

cancelling algorithms and machine learning.

2023 Part IV Project Group: 58

Authors: Edward Chan & Timothy Cabrera

Project Introduction

Speech enhancement is the application of algorithms to deconvolve a speech signal from a noise-corrupted signal. This is most commonly done through statistical estimation or neural networks that provide filter coefficients or apply non-analytic filtering methods.

However, for both cases, the estimation of intelligible speech from signals with very high noise levels remains a challenge without using a microphone array to utilise spatio-temporal information of the signal.

Many standard communication devices do not have access to these multi-channel microphone configurations. Therefore, this research project aims to contribute to the field of single-channel speech enhancement by investigating the potential of a proposed hybrid deep neural network and model-based approach for single-channel speech enhancement.

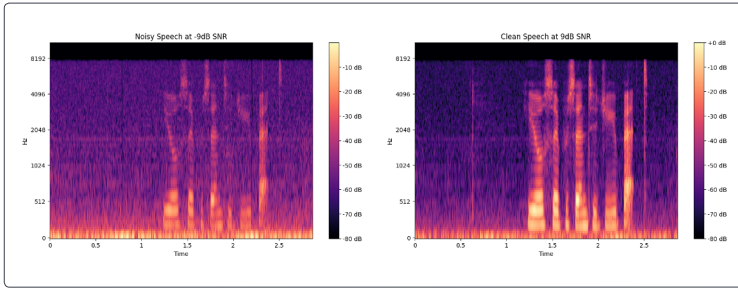
Project Goals

Our key project requirements are:

- A machine learning approach that can estimate noise and speech based on an observed mixture.
- A classifier that can distinguish between different signal-to-noise levels.
- The resulting enhanced signal from mixtures with SNR levels $< 0\text{dB}$ should be intelligible.

Background Theory

The signal-to-noise ratio (SNR) is a ratio that measures the presence of a signal relative to the noise in the signal. Figure 1 shows a spectrum of a mixture considered low SNR and high SNR.



Figures 1: Spectrograms representing a speech audio signal's low SNR and high SNR.

To enhance the signal, it is common practice to use a collection of pre-defined noise/speech spectra candidates (codebook), where each candidate is scaled such that the sum of all candidates will result in the desired spectrum of noise and speech, as shown in Figure 2. However, this only works for some types of mixtures.

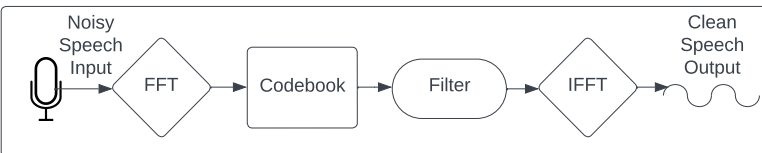


Figure 2: Conventional Speech Enhancement Algorithm [2]

Results

A confusion matrix of the testing data is used to distinguish between the different SNR levels from the CNN classifier model [3], as seen in Figure 4. The overall accuracy is 90.25%. Trained on 1200 utterances.

Actual label \ Predicted label	clean	9dB	6dB	3dB	0dB	-3dB	-6dB	-9dB
clean	148	2	0	0	0	0	0	0
9dB	2	121	27	0	0	0	0	0
6dB	0	10	138	2	0	0	0	0
3dB	0	0	10	136	4	0	0	0
0dB	0	0	0	4	139	7	0	0
-3dB	0	0	0	0	7	138	5	0
-6dB	0	0	0	0	0	7	122	21
-9dB	0	0	0	0	0	0	9	141

Figure 4: CNN classification model accuracy when testing using a confusion matrix.

A common assumption for speech enhancement is that the accurate estimation of noise will result in a better estimate of speech. Presented in Figure 5, is an example of the GAN accurately tracking the true underlying noise spectrum from a -9dB mixture.

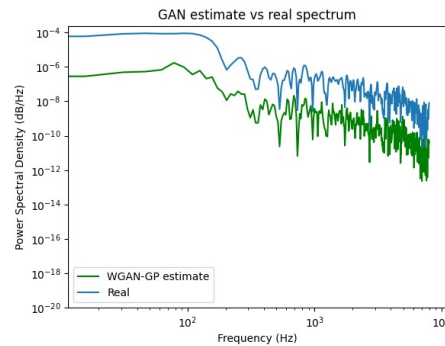


Figure 5: An example of the GAN noise spectrum vs. the real noise spectrum.

Discussion and Conclusion

The proposed design was able to accurately estimate the spectrum of noise at a range of different SNRs. In contrast, it could not reliably estimate the speech spectra at different SNRs. As such, a custom filter was designed such that only the noise spectra are needed for spectral attenuation. However, results show that it still results in distortion, and unintelligible speech.

Future considerations:

- Increase the training data for the GAN.
- Add more classifications, such as no speech with more background noises.
- Further, develop the GAN so that speech spectra can also be estimated reliably to be compatible with a Wiener filter.

Generative Adversarial Network Design & Implementation

Our system implementation enhances speech using a generative adversarial network, as shown in Figure 3.

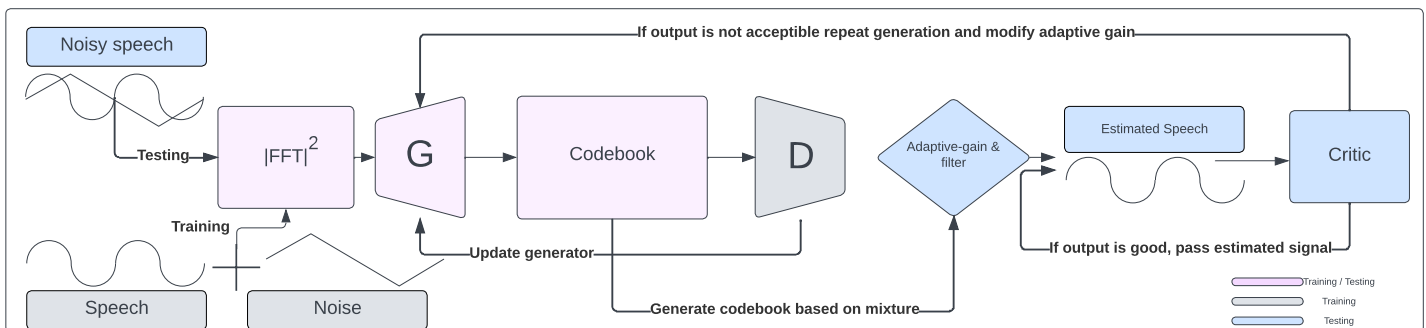


Figure 3: Flow chart of the Preposed Design

- Generative adversarial networks (GAN) are neural networks that are commonly used to generate realistic images through small input vectors [1].
- A GAN has been designed to take in the mixture's power spectral density and generate a codebook based on the observed mixture.
- The codebook is scaled through an adaptive gain filter in order to attenuate frequencies with noise.
- Trained on VCTK [5] (speech dataset) and DEMAND [4] (noise dataset) corpora. 300,000 utterances were used for training the GAN.
- STFT parameters: N-FFT = 1024, Window function = Hanning/Hann, Hop length = 128 samples, Samples per window = 512 samples, Zero-padding = left aligned.