

CONFORMER BASED AUTOMATIC SPEECH RECOGNITION WITH CTC

About Dataset

- Dataset: LibriSpeech “train-clean-100” subset
- Sampling rate: 16 kHz
- Average duration per utterance: ~10 seconds
- 100 hours of clean English speech
- Used for both training and testing

Baseline Model Setup

- Trained, full-precision (FP32) Conformer-CTC model as the reference.
- Evaluate on a subset of LibriSpeech ‘test-clean’ to record baseline metrics:
- Word Error Rate (WER) – measures recognition accuracy
- Model size – memory footprint
- Inference latency – time per input

Baseline ASR

- **Input:** Audio Input from the LibriSpeech dataset (Train/Test splits, 16 kHz).
- **Pre-processing:** Initial audio processing such as resampling and cropping.
- **Feature Extraction:** Preprocessing Mel Spectrogram Transformation.
- **Temporal Modeling:** The Downsample Convolutional Frontend followed by the CONFORMER ENCODER
- **Alignment:** The Classifier Head produces Frame Wise Probabilities which are aligned to the final transcription via CTC Greedy Decoding.
- **Training:** Adam optimizer, CTC loss.
- **Test-Output:** Final Text Transcription.

TRAINING



