

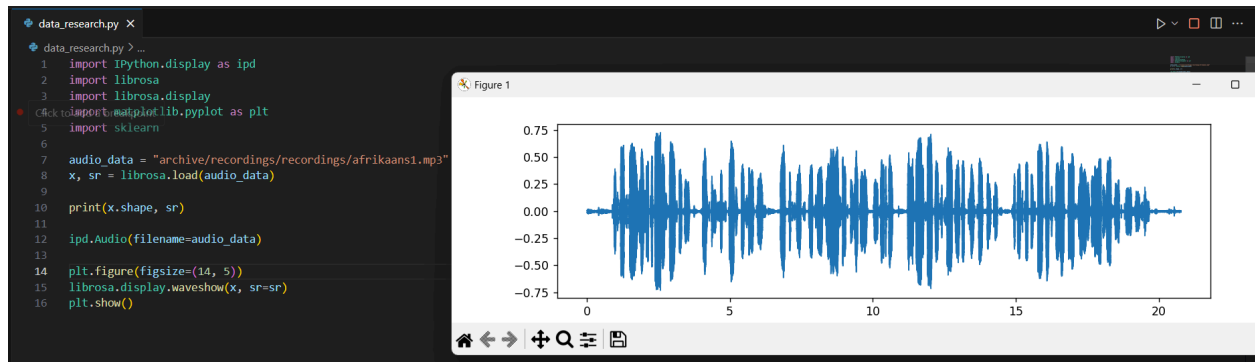
Biologically inspired artificial intelligence

Accent Detection

Patryk Hećko
Krzysztof Wojtyś

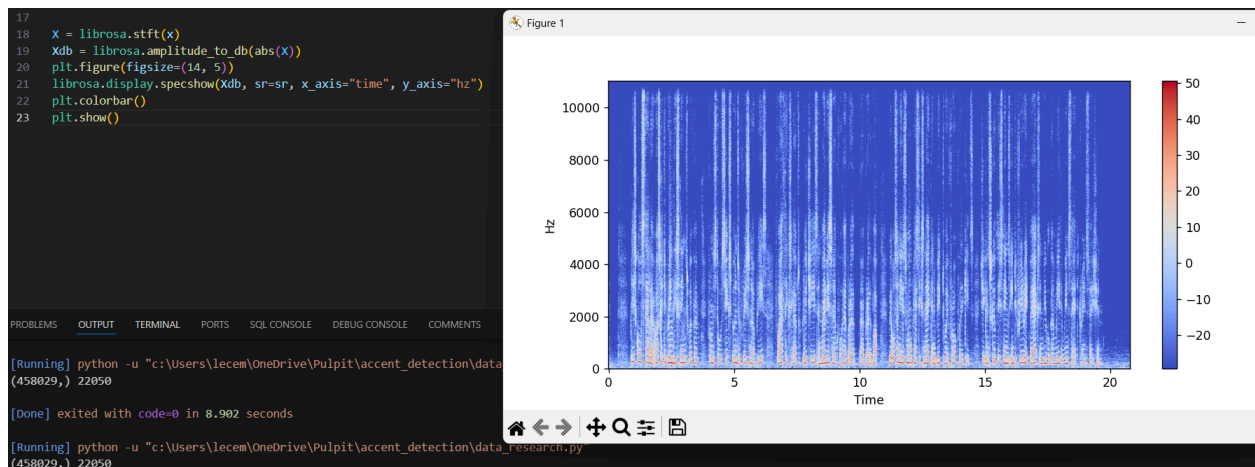
Informatyka sem. VI
Katowice

The first part of the project is examining the input audio recordings using librosa library to load the audio file and matplotlib to display the recording data.



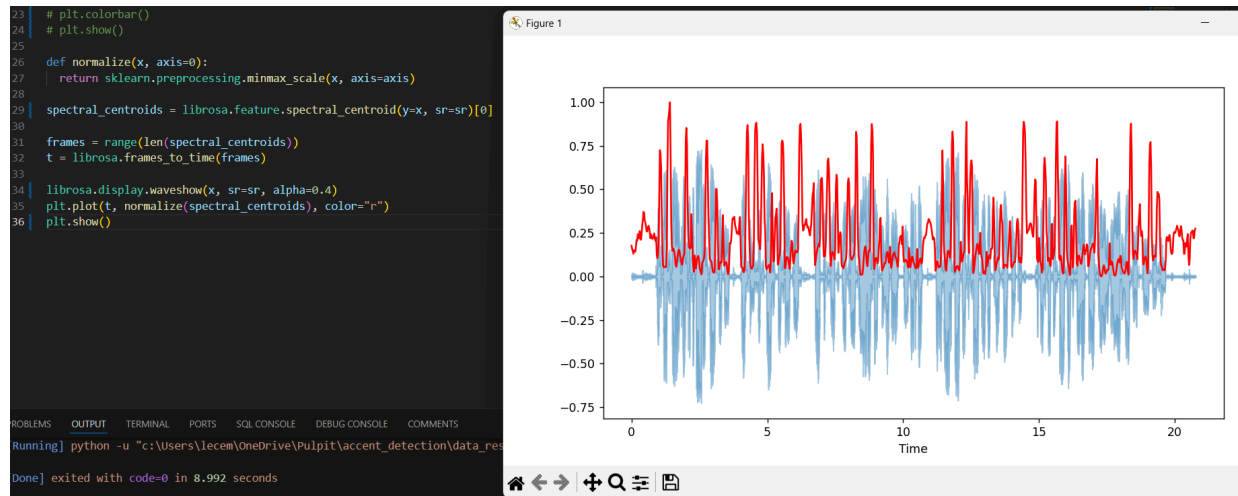
The given audio sample is sampled at 22050 Hz, dimensionality of 458029, it has 1 channel and duration of 21 s.

Next step was to plot a signal spectrogram to see the signal over time at various frequencies present in a particular waveform.

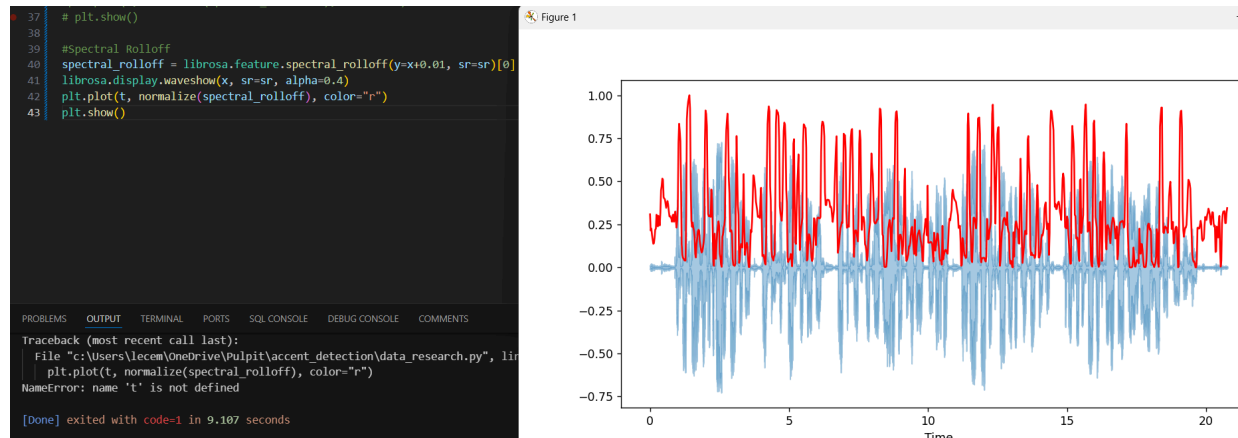


To best know what audio features will help us with our task, we extracted a few from the given audio signal.

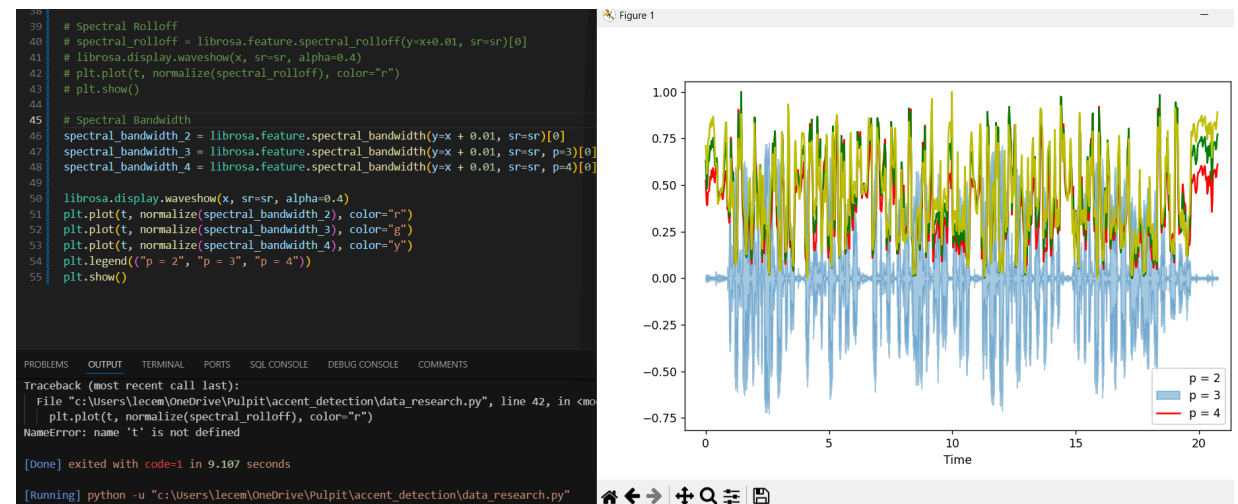
First feature to extract is spectral centroid, it indicates where the center of mass of the spectrum is located.



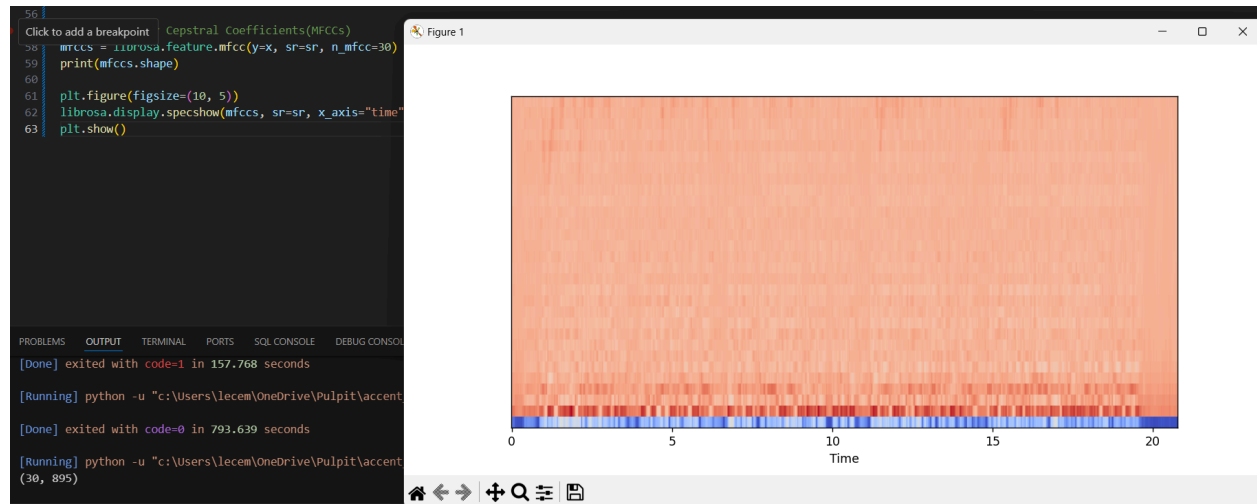
Next feature is spectral rolloff, which is a measure of the shape of the signal. It represents the frequency at which high frequencies decline to 0.



Next feature is spectral bandwidth, which is defined as the width of the band of light at one-half the peak maximum.



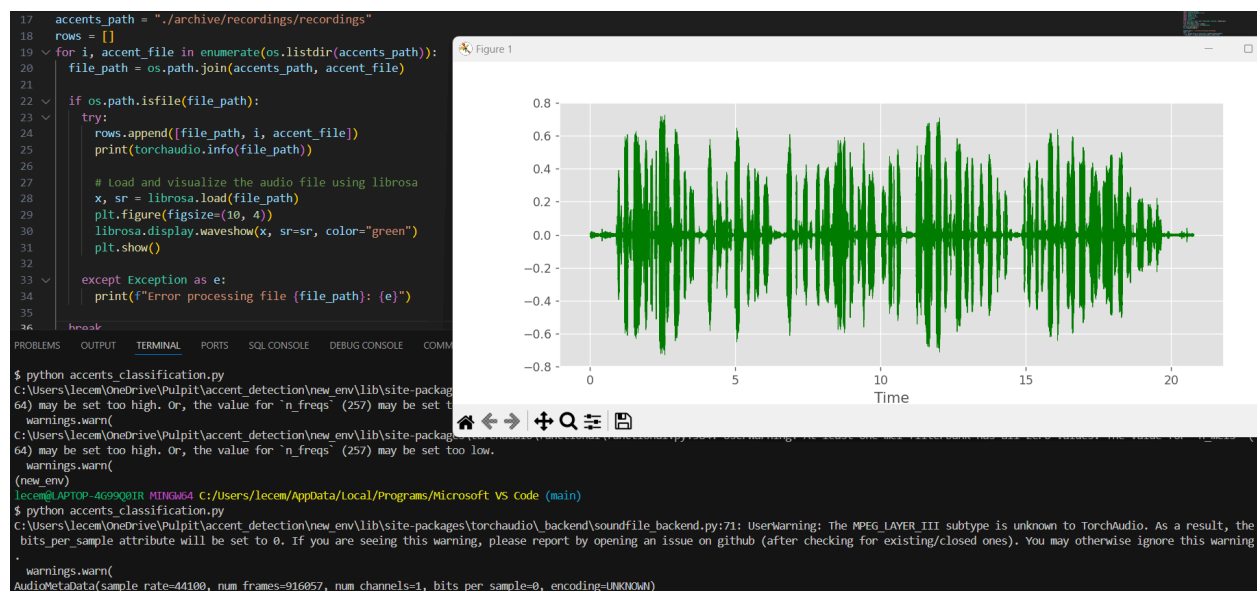
Last feature we extract is Mel-Frequency Cepstral Coefficients(MFCCs). The Mel frequency cepstral coefficients of a signal are a small set of features which concisely describe the overall shape of a spectral envelope. It models the characteristics of the human voice.



Mfcc computed 30 MFCCs over 30 frames

Because our dataset is collected human speech The Mel frequency cepstral coefficients suit it best and we will use these types of features for our further work.

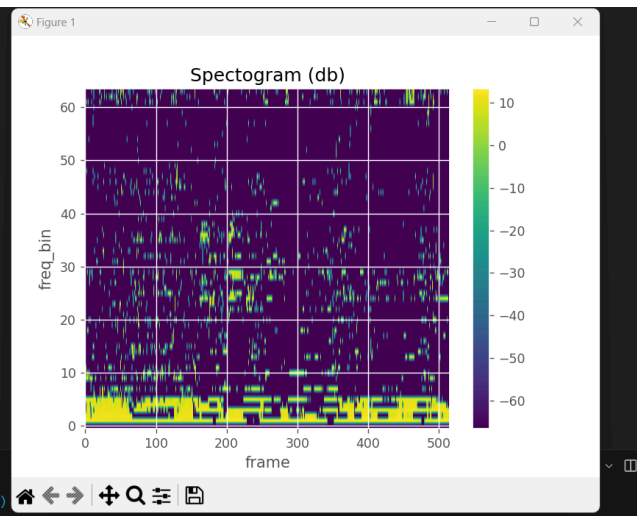
After examining the audio inputs we started implementing accents_classification file. First thing to implement was loading and displaying all recordings data.



After successfully retrieving recordings data we moved on to implement SpeechDataset class and plot_spectrogram function in order to create easy to distinguish frequency(frame) images.

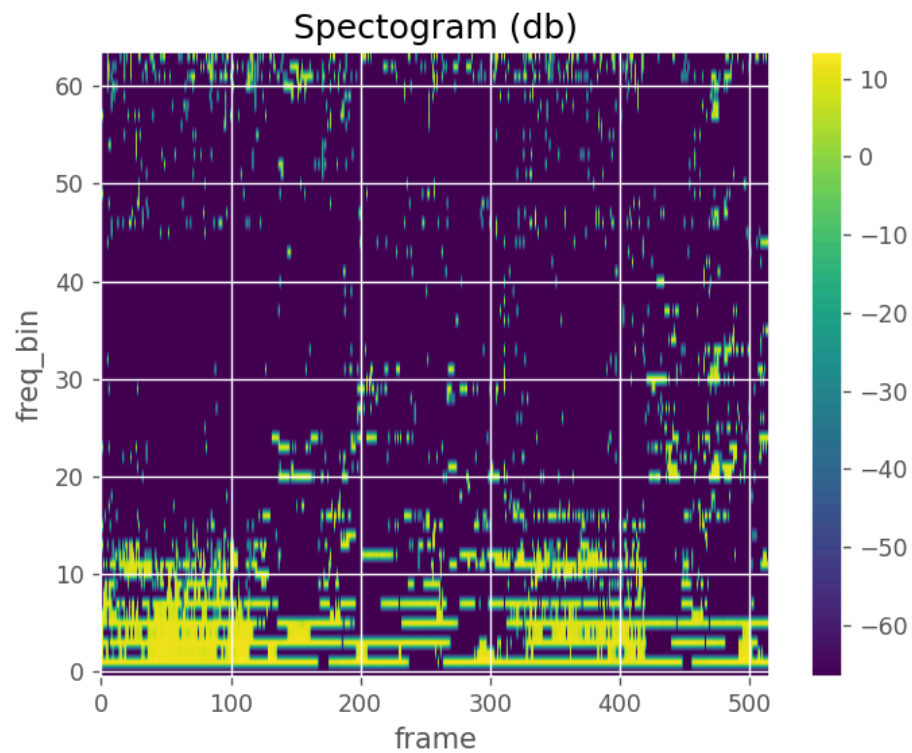
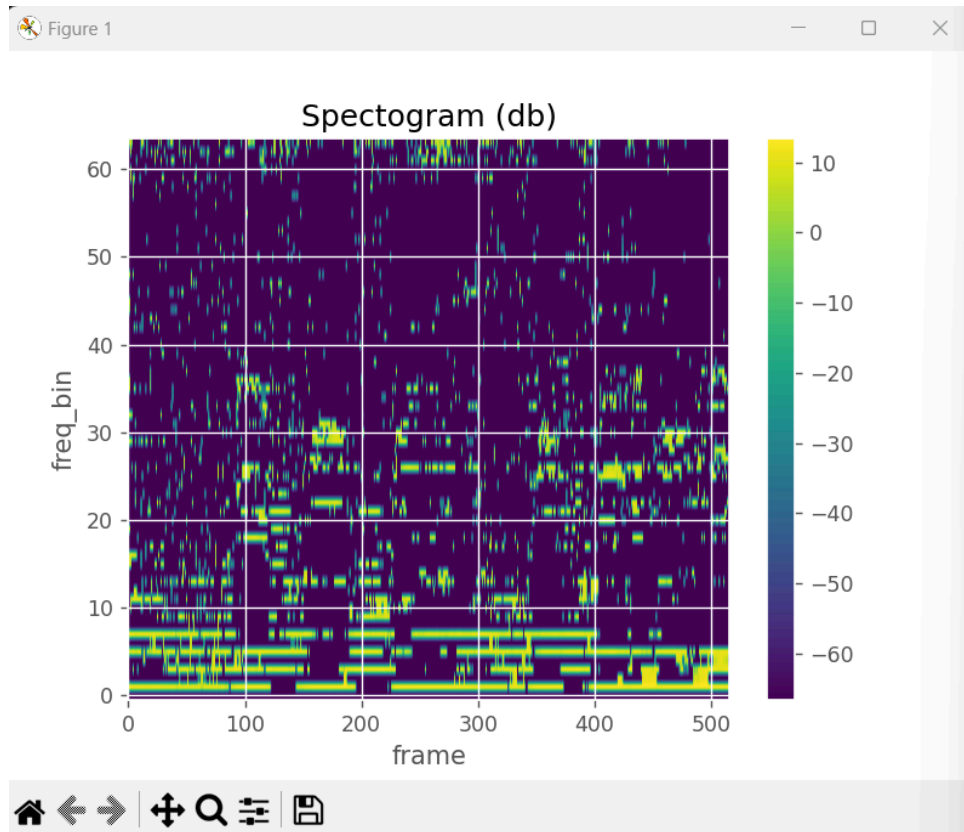
```
44 class SpeechDataset(Dataset):
45     def __init__(self, data_fr, data_path):
46         self.data_fr = data_fr
47         self.data_path = str(data_path)
48
49     def __len__(self):
50         return len(self.data_fr)
51
52     def __getitem__(self, idx):
53         audio_file = self.data_path + '/' + self.data_fr.loc[idx, "file_name"]
54         class_id = self.data_fr.loc[idx, "class_id"]
55         audio = preprocessor.load_audio(audio_file)
56         rechannel = preprocessor.double_channel(audio)
57         downsample = preprocessor.downsample(rechannel)
58         spectrogram = preprocessor.spectro_mfcc(downsample)
59         return spectrogram, class_id
60
61 speech_dataset = SpeechDataset(df, accents_path)
62
63 num_items = len(speech_dataset)
64 num_train = round(num_items * 0.7)
65 train_sample = num_items - num_train
66 train_ds, val_ds = random_split(speech_dataset, [num_train, train_sample])
67 # Training and validation data loaders
68 train_dl = DataLoader(train_ds, batch_size=10, shuffle=True)
69 val_dl = DataLoader(val_ds, batch_size=10, shuffle=False)
```

```
64 num_train = round(num_items * 0.7)
65 train_sample = num_items - num_train
66 train_ds, val_ds = random_split(speech_dataset, [num_train, train_sample])
67 # Training and validation data loaders
68 train_dl = DataLoader(train_ds, batch_size=10, shuffle=True)
69 val_dl = DataLoader(val_ds, batch_size=10, shuffle=False)
70
71 def plot_spectrogram(spec, ylabel="freq_bin", aspect="auto"):
72     fig, axs = plt.subplots(1, 1)
73     axs.set_title("Spectrogram (db)")
74     axs.set_xlabel("frame")
75     axs.set_ylabel(ylabel)
76     im = axs.imshow(librosa.power_to_db(spec), origin="lower", aspect=aspect)
77     fig.colorbar(im, ax=axs)
78     plt.show()
79
80 for i in range(5):
81     plot_spectrogram(train_ds[i][0][0])
```



PROBLEMS OUTPUT TERMINAL PORTS SQL CONSOLE DEBUG CONSOLE COMMENTS

```
lecem@LAPTOP-4G99Q0TR MINGW64 C:/Users/lecem/AppData/Local/Programs/Microsoft VS Code (main)
$ python accents_classification.py
```



As we can see the images are different based on the speaker.

After successfully displaying input audio as spectrograms we implemented a classifier model specifying our neural network data such as Convolution Layers, activation function and normalization. We decided to use the ReLu activation function simply because it is easier to compute than the sigmoid function and gives very satisfactory results.

```
83 class AudioClassifier(nn.Module):
84     def __init__(self):
85         super(AudioClassifier, self).__init__()
86         # Neural Network shape
87         self.conv = nn.Sequential([
88             nn.Conv2d(2, 8, kernel_size=(3, 3), stride=(2,2), padding=(1, 1)), # 2D Convolution layer
89             nn.ReLU(), # Rectified Linear Unit activation function
90             nn.BatchNorm2d(8), #normalization
91             nn.Conv2d(8, 16, kernel_size=(3, 3), stride=(2,2), padding=(1, 1)),
92             nn.ReLU(),
93             nn.BatchNorm2d(16),
94             nn.Conv2d(16, 32, kernel_size=(3, 3), stride=(2,2), padding=(1, 1)),
95             nn.ReLU(),
96             nn.BatchNorm2d(32),
97         ])
98         # Linear Classifier
99         self.ap = nn.AdaptiveAvgPool2d(output_size=1)
100         self.dropout = nn.Dropout(0.5)
101         # To establish in and out features
102         self.lin = nn.Linear(in_features=32, out_features=9)
103
104         # Forward propagation
105     def forward(self, inp_x):
106         inp_x = self.conv(inp_x)
107         inp_x = self.ap(inp_x)
108         inp_x = inp_x.view(inp_x.shape[0], -1)
109         inp_x = self.dropout(inp_x)
110         inp_x = self.lin(inp_x)
111         return inp_x
```

Given the model we displayed its data. We use PyTorch on gpu (cuda) rather than cpu because it is much faster given thousands of cores rather than just a few.

```

129 model = AudioClassifier()
130 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
131 model = model.to(device)
132 next(model.parameters()).device
133 summary(model, (2, 64, 258), 11)

```

PROBLEMS OUTPUT TERMINAL PORTS SQL CONSOLE DEBUG CONSOLE COMMENTS

lecem@LAPTOP-4G99Q0IR MINGW64 C:/Users/lecem/AppData/Local/Programs/Microsoft VS Code (main)

• \$ python accents_classification.py

```

-----
Layer (type)          Output Shape          Param #
-----
Conv2d-1              [11, 8, 32, 129]      152
ReLU-2                [11, 8, 32, 129]      0
BatchNorm2d-3         [11, 8, 32, 129]      16
Conv2d-4              [11, 16, 16, 65]      1,168
ReLU-5                [11, 16, 16, 65]      0
BatchNorm2d-6         [11, 16, 16, 65]      32
Conv2d-7              [11, 32, 8, 33]       4,640
ReLU-8                [11, 32, 8, 33]       0
BatchNorm2d-9         [11, 32, 8, 33]       64
AdaptiveAvgPool2d-10  [11, 32, 1, 1]        0
Dropout-11            [11, 32]              0
Linear-12              [11, 9]               297
=====
Total params: 6,369
Trainable params: 6,369
Non-trainable params: 0
-----
Input size (MB): 1.39
Forward/backward pass size (MB): 14.64
Params size (MB): 0.02
Estimated Total Size (MB): 16.05
-----
(new_env)

```

We also defined Accuracy Metric class in order for neural network to know the accuracy and be able to improve in the following propagations

```

class AccuracyMetric:
    def __init__(self):
        self.correct, self.total = None, None
        self.reset()

    def update(self, y_pred, y_true):
        self.correct += torch.sum(y_pred.argmax(-1) == y_true).item()
        self.total += y_true.size(0)

    def compute(self):
        return self.correct / self.total

    def reset(self):
        self.correct = 0
        self.total = 0

```


We then tried to train our model, but results were very dissatisfactory and iterations were way too slow, given our model works on gpu.

```

162     for epoch in range(epochs):
163         print(f"[INFO] Epoch: {epoch + 1}")
164         model.train()
165
166         batch_train_loss = []
167         batch_valid_loss = []
168
169         for X_batch, y_batch in tqdm(train_dl):
170             # single training step
171             model.zero_grad()
172             X_batch, y_batch = X_batch.to(device), y_batch.to(device)
173             y_predicted = model(X_batch)
174             v_predicted = v_predicted.to(torch.float)

```

PROBLEMS OUTPUT TERMINAL PORTS DEBUG CONSOLE COMMENTS SQL CONSOLE

```

return Variable.execution_engine.run_backward( # Calls into the C++ engine to run the backward pass
99%|██████████████████████████████████████████████████████████████████████████████| 99/100 [01:50<00:01, 1.13s/it]C:\Users\lece\OneDrive\Pulpit\ae
erWarning: Plan failed with a cudnnException: CUDNN_BACKEND_EXECUTION_PLAN_DESCRIPTOR: cudnnFinalize Descriptor Failed cu
rc\aten\native\cudnn\Conv_v8.cpp:919.)
return Variable.execution_engine.run_backward( # Calls into the C++ engine to run the backward pass
100%|██████████████████████████████████████████████████████████████████████████████| 100/100 [01:51<00:00, 1.11s/it]
100%|██████████████████████████████████████████████████████████████████████████████| 43/43 [00:51<00:00, 1.20s/it]
Train loss: 4.8609, Train accuracy: 0.2011 Validation loss: 4.1362, Validation accuracy: 0.2871
[INFO] Epoch: 2
100%|██████████████████████████████████████████████████████████████████████████████| 100/100 [01:50<00:00, 1.11s/it]
100%|██████████████████████████████████████████████████████████████████████████████| 43/43 [00:46<00:00, 1.09s/it]
Train loss: 4.0944, Train accuracy: 0.2645 Validation loss: 3.8365, Validation accuracy: 0.2855
[INFO] Epoch: 3
100%|██████████████████████████████████████████████████████████████████████████████| 100/100 [01:52<00:00, 1.13s/it]
100%|██████████████████████████████████████████████████████████████████████████████| 43/43 [00:47<00:00, 1.10s/it]
Train loss: 3.8929, Train accuracy: 0.2632 Validation loss: 3.7309, Validation accuracy: 0.2839
[INFO] Epoch: 4
: 0.2855
[INFO] Epoch: 5
100%|██████████████████████████████████████████████████████████████████████████████| 100/100 [01:54<00:00, 1.14s/it]
100%|██████████████████████████████████████████████████████████████████████████████| 43/43 [00:48<00:00, 1.12s/it]
Train loss: 3.7576, Train accuracy: 0.2632 Validation loss: 3.6917, Validation accuracy: 0.2839
(new env)

```

<https://towardsdatascience.com/audio-deep-learning-made-simple-part-1-state-of-the-art-techniques-da1d3dff2504>