

# Basic Speech Recognition with Pre-trained Wav2Vec2

## 1. Problem Statement and Objectives :

The objective of this work is to demonstrate a basic speech recognition pipeline using a pre-trained Wav2Vec2 model. Specifically, we aim to load a small subset of the LibriSpeech Automatic Speech Recognition (ASR) dataset, perform inference using the "facebook/wav2vec2-base-960h" model, and evaluate the model's performance on this data using the Word Error Rate (WER) metric. This serves as a foundational experiment to understand the capabilities of state-of-the-art pre-trained models for speech-to-text tasks.

## 2. Explanation of Experimental Setup and Methodology :

The experiment was conducted using the Python programming language with the Hugging Face Transformers and Datasets libraries. The following steps were involved:

**Dataset Loading:** A small subset (10 samples) of the "clean" split of the "train.100" configuration from the LibriSpeech ASR dataset was loaded using the `load_dataset` function in streaming mode. This allowed for efficient handling of the data.

**Model and Processor Loading:** The pre-trained Wav2Vec2 base model ("facebook/wav2vec2-base-960h") and its corresponding processor (Wav2Vec2Processor) were loaded using the `from_pretrained` method. The processor is essential for converting raw audio into the model's input format and decoding the model's output.

**Inference:** Each audio sample in the tiny dataset was processed individually. The audio waveform was converted to a PyTorch tensor and passed through the processor to obtain the input features. The pre-trained Wav2Vec2 model, set to evaluation mode (`model.eval()`), was then used to generate logits. The predicted token IDs were obtained by taking the `argmax` of the logits, and these IDs were subsequently decoded into text transcriptions using the processor's `batch_decode` method. The predicted transcriptions were converted to lowercase for consistent evaluation.

**Evaluation:** The Word Error Rate (WER) metric was used to evaluate the performance of the model. The `wer` metric from the Hugging Face Evaluate library was employed to calculate the WER between the predicted transcriptions and the ground truth reference texts from the dataset. Both sample-wise WER and an overall WER across the 10 samples were computed.

## 3. Results, Observations, and Analysis :

The sample-wise results, as observed in the executed code, generally showed a high degree of accuracy. Many of the predicted transcriptions were identical or very close to the original text, resulting in a low WER for those samples.

The overall Word Error Rate (WER) calculated across the 10 samples was 0.0083. This translates to an approximate accuracy of 99.17%.

#### Observations:

The pre-trained Wav2Vec2 model demonstrated an exceptional ability to transcribe the speech in this small subset of the LibriSpeech dataset with very few errors. This underscores the effectiveness of the pre-training process on a large and diverse corpus.

The high accuracy suggests that the audio quality in the "clean" subset is likely very good and aligns well with the data the Wav2Vec2 model was trained on.

The minimal errors observed might be due to subtle phonetic differences or very short, common words where misinterpretations are less likely.

#### Analysis:

The results of this basic experiment indicate that the pre-trained Wav2Vec2 model performs remarkably well on this small, clean subset of the LibriSpeech dataset, achieving a very low Word Error Rate and a high degree of accuracy. This highlights the strong capabilities of this model for speech recognition tasks without any fine-tuning. While these results are promising, it's important to note that the performance might vary on more challenging audio with background noise, different accents, or more complex vocabulary. Further evaluation on a larger and more diverse test set would provide a more comprehensive understanding of the model's generalization ability. This experiment serves as a successful initial demonstration of leveraging pre-trained models for speech recognition.