# Configuration Manual

MSc Research Project
Data Analytics

# Zeba Siddique

Student ID: x20227086

School of Computing
National College of Ireland

Supervisor:    Dr. Abubakr Siddig

| | |
|---|---|
| **Student Name:** | Zeba Siddique |
| **Student ID:** | x20227086 |
| **Programme:** | Data Analytics |
| **Year:** | 2022 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Dr. Abubakr Siddig |
| **Submission Due Date:** | 15/08/2022 |
| **Project Title:** | Configuration Manual |
| **Word Count:** | 1126 |
| **Page Count:** | 29 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | Zeba Siddique |
| **Date:** | 13th August 2022 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Zeba Siddique
x20227086

# 1 Introduction

This document is a configuration manual that contains all the necessary information required to achieve a child speech synthesis artefact. It comprises of the minimum and must-have system requirements for reproducing the thesis work. Along with the necessary system configurations and pre-requisites, this document with the help of code snippets describes the main blocks of the thesis. The document details the step-by-step instructions from data collection to generating results to executing the artefact.

# 2 Required Specifications

## 2.1 Hardware requirements

Figure 1 describes the hardware specifications (device and windows) used to carry out the thesis.



**Device Specifications**

HP Spectre x360 Convertible 14-ea0xxx

| | |
|---|---|
| Device name | DESKTOP-KHUORMC |
| Processor | 11th Gen Intel(R) Core (TM) i7-1165G7 @ 2.80GHz   2.80 GHz |
| Installed RAM | 16.0 GB (15.6 GB usable) |
| Device ID | B3992509-FEF9-4008-A80E-1FD4E8A47312 |
| Product ID | 00325-97258-57122-AAOEM |
| System type | 64-bit operating system, x64-based processor |
| Pen and touch | Pen and touch support with 10 touch points |

**Windows Specifications**

Windows 10

| | |
|---|---|
| Edition | Windows 10 Home |
| Version | 21H2 |
| Installed on | 23-Sep-21 |
| OS build | 19044.1889 |
| Experience | Windows Feature Experience Pack 120.2212.4180.0 |

Figure 1: Hardware Specifications

## 2.2 Software requirements

For execution of the artefacts, the below mentioned list of software must be installed on the system.

1. Anaconda Navigator for Windows (Version 4.12.0)

2. Python 3.9.13

3. Visual Studio Code (Version 1.70.0 [user setup])

4. VS Extensions:

    (a) Jupyter (v2022.7.1102252217)
    (b) Jupyter Keymap (v1.0.0)
    (c) Jupyter Notebook Renderers (v1.0.9)
    (d) Pylance (v2022.8.20)
    (e) Python (v2022.12.0)

5. Google Chrome (Version 104.0.5112.81)

## 2.3 Storage requirements/ Products/ Subscriptions

Following are the additional and essential requirements used to carry out the thesis:

1. Google Colaboratory Pro+

2. Google One/ Drive - 2TB storage

# 3 Data Collection

This research uses a freely available multi-speaker child speech dataset: My Science Tutor (MyST). This dataset was obtained through a shared drive after receiving necessary permissions from concerned authorities and agreeing to terms and conditions for the research license agreement. Access to MyST Corpus can be requested from the official website of the dataset [1]. The dataset contains audio references (.flac) and their transcripts (.trn) if any for each speaker based on the recording phase. The MyST corpus consists of 456 hours speech data from 1,371 students.

# 4 Data Cleaning and Pre-processing

Out of the total child speech data available in the MyST corpus only 45% of the audio references are transcribed. This thesis focuses on the audio references that have a transcript file. The data is sampled using the non-probability based purposive sampling to use the transcribed audio references for phase 2 sessions. Based on the works of Jain et al. (2022) on the same corpus, the MyST corpus is processed. The data cleaning and pre-processing steps undertaken for MyST corpus for this research are described below:

---

[1]MyST Corpus: `https://boulderlearning.com/request-the-myst-corpus/`

## 4.1 Import Libraries for preparing the MyST corpus

Figure 2 demonstrates the libraries to be imported to pre-process the MyST corpus. The figure also mentions the different .py files that facilitate the cleaning and pre-processing of MyST corpus.

```
MyST Preprocessing > DataPrep > libs > 🐍 flags.py > ...
    1    from enum import Enum
    2

MyST Preprocessing > DataPrep > 🐍 extract_readings.py > ...
    1    import asyncio
    2    import logging
    3    from typing import Any
    4
    5    from libs.flags import Dataset
    6    from libs.data_placement import DataSetChanges
    7

MyST Preprocessing > DataPrep > libs > 🐍 data_placement.py > ...
    1    import os
    2    import string
    3    import librosa
    4    import matplotlib.pyplot as plt
    5    from multimethod import RETURN
    6    import numpy as np
    7    import pandas as pd
    8    import soundfile as sf
    9    import librosa.display
   10
   11    from glob import glob
   12    from typing import Dict, List
   13    from libs.flags import Dataset
   14    from asyncio.log import logger
   15
```

Figure 2: Required Python Libraries for MyST corpus

## 4.2 Creating Constants

Figure 3 demonstrates the constants used during the cleaning and pre-processing of the MyST corpus.

```
MyST Preprocessing > DataPrep > libs > 🐍 flags.py > ...
    1    from enum import Enum
    2
    3    class Dataset(Enum):
    4        NAME = "My Science Tutor"
    5        URL = ""
    6        ACTUAL_PATH  = r'C:\Users\zzeba\Documents\Zeba\NCI\Sem3\Thesis\myst-v0.4.2\data'
    7        MODIFIED_PATH = r'C:\Users\zzeba\Documents\Zeba\NCI\Sem3\Thesis\MySTDataset'
    8        DEV_RUN = "development"
    9        TRAIN_RUN = "train"
   10        TEST_RUN = "test"
   11        RECORDING_CSV_FILE_PATH = r'C:\Users\zzeba\Desktop\Child TTS - MyST\MyST Preprocessing\DataPrep\InitialAudioDetails.csv'
   12        METADATA_FILE_PATH = r'MyST Preprocessing\MyST_DataSet\metadata.txt'
   13
```

Figure 3: Defining Constants for MyST corpus

## 4.3 Data Cleaning

Figure 4 demonstrates the execution point for cleaning the MyST corpus.

```
MyST Preprocessing > DataPrep > 🐍 extract_readings.py > ...
27    async def extract_readings(run_type: str) -> None:
28        """Pre-process the MyST Speech corpus
29        """
30        try:
31            await DataSetChanges._read_dataset(run_type)
32        except:
33            logger.exception('Get Dataset Details: %s')
34
35    if __name__ == '__main__':
36        params_myst = {
37            'run_type' : 'development'
38        }
39        handler(params_myst, None)
40
41
MyST Preprocessing > DataPrep > libs > 🐍 data_placement.py > ♦ _extract_recording_details
16    class DataSetChanges():
17
18        async def _read_dataset(run_type: str) -> None:
19            """Navigate to the file path and wait for processing to finish"""
20            await _extract_recording_details() # step 1
21            recording_df = await _read_recording_details_from_csv() # step 2
22            await _remove_silence_from_audio_files(recording_df=recording_df) # step 3
23
24
```

Figure 4: Code Execution Point

Figure 5 demonstrates the selection of transcribed audio references that are 10-15 seconds in length and free from dis-fluency markers. The punctuation in the transcripts are replaced. Additionally, the details of the selected transcribed audio files are saved in a CSV file for further processing

```
MyST Preprocessing > DataPrep > libs > 🐍 data_placement.py > ...
25    async def _extract_recording_details():
26        """ Get the valid audio references and transcripts from the
27        MyST corpus
28        """
29        duration_of_recordings = []
30        root_directory = Dataset.ACTUAL_PATH.value
31        batch_counter = 1
32
33        speech_data_details = []
34
35        for dir_path, sub_directories, files in os.walk(root_directory):
36            for filename in files:
37                is_valid_session = await _check_for_phase_two_sessions(directory_path=dir_path)
38                if is_valid_session:
39                    if filename.endswith(".flac"):
40                        flac_file = os.path.join(dir_path,filename)
41                        trn_file = flac_file.replace(".flac", ".trn")
42                        is_valid_transcript, transcript_content = await _check_for_valid_transcript(trn_file=trn_file)
43                        if is_valid_transcript :
44                            speech_data = await _get_data_details_in_row(
45                                flac_file=flac_file, trn_file=trn_file,
46                                transcript_content=transcript_content
47                            )
48                            if speech_data != []:
49                                speech_data_details.append(speech_data)
50                                if (len(speech_data_details) % 100 == 0):
51                                    await _write_speech_details_to_file(speech_data_details)
52                                    speech_data_details = []
53                                    print("Batch completed : ", batch_counter)
54                                    batch_counter += 1
55                        else:
56                            continue
57        if len(speech_data_details) > 0:
58            await _write_speech_details_to_file(speech_data_details)
59            speech_data_details = []
60            print("Batch completed as last : ", batch_counter)
61
62        print ("function finished execution")
```

Figure 5: Selection of Audio-Transcripts from MyST corpus

Figures 6, 7, and 8 demonstrates the transcript cleaning, audio selection, and creation of CSV file respectively.

```python
64    async def _check_for_phase_two_sessions(directory_path):
65        """Pick audio references that are recorded
66        for Phase 2 sessions of the MyST Corpus
67        Purposive sampling: Non-probability"""
68        isValid = False
69        phase_two_sessions = ['_EE_', '_MX_', '_SMP_', '_SRL_', '_LS_']
70        for session_name in phase_two_sessions:
71            if session_name in directory_path:
72                isValid = True
73        return isValid
74
75    async def _check_for_valid_transcript(trn_file):
76        """ Check for transcripts that do not have any
77        disfluency markers. Return False if a the transcript contains
78        disfluency else remove the punctuations from the transcript
79        if exist and return True
80        """
81        isValid = True
82        transcript_content = None
83        disfluency_markers = ['<NO_SIGNAL>',
84        '<SILENCE>', '<BREATH>', '<LAUGH>', '<COUGH>',
85        '<NOISE>', '<SIDE_SPEECH>', '<SNIFF>', '<ECHO>',
86        '<DISCARD>', '<INDISCERNIBLE>', '(*)', '(())']
87        if os.path.exists(trn_file):
88            with open(trn_file, 'r', encoding='ascii', errors='ignore') as transcript_file:
89                transcript_content = transcript_file.read()
90                transcript_content = transcript_content.upper()
91                #print(transcript_content)
92                if any(ele in transcript_content for ele in disfluency_markers):
93                    isValid = False
94                    transcript_content = None
95        else:
96            isValid = False
97
98        if isValid :
99            transcript_content = ''.join([i for i in transcript_content if i not in string.punctuation])
100            #print("Removed Punctuation")
101        return isValid, transcript_content
```

Figure 6: Purposive Sampling and Transcript Cleaning

6

```python
104  async def _get_data_details_in_row(flac_file : str, trn_file : str, transcript_content : str):
105      """ Get valid audio reference and transcript file names along
106      with additional details in an array
107      """
108      speech_details = []
109
110      samples, sample_rate = librosa.load(flac_file)
111      duration_of_recording= round(librosa.get_duration(y = samples, sr = sample_rate),2)
112      if duration_of_recording >= 10 and duration_of_recording <= 15:
113          speech_details = [
114              _get_speaker_id(file_name=flac_file),
115              duration_of_recording,
116              flac_file,
117              trn_file,
118              transcript_content.strip(),
119              sample_rate
120          ]
121      return speech_details
122
123  def _rows_to_dataframe(speech_data_details : List) -> pd.DataFrame:
124      """Convert the array into a pandas dataframe"""
125      speech_details_df = pd.DataFrame(
126          speech_data_details,
127          columns=[
128              'SpeakerID',
129              'AudioLength',
130              'AudioFilePath',
131              'TranscriptFilePath',
132              'TranscriptText',
133              'AudioSampleRate'
134          ]
135      )
136
137      return speech_details_df
138
139
140  def _get_speaker_id(file_name : str):
141      """Get the speaker ID from a file name"""
142      return ((file_name.split("myst_"))[1]).split("_")[0]
```

Figure 7: Selection of Audio and Creation of DataFrame

```python
143
144  async def _write_speech_details_to_file(speech_data_details):
145      """Put the details of the valid files in a text file"""
146      speech_data_dataframe = _rows_to_dataframe(speech_data_details)
147      file_name = r'C:\Users\zzeba\Desktop\MyST Preprocessing\DataPrep\InitialAudioDetails.csv'
148      if not os.path.exists(file_name):
149          speech_data_dataframe.to_csv('DataPrep/InitialAudioDetails.csv',
150          index=False, encoding='utf-8-sig'
151          )
152          print("Created CSV file")
153      else:
154          speech_data_dataframe.to_csv('DataPrep/InitialAudioDetails.csv',
155          mode='a', header=False, index=False, encoding='utf-8-sig'
156          )
157          print("Appended to CSV file")
```

Figure 8: MyST CSV File Creation

## 4.4 Data Pre-processing

For performing data pre-processing on the MyST corpus, the CSV file is retrieved to obtain the cleaned audio references. The following figures 9 and 10 demonstrate the steps carried on the MyST corpus.

```python
async def _read_recording_details_from_csv():
    """Read the CSV file consisting details of recordings and return a pandas dataframe"""
    if os.path.exists(Dataset.RECORDING_CSV_FILE_PATH.value):
        recording_df = pd.read_csv(Dataset.RECORDING_CSV_FILE_PATH.value)
    return recording_df

async def _remove_silence_from_audio_files(recording_df):
    """Remove all the silence between the audio files"""
    if len(recording_df) > 0:
        try:
            for row in recording_df.itertuples(index=True, name='Pandas'):
                if os.path.exists(row.AudioFilePath):
                    wav_file_name = (row.AudioFilePath.split("\\")[-1]).split(".flac")[0]
                    wav_file_name_with_extension = wav_file_name + ".wav"
                    wav_file_path = r'MyST Preprocessing\MyST_DataSet\wavs' + "\\" + wav_file_name_with_extension

                    audio_file, sample_rate = librosa.load(row.AudioFilePath, sr= row.AudioSampleRate, mono=True)
                    # remove the silence from between the voice signals
                    clips = librosa.effects.split(audio_file, top_db=100)

                    # combine audio clips without silence
                    wav_data = []
                    if(len(clips) > 0):
                        for clip in clips:
                            audio_data = audio_file[clip[0]: clip[1]]
                            wav_data.extend(audio_data)
                        # save the wav file in the folder
                        sf.write(wav_file_path, wav_data, row.AudioSampleRate)
                    else:
                        print("no clips")
                        sf.write(wav_file_path, audio_file, row.AudioSampleRate)
                    # write the audio, transcript, and speaker details into a txt file
                    await _write_transcript_to_metadata(row, wav_file_name)
        except:
            logger.exception('Error in _remove_silence_from_audio_files: %s')
        finally:
            logger.exception('Finished executing _remove_silence_from_audio_files')
```

Figure 9: MyST Data Pre-processing

```
198          """Write details of the audio recording to a txt file"""
199          file_path = r'MyST Preprocessing\MyST_DataSet\metadata.txt'
200          audio_info = audio_file_name + "|" + audio_details.TranscriptText + "|" + str(audio_details.SpeakerID) + "\n"
201          file_mode = ""
202          try:
203              if not os.path.exists(file_path):
204                  file_mode = "w"
205              else:
206                  file_mode = "a"
207              with open(file_path, mode = file_mode) as f:
208                  f.write(audio_info)
209          except:
210              logger.exception('Error in _write_transcript_to_metadata: %s')
```

Figure 10: MyST Data Metadata

Figure 11 demonstrates the distribution of duration of the pre-processed audio reference files.

```
244      async def _plot_audio_length_distribution(speech_data_dataframe : pd.DataFrame):
245          # Plot audio duraiton distribution
246          n, bins, patches = plt.hist(x=speech_data_dataframe['AudioLength'], bins='auto', color='#0504aa',
247                                      alpha=0.7, rwidth=0.85)
248          plt.grid(axis='y', alpha=0.75)
249          plt.xlabel('Duration in seconds')
250          plt.ylabel('Number of Audio')
251          plt.title('Distribution of audio duration')
252          maxfreq = n.max()
253          # Set a clean upper y-axis limit.
254          plt.ylim(ymax=np.ceil(maxfreq / 10) * 10 if maxfreq % 10 else maxfreq + 10)
255          plt.show()
```



Figure 11: Pre-processed Audio Distribution

9

# 5  Implementation of the HiFi-GAN vocoder

This section describes the implementation of the HiFi-GAN vocoder that uses the python TTS library (for training from scratch) [2] and the coqui-ai TTS Github repository (for fine-tuning) [3]. The rest of the artefeacts are executed on Google Colab Pro+ with hardware accelerator selected as 'GPU', with High Run-time and background execution enabled. Note: The pre-processed MyST corpus having folder structure as

- Folder: DataSetName
    - Folder: wavs (containing all pre-processed files)
    - Folder: wavs1 (containing 10% files from the wavs folder)
    - File: metadata.txt

## 5.1  Run pre-requisites



Figure 12: Colab Pre-requisites for HiFi-GAN

---

[2]Python TTS library `https://pypi.org/project/TTS/`
[3]Python TTS library `https://github.com/coqui-ai/TTS`

Figure 13: Mount Google Drive

## 5.2 Import Libraries for implementing the HiFi-GAN vocoder

The following Figure demonstrates all the necessary libraries required to implement the HiFi-GAN vocoder to be trained on the MyST corpus.



Figure 14: Import Libraries for HiFi-GAN vocoder

## 5.3 Training the HiFi-GAN vocoder on MyST corpus

The implementation of the HiFi-GAN vocoder for this research uses the HiFi-GAN generator [4], discriminator [5] and the training config file [6]. The modifications in the configurations are demonstrated in the subsequent figures.



Figure 15: Modifying the parameters for HiFi-GAN vocoder

---

[4]HiFi-GAN Generator https://github.com/coqui-ai/TTS/blob/dev/TTS/vocoder/models/hifigan_generator.py

[5]HiFi-GAN Discriminator https://github.com/coqui-ai/TTS/blob/dev/TTS/vocoder/models/hifigan_discriminator.py

[6]HiFi-GAN Config https://github.com/coqui-ai/TTS/blob/dev/TTS/vocoder/configs/hifigan_config.py

Figure 16: Initiate Audio Processor for HiFi-GAN vocoder



Figure 17: Initiate Trainer Class for HiFi-GAN vocoder

The following figure displays the HiFi-GAN vocoder trained from scratch on pre-processed MyST corpus for 500 epochs

```
Step 12: Start the training of the HiFi-GAN vocoder model from scratch.

Note that: Observe the evaluation results to understand the loss trends of the HiFi-GAN vocoder model

[ ] trainer.fit()
        | > avg_D_mse_gan_fake_loss: 0.00518 (-0.00296)
        | > avg_loss_0: 0.45660 (+0.03587)
        | > avg_G_l1_spec_loss: 0.68739 (-0.01777)
        | > avg_G_mse_fake_loss: 0.53558 (+0.09378)
        | > avg_G_feat_match_loss: 0.04207 (+0.00070)
        | > avg_G_gen_loss: 30.93272 (-0.79960)
        | > avg_G_adv_loss: 5.07904 (+0.16968)
        | > avg_loss_1: 36.01175 (-0.62995)


    > EPOCH: 498/500
    --> /content/gdrive/MyDrive/Zeba_TTS/TTS_CP_HiFiGan/run-August-06-2022_04+27PM-0000000
    /usr/local/lib/python3.7/dist-packages/torch/utils/data/dataloader.py:560: UserWarning: This DataLoade
      cpuset_checked))

    > TRAINING (2022-08-07 06:19:00)

    > EVALUATION


    --> EVAL PERFORMANCE
        | > avg_loader_time: 0.00517 (+0.00008)
        | > avg_D_mse_gan_loss: 0.40273 (-0.05387)
        | > avg_D_mse_gan_real_loss: 0.04795 (-0.10758)
        | > avg_D_mse_gan_fake_loss: 0.02291 (+0.01773)
        | > avg_loss_0: 0.40273 (-0.05387)
        | > avg_G_l1_spec_loss: 0.68267 (-0.00472)
        | > avg_G_mse_fake_loss: 0.34242 (-0.19316)
        | > avg_G_feat_match_loss: 0.04707 (+0.00500)
        | > avg_G_gen_loss: 30.72021 (-0.21251)
        | > avg_G_adv_loss: 5.42630 (+0.34726)
        | > avg_loss_1: 36.14650 (+0.13475)


    > EPOCH: 499/500
    --> /content/gdrive/MyDrive/Zeba_TTS/TTS_CP_HiFiGan/run-August-06-2022_04+27PM-0000000
    /usr/local/lib/python3.7/dist-packages/torch/utils/data/dataloader.py:560: UserWarning: This DataLoade
      cpuset_checked))

    > TRAINING (2022-08-07 06:20:40)

    > EVALUATION


    --> EVAL PERFORMANCE
        | > avg_loader_time: 0.00519 (+0.00002)
        | > avg_D_mse_gan_loss: 0.41995 (+0.01722)
```

Figure 18: Training the HiFi-GAN vocoder from scratch on the MyST corpus

## 5.4 Generating Audio References from the trained HiFi-GAN vocoder

The audio references were generated as shown in Figure 19 from the HiFi-GAN vocoder model trained on MyST corpus for 100 and 200 epochs for preliminary analysis. Figure 20 demonstrates the use of IPython library to hear the audio file generated from the vocoder

14

Figure 19: Generating audio references from the HiFi-GAN vocoder for 100 and 200 epochs



Figure 20: Playing the audio generated from HiFi-GAN vocoder for 100 and 200 epochs

As the audio generated from the HiFi-GAN vocoder trained from the scratch was purely metallic, the transfer learning approach was adopted.

# 6 Transfer Learning approach for the HiFi-GAN vocoder

This section describes the implementation of the transfer learning approach of the HiFi-GAN vocoder. A pre-trained HiFi-GAN model was used to be further fine-tuned on the MyST dataset. The generator and the discriminator is the same as used in 5. The configuration parameters changed are described in the subsequent figures.

## 6.1 Import/ Install repository for implementing the transfer learning for the pre-trained HiFi-GAN vocoder



Figure 21: Cloning the TTS git repository

## 6.2 Modify the hyper-parameters in the HiFi-GAN vocoder configuration file



Figure 22: Modifying the parameters for fine-tuning the HiFi-GAN vocoder

## 6.3   Start Fine-tuning and Training the HiFi-GAN Vocoder



Figure 23: Fine-tune and train the HiFi-GAN vocoder on MyST corpus

The following figure demonstrates the repetition of the fine-tuning and training process for desired number of epochs (here 200).

```
Step 18: Repeat the training process for various epochs by changing the --restore_path parameter for the latest checkpoint available in each run

[ ]  !CUDA_VISIBLE_DEVICES="0" python "/content/TTS/recipes/ljspeech/hifigan/train_hifigan.py" \
        --output_path="/content/gdrive/MyDrive/Zeba_HiFi/FineTune_HiFi"\
        --restore_path="/content/gdrive/MyDrive/Zeba_HiFi/FineTune_HiFi/run-August-10-2022_12+38PM-d46fbc24/best_model_980.pth"\
        --save_all_best = True

   --> EVAL PERFORMANCE
     | > avg_loader_time: 0.00421 (-0.00006)
     | > avg_D_mse_gan_loss: 0.36248 (-0.01204)
     | > avg_D_mse_gan_real_loss: 0.04637 (-0.05039)
     | > avg_D_mse_gan_fake_loss: 0.01144 (+0.00312)
     | > avg_loss_0: 0.36248 (-0.01204)
     | > avg_G_l1_spec_loss: 0.67334 (-0.01766)
     | > avg_G_mse_fake_loss: 0.50474 (-0.01965)
     | > avg_G_feat_match_loss: 0.05486 (+0.00170)
     | > avg_G_gen_loss: 30.30023 (-0.79476)
     | > avg_G_adv_loss: 6.42967 (+0.16438)
     | > avg_loss_1: 36.72990 (-0.63039)


!CUDA_VISIBLE_DEVICES="0" python "/content/TTS/recipes/ljspeech/hifigan/train_hifigan.py" \
   --output_path="/content/gdrive/MyDrive/Zeba_HiFi/FineTune_HiFi"\
   --restore_path="/content/gdrive/MyDrive/Zeba_HiFi/FineTune_HiFi/run-August-10-2022_01+49PM-d46fbc24/best_model_1111.pth"\
   --save_all_best = True

     | > D_mse_gan_loss: 0.31493  (0.32491)
     | > D_mse_gan_real_loss: 0.04014  (0.03069)
     | > D_mse_gan_fake_loss: 0.00761  (0.02420)
     | > loss_0: 0.31493  (0.32491)
     | > grad_norm_0: 4.23894  (5.86320)
     | > G_l1_spec_loss: 0.53196  (0.52914)
     | > G_mse_fake_loss: 0.53478  (0.50658)
     | > G_feat_match_loss: 0.06495  (0.06166)
     | > G_gen_loss: 23.93812  (23.81123)
     | > G_adv_loss: 7.54901  (7.16551)
     | > loss_1: 31.48713  (30.97674)
     | > grad_norm_1: 1360.89795  (880.02283)
     | > current_lr_0: 0.00009
     | > current_lr_1: 0.00009
     | > step_time: 2.07910  (2.13160)
     | > loader_time: 0.00350  (0.00204)


 > CHECKPOINT : /content/gdrive/MyDrive/Zeba_HiFi/FineTune_HiFi/run-August-10-2022_03+03PM-d46fbc24/checkpoint_2100.pth

 > EVALUATION

   --> EVAL PERFORMANCE
     | > avg_loader_time: 0.00453 (-0.00001)
     | > avg_D_mse_gan_loss: 0.33334 (-0.03052)
     | > avg_D_mse_gan_real_loss: 0.04228 (-0.03455)
     | > avg_D_mse_gan_fake_loss: 0.00998 (+0.00377)
     | > avg_loss_0: 0.33334 (-0.03052)
     | > avg_G_l1_spec_loss: 0.67433 (-0.01866)
     | > avg_G_mse_fake_loss: 0.46126 (-0.01942)
```

Figure 24: Restoring the training process for fine-tuning the HiFi-GAN vocoder

# 7 Synthesizing voice from the proposed model

This section describes the implementation of the proposed model (Speaker encoder - Tacotron 2 - HiFi-GAN) . The implementations uses the implementation of the speaker

encoder and Tacotron model [7] and the pre-trained model weights from [8] and [9]. Following figures demonstrate the steps to synthesize voice from the proposed model.

## 7.1 Import/ Install repository for implementing the proposed model that contains transfer learning for the pre-trained HiFi-GAN vocoder



Figure 25: Cloning the git repository and load default models

---

[7]Speaker encoder and Tacotron 2 Implementation `https://github.com/CorentinJ/Real-Time-Voice-Cloning`

[8]Speaker encoder checkpoints `https://drive.google.com/drive/folders/1FuAY2XXcUOvLVo1f9QYQjhs_g9eURbio`

[9]Acoustic model checkpoints `https://drive.google.com/drive/folders/1wcxVnJ5mQZNdl1r_aLzY86iIAgRm4hQH`

Figure 26: Download pre-trained models and use HiFi-GAN pre-trained model

## 7.2 Synthesize voice using pre-trained HiFi-GAN vocoder in the proposed model

The following figures demonstrate loading the pre-trained HiFi-GAN vocoder model and synthesize voice from the proposed model.



Figure 27: Load pre-trained HiFi-GAN vocoder model and upload an audio file to synthesize

As shown in the figure below, enter the desired text to generate the synthesized voice. Additionally, save the input and the synthesized file for evaluation

Figure 28: Enter desired text to generate the synthesized voice

## 7.3 Synthesize voice using fine-tuned HiFi-GAN vocoder for 100 epochs in the proposed model

The following figures demonstrate the generation of synthesized voice for the fine-tuned HiFi-GAN vocoder model for 100 epochs.



Figure 29: Load fine-tuned HiFi-GAN vocoder model (100 epochs) and upload an audio file to synthesize

As shown in the figure below, enter the desired text to generate the synthesized voice. Additionally, save the input and the synthesized file for evaluation

Figure 30: Enter desired text to generate the synthesized voice

## 7.4 Synthesize voice using fine-tuned HiFi-GAN vocoder for 200 epochs in the proposed model

The following figures demonstrate the generation of synthesized voice for the fine-tuned HiFi-GAN vocoder model for 200 epochs.



Figure 31: Load fine-tuned HiFi-GAN vocoder model (200 epochs) and upload an audio file to synthesize

As shown in the figure below, enter the desired text to generate the synthesized voice. Additionally, save the input and the synthesized file for evaluation

Figure 32: Enter desired text to generate the synthesized voice

# 8 Evaluation: Generate MOSNet scores

This section describes the generation of MOSNet scores for the original and the synthesized audio files produced in 7. The MOSNet for this research is implemented from [10]. The following figure demonstrates the generation of MOSNet scores for this research.

## 8.1 Import/ Install repository for implementing and generating MOSNet scores



Figure 33: Cloning the git repository, change directory, and install requirements

---

[10]MOSNet Implementation `https://github.com/lochenchou/MOSNet.git`

## 8.2 Get Valid Directories for MOSNet to run



```
[ ]  """Loop through each speaker folders to get their MOSNet scores
     using the MOSNet github repo code
     """
     root_directory = "/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices"
     previous_directory = ""

     for dir_path, sub_directories, files in os.walk(root_directory):
       for filename in files:
         if filename.endswith(".wav") and ("Pretrained" in dir_path):
           if previous_directory == "" or previous_directory != dir_path:
             previous_directory = dir_path
             print("Valid directories :", previous_directory)
```
```
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_0/002013_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_0/002033_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002026_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002030_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002102_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002109_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002113_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002269_Pretrained
Valid directories : /content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002274_Pretrained
```

Figure 34: Get Valid Directories for MOSNet

## 8.3 Generate MOSNet scores for synthesized voices from pre-trained HiFi-GAN vocoder



```
[ ]  # HiFi-GAN purely pre-trained model
     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_0/002013_Pretrained'
     print("\n")
     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_0/002033_Pretrained'
```
```
Loading model weights
CNN_BLSTM init
Start evaluating 2 waveforms...
100% 2/2 [00:04<00:00,  2.07s/it]
Average: 2.3395


Loading model weights
CNN_BLSTM init
Start evaluating 2 waveforms...
100% 2/2 [00:03<00:00,  1.99s/it]
Average: 2.3425000000000002
```

Figure 35: MOSNet scores for synthesized voices from pre-trained HiFi-GAN vocoder

## 8.4 Generate MOSNet scores for synthesized voices from fine-tuned HiFi-GAN vocoder (100 epochs)

```
[ ]  # HiFi-GAN fine-tuned model for 100 epochs
     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002026_Pretrained'

     print("\n")

     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002030_Pretrained'

     print("\n")

     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_100/002102_Pretrained'


     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  1.71s/it]
     Average: 2.904

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:04<00:00,  2.08s/it]
     Average: 2.6870000000000003

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  1.93s/it]
     Average: 2.6405
```

Figure 36: MOSNet scores for synthesized voices from fine-tuned HiFi-GAN vocoder (100 epochs)

## 8.5 Generate MOSNet scores for synthesized voices from fine-tuned HiFi-GAN vocoder (200 epochs)



```
[ ]  # HiFi-GAN fine-tuned model for 200 epochs
     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002109_Pretrained'

     print("\n")

     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002113_Pretrained'

     print("\n")

     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002269_Pretrained'

     print("\n")

     !python ./custom_test.py \
     --rootdir '/content/gdrive/MyDrive/Zeba_TTS/Synthesized_Voices/Pretrained_200/002274_Pretrained'

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  1.77s/it]
     Average: 2.928

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  1.75s/it]
     Average: 2.9215

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  1.76s/it]
     Average: 2.721

     Loading model weights
     CNN_BLSTM init
     Start evaluating 2 waveforms...
     100% 2/2 [00:03<00:00,  2.00s/it]
     Average: 2.389999999999997
```

Figure 37: MOSNet scores for synthesized voices from fine-tuned HiFi-GAN vocoder (200 epochs)

## 9 Evaluation: Plot Visualization

This section describes the generation of mel spectrograms and waveforms for the original and synthesized voices produced in 7. The following figures demonstrate the visualizations done in this research.

## 9.1 Plot Mel Spectrograms and Waveforms



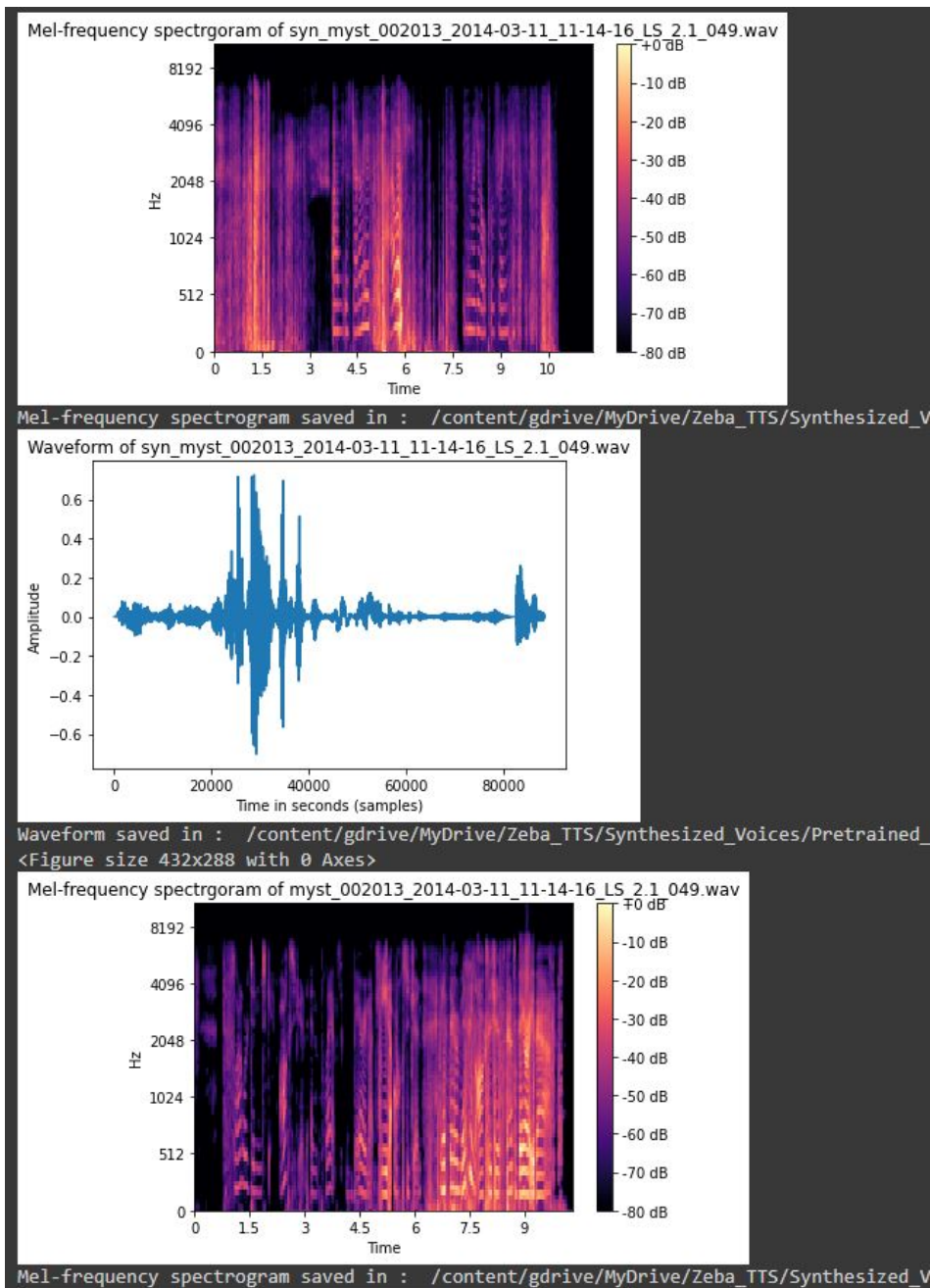Figure 38: Generate and Store Mel Spectrograms and Waveforms

Figure 39: Produced Mel Spectrograms and Waveforms

# References

Jain, R., Yiwere, M. Y., Bigioi, D., Corcoran, P. and Cucu, H. (2022). A text-to-speech pipeline, evaluation methodology, and initial fine-tuning results for child speech synthesis, *IEEE Access* **10**: 47628–47642.