

Speech Emotion Recognition (SER)

Model: Long Short-Term Memory (LSTM)

Problem Statement:

The aim of this project is to construct and employ an LSTM classification model for SER.

GOAL:

Accurately assess the emotional state of speakers in audio recordings.

HOW:

Using the librosa library standardize and extract features voice signals from a dataset created for SER modeling and then build, train and test the SER model.

DELIVERABLE:

MVP - A functioning LSTM model

What is LSTM

Data:

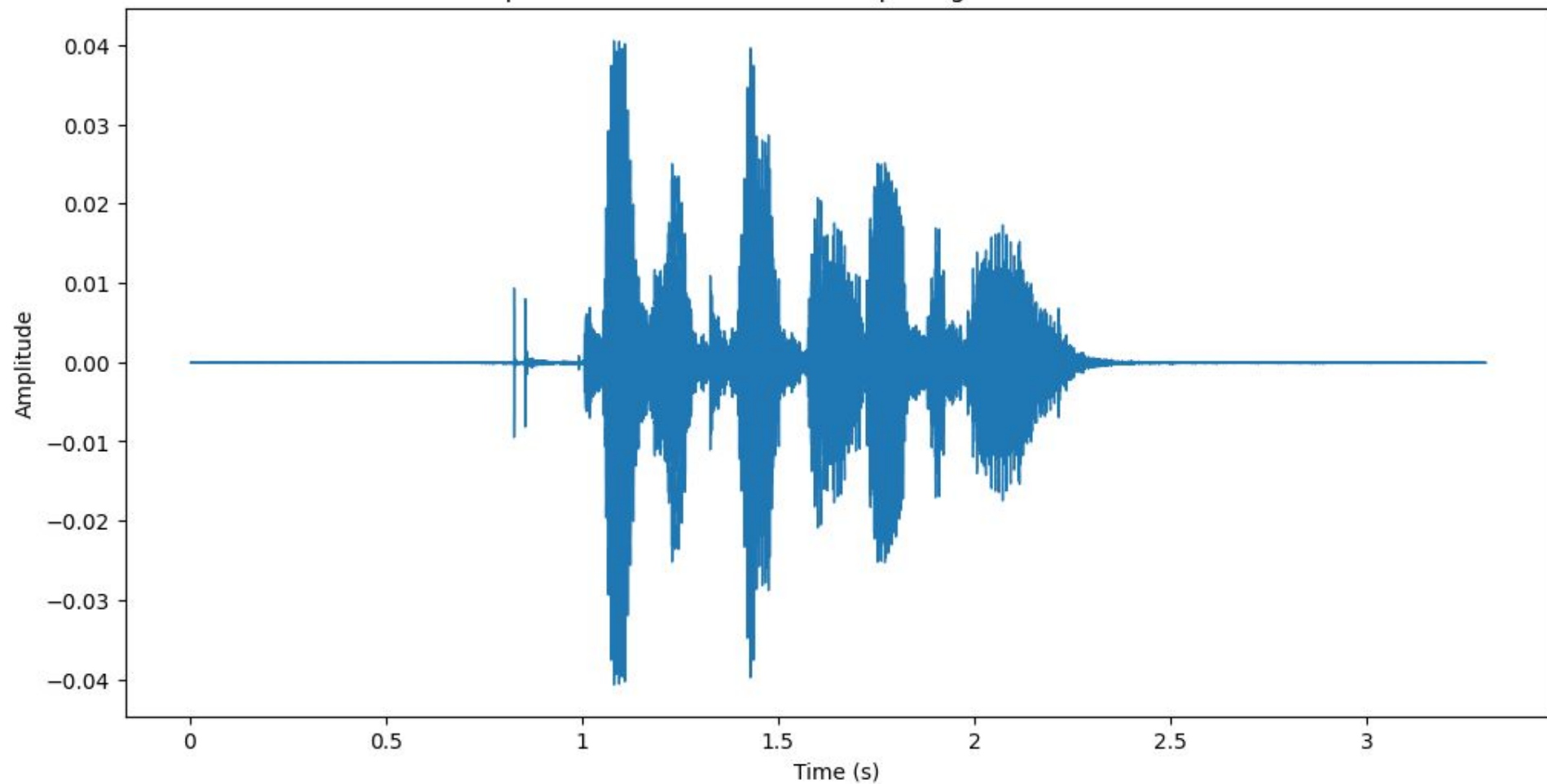
[Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS)](<https://www.kaggle.com/datasets/uwrfkaggler/ravdess-emotional-speech-audio>)

[Toronto emotional speech set (TESS)](<https://www.kaggle.com/datasets/ejlok1/toronto-emotional-speech-set-tess>)

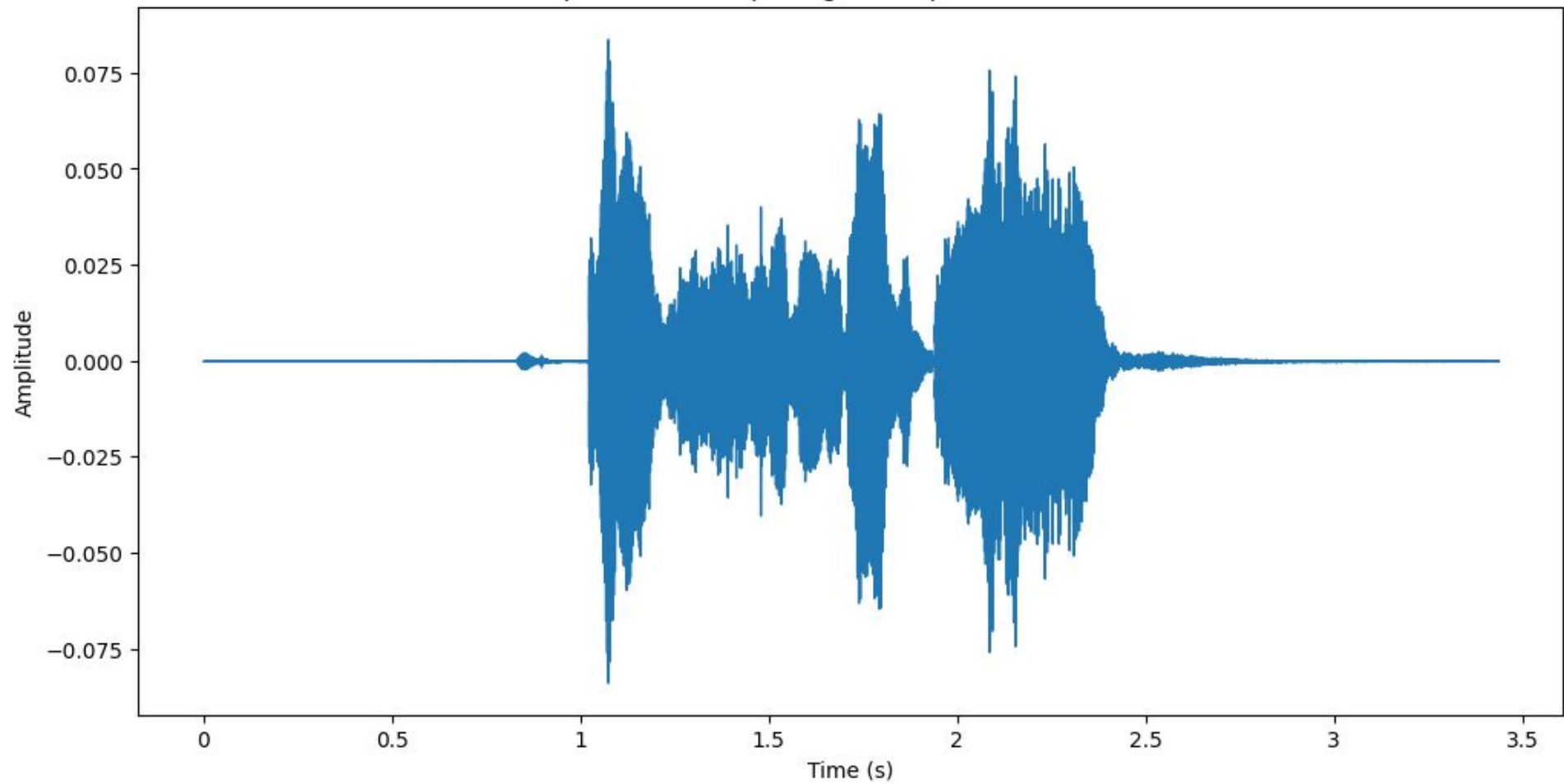
[Surrey Audio-Visual Expressed Emotion (SAVEE)](<https://www.kaggle.com/datasets/ejlok1/surrey-audiovisual-expressed-emotion-savee>)

[Crowd Sourced Emotional Multimodal Actors Dataset (CREMA-D)](<https://www.kaggle.com/datasets/ejlok1/cremad>)

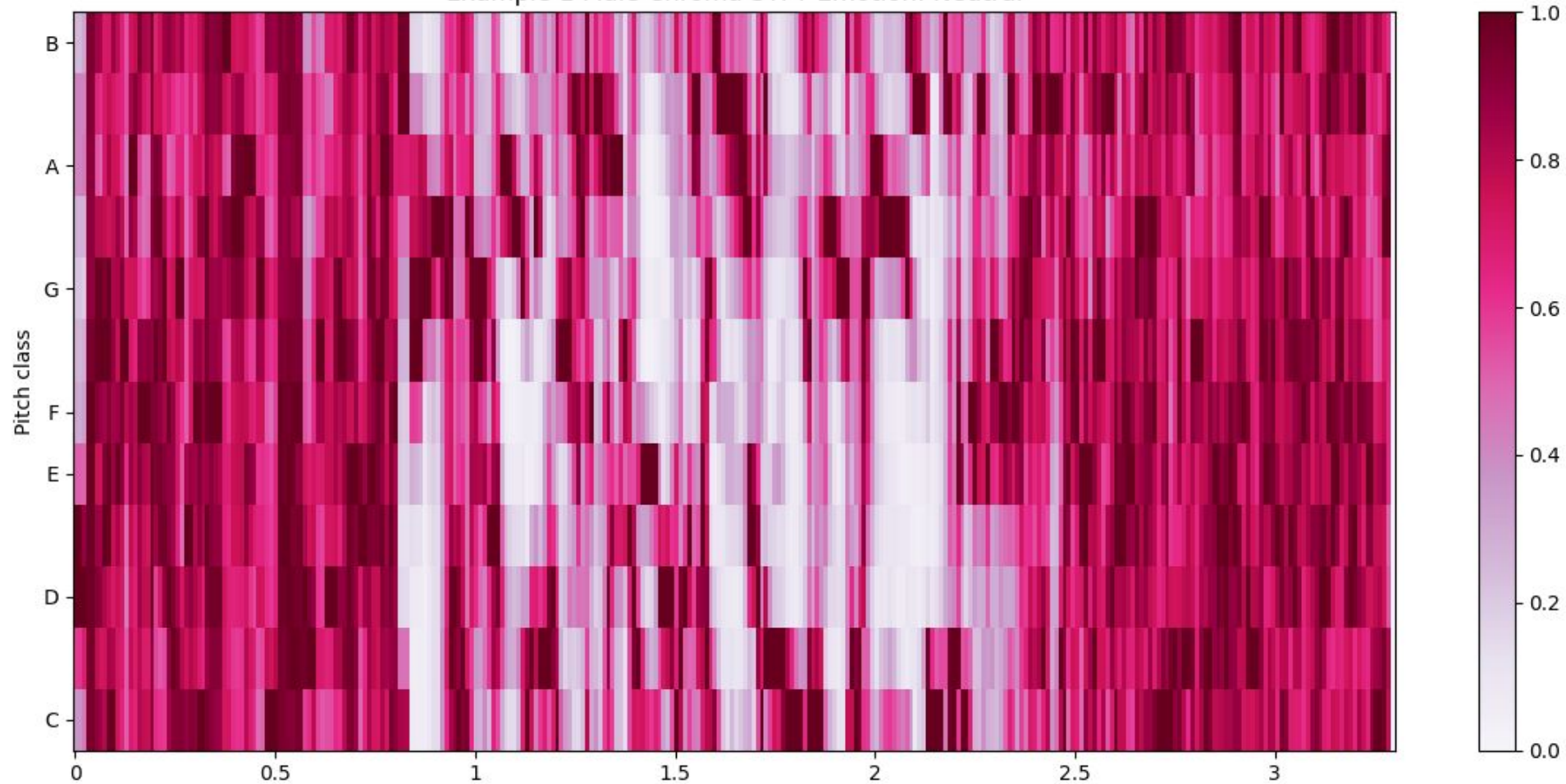
Example 1 Audio Waveform Male Spectrogram Emotion: Neutral



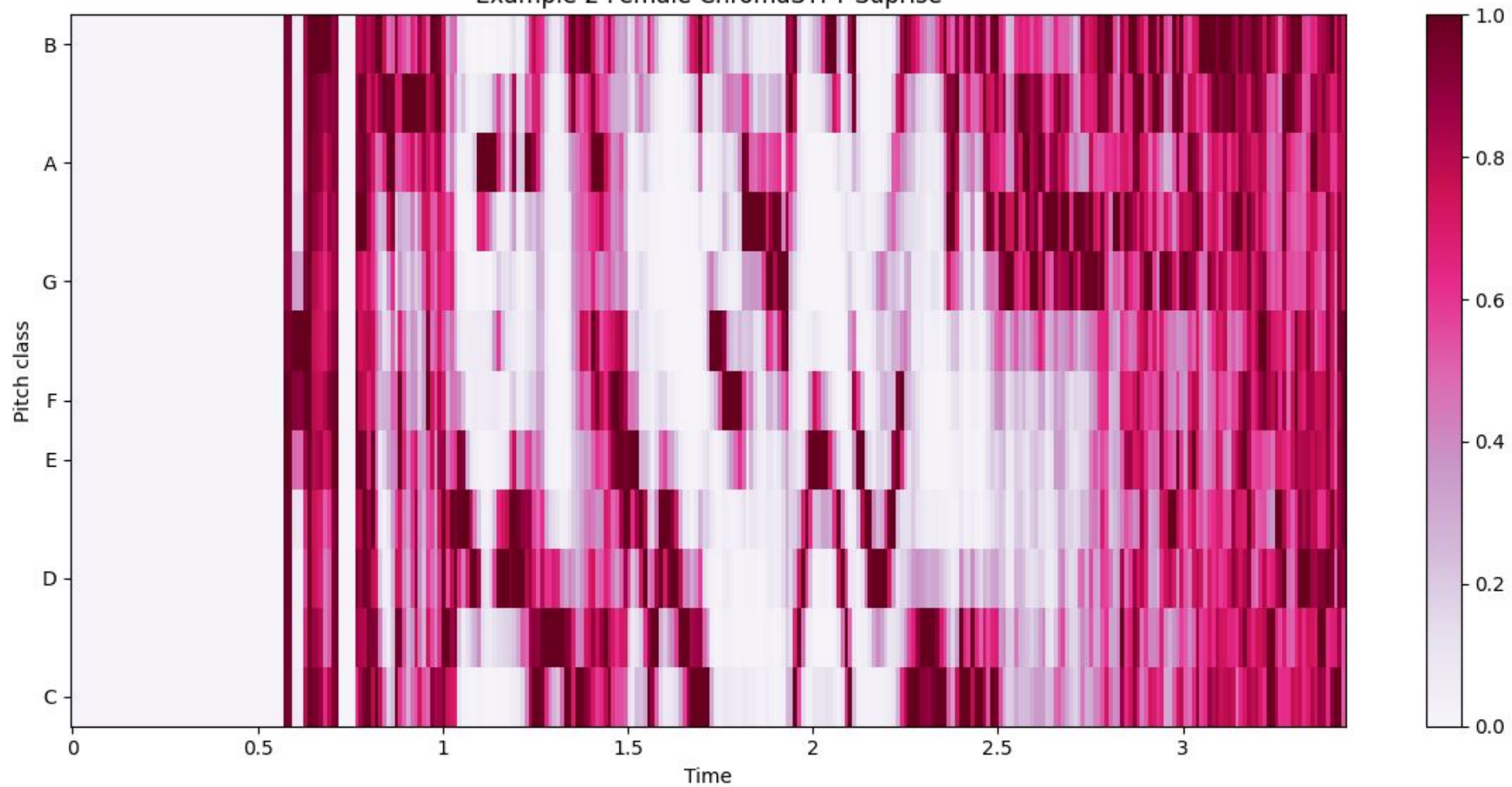
Example 2 Female Spectrogram Surprise Audio Waveform



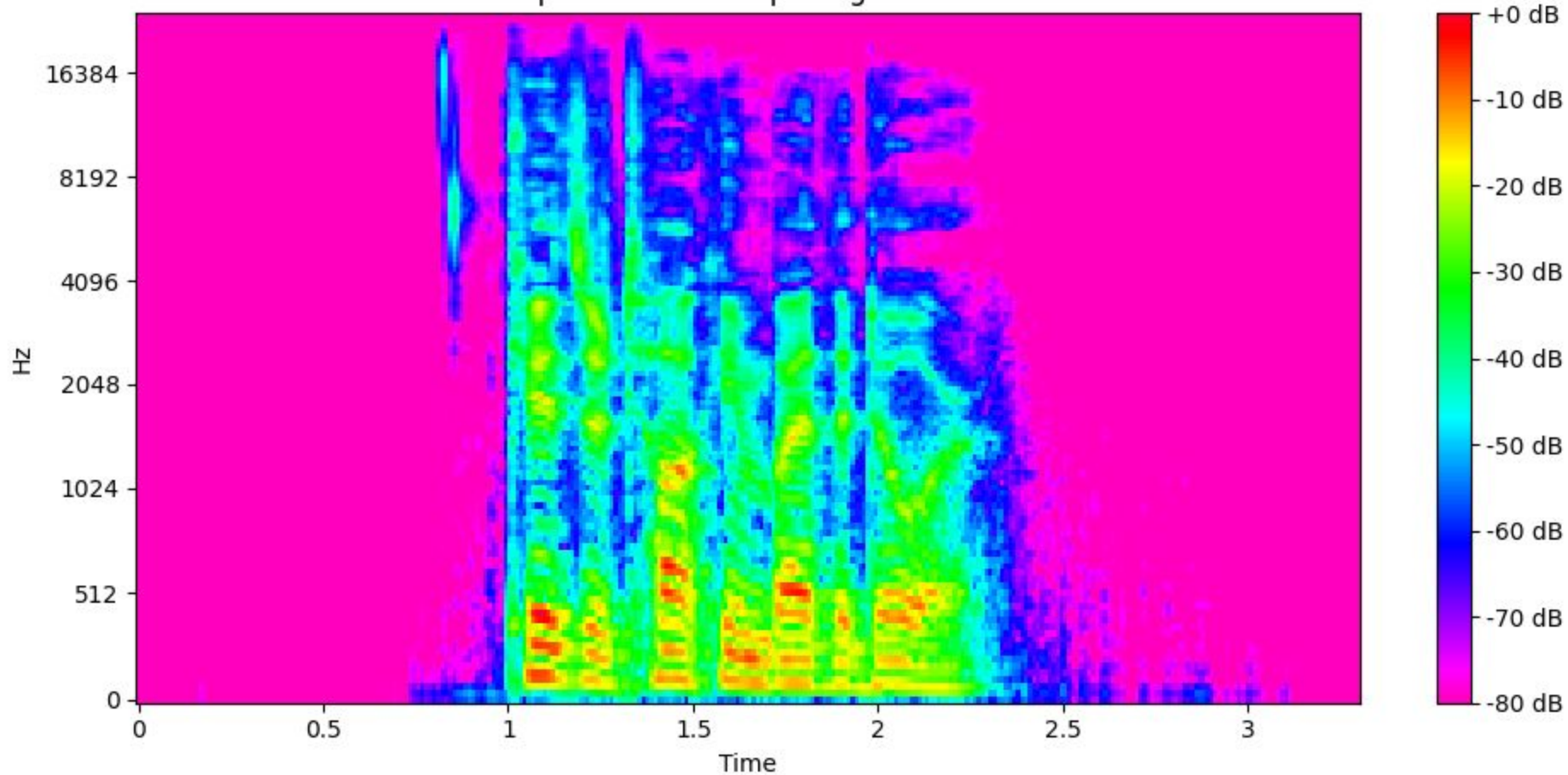
Example 1 Male Chroma STFT Emotion: Neutral



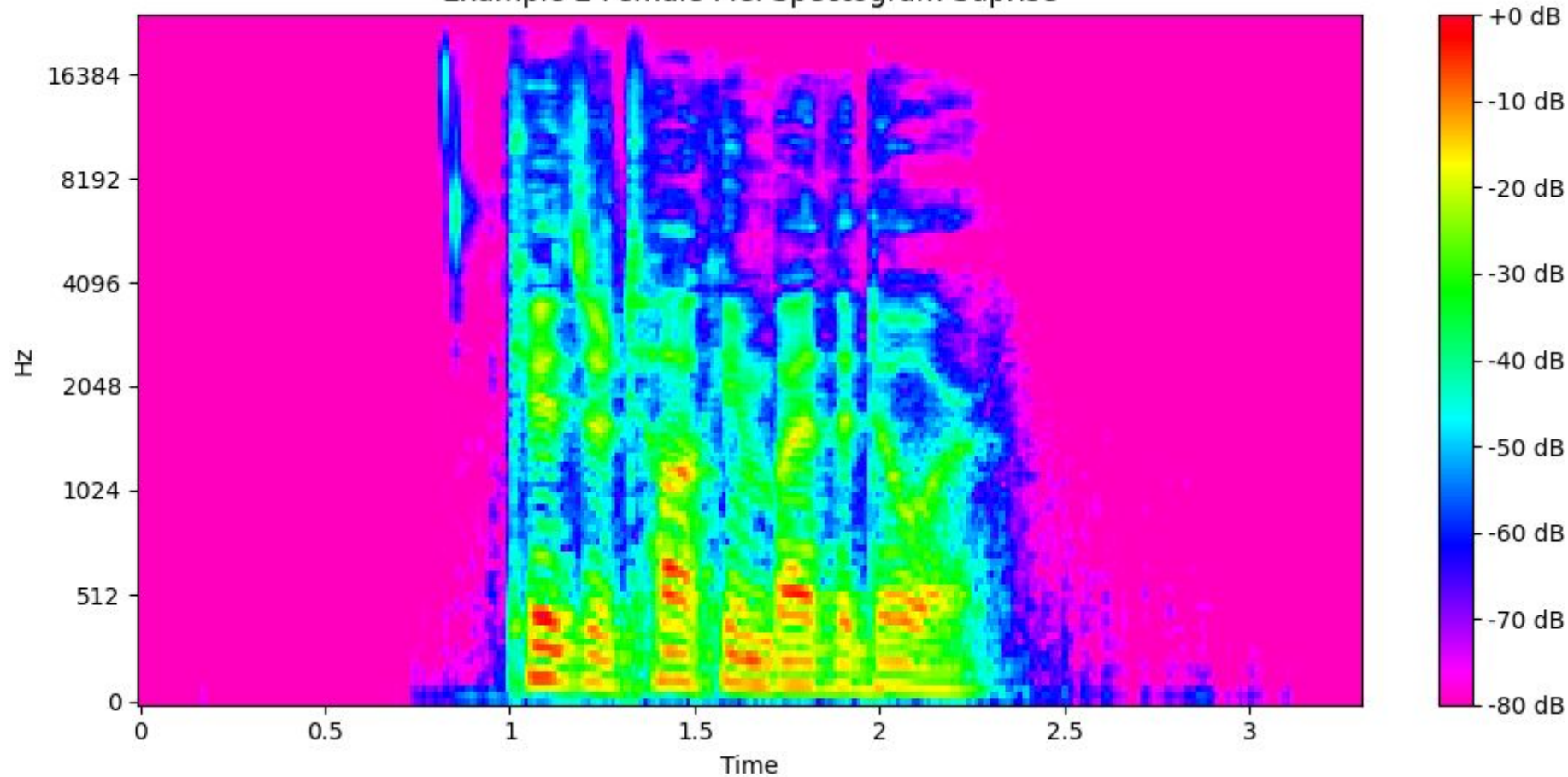
Example 2 Female ChromaSTFT Surprise



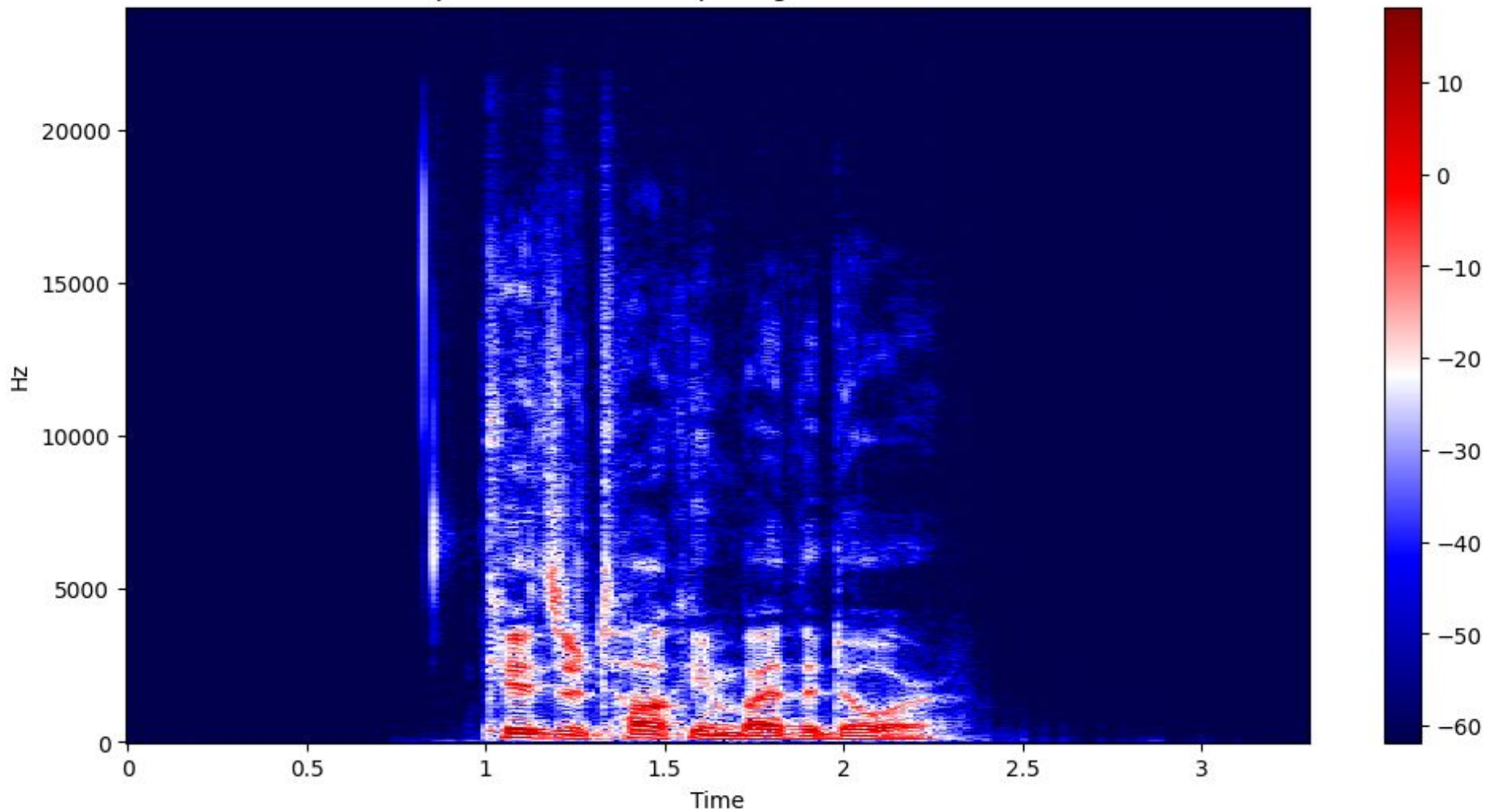
Example 1 Male Mel Spectrogram Neutral



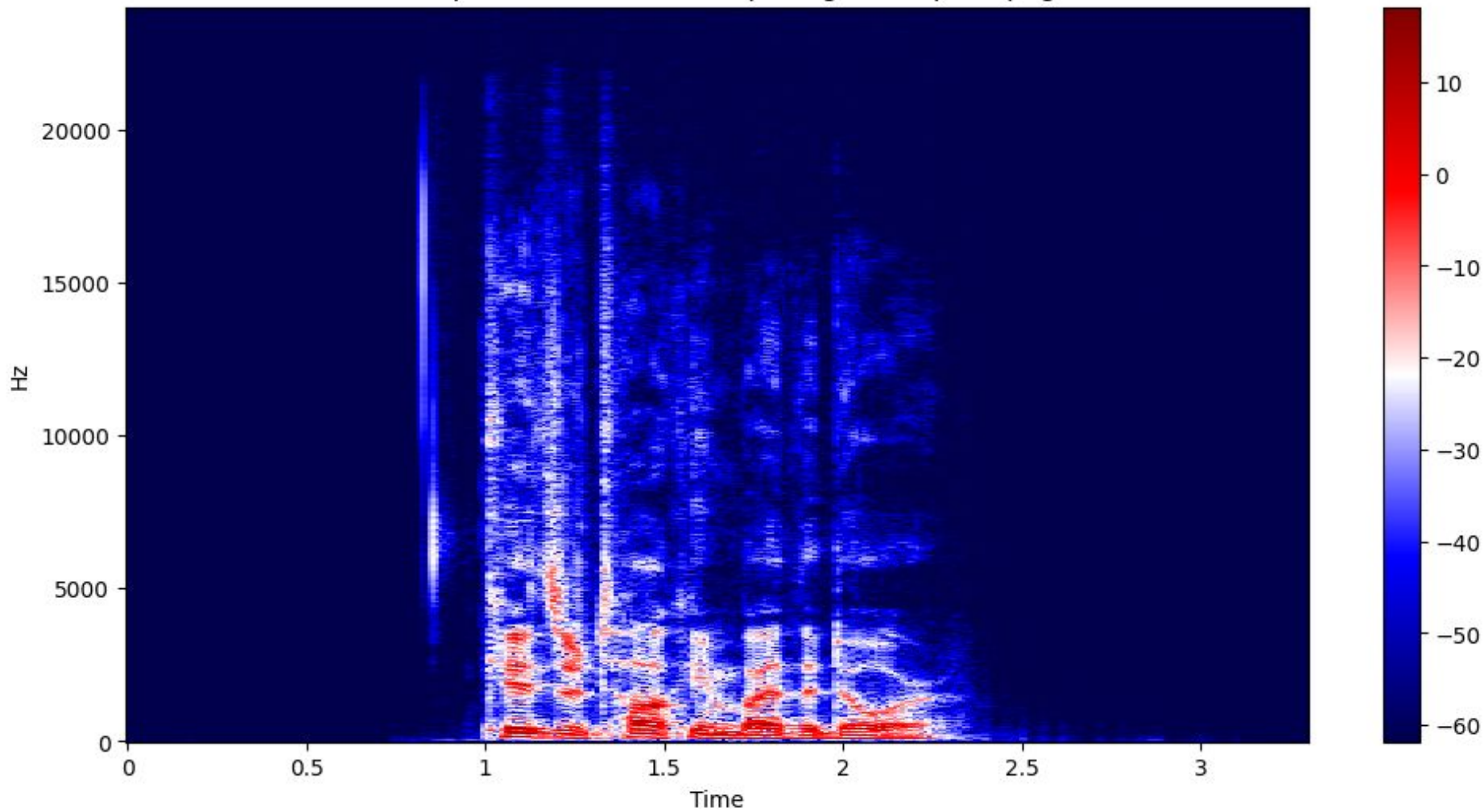
Example 2 Female Mel Spectrogram Surprise



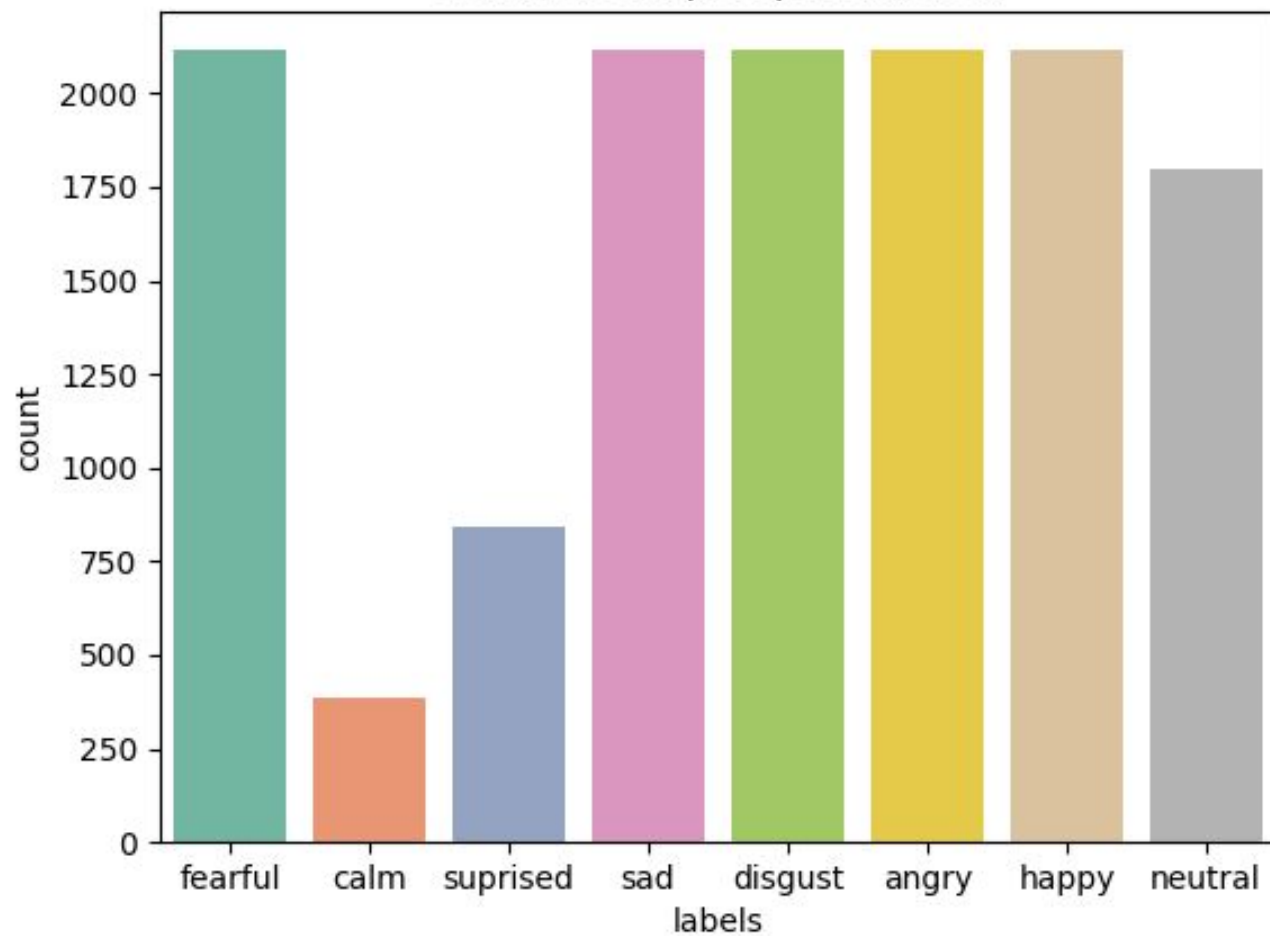
Example 1 Male Fourier Spectrogram Emotion: Neutral



Example 2 Female Fourier Spectrogram Surprise.png

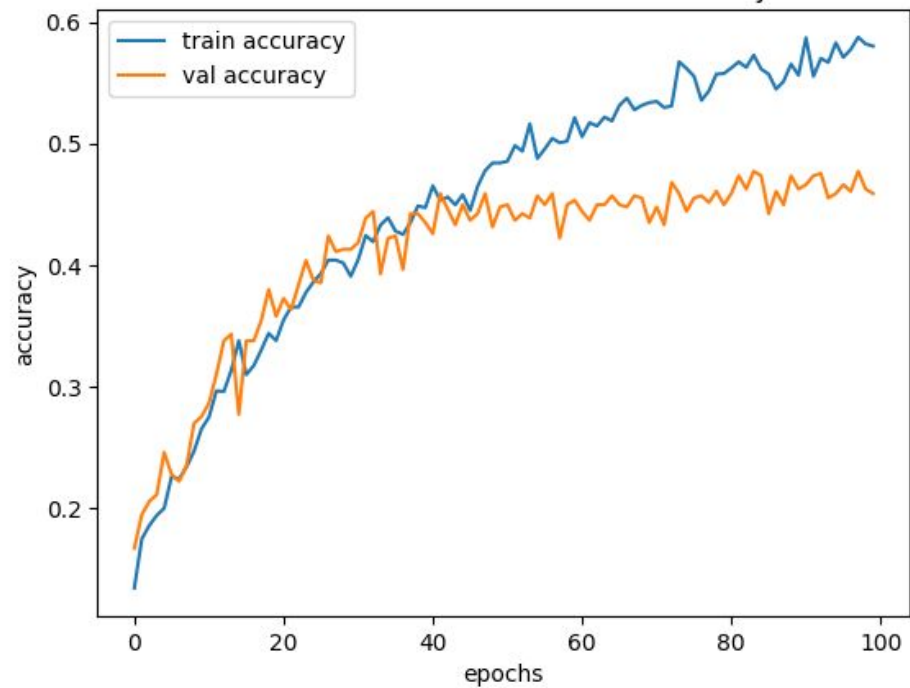


Count of Samples per Emotion

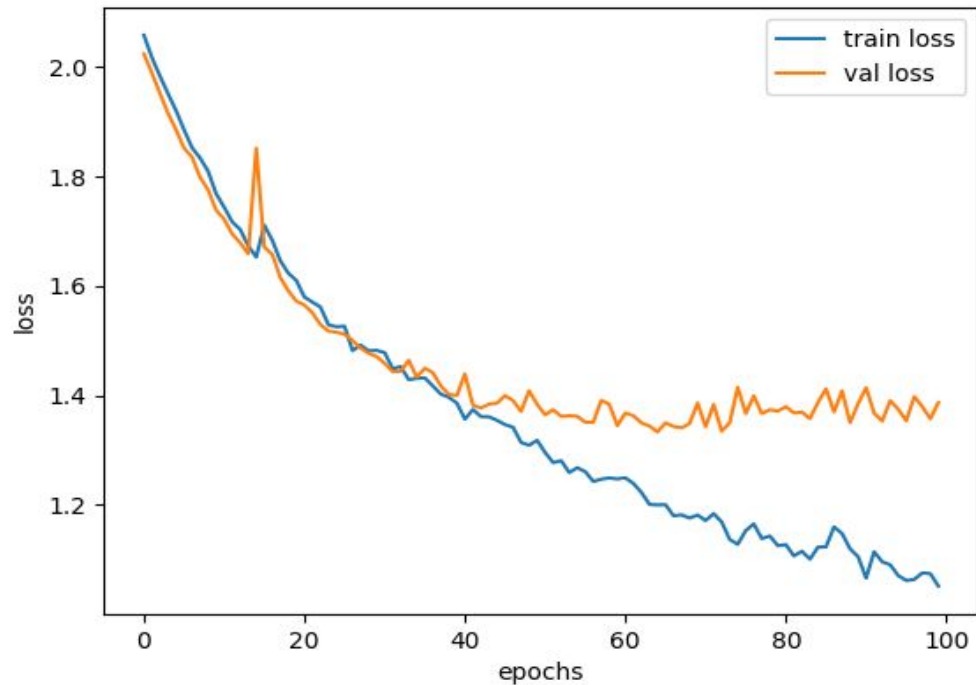


Performance

LSTM1 train versus validation accuracy

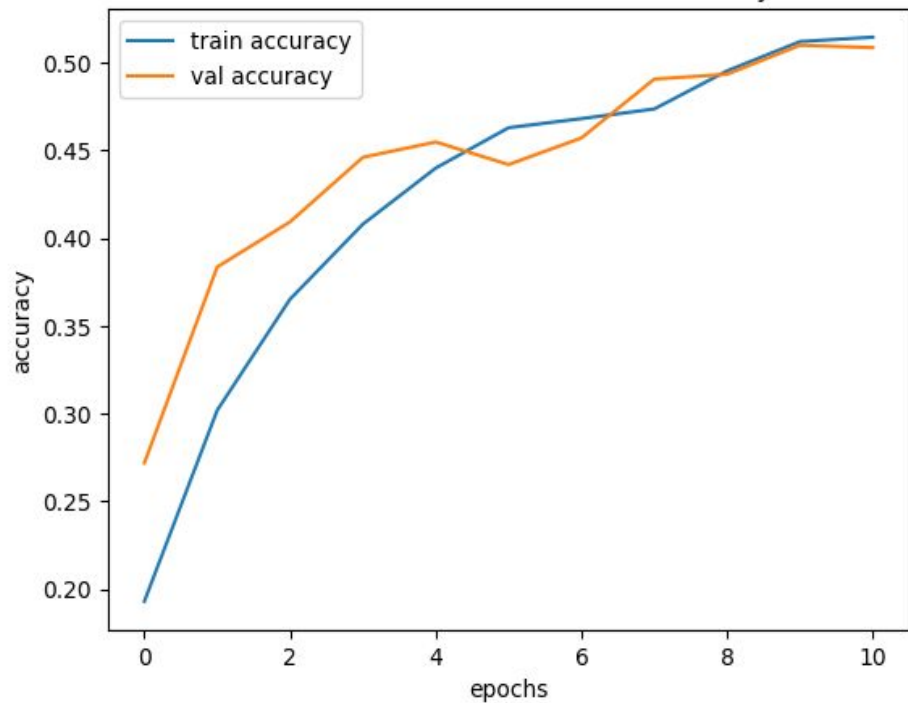


LSTM1 train versus validation loss

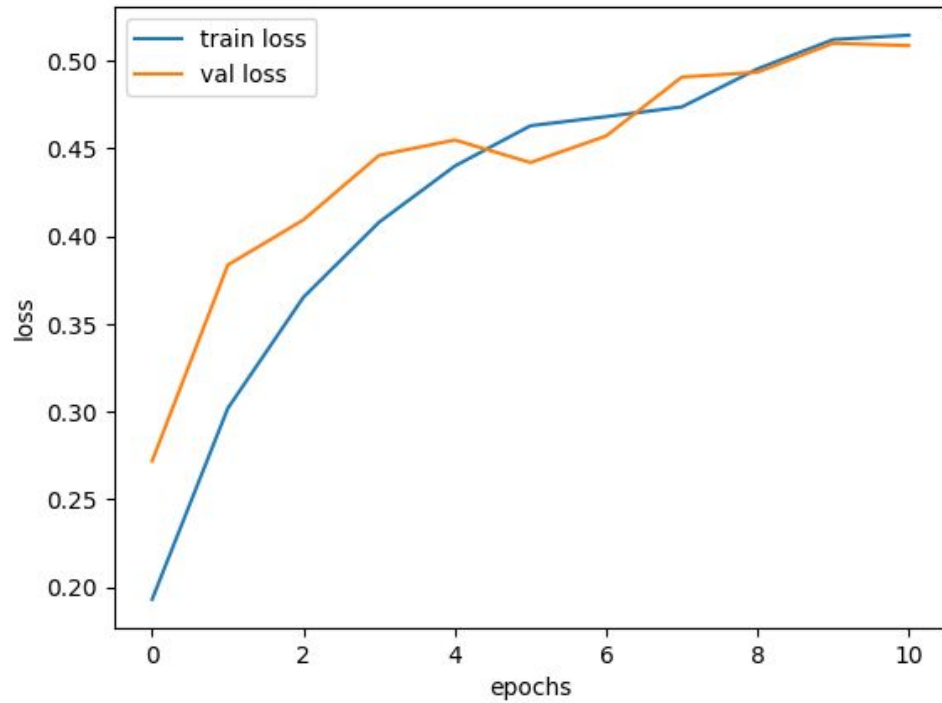


Performance

LSTM2 train versus validation accuracy

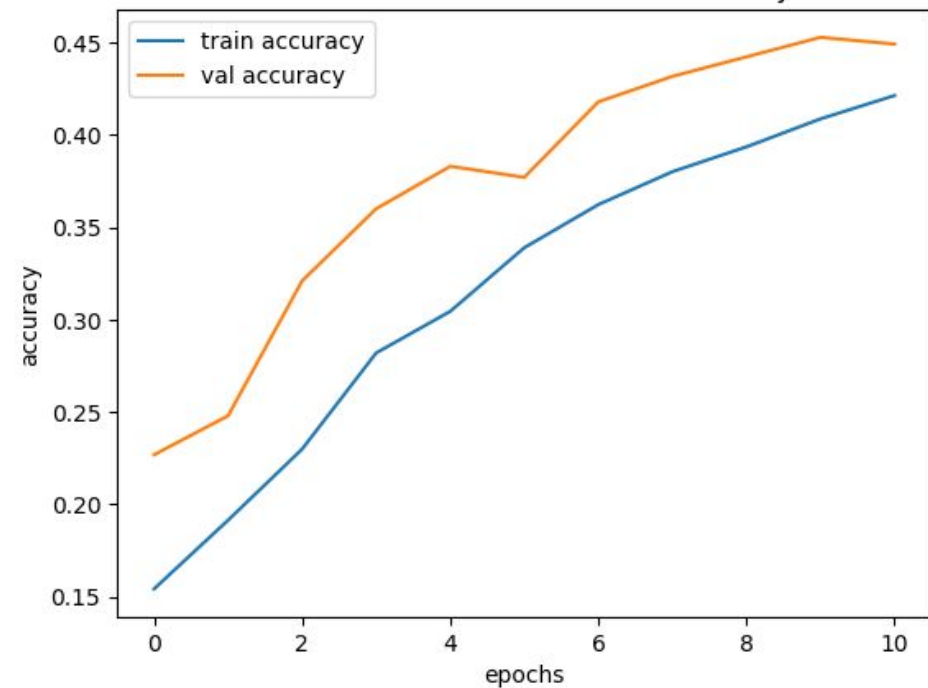


LSTM2 train versus validation loss

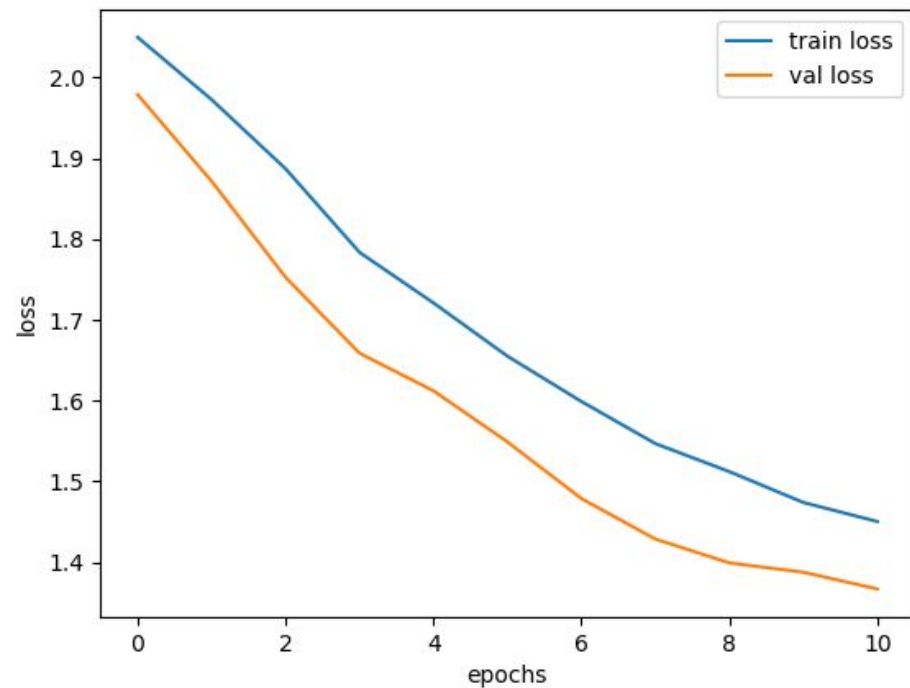


Performance

LSTM3 train versus validation accuracy

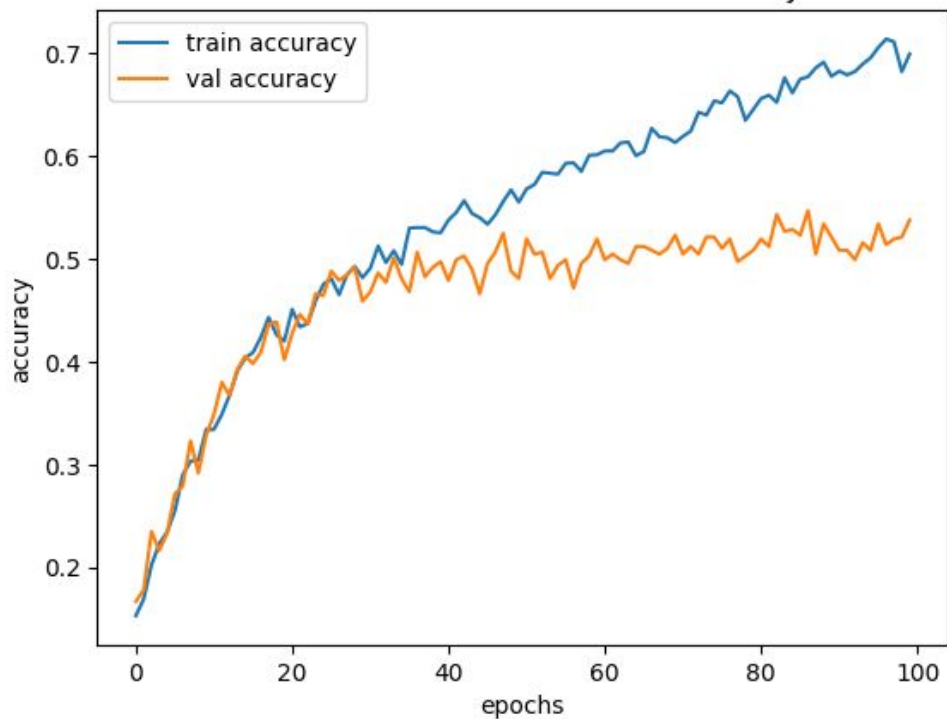


LSTM3 train versus validation loss

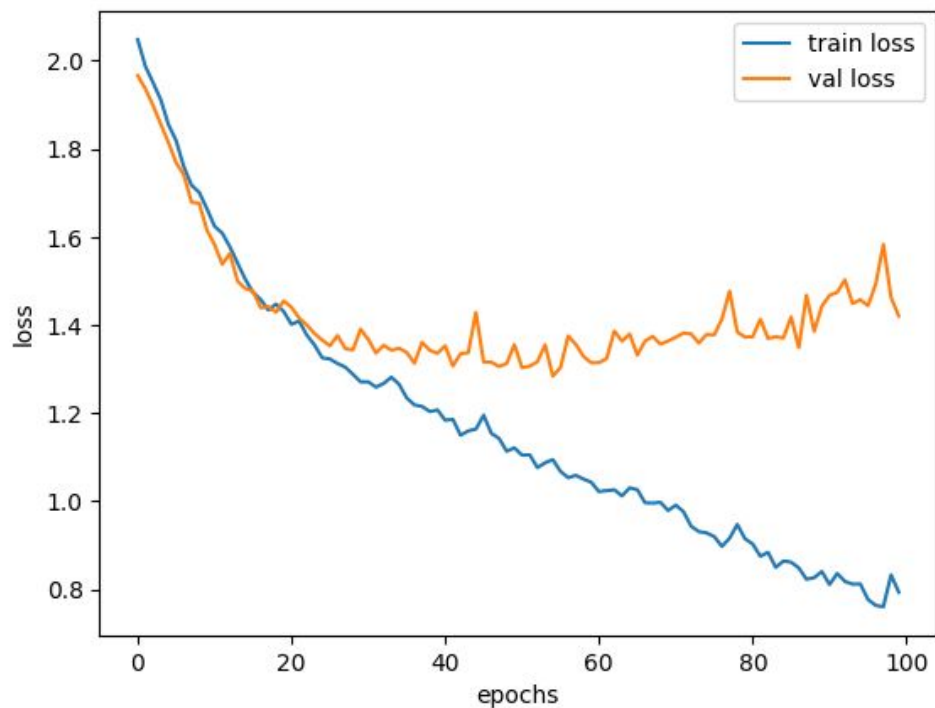


Performance

LSTM4 train versus validation accuracy



LSTM4 train versus validation loss



Outcomes

- Best model:

model_LSTM, 'first_lstm_model.h5'

- Best score:

accuracy: 0.6710 - loss: 0.8293 - val_accuracy: 0.6022 - val_loss: 1.1673

Future work:

- address class imbalances
- continue to tune and iterate or at least identify performance decline.
- Access other possible standardisations techniques, feature extraction and tuning methods.
- Utilize pretrained models such as Whisper, WavLM, and Wav2Vec 2.0, which can be fine-tuned for SER tasks.

Sources: