

Task 4

Compare your finetuned model from task 3 to be compared against the results from task 2a for the cv-valid-dev mp3 dataset. Describe your observations and propose a series of steps (including datasets and experiments) to improve the accuracy. Your answer can be saved as training-report.pdf under the main repository

Please find the answer at the end of this document.

```
In [1]:
import os
import pandas as pd
import torch
from transformers import Wav2Vec2ForCTC, Wav2Vec2Processor
import librosa
import numpy as np
import evaluate
from tqdm import tqdm
import matplotlib.pyplot as plt

In [2]:
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

In [3]:
# Load the model and tokenizer
model_name = "facebook/wav2vec2-base-960h"
model = Wav2Vec2ForCTC.from_pretrained(model_name).to("cuda" if torch.cuda.is_available() else "cpu")
processor = Wav2Vec2Processor.from_pretrained(model_name)

Some weights of Wav2Vec2ForCTC were not initialized from the model checkpoint at facebook/wav2vec2-base-960h and are newly initialized: ['wav2vec2.masked_sp
ec_embed']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

In [4]:
# Path to the folder that contains the mp3 files and the CSV file for cv-valid-dev
cv_folder = 'C:/Users/lingy/OneDrive/Desktop/HTX/HTX_techtest/common_voice' # Update this to your local path
cv_valid_dev_audio_folder = os.path.join(cv_folder, 'cv-valid-dev')
cv_valid_dev_csv_file = os.path.join(cv_folder, 'cv-valid-dev.csv')

# Load the CSV file into a dataframe
cv_valid_dev_df = pd.read_csv(cv_valid_dev_csv_file)

# Update the 'filename' column to include the full path to the audio files
cv_valid_dev_df['filename'] = cv_valid_dev_df['filename'].apply(lambda x: os.path.join(cv_valid_dev_audio_folder, x))

In [5]:
def preprocess_audio(file_path, target_sr=16000):
    # Load audio file using librosa
    audio, sr = librosa.load(file_path, sr=None)

    # Resample audio if needed
    if sr != target_sr:
        audio = librosa.resample(audio, orig_sr=sr, target_sr=target_sr)

    # Normalize the audio to [-1, 1]
    audio = audio / np.max(np.abs(audio))

    # Convert to tensor for Wav2Vec2 processor
    input_values = processor(audio, sampling_rate=target_sr, return_tensors="pt").input_values
    input_values = input_values.squeeze(0)
    input_values = input_values.unsqueeze(0)
    return input_values # Remove batch dimension

In [6]:
```

```
In [ ]:  
# Create a new column with transcriptions  
cv_valid_dev_df['generated_text_fine-tuned'] = cv_valid_dev_df['filename'].apply(lambda x: transcribe_audio(x))  
  
# Save the updated dataframe to a new CSV file  
cv_valid_dev_df['generated_text_fine-tuned'] = cv_valid_dev_df['generated_text_fine-tuned'].str.lower()  
cv_valid_dev_df.to_csv(cv_valid_dev_csv_file, index=False)  
In [11]:  
# Initialize the WER metric from the evaluate library  
wer_metric = evaluate.load("wer")  
In [12]:  
# Define a function to compute WER for the text generated from "fine-tuned" model with the groundtruth for cv_valid_dev dataset  
def evaluate_model_generated_text_fine_tuned(df):  
    predictions = []  
    references = []  
  
    for idx, row in tqdm(df.iterrows(), total=len(df)):  
        # Get the file path and ground truth transcription  
        predicted_text = row['generated_text_fine-tuned']  
        ground_truth = row['text']  
  
        # Append predictions and references for evaluation  
        predictions.append(predicted_text)  
        references.append(ground_truth)  
  
    # Compute the Word Error Rate (WER)  
    wer_score = wer_metric.compute(predictions=predictions, references=references)  
    return wer_score  
In [13]:  
# Evaluate the text generated from "fine-tuned" model with the groundtruth for cv_valid_dev dataset  
wer_score = evaluate_model_generated_text_fine_tuned(cv_valid_dev_df)  
  
# Log the overall performance  
print(f"Word Error Rate (WER) on the cv_valid_dev set (groundtruth VS fine-tuned): {wer_score * 100:.2f}%")  
100% ██████████ 4076/4076 [00:00<00:00, 208  
3it/s]  
Word Error Rate (WER) on the cv_valid_dev set (groundtruth VS fine-tuned): 14.39%  
In [19]:
```

In []: