<https://github.com/pytorch/audio>

⁴<https://pytorch.org/blog/accelerating-large-language-models/>

vocabulary dimensions, it can achieve much higher decoding throughputs than the CPU-based implementation, which we demonstrate in Section 5.2. To our knowledge, the implementation is the first and only publicly available CUDA-compatible CTC decoder.

Conformer. Conformer is a transformer-based acoustic model architecture that has achieved state-of-the-art results for ASR [9, 22]. We have developed a PyTorch-based implementation of Conformer and published an RNN-Transducer ASR training recipe that uses it. Using the recipe, we produced a model that achieves word-error-rate (WER) parity with comparable open-source implementations, which will be discussed in Section 5.3

Emformer. Emformer is a streaming-capable efficient memory transformer-based acoustic model [10]. For on-device streaming ASR applications, it has demonstrated state-of-the-art performance balancing word error rate, latency, and model size. Moreover, because it applies a novel parallel block processing scheme for training, it can be trained very efficiently. We have introduced an implementation of Emformer matching that described in [10] along with an Emformer transducer ASR training recipe and pre-trained inference pipeline. Our implementation is the first to be publicly available, and it has been adopted and extended by icefall⁵.

Streaming AV-ASR. AV-ASR involves transcribing text from audio and video. The vast majority of work to date [23, 24, 25] has focused on developing non-streaming AV-ASR models; studies on streaming AV-ASR, i.e. transcribing text from audio and video streams in real time, are comparatively limited [26]. Auto-AVSR [25] is an effective approach to scale up audio-visual data, which enables training more accurate and robust speech recognition systems. We extend Auto-AVSR to real-time AV-ASR and provide an example Emformer transducer training pipeline that incorporates audio-visual input. As far as we know, the AV-ASR model is the first to be capable of real-time inference on CPU.

Advanced media I/O. We have added advanced media processing capabilities to TorchAudio. Class `StreamReader` can decode not only audio but also images and videos to PyTorch tensors. Similarly, `StreamWriter` can encode tensors as audio, images, and videos. Both support streaming processing as well as applying transforms such as resampling and resizing. They are capable of interfacing with numerous sources and destinations, including file paths and objects, network locations, and devices, e.g. microphones and webcams. Using these features, one can for instance stream audio chunk by chunk from a remote video file and process the corresponding tensors in an online fashion.

We convey the simplicity and versatility of the API via code samples. Figure 1 instantiates `StreamReader` specifying the data source to be a network location, configures output audio and video streams, and iterates over tensors representing chunks of audio and video streamed from the output. Appendix A provides additional examples that illustrate how to read from media devices and write to a Real-Time Messaging Protocol server.

Furthermore, `StreamReader` and `StreamWriter` can leverage hardware video processors available on NVIDIA GPUs to greatly accelerate decoding and encoding.

⁵<https://github.com/k2-fsa/icefall/tree/master/egs/librispeech/ASR>

```
# Stream video and audio from remote location
r = torchaudio.io.StreamReader(
    src="https://example.com/video.mp4")
r.add_basic_audio_stream(
    frames_per_chunk=1600,
    sample_rate=8000,
)
r.add_basic_video_stream(
    frames_per_chunk=5,
    frame_rate=25,      # change frame rate
    width=224,          # change width and height
    height=196,
)
for (audio, video,) in r.stream():
    # audio.shape == [frames, channels]
    # video.shape ==
    #     [frames, channels, height, width]
```

Fig. 1. `StreamReader` usage example.

CTC-based forced alignment. We have added support for forced alignment generation, which computes frame-level alignments between audio and transcripts using a CTC-trained neural network model [27]. The function `forced_align` is compatible with both CPU [20, 28] and GPU [29], providing flexibility to users. The GPU implementation is highly scalable and enables efficient processing of long audio files, and represents the first publicly available GPU-based solution for computing forced alignments.

We provide a tutorial that demonstrates how to effectively use the API. The tutorial also explains how to perform forced alignment for more than 1000 languages using the CTC-based alignment model from the Massively Multilingual Speech (MMS) project [29].

Multi-channel speech enhancement. Multi-channel speech enhancement aims to remove noise and interfering speech from multi-channel audio by leveraging spatial properties. Relative to single-channel speech enhancement, multi-channel speech enhancement can produce higher-quality outputs and further enhance the performance of downstream tasks such as ASR [30].

Estimating time-frequency masks and applying them to Minimum Variance Distortionless Response (MVDR) beamforming is an established technique capable of robustly improving multi-channel speech enhancement [31, 32, 33]. To support such work, we have implemented a time-frequency mask prediction network and an MVDR beamforming module along with a corresponding training recipe in TorchAudio.

Reference-less speech assessment. Speech assessment is essential for developing speech enhancement systems. Existing metrics require either human listening tests, e.g. Mean Opinion Score (MOS), which are expensive and unscalable, or reference clean speech, e.g. Short-Time Objective Intelligibility (STOI), Perceptual Evaluation of Speech Quality (PESQ), scale-invariant signal-to-distortion ratio (Si-SDR), which are impractical for real-world usage.

To address the limitations of such metrics, we have introduced TorchAudio-Squim [34] — TorchAudio-Speech Quality and Intelligibility Measures — which comprises two neural network based models: one for predicting objective metrics (STOI, wideband PESQ, Si-SDR), and one for predicting subjective metrics (MOS), without reference clean speech. Broadly speaking, this speech assessment feature establishes a protocol for evaluating speech enhancement

Subset	Greedy	Greedy in [19]	4-gram	4-gram in [5]
test-clean	10.9	10.5	4.4	4.3
test-other	17.8	17.6	9.5	9.4

Table 2. WER (%) results of HuBERT fine-tuning model on test subsets of the LibriSpeech dataset. Greedy refers to greedy search decoding and 4-gram to beam search decoding with a 4-gram language model.

without needing any reference signals. We present a case study of its effectiveness in Section 5.6.

5. EMPIRICAL EVALUATIONS

We demonstrate the utility of TorchAudio’s new features via studies.

5.1. Self-supervised learning

For the HuBERT recipes, we follow the methodology described in [5] of first pre-training a model and then fine-tuning it. To pre-train the model, we run two iterations of training. The first iteration trains a HuBERT model on the 960-hour LibriSpeech dataset for 250K steps, with the output targets being clusters mapped from masked frames by a 100-cluster k-means model trained on MFCC features extracted from the dataset. The second iteration trains another HuBERT model on the dataset for 400K steps, with the output targets being clusters assigned by a 500-cluster k-means model trained on intermediate feature representations generated by the first iteration’s model. Then, we fine-tune this final pre-trained model on the 10-hour LibriLight dataset with CTC loss.

Table 2 shows WERs produced in [5] and [19] by evaluating the fine-tuned "Base" model described in the original publication [5] on LibriSpeech’s test subsets, alongside WERs produced using the same model trained via our recipe and the same decoding strategies. The results validate that our HuBERT training recipes are capable of producing models of quality similar to those described in [5] and [19]. These along with the aforementioned modularity and usability benefits make the models and training recipes particularly promising for users to build upon. Indeed, Chen et al. [19] adopt TorchAudio’s HuBERT implementation and fine-tuning recipe applying slightly different training approaches, e.g. different k-means training strategies and mixed-precision training with brain floating-point (bfloat16), and achieve better performance than the original (7.6% and 7.4% relative WER improvement on test-clean and test-other subsets) while consuming far fewer GPU hours.

5.2. CTC decoder

CPU CTC decoder. The experiments in Figure 2 are conducted on LibriSpeech’s test-other set on a Wav2Vec 2.0 base model trained on 100 hours of audio. Decoding uses the official KenLM 4-gram LM and takes place on a single thread Intel® Xeon® E5-2696 v3 CPU. Because different decoder libraries support different parameters and have different underlying implementations, we first do a sweep for each decoder library for its baseline hyperparameters, and then run decoding with increasing beam sizes for additional data points. We display the wall-clock-time-WER relationship with pyctcdecode⁶ and NeMo [17], where the time in seconds is for decoding the entirety of the test-other dataset. The results show that, for a given target

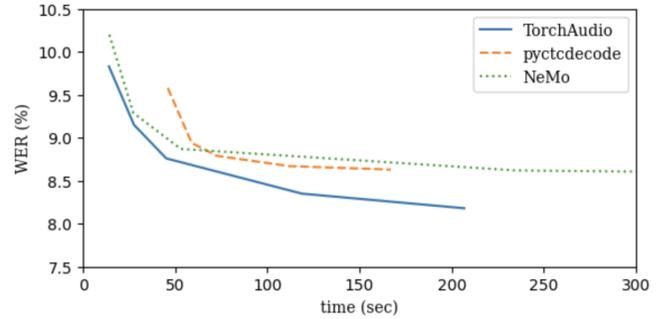


Fig. 2. Comparison of WER (%) vs. time (s) for CPU beam search decoding for TorchAudio, pyctcdecode, and NeMo.

Decoder	Configuration	WER (%)	Oracle WER (%)	Time (s)
CUDA	a = 0.95	5.81	4.11	2.57
CUDA	a = 1.0	5.81	4.09	6.24
CPU	b = 10	5.86	4.30	28.61
CPU	b = 500	5.86	4.30	791.80

Table 3. Comparison of decoding performance for TorchAudio’s CUDA and CPU CTC decoders. a is the blank frame skip threshold and b the max tokens per step.

WER, TorchAudio’s decoder runs in less time than the baselines. TorchAudio also supports a wider variety of customizable parameters for better hyperparameter tuning and overall WERs.

CUDA CTC decoder. The experiments in Table 3 are conducted on LibriSpeech’s test-other set using a single V100 GPU and Intel® Xeon® E5-2698 v4 CPU. For both recipes, a batch size of 4 and a beam size of 10 were applied. The CUDA CTC Decoder uses a CUDA kernel to implement the blank collapse method in [35]. By setting the blank frame skip threshold to 0.95, the decoding speed can be increased by 2.4 times without sacrificing accuracy. Since the CPU decoder does not support blank collapsing, the CPU decoder’s effective blank frame skip threshold is 1.0. For comparability’s sake, then, we include results for the CUDA decoder configured with a blank frame skip threshold of 1.0. By way of CUDA’s parallelism, the CUDA decoder allows for performing beam search on all tokens at every step. Thus, the CUDA decoder’s effective max tokens per step is the vocabulary size, which is 500 in this experiment. Accordingly, we include results for the CPU decoder configured with a max tokens per step of 500 to mimic the CUDA decoder’s behavior. Our experimental results show that, compared to the CPU decoder, the CUDA decoder achieves a lower WER and N-best oracle WER while increasing decoding speed by a factor of roughly 10.

5.3. Conformer

Model architecture. Rather than pursuing state-of-the-art performance, our primary goal is to validate TorchAudio’s implementations of Conformer, RNN-T loss, and data operations. Accordingly, we adopt an architecture similar to that used in the baseline Conformer

⁶<https://github.com/kensho-technologies/pyctcdecode>

Toolkit	# param.	Vocab. size	WER (%)	
			test-clean	test-other
ESPnet	91.8M/94.3M*	5000	3.1/2.8*	7.4/6.6*
TorchAudio	87.4M	1024	2.89	7.08

Table 4. Comparison of Conformer transducer recipe model performance between ESPnet and TorchAudio. *With CTC auxiliary loss.

transducer recipes in the ESPnet⁷ and icefall⁸ toolkits. To ensure comparability, we configure the model architecture to be as similar as possible to those of the baselines. As in the baselines, the encoder has a 512-d output, a subsampling rate of 4, a kernel size of 31, and 8 attention heads with a 2048-d feed-forward unit. In total, our model has 87.4M parameters, with the encoder owning 92% of them. In contrast, ESPnet’s has 94.3M parameters, while icefall’s has 84.0M. The differences in parameter count stem mostly from small differences in model architecture, which empirically do not significantly impact performance. For instance, whereas the baselines’ encoders use positional embeddings, ours does not.

Training strategy. Building upon TorchAudio’s base LibriSpeech Conformer transducer recipe, we create two training recipes, with one reproducing ESPnet’s recipe and the other icefall’s. Both apply online speed perturbation with factors uniformly sampled from {0.9, 1.0, 1.1}. The former applies SpecAugment [36] with parameters $(T, m_T, F, m_F) = (40, 2, 30, 2)$ and omits additive noise. The latter applies SpecAugment with parameters (100, 10, 27, 2) and includes additive noise, which entails sampling a waveform from MUSAN’s “noise” and “music” subsets [37] and adding it to a training sample with probability 0.5 and signal-to-noise ratio (SNR) uniformly sampled from (15, 30) dB.

Both recipes use the Adam optimizer with weight decay factor $2e-6$. Our learning rate scheduler is similar to Noam [38] in warmup (up to 40 epochs) and annealing (starting from epoch 120 with factor 0.96) steps, with the addition of an 80-epoch plateau at value $8e-4$.

Results. With comparable model architectures and training setups, our Conformer transducer recipe performs similarly to or better than ESPnet’s and icefall’s.

Table 4 shows that the performance of our recipe lies between those of the two ESPnet baselines. Compared with the baseline without CTC auxiliary loss, our recipe produces a model that achieves a 4.3%/7.8% relative improvement on test-other/clean. We note, however, that including the auxiliary loss allows the ESPnet recipe to achieve a 10.8%/9.7% relative improvement on test-other/test-clean. With the same SpecAugment policy and usage of additive noise, our model performs similarly to icefall’s on test-clean and outperforms it by 9.1% on test-other (Table 5).

5.4. Streaming AV-ASR

Datasets⁹. In this study, we use the LRS3 dataset [39], which consists of 151,819 TED Talk video clips totaling 438 hours. Following [25], we also include English-speaking videos from AVSpeech (1,323 hours) [40] and VoxCeleb2 (1,307 hours) [41] as additional

⁷<https://github.com/espnet/espnet/tree/master/egs2/librispeech/asr1#conformer-rnn-transducer>

⁸<https://github.com/k2-fsa/icefall/blob/master/egs/librispeech/ASR/RESULTS.md#2022-04-19>

⁹All data collection and processing performed at Imperial College London.

Toolkit	# param.	Vocab. size	WER (%)	
			test-clean	test-other
icefall	84.0M	500	2.59	6.15
TorchAudio	87.4M	1024	2.54	5.59

Table 5. Comparison of Conformer transducer recipe model performance between icefall and TorchAudio.

training data along with automatically-generated transcriptions. Our model is fed raw audio waveforms and face region of interests (ROIs). We do not use mouth ROIs as in [25, 42, 43] or facial landmarks or attributes during both training and testing.

Model architecture and training. We consider two configurations: Small with 12 Emformer blocks and Large with 28, with 34.9M and 383.3M parameters, respectively. Each AV-ASR model composes frontend encoders, a fusion module, an Emformer encoder, and a transducer model. We use convolutional frontends [25] to extract features from raw audio waveforms and facial images. The features are concatenated to form 1024-d features, which are then passed through a two-layer multi-layer perceptron and an Emformer transducer model [10]. The entire network is trained using RNN-T loss.

Results. Non-streaming evaluation results on LRS3 are presented in Table 6. Our audio-visual model with an algorithmic latency [10] of 800 ms (160 ms + 1280 ms × 0.5) yields a WER of 1.3%, which is on par with those achieved by state-of-the-art offline models such as AV-HuBERT, RAVEn, and Auto-AVSR. We also perform streaming evaluation adding babble acoustic noise to the raw audio waveforms at various signal-to-noise ratios. With increasing noise level, the performance advantage of our audio-visual model over our audio-only model grows (Table 7), indicating that incorporating visual data improves noise robustness. Furthermore, we measure real-time factors (RTFs) using a laptop with an Intel® Core™ i7-12700 CPU running at 2.70 GHz and an NVIDIA 3070 GeForce RTX 3070 Ti GPU. To the best of our knowledge, this is the first AV-ASR model that reports RTFs on the LRS3 benchmark. The Small model achieves a WER of 2.6% and an RTF of 0.87 on CPU (Table 8), demonstrating its potential for real-time on-device inference applications.

Method	Total Hours	WER (%)
ViT3D-CM [23]	90,000	1.6
AV-HuBERT [43]	1,759	1.4
RAVEn [42]	1,759	1.4
Auto-AVSR [25]	3,448	0.9
Ours	3,068	1.3

Table 6. Non-streaming evaluation results for audio-visual models on the LRS3 dataset.

5.5. Multi-channel speech enhancement

Datasets. To validate the efficacy of TorchAudio’s MVDR beamforming module, we use the L3DAS22 3D speech enhancement task (Task1) dataset [44] which contains 80 and 6 hours of audio for training and development, respectively. Each sample in the dataset comprises a far-field mixture recorded by two four-channel ambisonic

Type	∞ dB	10 dB	5 dB	0 dB	-5 dB	-10 dB
A	1.6	1.8	3.2	10.9	27.9	55.5
A+V	1.6	1.7	2.1	6.2	11.7	27.6

Table 7. Streaming evaluation WER (%) results at various signal-to-noise ratios for our audio-only (A) and audio-visual (A+V) models on the LRS3 dataset under 0.80-second latency constraints.

Model	Device	WER (%)	RTF
Large	GPU	1.6	0.35
Small	GPU	2.6	0.33
	CPU		0.87

Table 8. Impact of AV-ASR model size and device on WER and RTF. Note that the RTF calculation includes the pre-processing step wherein the Ultra-Lightweight Face Detection Slim 320 model is used to generate face bounding boxes.

microphone arrays and the corresponding target dry clean speech and transcript.

Model architecture and training. Experiments are conducted following the mask-based MVDR beamforming methodology described in [31]. First, a Conv-TasNet-based mask network is applied to compute the complex-valued spectrum and estimate the time-frequency masks for speech and noise. The mask network consists of a short-time Fourier transform (STFT) layer and a Conv-TasNet model with its feature encoder and decoder removed. Then, the MVDR module is applied to the masks and multi-channel spectrum to produce the beamforming weights. Finally, the beamforming weights are multiplied with the multi-channel STFT to produce a single-channel enhanced STFT from which the enhanced waveform is derived via inverse STFT. We use Ci-SDR [45] as the loss function since dry clean signals are generally not aligned with multi-channel inputs in real-world scenarios. Model configurations and training details can be found in [46].

Results. We evaluate the impact of the mask-based MVDR beamforming model alongside various baselines on downstream ASR performance. First, we evaluate each model on the test set of the L3DAS22 dataset to generate the corresponding enhanced speech. Then, we evaluate the Conformer transducer model presented in Section 5.3 on the enhanced speech and compute the WER between the generated transcriptions and the true transcriptions. Separately, we also evaluate a Wav2Vec-2.0-based ASR model on the enhanced speech and dry clean speech and compute the WER between the two sets of generated transcriptions, per the L3DAS22 Challenge’s WER metric. The results (Table 9) imply that the mask-based MVDR model significantly improves ASR performance compared to other methods, validating the efficacy of TorchAudio’s MVDR module.

5.6. Reference-less speech assessment

As discussed in Section 5.5, it can be challenging to compute signal-level speech enhancement metrics (e.g., Si-SDR) in real-world scenarios since obtaining aligned dry clean signals is difficult. Here, we

Model	WER (%)	WER* (%)
MIMO-UNet [44] (baseline)	9.4	25.0
FasNet [47]	-	14.2
DCCRN [48]	-	18.8
Demucs v3 [49]	-	15.3
Mask-based MVDR	3.5	5.6
Noisy Mixture	11.5	46.7
Dry Clean	2.6	0.0

Table 9. WER results for the mask-based MVDR beamforming model and other baselines on the test set of the L3DAS22 dataset. WER corresponds to the WER computed for the Conformer transducer model and WER* to the L3DAS22 Challenge’s WER metric.

Model	Real metrics			TorchAudio-Squim		
	WER	PESQ	Ci-SDR	STOI	PESQ	Si-SDR
MIMO-UNet [44]	9.4	1.93	8.26	0.90	1.83	10.33
mask-based MVDR	3.5	2.46	19.00	0.92	2.25	15.95
Noisy Mixture	11.5	1.21	1.87	0.67	1.22	-1.70
Dry Clean	2.6	4.64	∞	0.99	3.56	23.93

Table 10. WER, PESQ, Ci-SDR, and SQUIM metric results for the mask-based MVDR beamforming model and other baselines on the test set of the L3DAS22 dataset.

conduct a case study of TorchAudio-Squim’s utility in evaluating enhanced signals assuming such scenarios. Using TorchAudio-Squim, we estimate STOI, PESQ, and Si-SDR for the enhanced speech generated in Section 5.5. Table 10 suggests that the scores predicted by TorchAudio-Squim are consistent with actual speech quality and intelligibility.

By jointly leveraging TorchAudio’s mask-based MVDR beamforming model, Conformer transducer model, and TorchAudio-Squim, we show that one can perform multi-channel speech enhancement, ASR, and speech quality assessment all within TorchAudio.

6. CONCLUSION

TorchAudio 2.1 offers many compelling audio and speech machine learning components. Not only are its components well-designed and easy to use, but they are also effective and performant, as corroborated by our empirical results. Consequently, the library establishes a sound basis for future work in alignment with its ultimate goal of accelerating the advancement of audio technologies, and we look forward to seeing what its incredible community of users will achieve with it next.

7. ACKNOWLEDGEMENTS

We thank all of TorchAudio’s contributors on Github, including Grigory Sizov, Joel Frank, Kuba Rad, Kyle Finn, and Piotr Bialecki. We thank Andrey Talman, Danil Baibak, Eli Uriegas, Nikita Shulga, and Omkar Salpekar for helping with TorchAudio’s releases. We thank Abdelrahman Mohamed, Buye Xu, Daniel Povey, Didi Zhang, Donny Greenberg, Hung-yi Lee, Laurence Rouesnel, Matt D’Zmura, Mei-Yuh Hwang, Soumith Chintala, Thomas Lunner, Wei-Ning Hsu, and Xin Lei for the many valuable discussions.

8. REFERENCES

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, et al., “TensorFlow: Large-scale machine learning on heterogeneous systems,” 2015, Software available from tensorflow.org.
- [2] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, et al., “Pytorch: An imperative style, high-performance deep learning library,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [3] Yao-Yuan Yang, Moto Hira, Zhaoheng Ni, Anjali Chourdia, Artyom Astafurov, et al., “Torchaudio: Building blocks for audio and speech processing,” *ICASSP*, pp. 6982–6986, 2022.
- [4] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [5] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, et al., “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.
- [6] Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdel rahman Mohamed, and Michael Auli, “Unsupervised cross-lingual representation learning for speech recognition,” in *Interspeech*, 2020.
- [7] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 16, no. 6, pp. 1505–1518, 2022.
- [8] Alex Graves, Santiago Fernández, Faustino J. Gomez, and Jürgen Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” *ICML*, 2006.
- [9] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, et al., “Conformer: Convolution-augmented Transformer for Speech Recognition,” in *Interspeech*, 2020, pp. 5036–5040.
- [10] Yangyang Shi, Yongqiang Wang, Chunyang Wu, Ching-Feng Yeh, Julian Chan, et al., “Emformer: Efficient memory transformer based acoustic model for low latency streaming speech recognition,” in *ICASSP*, 2021, pp. 6783–6787.
- [11] Brian McFee, Colin Raffel, Dawen Liang, Daniel PW Ellis, Matt McVicar, et al., “librosa: Audio and music signal analysis in python,” in *SciPy*. Citeseer, 2015, vol. 8, pp. 18–25.
- [12] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, et al., “Array programming with NumPy,” *Nature*, vol. 585, no. 7825, pp. 357–362, Sept. 2020.
- [13] Piotr Żelasko, Daniel Povey, Jan Trmal, and Sanjeev Khudanpur, “Lhotse: a speech data representation library for the modern deep learning ecosystem,” *ArXiv*, vol. abs/2110.12561, 2021.
- [14] Shinji Watanabe, Takaaki Hori, Shigeki Karita, Tomoki Hayashi, Jiro Nishitoba, et al., “ESPnet: End-to-end speech processing toolkit,” in *Interspeech*, 2018, pp. 2207–2211.
- [15] Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, et al., “Speechbrain: A general-purpose speech toolkit,” *arXiv preprint arXiv:2106.04624*, 2021.
- [16] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, et al., “fairseq: A fast, extensible toolkit for sequence modeling,” in *North American Chapter of the Association for Computational Linguistics*, 2019, pp. 48–53.
- [17] Oleksii Kuchaiev, Jason Li, Huyen Nguyen, Oleksii Hrinchuk, Ryan Leary, et al., “Nemo: a toolkit for building ai applications using neural modules,” *arXiv preprint arXiv:1909.09577*, 2019.
- [18] Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang, Cheng-I Jeff Lai, Kushal Lakhotia, et al., “Superb: Speech processing universal performance benchmark,” *arXiv preprint arXiv:2105.01051*, 2021.
- [19] William Chen, Xuankai Chang, Yifan Peng, Zhaoheng Ni, Soumi Maiti, et al., “Reducing barriers to self-supervised learning: Hubert pre-training with academic compute,” *arXiv preprint arXiv:2306.06672*, 2023.
- [20] Jacob D Kahn, Vineel Pratap, Tatiana Likhomanenko, Qiantong Xu, Awni Hannun, et al., “Flashlight: Enabling innovation in tools for machine learning,” in *ICML*. PMLR, 2022, pp. 10557–10574.
- [21] Kenneth Heafield, “Kenlm: Faster and smaller language model queries,” in *Proceedings of the Sixth Workshop on Statistical Machine Translation, USA, 2011, WMT ’11*, p. 187–197, Association for Computational Linguistics.
- [22] Pengcheng Guo, Florian Boyer, Xuankai Chang, Tomoki Hayashi, Yosuke Higuchi, et al., “Recent developments on espnet toolkit boosted by conformer,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5874–5878.
- [23] Dmitriy Serdyuk, Otavio Braga, and Olivier Siohan, “Transformer-based video front-ends for audio-visual speech recognition for single and multi-person video,” in *Interspeech*, 2022, pp. 2833–2837.
- [24] Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdelrahman Mohamed, “Learning audio-visual speech representation by masked multi-modal cluster prediction,” in *ICLR*, 2022.
- [25] Pingchuan Ma, Alexandros Haliassos, Adriana Fernandez-Lopez, Honglie Chen, Stavros Petridis, et al., “Auto-avr: Audio-visual speech recognition with automatic labels,” in *ICASSP*, 2023, pp. 1–5.
- [26] Pingchuan Ma, Niko Moritz, Stavros Petridis, Christian Fuegen, and Maja Pantic, “Streaming audio-visual speech recognition with alignment regularization,” in *Interspeech*, 2023.
- [27] Ludwig Kürzinger, Dominik Winkelbauer, Lujun Li, Tobias Watzel, and Gerhard Rigoll, “Ctc-segmentation of large corpora for german end-to-end speech recognition,” in *International Conference on Speech and Computer*. Springer, 2020, pp. 267–278.
- [28] Vineel Pratap, Awni Hannun, Qiantong Xu, Jeff Cai, Jacob Kahn, et al., “Wav2letter++: A fast open-source speech recognition system,” in *ICASSP*. IEEE, 2019, pp. 6460–6464.
- [29] Vineel Pratap, Andros Tjandra, Bowen Shi, Paden Tomasello, Arun Babu, et al., “Scaling speech technology to 1,000+ languages,” 2023.
- [30] Reinhold Haeb-Umbach, Shinji Watanabe, Tomohiro Nakatani, Michiel Bacchiani, Bjorn Hoffmeister, et al., “Speech processing for digital home assistants: Combining signal processing with deep-learning techniques,” *IEEE Signal processing magazine*, vol. 36, no. 6, pp. 111–124, 2019.
- [31] Shoko Araki, Masahiro Okada, Takuya Higuchi, Atsunori Ogawa, and Tomohiro Nakatani, “Spatial correlation model based observation vector clustering and mvdr beamforming for meeting recognition,” in *ICASSP*. IEEE, 2016, pp. 385–389.
- [32] Hakan Erdogan, John R Hershey, Shinji Watanabe, Michael I Mandel, and Jonathan Le Roux, “Improved mvdr beamforming using single-channel mask prediction networks,” in *Interspeech*, 2016, pp. 1981–1985.
- [33] Zhuohuang Zhang, Yong Xu, Meng Yu, Shi-Xiong Zhang, Lianwu Chen, et al., “Adl-mvdr: All deep learning mvdr beamformer for target speech separation,” in *ICASSP*. IEEE, 2021, pp. 6089–6093.
- [34] Anurag Kumar, Ke Tan, Zhaoheng Ni, Pranay Manocha, Xiaohui Zhang, et al., “Torchaudio-squim: Reference-less speech quality and intelligibility measures in torchaudio,” in *ICASSP*. IEEE, 2023, pp. 1–5.
- [35] Minkyu Jung, Ohhyeok Kwon, Seunghyun Seo, and Soonshin Seo, “Blank collapse: Compressing ctc emission for the faster decoding,” *arXiv preprint arXiv:2210.17017*, 2022.
- [36] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, et al., “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Interspeech*. sep 2019, ISCA.
- [37] David Snyder, Guoguo Chen, and Daniel Povey, “MUSAN: A Music, Speech, and Noise Corpus,” 2015, arXiv:1510.08484v1.
- [38] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, et al., “Attention is all you need,” in *NeurIPS*, 2017.

- [39] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman, “LRS3-TED: a large-scale dataset for visual speech recognition,” *arXiv preprint arXiv:1809.00496*, 2018.
- [40] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, et al., “Looking to listen at the cocktail party: a speaker-independent audio-visual model for speech separation,” *ACM Transactions on Graphics*, vol. 37, no. 4, pp. 112:1–112:11, 2018.
- [41] J Chung, A Nagrani, and A Zisserman, “Voxceleb2: Deep speaker recognition,” *Interspeech*, 2018.
- [42] Alexandros Haliassos, Pingchuan Ma, Rodrigo Mira, Stavros Petridis, and Maja Pantic, “Jointly learning visual and auditory speech representations from raw data,” in *ICLR*, 2023.
- [43] Bowen Shi, Wei-Ning Hsu, and Abdelrahman Mohamed, “Robust self-supervised audio-visual speech recognition,” in *Proceedings of Interspeech*, 2023, pp. 2118–2122.
- [44] Eric Guizzo, Christian Marinoni, Marco Pennese, Xinlei Ren, Xiguang Zheng, et al., “L3das22 challenge: Learning 3d audio sources in a real office environment,” in *ICASSP*. IEEE, 2022, pp. 9186–9190.
- [45] Christoph Boeddeker, Wangyou Zhang, Tomohiro Nakatani, Keisuke Kinoshita, Tsubasa Ochiai, et al., “Convolutional transfer function invariant sdr training criteria for multi-channel reverberant speech separation,” in *ICASSP*. IEEE, 2021, pp. 8428–8432.
- [46] Yen-Ju Lu, Samuele Cornell, Xuankai Chang, Wangyou Zhang, Chenda Li, et al., “Towards low-distortion multi-channel speech enhancement: The espnet-se submission to the l3das22 challenge,” in *ICASSP*. IEEE, 2022, pp. 9201–9205.
- [47] Yi Luo, Cong Han, Nima Mesgarani, Enea Ceolini, and Shih-Chii Liu, “Fasnet: Low-latency adaptive beamforming for multi-microphone audio processing,” in *ASRU*. IEEE, 2019, pp. 260–267.
- [48] Yanxin Hu, Yun Liu, Shubo Lv, Mengtao Xing, Shimin Zhang, et al., “Dccrn: Deep complex convolution recurrent network for phase-aware speech enhancement,” *Interspeech*, 2020.
- [49] Alexandre Défossez, “Hybrid spectrogram and waveform source separation,” *arXiv preprint arXiv:2111.03600*, 2021.

A. USAGE EXAMPLES

```
# Stream audio and video from media devices
r = torchaudio.io.StreamReader(
    # Use cameras and mics in macOS
    format="avfoundation",
    # Pick the first video device and
    # third audio device
    src="0:2",
    # Configure the frame rate of the camera
    option={"framerate": "30"},
)
r.add_basic_audio_stream(
    frames_per_chunk=1600,
    sample_rate=8000,
)
r.add_basic_video_stream(
    frames_per_chunk=5,
    frame_rate=25,    # change frame rate
    width=224,       # change width and height
    height=196,
)
for (audio, video,) in r.stream():
    # audio.shape == [frames, channels]
    # video.shape ==
    #     [frames, channels, height, width]
```

Fig. 3. Sample code that uses `StreamReader` to process streaming input from media devices.

```
# Live stream audio and video using
# Real-Time Messaging Protocol
w = torchaudio.io.StreamWriter(
    dst="rtmp://localhost:80", format="flv")
w.add_audio_stream(
    sample_rate=8000, num_channels=1)
w.add_video_stream(
    frame_rate=30, width=128, height=96)

with w.open(option={"listen": "1"}):
    for (video_chunk, audio_chunk) in generator():
        w.write_audio_chunk(0, audio_chunk)
        w.write_video_chunk(1, video_chunk)
```

Fig. 4. Sample code that uses `StreamWriter` to live-stream audio-visual data via Real-Time Messaging Protocol server.