# A System for Bandwidth Extension of Narrow-band Speech

Georgios Manos

Bachelor's Thesis
Department of Computer Science
University of Crete

ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
UNIVERSITY OF CRETE

Advisor: Prof. Yannis Stylianou
Supervisor: Dr George P. Kafentzis

December 5, 2024

# Contents

# List of Figures

# List of Tables

# Abstract

This work implements a systematic approach to expand the bandwidth of a narrowband (NB) signal, specifically designed for signals like speech. It utilizes a parametric technique based on a discrete acoustic tube model (DATM) that does not necessitate prior training. The process involves several steps: computing narrowband linear predictive coefficients (LPCs) from the incoming narrowband speech signal, recursively calculating narrowband partial correlation coefficients (PARCORS), deriving $M$ area coefficients from these partial correlation coefficients, and then extracting $M$ area coefficients through interpolation. Subsequently, it computes wideband (WB) PARCOR values based on the $M$ area coefficients and derives wideband LPCs from these wideband PARCORS. This methodology then synthesizes a wideband signal using these LPCs alongside a wideband excitation signal. The resulting synthesized signal undergoes highpass filtering to generate a highband signal, which is combined with the original narrowband signal to produce the desired wideband signal. In an optimized version of this approach, it involves converting $M$ area coefficients into logarithmic representations (log-area coefficients), extracting $M$ log-area coefficients using a shifted interpolation technique, and reverting these extracted log-area coefficients back to the original $M$ area coefficients before generating the wideband PARCORS.

This study assesses six speech samples sourced from the TIMIT database. Employing LibROSA, we downsample the speech signal to 8kHz and then implement a 2x bandwidth expansion to restore it to 16kHz. The evaluation involves three methods: auditory assessment, visual examination of waveforms and spectrograms, and the application of three metrics - Mean Squared Distance (MSD) between their log-magnitude spectra, Spectral Convergence (SC), and L1-norm Mel-Spectrogram difference. Results show that the synthesized WB signal is closer to the original WB signal.

# Acknowledgements

This work was advised by Professor Yannis Stylianou and supervised by Dr George Kafentzis. Meeting and collaborating with them was definitely my most valuable experience that I earned during my undergraduate studies. Their contribution not only to this work but in general on my academic career played a pivotal role and shaped most of my decisions, while it helped me build many of the soft skills that got me to where I am today.

The way I met them was of course through their courses. Their courses are recognized among every student to be of the highest quality, along with the lectures. Aside from making the material simplified and intuitive, they were always able to present the content in a fun way, with many jokes surrounding them. In the complicated world of signal processing, they were able to spark my curiosity to pursue the field further, in an attempt to understand concepts that I had never seen before. This took me through CS370 Digital Signal Processing until the post-graduate course of CS578 Speech Processing, courses which I believe even up to date (December 2024) were by far of the highest quality that I have ever had the chance to attend. Courses which made me fall in love with the field, and also were the most influential factors for my teaching values that I hold.

Professor Yannis taught me to persevere and never give up. Dr George taught me how to align my needs and my goals in this journey to keep pushing in the right direction. I would like to wholeheartedly thank them for their key contribution to not only my career, but my personality and my core values. I really am hoping that our paths meet again in the future.

# Chapter 1

# Introduction

In this work, we explore the system[1] invented by David Malah and Richard Vandervoot Cox for extending the bandwidth of a signal. The goal of this work is to implement a version of the system as suggested by the patent in a modular way such that one can modify its components and the system's configuration, while also assessing its performance both qualitative and quantitative. The code is written in Python and is available on GitHub. Figure 1.2 and 1.1 were taken directly from the patent, while some equations and implementation ideas were taken from CS578 Speech Processing course.

## 1.1 Short System Description

This is a system for extending the bandwidth of a narrowband signal such as a speech signal. The method applies a parametric approach to bandwidth extension but does not require training. The parametric representation relates to a discrete acoustic tube model (DATM). The method comprises computing narrowband linear predictive coefficients (LPCs) from a received narrowband speech signal, computing narrowband partial correlation coefficients (parcors) using recursion, computing $M_{nb}$ area coefficients from the partial correlation coefficient, and extracting $M_{wb}$ area coefficients using interpolation. Wideband parcors are computed from the $M_{wb}$ area coefficients and wideband LPCs are computed from the wideband parcors. The method further comprises synthesizing a wideband signal using the wideband LPCs and a wideband excitation signal, highpass filtering the synthesized wideband signal, and combining the highband signal with the original narrowband signal to generate a wideband signal. As mentioned by the inventors, the preferred variation of the invention is implemented here, where the $M_{nb}$ area coefficients are converted to log-area coefficients for the purpose of extracting, through shifted interpolation, $M_{wb}$ log-area coefficients. The $M_{wb}$ log-area coefficients are then converted to $M_{wb}$ area coefficients before generating the wideband parcors.

The system is described in Figure 1.1, and explained in depth in Section 2

## 1.2 Experiments Setup

The present implementation of the patent follows the flow diagram as shown in Figure 1.2. The bandwidth extension is performed on a frame-by-frame basis. The inventors mention that some of the parameter values discussed are merely default values used in simulation. In this work, we mostly follow their suggested default values, while in section 3, we also discuss the effect of some of those parameters.

This work is evaluated on 6 speech samples taken from the TIMIT database. Also, using librosa, we are downsampling the speech signal to 8khz and applying 2x bandwidth expansion (back to 16khz).

The results are evaluated in 3 ways; aurally, visually through waveforms and spectrograms, and using 3 metrics, Mean Squared Distance (MSD) between their log magnitude spectrum, Spectral Convergence (SC), L1-norm Mel Spectrogram difference.

---

[1]You can find the system patent here
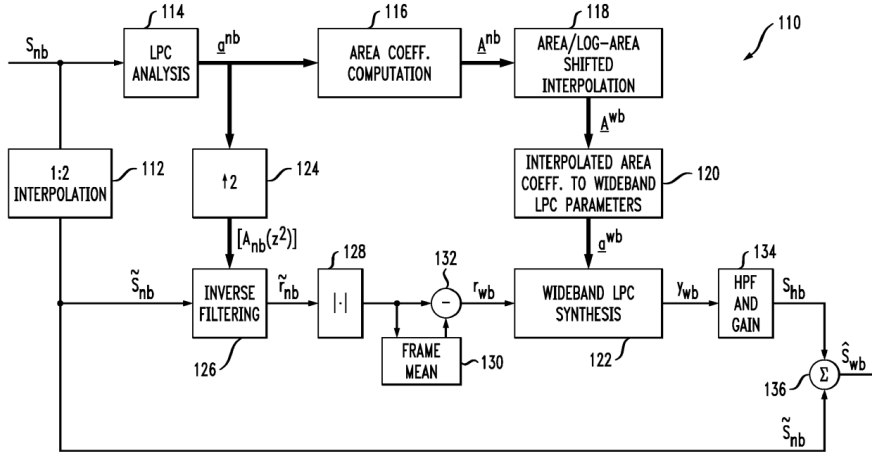
FIG. 8

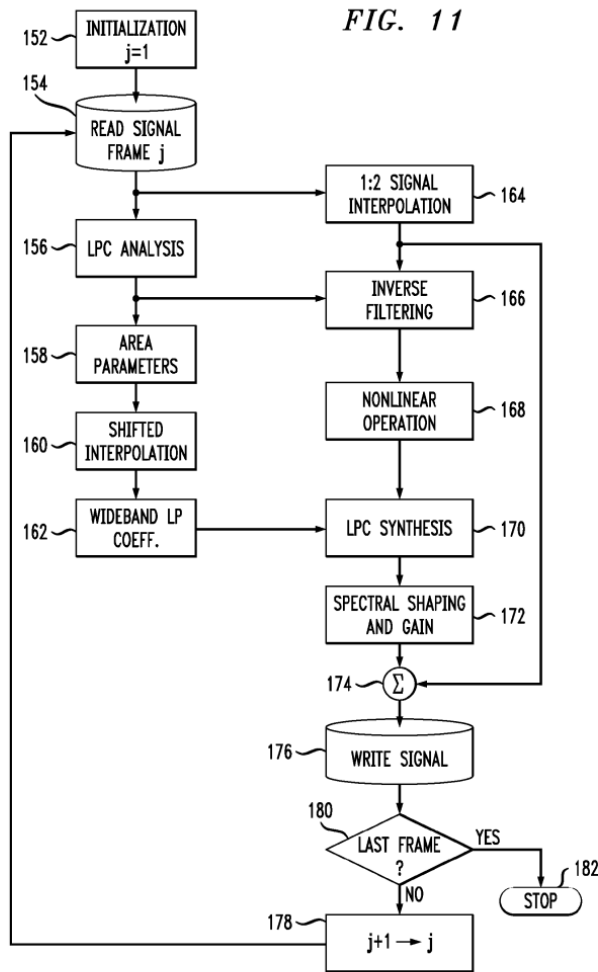Figure 1.1: System description

FIG. 11

Figure 1.2: Algorithm flowchart

# Chapter 2

# A System for Bandwidth Extension of Speech

In this chapter, we discuss the system presented in patent by Malah et al.

## 2.1 Parameters

During the initialization (**152**), the following parameters are established: Input signal frame length = N (60 ms), Frame update step = N / 2, Number of narrowband DATM sections M (16), Sampling Frequency (in Hz)=$f_s^{nb}$ (8000), Input signal upper cutoff frequency in Hz=$F_c$ (3900, as the input was a microphone; inventors suggest using 3600 for MIRS input and 3400 for IRS telephone speech) along with the order (20) of the high-pass butterworth filter (**172**), and R(0) modification parameter $\delta$ (0.01). Also, the shift parameter (0.25) for the area shifted interpolation (**160**) applied to the Area Parameters is configured, and the pre-emphasis coefficient (0.97) applied prior to the LPC analysis. The inventors mention that the values may vary depending on the source characteristics and application.

## 2.2 Signal Analysis

A speech signal in wav format is read from disk. For frame j, the signal undergoes an LPC analysis (**156**) that comprises of the following steps: computing a correlation coefficient $\rho_1$, apply a Hann window of length N to the signal frame, , pre-emphasizing the windowed signal frame using $(1 - \rho_1 z^{-1}$, computing the M+1 autocorrelation coefficients R(0), R(1), ..., R(M), modifying R(0) by a factor $(1+\delta)$, and applying the Levinson-Durbin recursion[1] to find LP coefficients $\underline{\alpha}^{nb}$ and reflection coeffcients $\underline{k}^{nb}$. Parcors are then computed through reflection coefficients using the relationship $\underline{k}^{nb} = -\underline{r}^{nb}$. Finally, the gain G is computed using equation 2.1

$$G = \sqrt{\sum_i^{M+1} \underline{\alpha}_i^{nb} R(i)} \tag{2.1}$$

## 2.3 Computing Area Parameters

Next, the area parameters are computed (**158**) according to an important aspect of the invention. Computation of these parameters comprises computing M area coefficients via equation 2.3 that describes the parameters of the discrete acoustic tube model (DATM), and computing M log-area coefficients. The relationship between the LP model parameters and the area parameters of the DATM are given by the backward recursion of equation 2.3, where $A_1$ corresponds to the cross-section at the lips and $A_{M^{nb+1}}$ corresponds to the cross-section at the glottis opening. Computing the M log-area coefficients is mentioned to be an optional step but a preferred one and thus followed in this implementation. The computed area or log-area coefficients are shift-interpolated (**160**) by a desired factor with a proper sample shift. Cubic spline is applied by default as interpolation method.

$$r_i^{wb} = \frac{A_i^{wb} - A_{i+1}^{Wb}}{A_i^{wb} + A_{i+1}^{Wb}}, \ i = 1, 2, ..., M_{wb} \tag{2.2}$$

The next step relates to calculating wideband LP coefficients (**162**) and comprises computing the wideband parcors from interpolated area coefficients via equation 2.2 and computing wideband LP coefficients, $\alpha^{wb}$, by

---

[1]A thank you to Dr. George Kafentzis for providing a python implementation of the algorithm

applying the Step-Down Recursion to the wideband parcors. As log-area coefficients were used, exponentiation is applied to obtain the interpolated area coefficients.

$$A_i = \frac{1 + r_i}{1 - r_i} A_{i+1}; \ i = M_{nb}, M_{nb} - 1, ..., 1 \tag{2.3}$$

## 2.4    Producing Wideband Signal

Returning now to the branch from the output of step 154, step 164 relates to signal interpolation. Step 164 comprises interpolating the narrowband input signal, $S_{nb}$, by a factor (i.e. 2, upsampling and lowpass filtering). This step results in a narrowband interpolated signal $\tilde{S_{nb}}$. The signal $\tilde{S_{nb}}$ is inverse filtered (166) using, for example, a transfer function of $A_{nb}(z^2)$ having the coefficients shown in equation 2.4, resulting in a narrow band residual signal $\tilde{r_{nb}}$ sampled at the interpolated-signal rate.

$$\underline{\alpha}^{nb} \uparrow 2 = \{1, 0, \alpha_1^{nb}, 0, \alpha_2^{nb}, 0, ..., \alpha_{M^{nb}-1}^{nb}, 0, \alpha_{M^{nb}}^{nb}\} \tag{2.4}$$

Next, a non-linear operation is applied to the signal output from the inverse filter. The operation comprises fullwave rectification (absolute value) of residual signal $\tilde{r}_{nb}$ (168). The inventors mention that other nonlinear operators may also optionally be applied. Also, other potential elements associated with step **168** may comprise computing frame mean and subtracting it from the rectified signal (as shown in Fig 1.1), generating a zero-mean wideband excitation signal $r_{wb}$; optional compensation of spectral tilt due to signal rectification via LPC analysis of the rectified signal and inverse filtering. The preferred setting here is no spectral tilt compensation.

Next, the highband signal must be generated before being added (**174**) to the original narrowband signal. This step comprises exciting a wideband LPC synthesis filter (**170**) (with coefficients $\underline{\alpha}^{wb}$ by the generated wideband excitation signal $r_{wb}$, resulting in a wideband signal $y_{wb}$. Fixed or adaptive de-emphasis are optional, but the default and preferred setting is no de-emphasis. The resulting wideband signal $y_{wb}$ may be used as the output signal or may undergo further processing. In this implementation, the further processing involves the following steps: the wideband signal $y_{wb}$ is highpass filtered (**172**) using a butterworth digital filter of order 16 with cutoff frequency $f_c = 3900 hz$ to generate a highband signal and the gain is also applied here. The inventors mention using a fixed gain value (e.g. 2) here instead or adaptive gain matching, but we used the previously computed gain value. The resulting signal is $S_{hb}$ (as shown in Fig 1.1). The butterworth filter was applied using *scipy.signal.filtfilt* as it is zero-phase filtering, which doesn't shift the signal as it filters.
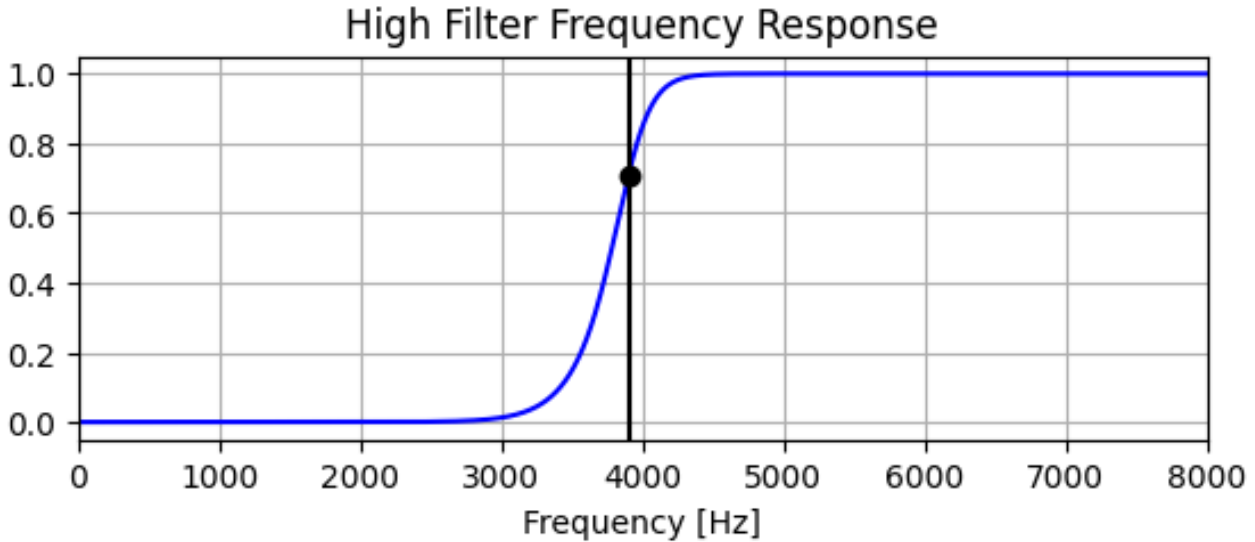
The Buttworth filter is presented in Figure 2.1.



Figure 2.1: Butterworth Filter Frequency Response

## 2.5    Signal Synthesis

Finally, the output wideband signal is generated. This step comprises generating the output wideband speech signal by summing (**174**) the generated highband signal $S_{hb}$ with the narrowband interpolated input signal $\tilde{S}_{nb}$. The resulting summed signal is stored in the main buffer. The output signal frame (of 2N samples) can either be

overlap-added (OLA, with a half-frame shift of N samples) to a signal buffer, or because $\tilde{S}_{nb}$ is an interpolated original signal, the center half-frame (N samples out of 2N) is extracted and concatenated with previous output stored in a signal buffer. OLA is used in this implementation.

# Chapter 3

# Experiments

In this chapter, the system's results are presented for speech file *arctic_bdl1_snd_norm.wav*. All speech files were originally sampled at 16khz, and were downsampled to 8khz prior to the system's input. The system's output is then compared to the narrowband signal (8khz), the original wideband signal (16khz) - which will be considered as ground truth, as well as a signal which was simply upsampled to 16khz without bandwidth expansion.

This report also includes results for the rest of the speechfiles in the database, presented in section 5.

## 3.1 Time Domain Waveforms

The signals' waveforms are presented in Figure 3.1. It includes the waveforms for the original signal at 16khz *"Orig (WB)"*, the signal that was upsampled using High-quality FFT-based bandlimited interpolation *"Interpolated"* and finally the bandwidth expansion system output *"Reconstructed"*.
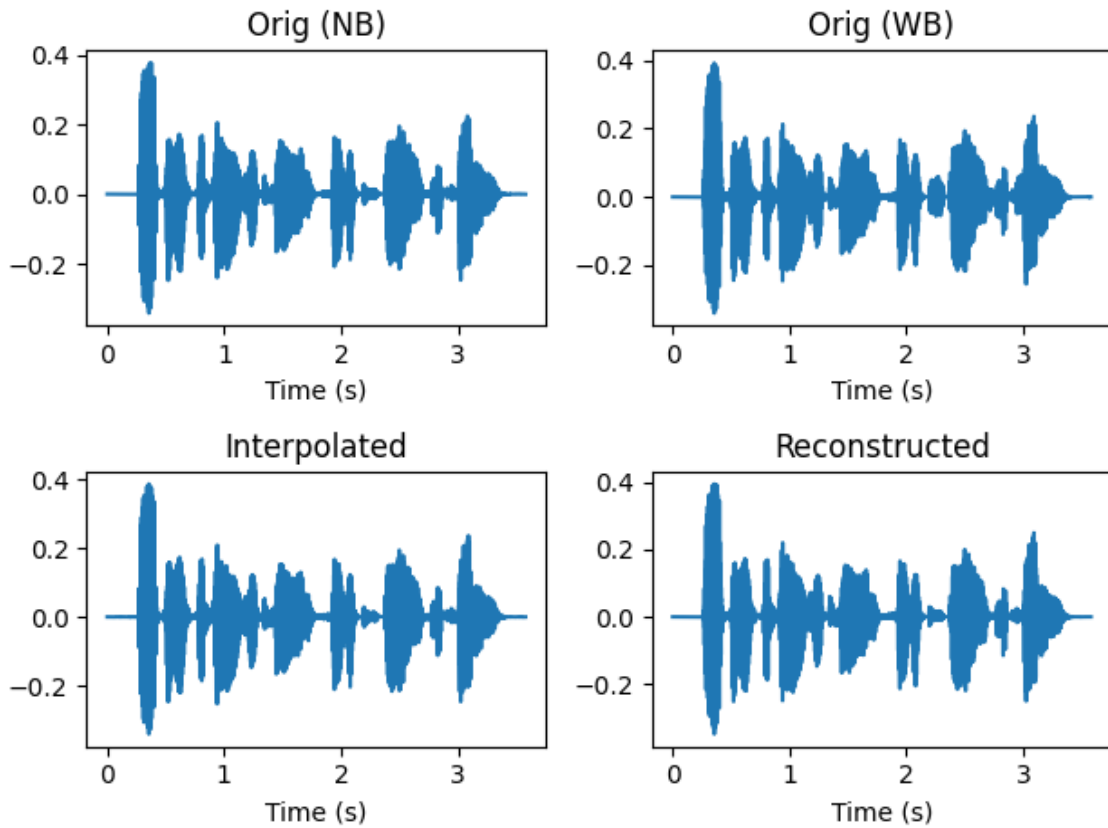


Figure 3.1: Signal Waveforms for each case: Original Narrowband, Original Wideband, Interpolated and Reconstructed signal with BE.

One may observe that all 3 waveforms overlap almost totally, with the reconstructed signal almost completely covering the interpolated one, while only on a few segments the original one differs. However, the
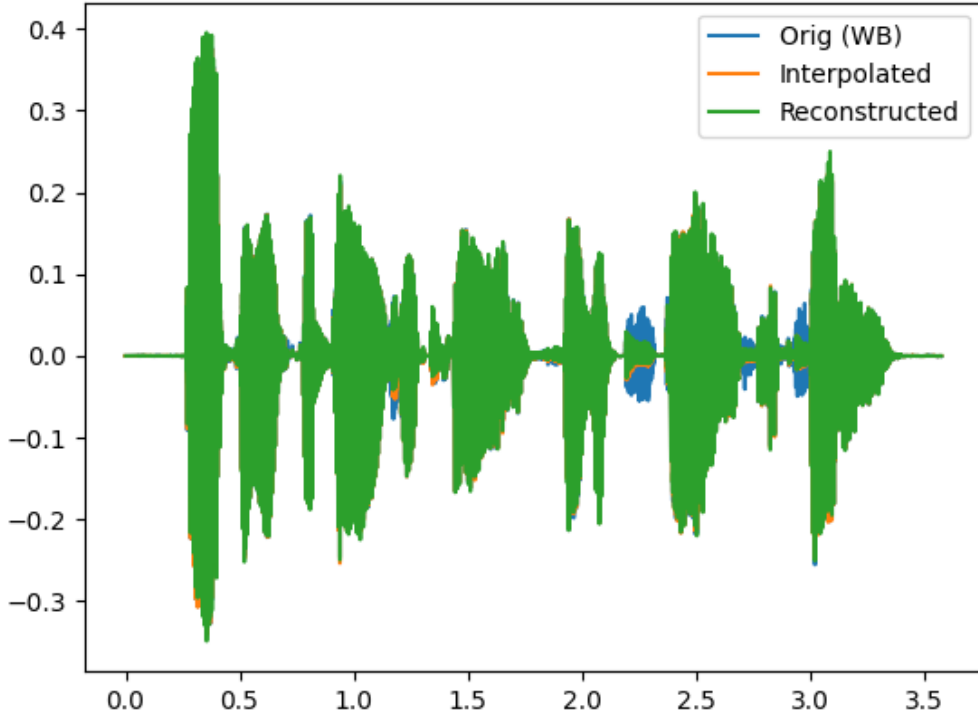
Figure 3.2: Signal Waveforms Merged

auditory system decomposes sound into elementary time-frequency representations, where there are several acoustic cues for speech perception, including Formants, Harmonics, voice pitch and syllable rate.

Aurally, as per to the author's subjective opinion, the narrowband signal has poor resolution, resulting to a telephone sounding like speech, with small-but-existing impact to the speech intelligibility. The wideband original signal clearly has a higher resolution. One may also observe that aural difference and higher resolution on the SBE signal; however, the voice sounds a little saturated. The interpolated signal is almost identical aurally with the narrowband signal, and that is, as shown in Section 3.3, due to the lack of information on the higher band.

## 3.2   Metrics

The signals are also compared in 3 different metrics against the ground truth (i.e. 16khz original sampled signal); L1 Norm of Magnitude Spectrum Difference (MSD), Spectral Convergence (SC), and L1-Norm of Mel Spectrogram difference. The results are presented in Table 3.1.

Each metric is presented twice; once comparing the interpolated signal (*Interpolated*) and the second time the bandwidth expanded (*SBE*) signal to the Wideband original signal.

Magnitude Spectrum Difference quantifies the maximum frame magnitude spectrum absolute difference between the 2 signals (lower is better). The mathematical formulation is given in Equation 3.2, where STFT is the discrete Short Time Fourier Transform of the ground truth speech signal, and $\|A\|_1$ denotes the 1st norm of a matrix A:

$$\|A\|_1 = \max_{0 \le j \le m} \sum_{i=1}^{n} |\alpha_{i,j}| \tag{3.1}$$

The results on the table shows that on every case, the MSD of the simply interpolated signal is higher than the bandwidth expanded (SBE) one, and thus on average, the SBE signal is closer to the ground truth in the frequency domain.

$$MSD = \||STFT(s[n])| - |STFT(\hat{s}[n])|\|_1 \tag{3.2}$$

Spectral Convergence emphasizes spectral peaks and other spectral components when comparing 2 signals. For this purpose, it became popular as a Loss Function when training Neural Networks on speech signals. The

SC function is defined in Equation 3.4 where $\|A\|_F$ denotes the *Frobenius* norm of a matrix A:

$$\|A\|_F = \sqrt{\sum_{i=0}^{N-1} \sum_{j=0}^{K-1} |\alpha_{i,j}|^2} \tag{3.3}$$

. In this case, we can see that exactly for half samples the SBE performs better than the interpolated signal.

$$L_{sc}(s[n], \hat{s}[n]) = \frac{\||STFT(s[n])| - |STFT(\hat{s}[n])|\|_F}{\||STFT(s[n])|\|_F} \tag{3.4}$$

Humans perceive sound in a logarithmic scale rather than a linear scale. The Mel Scale was developed to take this into account by conducting experiments with a large number of listeners. It is a scale of pitches, such that each unit is judged by listeners to be equal in pitch distance from the next.

A Mel Spectrogram makes two important changes relative to a regular Spectrogram that plots Frequency vs Time; It uses the Mel Scale instead of Frequency on the y-axis and it uses the Decibel Scale instead of Amplitude to indicate colors. Therefore, L1-Norm of the difference of 2 Mel Spectrogram quantifies the maximum distance on the Mel Scale between 2 frames' spectrograms of the respective 2 signals. Also in this case, we can see that on exactly half of the speech samples, the SBE signal has lower score than the interpolated one. Its also worth noting here that although SBE had better performance on the MSD metric, it appears that it is not the case in the Mel Spectrogram distance. This is probably for 2 reasons: first, the SBE voice has that saturation effect on it, with major impact on fricatives, such as 's' and 'sh'. Secondly, the mel scale is linear on the lower band and logarithmic on the higher band. Thus, any difference in the lower band will have larger impact than any difference in the higher band, and our system's goal is to add high-band information to our interpolated signal.

| Filename | MSD Interpolated | MSD SBE | SC Interpolated | SC SBE | Mel Diff Interpolated | Mel Diff SBE |
|---|---|---|---|---|---|---|
| arctic_bdl1_snd_norm | 0.105 | **0.095** | **0.090** | 0.094 | **2.257** | 2.467 |
| arctic_slt1_snd_norm | 0.042 | **0.036** | 0.042 | **0.041** | **0.182** | 1.186 |
| fjlg0_si1889 | 0.107 | **0.100** | 0.221 | **0.218** | 1.326 | **1.317** |
| mdns0_si873 | 0.134 | **0.118** | **0.171** | 0.175 | 3.343 | **3.291** |
| mdns0_sa1 | 0.165 | **0.150** | 0.211 | **0.208** | 8.441 | **8.344** |
| falk0_sa1 | 0.090 | **0.087** | **0.159** | 0.162 | **0.748** | 1.245 |

Table 3.1: Metrics Table over all Speech Signals

The 6 samples include 6 separate speakers, 3 male and 3 female ones. The 2nd, 3rd and last row of the table correspond to the female speakers, while the other rows correspond to male speakers. Also, the first 2 and last 2 rows correspond to the same phrase spoken by 2 different speakers (1 male 1 female in each case). One may observe that on average, the metrics show that the system approaches better the wideband speech for female voices than the male ones. Further testing is needed to confirm and see the limits of the system.

## 3.3 Spectrograms

The spectrograms of the 4 different signal types are presented in Figure 3.3. As we can see, the narrowband signal has no spectral information above 4khz (i.e. $0.5 \cdot F_s$, where $F_s = 8$khz) whatsoever, while the wideband signal goes up to 8khz. Respectively, while the interpolated signal has a sampling rate $F_s = 16$khz, we see no significant information on the highband zone $(4-8\text{khz})$ while the signal has mostly information on the lowband zone $(0-4\text{khz})$. On the BE signal, the highband zone has much more spectral information.

The spectrograms produced are a result of a narrowband analysis, as a long window was used (60ms). Narrowband spectrogram gives good frequency resolution as the harmonics are effectively resolved (horizontal striations on the spectrogram). However, it also gives poor time resolution, because the long analysis window covers several pitch periods and thus is unable to reveal fine periodicity changes over time. Since this work focuses on the frequency resolution of a signal when applying the artifical bandwidth expansion method, narrowband analysis was preferred. It should be noted that colors in spectrograms have a meaning; intense yellow color corresponds to high magnitude values (high energy), whereas green or blue color correspond to lower magnitude values (low energy).

To further expand our understanding of the system internally, we can look on the magnitude spectrums of the signal frames. In Figure 3.4, we can see how the highband signal is synthesized, and in Figure 3.5 we can visually compare the synthesized signal with the original wideband.
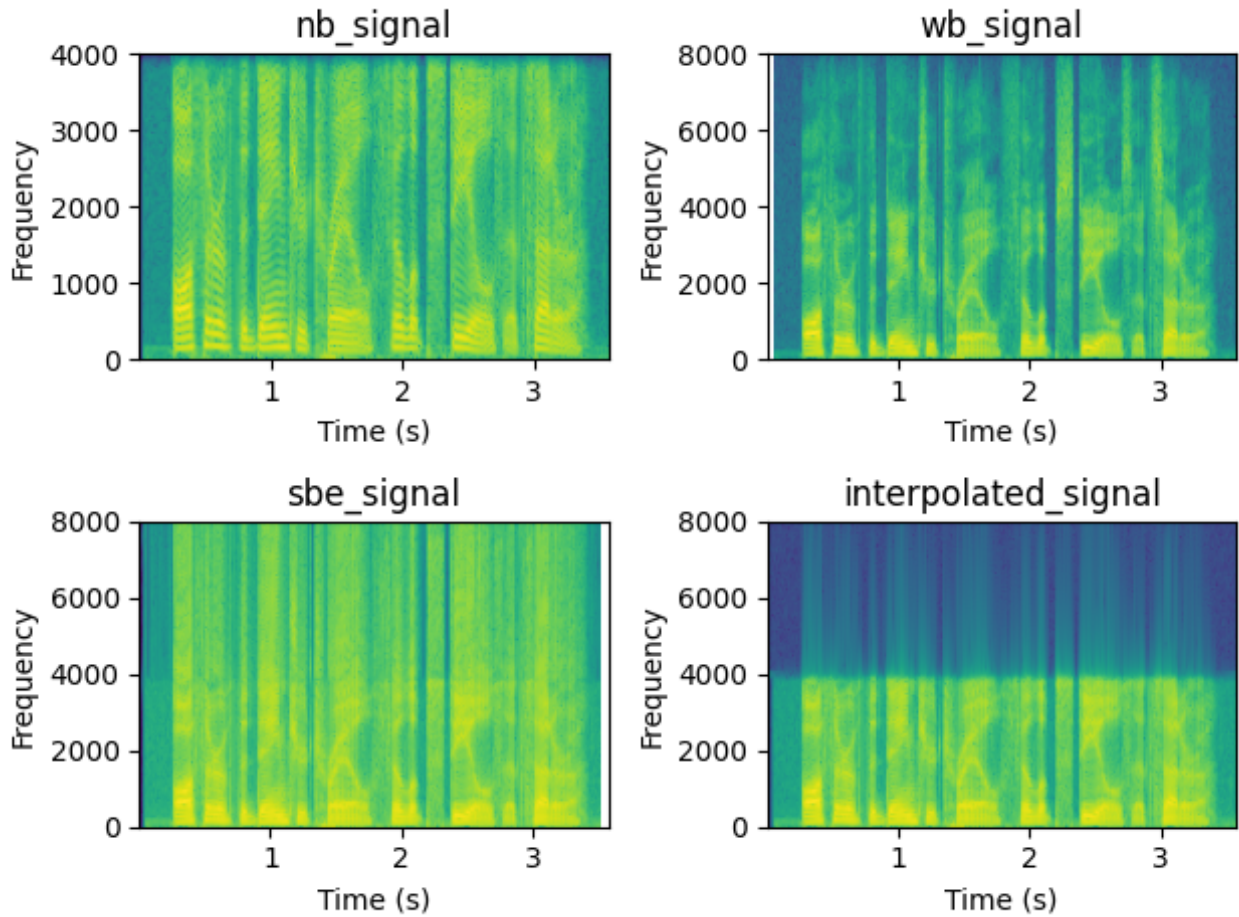
Figure 3.3: Spectrograms for each case: 'nb_signal' corresponds to the downsampled signal spectrogram (8khz), 'wb_signal' for the original signal (16khz), 'sbe_signal' for the reconstructed signal with BE, 'interpolated_signal' for the reconstructed signal with interpolation.

We can see that while on the lower band (0 - 3800 Hz) the signals match perfectly, on the higher band there is a significant difference. To see how well that generalizes, we can compute the mean absolute difference between them for every frame and gather the results in a histogram. However, for better clarity, we can do that in 3 separate histograms: 1 for the lower band (0 - 3kHz), another for the middle band (3-5kHz), where also the high-pass filter cutoff frequency resides, and finally a last one for the higher band (5-8kHz). The result is presented in 3.6.

Figure 3.4: FFT magnitudes for a random signal frame. From top to bottom; The output of LPC Synthesis (**170**), the highband signal after applying high-pass butterworth filter to LPC output, the interpolated signal and finally the reconstructed wideband signal as the sum of the interpolated and highband signals

Figure 3.5: The log Magnitudes of the BE signal (patent) and the ground truth (original wideband) on the same plot.

Figure 3.6: Histograms representing the mean absolute difference between the spectrograms in the log space of the original wideband signal and the SBE signal.

# Chapter 4

# Conclusion and Future Work

This invention focuses on a novel bandwidth extension approach in the category of parametric methods that do not require training. Almost all parametric techniques use an LPC synthesis filter for wideband signal generation (typically an intermediate wideband signal which is further highpass filtered), by exciting it with an appropriate wideband excitation signal.

In Section 3.2, we compared the resulting signals using various metrics. In Section 3.3, we compared the signals visually through their spectrograms. The metrics showed that the SBE signal is closer to the wideband signal than the interpolated one, but also, the SBE signal - regardless of the voice saturation - appears to have the wideband resolution, as well as better frequency information on the highband (4-8khz), compared to the interpolated signal, resulting in a signal with better aural perception.

As a future work, the most important part is to find out where that voice saturation comes from. This may be a result either of redesigning the system's components (i.e. better implementation of specific algorithms) or changing the system's configuration and parameters (i.e. switch to wideband analysis, change cutoff frequency, LPC order).

Also, we only applied a bandwidth expansion of a factor 2 (8khz to 16khz). The system should further be tested for other scales (e.g. 16khz to 32khz, 16khz to 44khz, 8khz to 32khz etc) and compare the results on the same basis.

Finally, we showed that the system performs better for the female voices rather than the male ones. Aurally, there appears to be an impact on the fricatives in a spoken sentence by the downsampling, that the system couldn't restore. Further experimentation will be needed to confirm those hypotheses and understand why that is the case.

Reported bandwidth extension methods can be classified into two types - parametric and non-parametric. Non-parametric methods usually convert directly the received narrowband speech signal into a wideband signal using simple techniques like spectral folding and non-linear processing. This invention aims to find a relationship between the narrowband and wideband speech parameters through a parametric method that doesn't require training. Other approaches may even be tried for that goal, such as neural-net-based methods and statistical methods.The inventors mention that the main advantages of a non-parametric approach are its relatively low complexity and its robustness, stemming from the fact that no model needs to be defined and, consequently, no parameters need to be extracted and no training is needed. These characteristics, however, typically result in lower quality when compared with parametric methods.

# Chapter 5

# Appendix

## 5.1 arctic_slt1_snd_norm



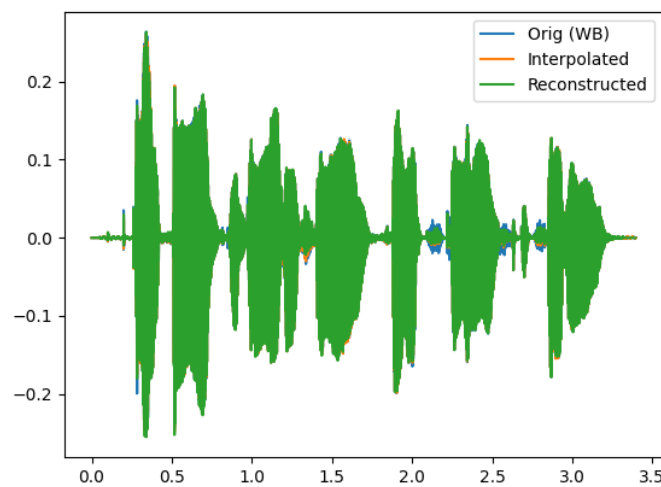Figure 5.1: arctic_slt1_snd_norm Signal Waveforms



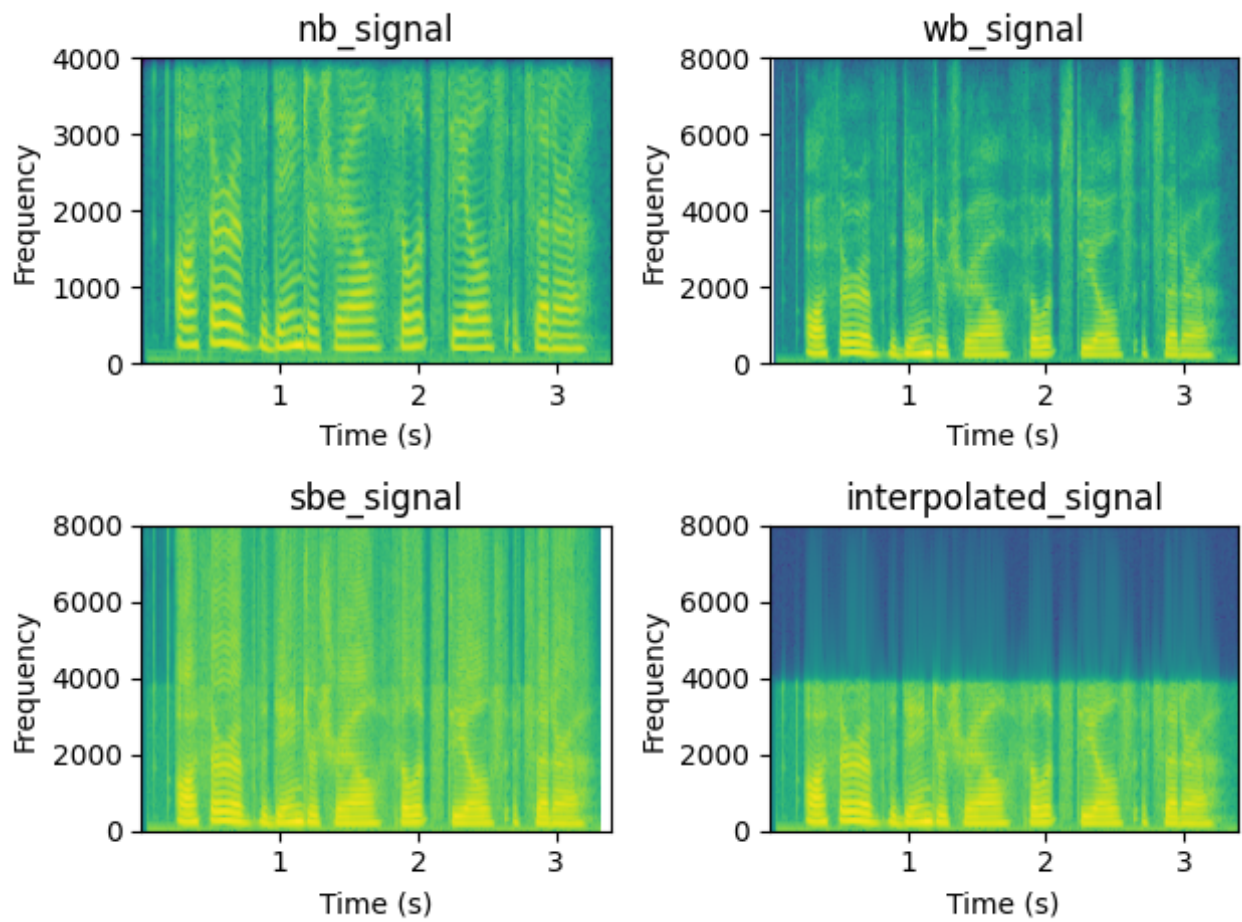Figure 5.2: arctic_slt1_snd_norm Signal Waveforms Merged

Figure 5.3: arctic_slt1_snd_norm Narrowband vs Wideband vs Interpolated vs SBE
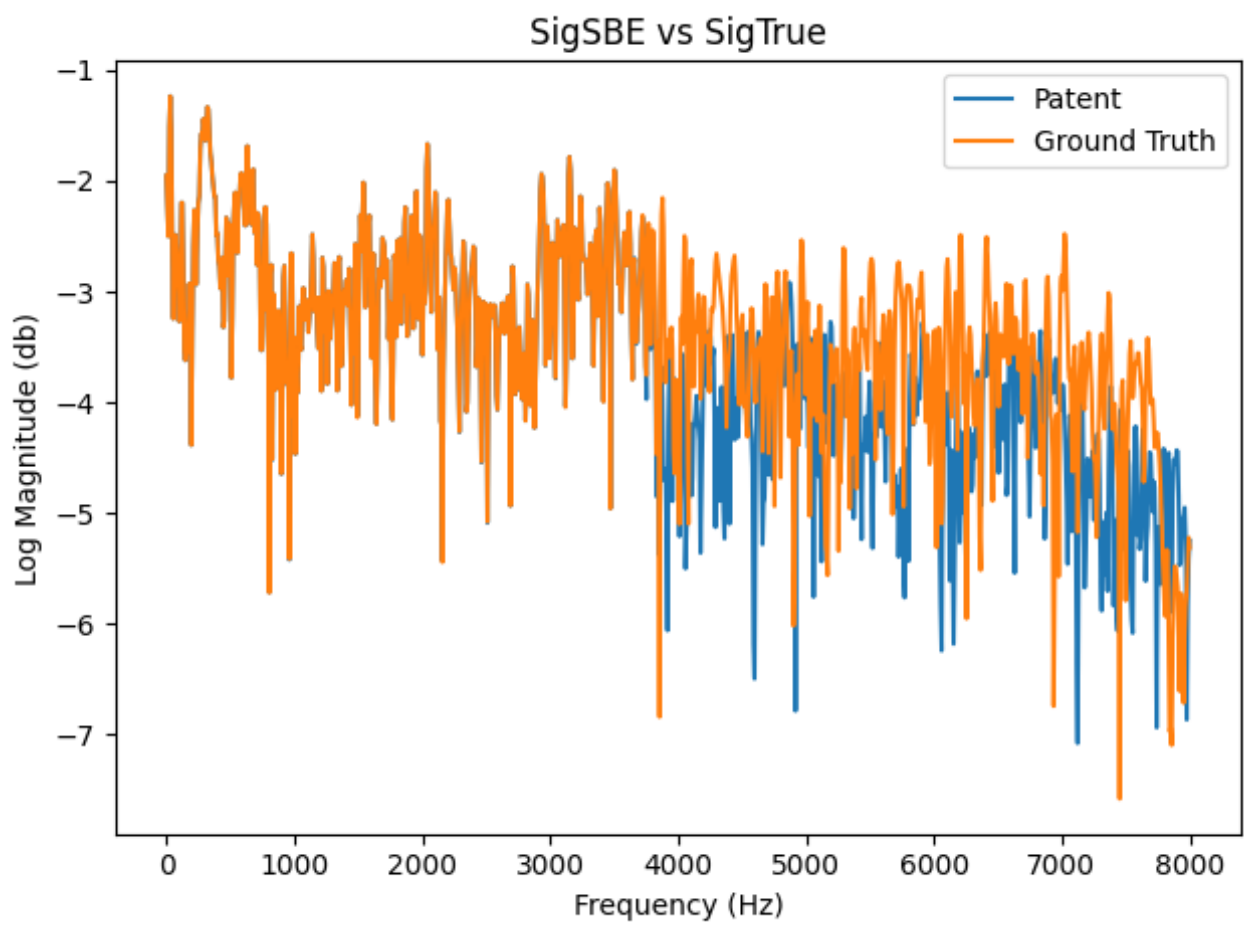
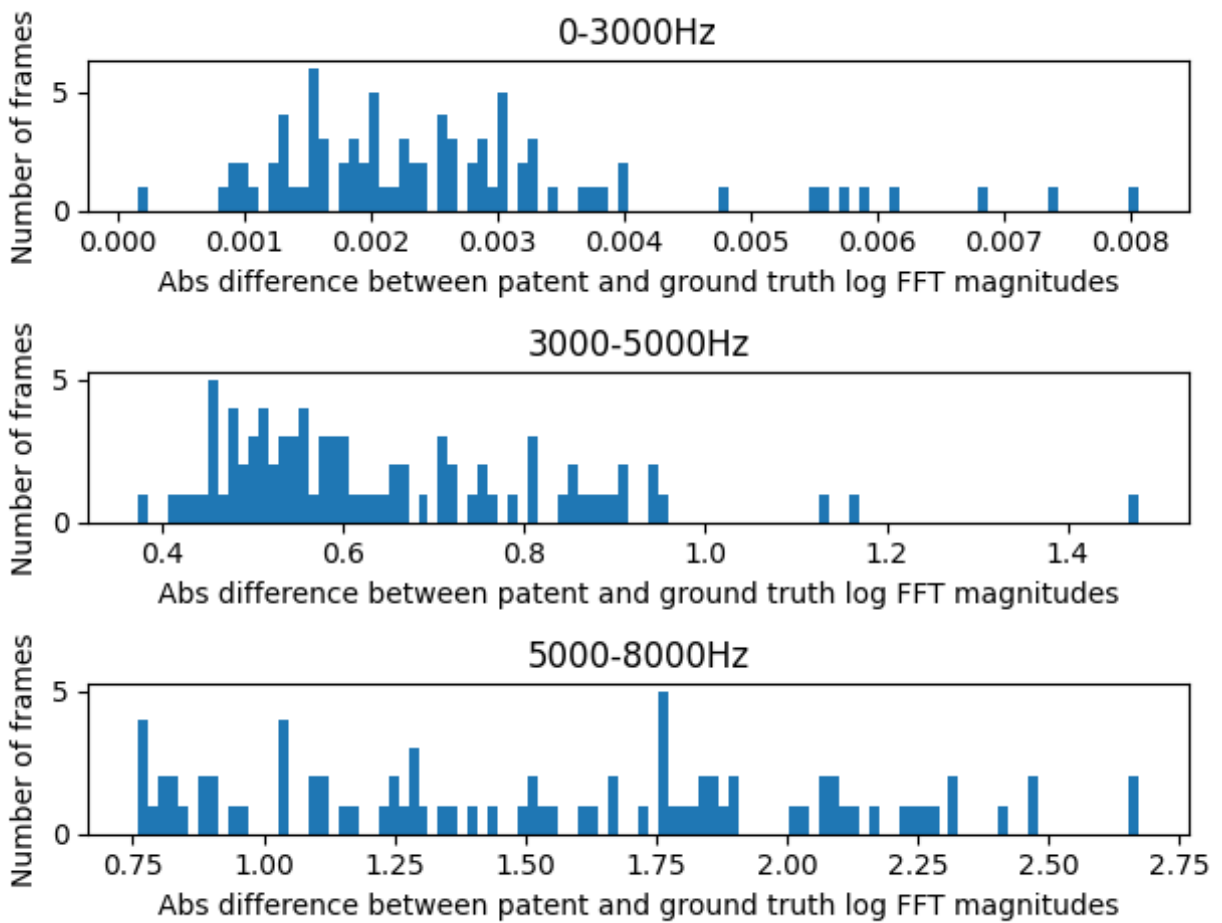Figure 5.4: arctic_slt1_snd_norm Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.5: arctic_slt1_snd_norm Mean Absolute Difference Histograms per band
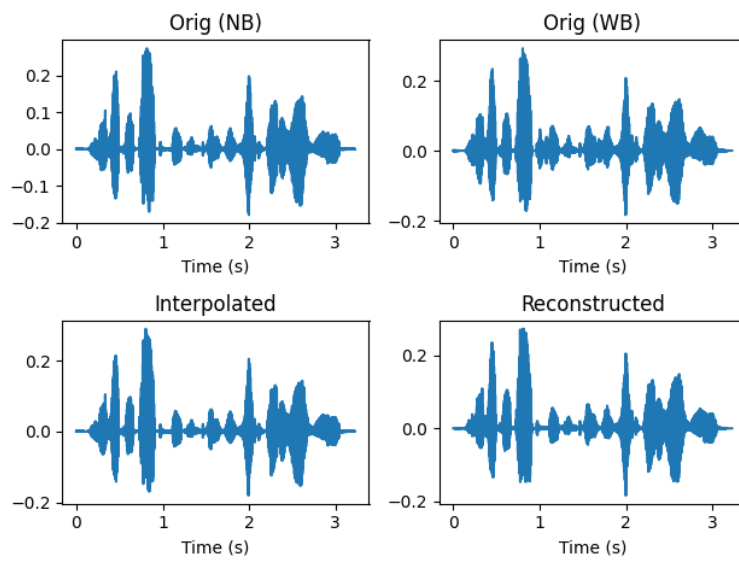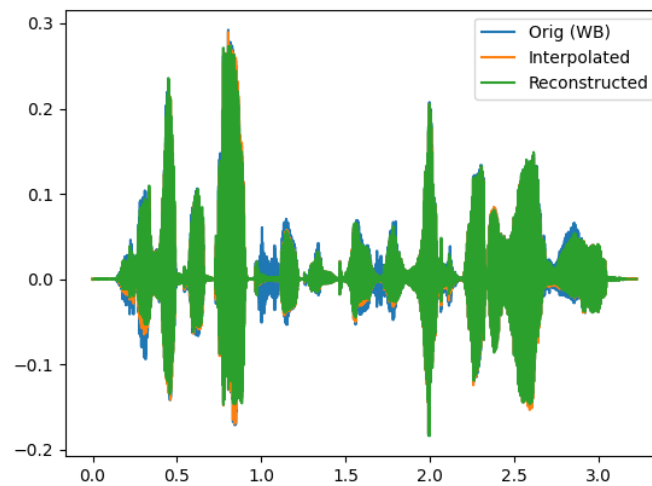
## 5.2    falk0_sa1
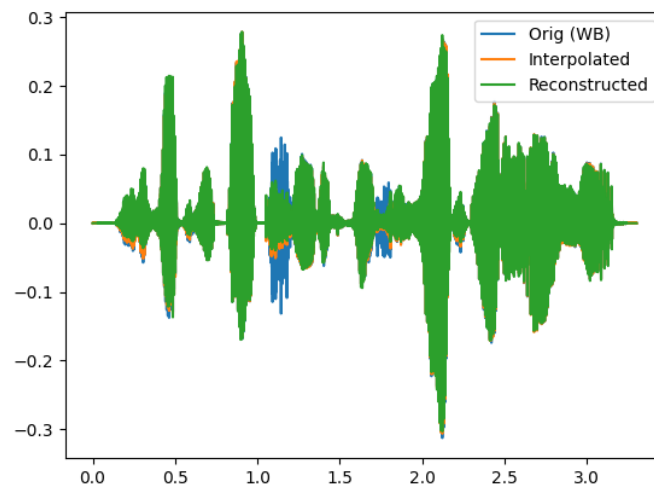


Figure 5.6: falk0_sa1 Signal Waveforms
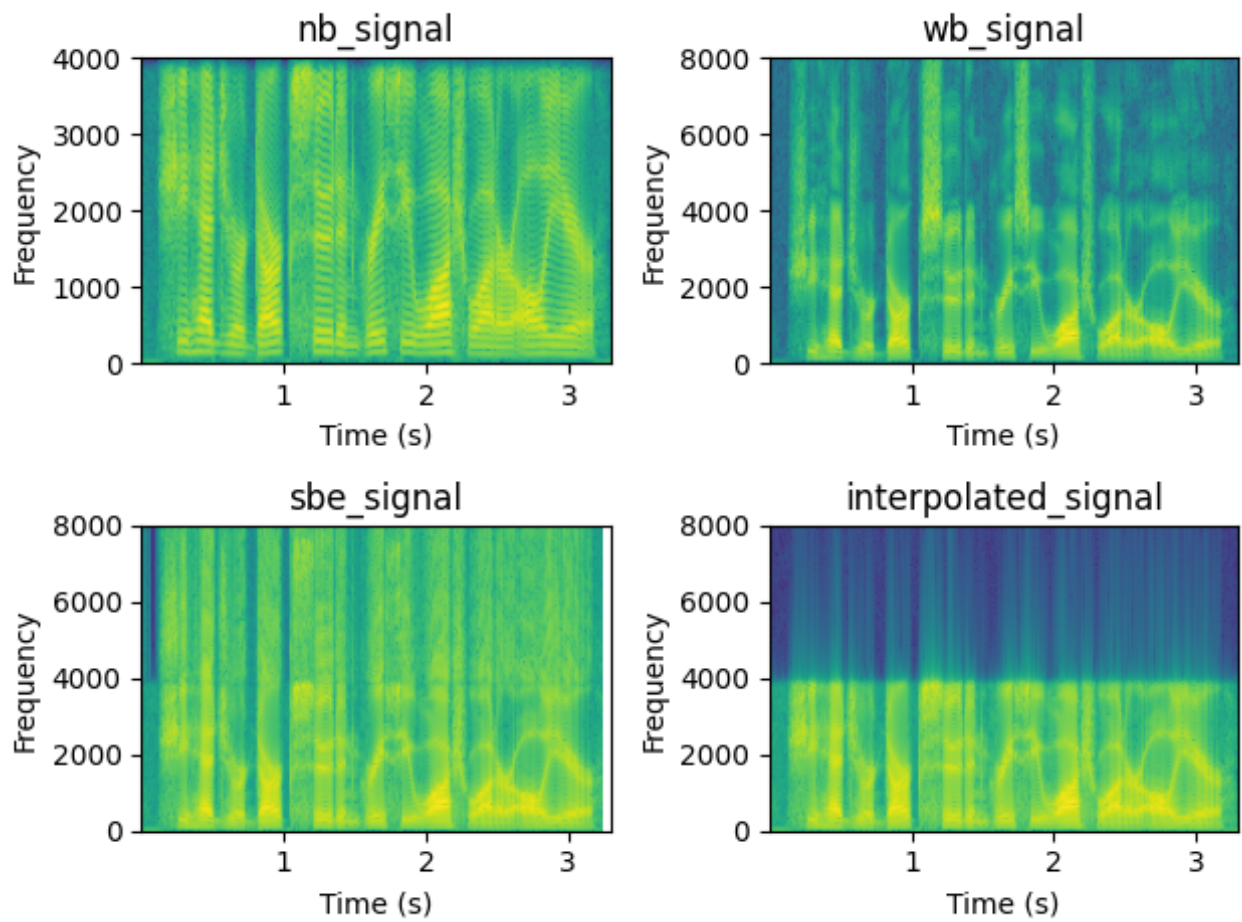


Figure 5.7: falk0_sa1 Signal Waveforms Merged

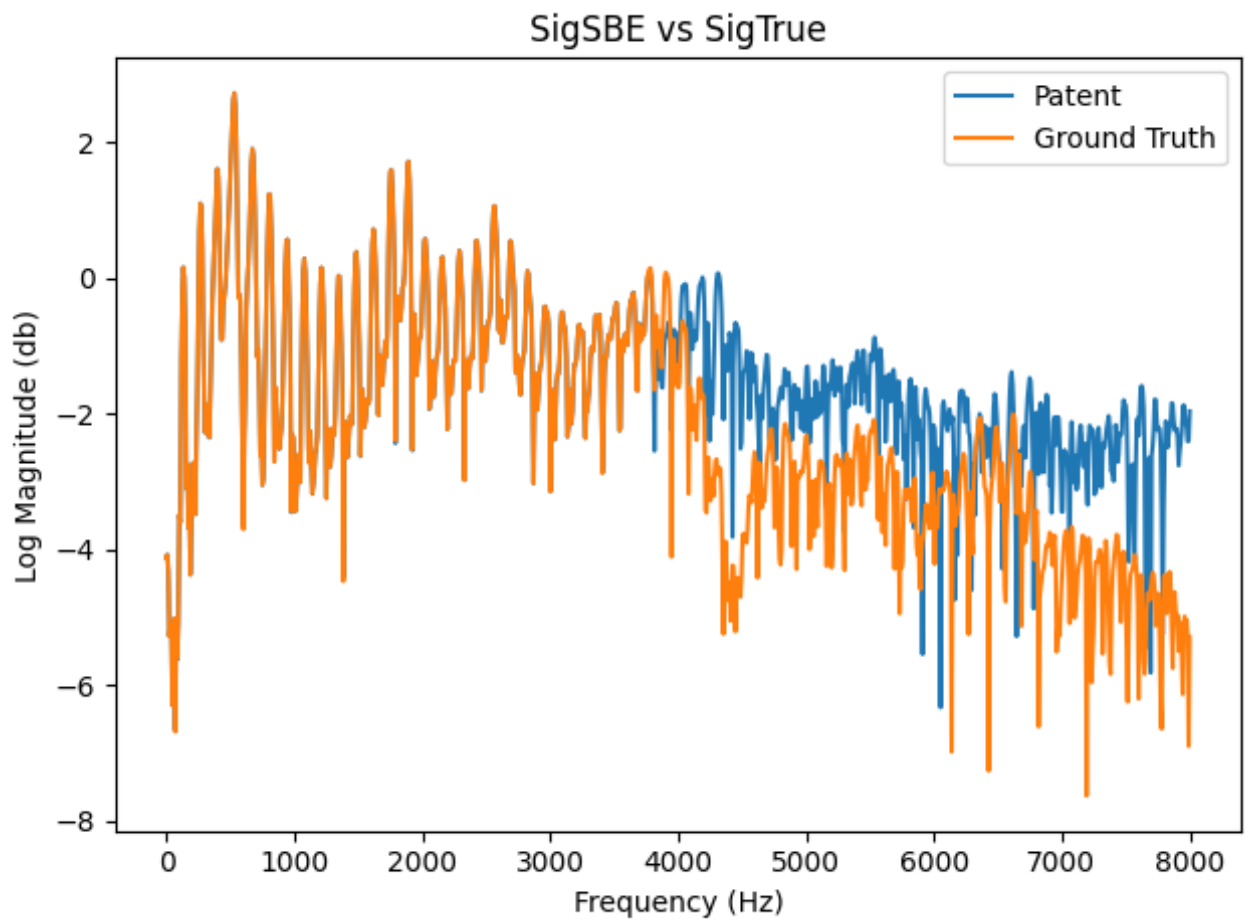Figure 5.8: falk0_sa1 Narrowband vs Wideband vs Interpolated vs SBE

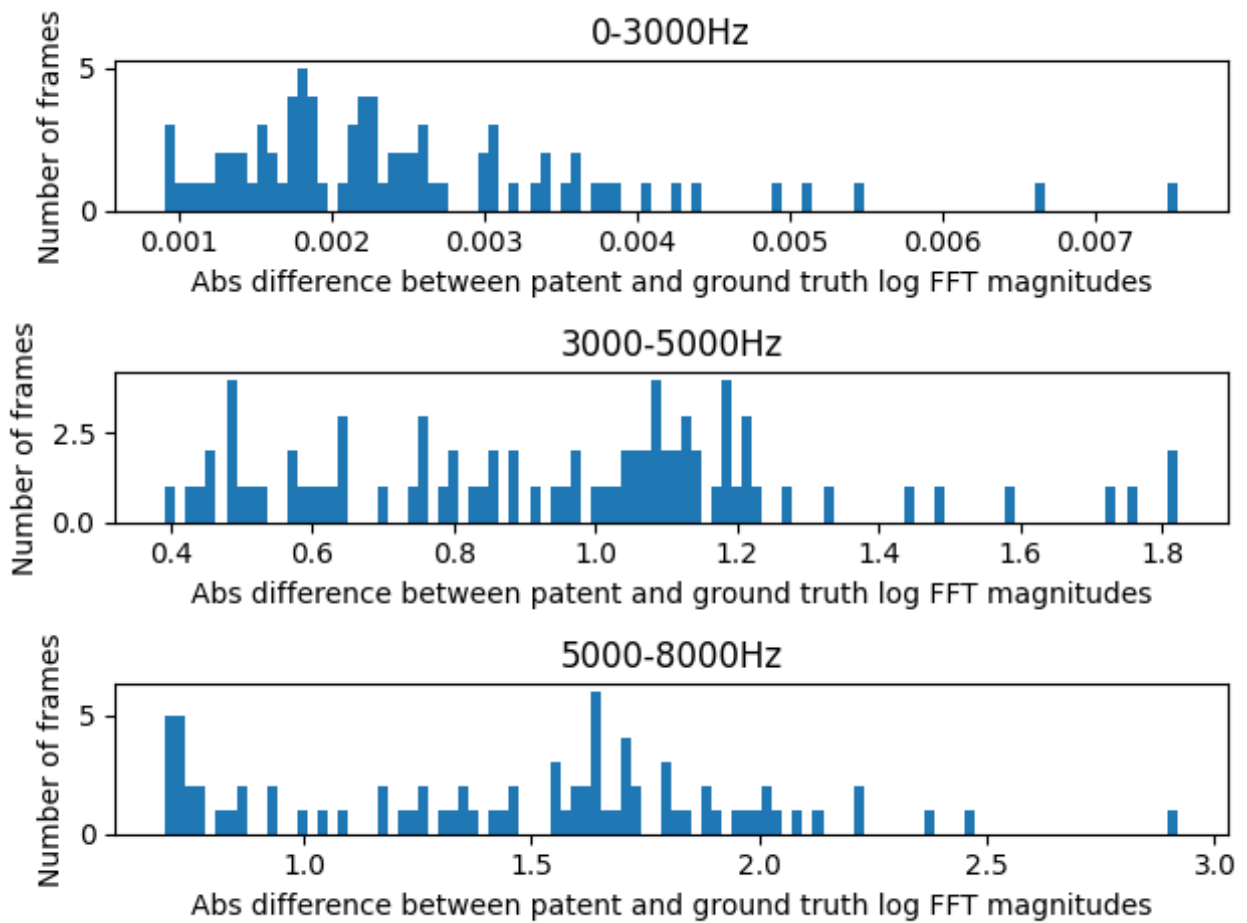Figure 5.9: falk0_sa1 Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.10: falk0_sa1 Mean Absolute Difference Histograms per band

## 5.3   fjlg0_si1889



Figure 5.11: fjlg0_si1889 Signal Waveforms



Figure 5.12: fjlg0_si1889 Signal Waveforms Merged

Figure 5.13: fjlg0_si1889 Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.14: fjlg0_si1889 Narrowband vs Wideband vs Interpolated vs SBE

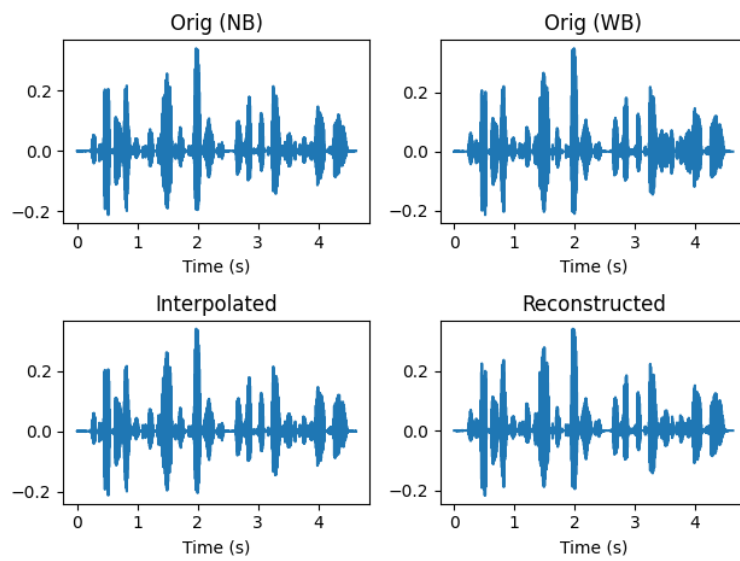Figure 5.15: fjlg0_si1889 Mean Absolute Difference Histograms per band

## 5.4   mdns0_sa1



Figure 5.16: mdns0_sa1 Signal Waveforms



Figure 5.17: mdns0_sa1 Signal Waveforms Merged

Figure 5.18: mdns0_sa1 Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.19: mdns0_sa1 Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.20: mdns0_sa1 Mean Absolute Difference Histograms per band

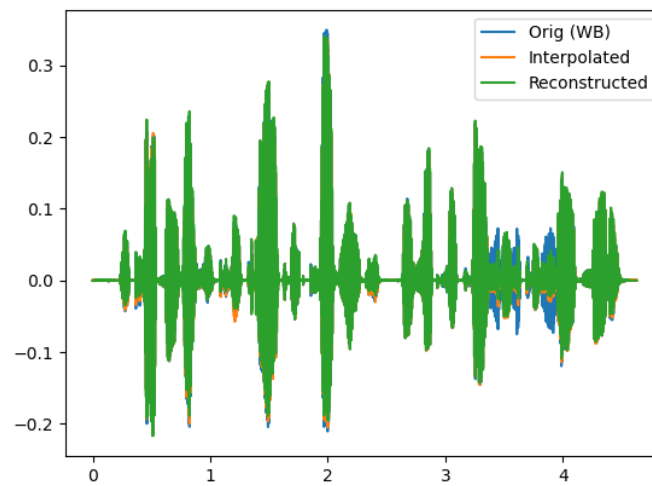# 5.5 mdns0_si873



Figure 5.21: mdns0_si873 Signal Waveforms
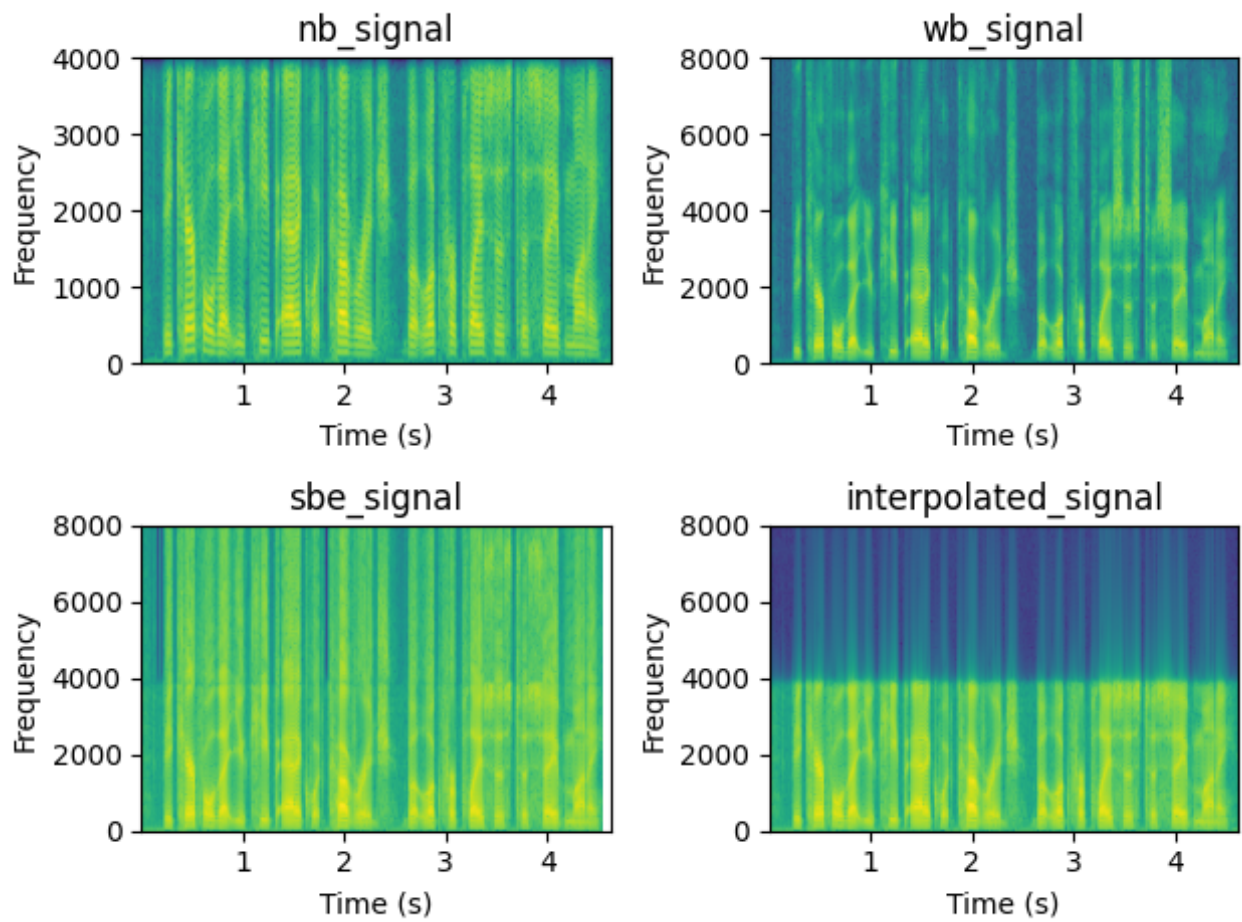


Figure 5.22: mdns0_si873 Signal Waveforms Merged

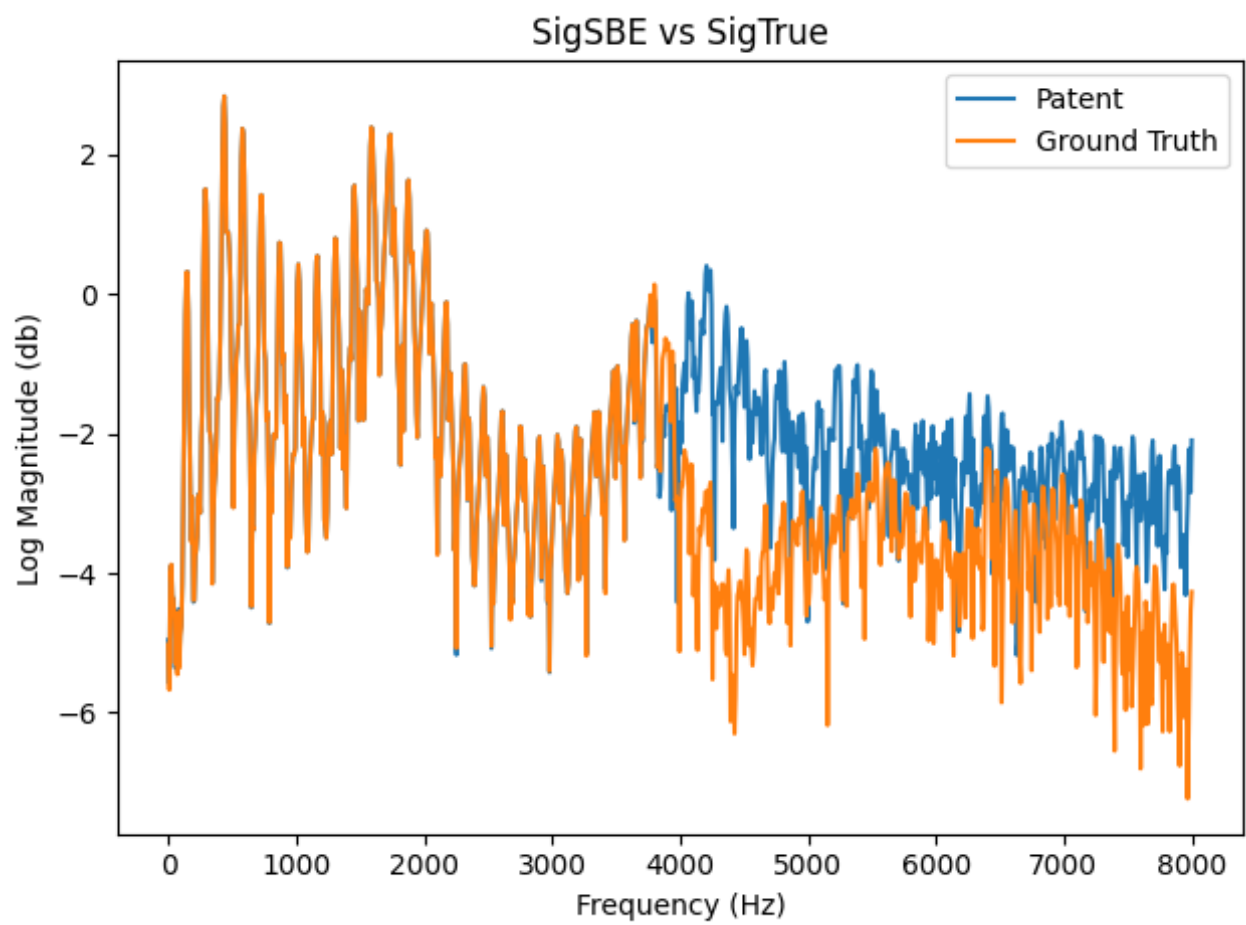Figure 5.23: mdns0_si873 Narrowband vs Wideband vs Interpolated vs SBE

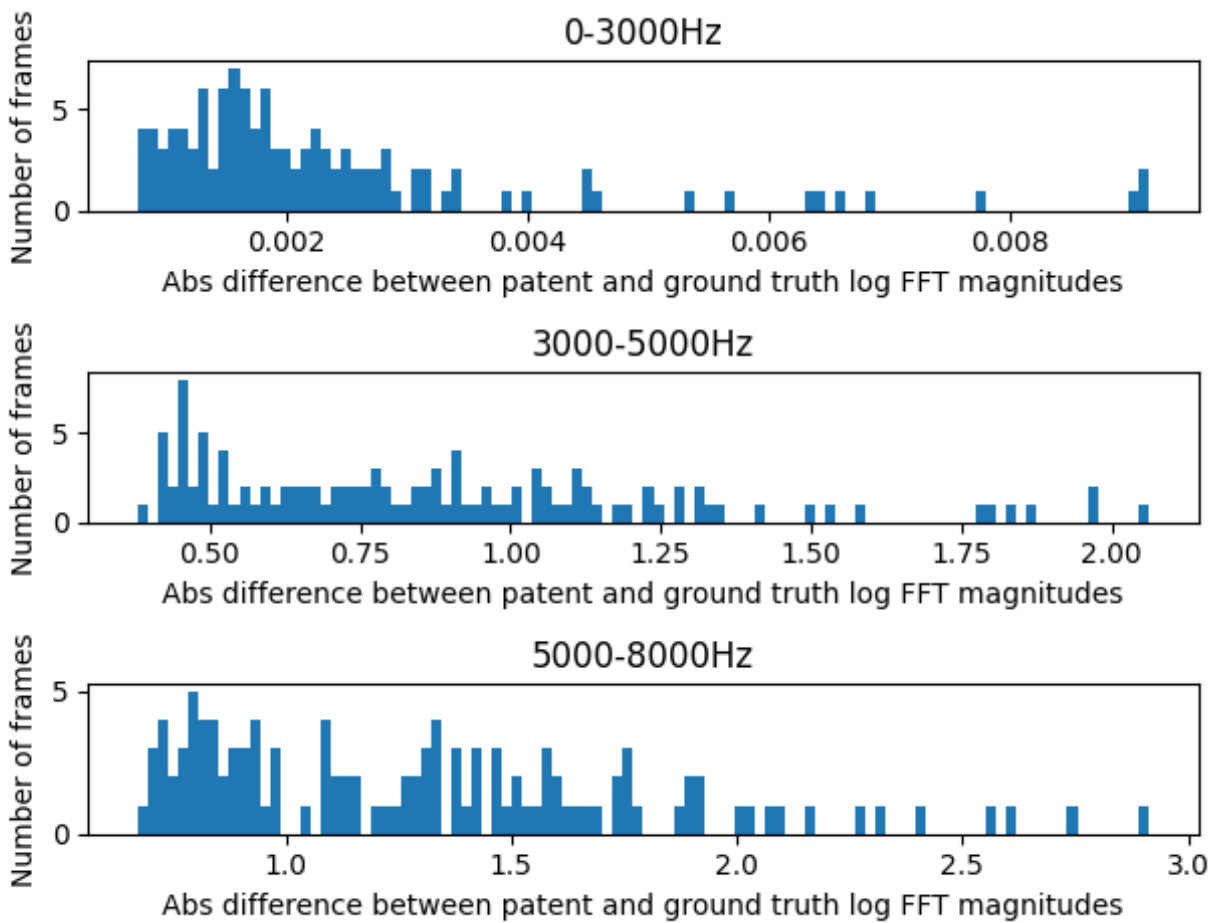Figure 5.24: mdns0_si873 Narrowband vs Wideband vs Interpolated vs SBE

Figure 5.25: mdns0_si873 Mean Absolute Difference Histograms per band