# TWO CHANNEL AUDIO ZOOMING SYSTEM FOR SMARTPHONE

*A. khandelwal, E.B. Goud, Y. Chand, L. Kumar, S. Prasad*

*N. Agarwala, R. Singh*

Electrical Engineering Department
Indian Institute of Technology, Delhi
anant.iitd.2085@gmail.com
{ee172241/39,lkumar,sprasad}@ee.iitd.ac.in

SRI Noida
ritesh.s7@samsung.com
n.agarwala@samsung.com

## ABSTRACT

In this paper, two microphone based systems for audio zooming is proposed for the first time. The audio zooming application allows sound capture and enhancement from the front direction while attenuating interfering sources from all other directions. The complete audio zooming system utilizes beamforming based target extraction. In particular, Minimum Power Distortionless Response (MPDR) beamformer and Griffith Jim Beamformer (GJBF) are explored. This is followed by block thresholding for residual noise and interference suppression, and zooming effect creation. A number of simulation and real life experiments using Samsung smartphone (Samsung Galaxy A5) were conducted. Objective and subjective measures confirm the rich user experience.

## 1. INTRODUCTION

Portable devices for communications like smartphones have become an inseparable part of life. The increasing dependency on the smartphones is due to numerous features they support, ranging from health and convenience to entertainment. One such useful feature being developed is audio zooming [1, 2] where sound from desired direction is enhanced while suppressing interferences from all other directions. This is desirable while trying to listen to a sound in the presence of one or more noise and interfering sources. Practical examples of such environments include that of a railway station, a stadium, classroom and market place. A pictorial depiction of the audio-zooming application for two microphone based smartphone is presented in Figure 1. The evolution of compact device technology and computational power have resulted in use of multiple microphones in a smartphone to exploit the spatial diversity. Many smartphones today utilize two or more microphones. Apple iPhone-5[1] makes use of three microphones for beamforming and noise cancellation. Audio zooming has also been reported recently in some smartphones [2]. However, a significant improvement is required in the presence of severe noise, reverberation and multiple interferences. To the best of our knowledge, the only scientific publication for audio zooming in smartphone is [3] that utilizes MVDR beamforming with a linear array of four microphones.

As most of the current smartphones have two microphones, the possibility of real time audio zooming with smartphone having two microphones is explored in this paper. The complete audio zooming system consists of two blocks as shown in Figure 2. For the beam-

---

[1]https://www.idownloadblog.com/2012/09/12/iphone-5-three-mics/
[2]https://www.youtube.com/watch?v=zzTUAcZ8FRQ, https://www.youtube.com/watch?v=YNh4snIzmq4

forming block, the simplest two channel frequency and time domain beamforming is investigated due to limited degree of freedom available. In particular, Minimum Power Distortionless Response (MPDR) beamformer [4] and Griffith Jim Beamformer (GJBF) [5] are explored in frequency and time domain respectively. The two channel beamformer extracts the target source. This spatial filtering causes some suppression of the interference, but a significant residual interference may be present due to the use of only two microphones. A novel block-thresholding [6] based post-filtering is formulated for creating the audio zooming effect. The additional novelty of the work is in exploration of two channel based audio zooming system deployable on smartphone. The filter length in proposed time domain GJBF can be estimated dynamically based on different environmental condition.
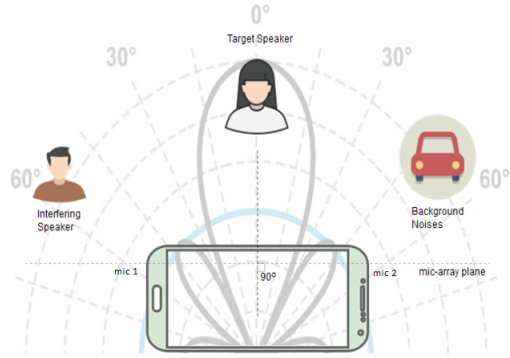


**Fig. 1**: Two microphone smartphone based audio zooming: schematic depiction

## 2. THE PROPOSED AUDIO ZOOMING SYSTEM

We consider a smartphone with two identical and omni-directional microphones located at the top and the bottom. The target source to be acoustically zoomed in, is made incident on the smartphone normal to the plane containing the microphones, as shown in Figure 1. The sources incident from other directions are assumed to be interferences for the audio zooming application. The target and the interference are assumed to be in the far-field region. The proposed audio zooming systems consist of beamforming followed by postprocessing. In particular, two audio zooming systems based on frequency and time domain beamforming, have been proposed and their performance have been analyzed . Motivation for using time domain based Griffith Jim Beamformer (GJBF) comes from the fact that it
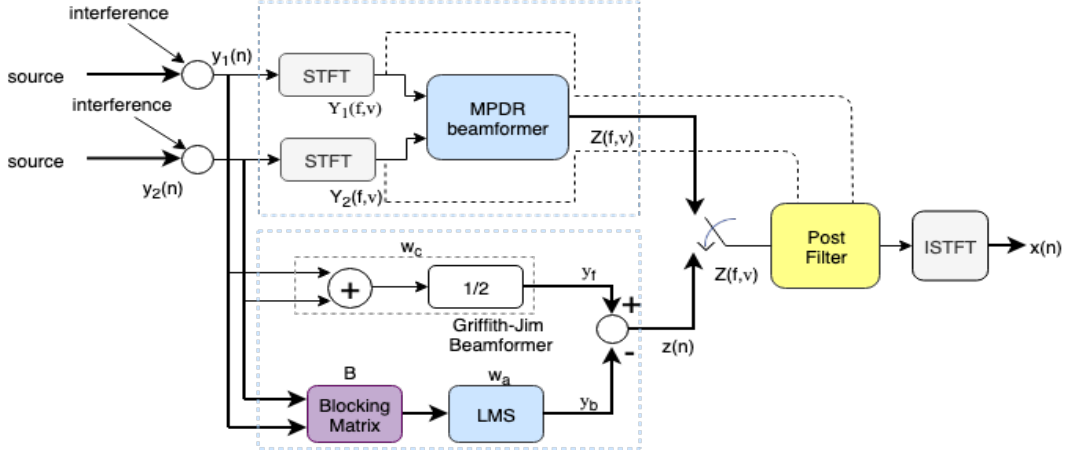
**Fig. 2**: Two stage audio zooming system using modified Griffith-Jim and MPDR beamformer

can control the level of interference at the output without the directional information of interferences [5]. In the ensuing Section, two channel based MPDR beamformer is presented, followed by modified time domain GJBF.

### 2.1. Two Channel MPDR Beamformer

A general wideband array data model can be written in STFT domain as

$$\mathbf{Y}(f,\nu) = \mathbf{D}(\Psi,\nu)\mathbf{S}(f,\nu) + \mathbf{N}(f,\nu) \tag{1}$$

where $\mathbf{Y}(f,\nu)$ is the received array signal, and $\mathbf{N}(f,\nu)$ is zero mean, uncorrelated sensor noise. For $M$ microphones and $L$ sources, $\mathbf{D}(\Psi,f)$ is $M \times L$ array manifold given by

$$\mathbf{D}(\Psi,f) = [\mathbf{d}(\Psi_1,f), \mathbf{d}(\Psi_2,f), \dots, \mathbf{d}(\Psi_L,f)] \tag{2}$$

where $\mathbf{d}(\Psi_l,f)$ represents the steering vector for $l^{th}$ source given by

$$\mathbf{d}(\Psi_l,f) = \begin{bmatrix} e^{-j2\pi f\tau_{1l}} & e^{-j2\pi f\tau_{2l}} & \dots & e^{-j2\pi f\tau_{Ml}} \end{bmatrix}^T \tag{3}$$

$\Psi_l = (\theta_l, \phi_l)$ is the incident direction of the $l^{th}$ source. Here $\tau_{ml}$ represents the time delay of arrival of the $l^{th}$ signal at the $m^{th}$ microphone with respect to a reference microphone. MPDR beamforming problem is equivalent to minimizing the output power with a distortionless response in the target direction given as

$$\min_{W} \mathbf{W}^H \mathbf{R_Y} \mathbf{W} \quad subject\ to \quad \mathbf{W}^H \mathbf{d}(\Psi_l,f) = 1 \tag{4}$$

The solution to the above optimization problem under diagonal loading is given by

$$\hat{\mathbf{W}}_{(\Psi_l,f)} = \frac{(\hat{\mathbf{R}}_{Y,f} + \alpha\mathbf{I})^{-1}\mathbf{d}(\Psi_d,f)}{\mathbf{d}^H(\Psi_d,f)(\hat{\mathbf{R}}_{Y,f} + \alpha\mathbf{I})^{-1}\mathbf{d}(\Psi_d,f)} \tag{5}$$

where $\alpha$ is the diagonal loading factor. The co-variance matrix is estimated as

$$R_{Y,f} = \frac{1}{K} \sum_{k=0}^{k=K} Y(f,k) * Y^*(f,k) \tag{6}$$

where K is total number of time frames.

### 2.2. Modified Griffiths-Jim Adaptive Beamformer

Audio zooming application is additionally, explored using two channel time domain beamformer. The target is assumed to be incident from broadside, resulting in identical delays at the two microphones. The $n^{th}$ snapshot of the received signal at the $m^{th}$ microphone is written as

$$y_m(n) = s(n) + n_m(n),\ n = \{0, 1, \cdots, N_s-1\}, m = \{1, 2\} \tag{7}$$

where $s(n)$ is the target signal to be zoomed in, and $n_m(n)$ is the total noise and interferences.

As the target signal undergoes identical delays at the two microphones, no phase adjustment is required herein when compared to the original GJBF [5]. The constrained weight $w_c = \frac{1}{2}[1,1]^T$ for the target will result in simple addition of the two channel signal providing signal plus interference $y_f$ in the upper branch of GJBF in Figure 2. The blocking matrix $B = [1,-1]^T$ subtracts the two channel data thus it does not allow the signal from the constraint direction in the lower branch of GJBF. The adaptive weight $w_a$ in the lower branch is chosen to estimate the signal at the output of $w_c$ as a linear combination of the data at the output of the blocking matrix $B$. As blocking matrix does not allow the signal from the constraint direction, the signal estimated by $w_a$ is the interference close to interference present at the output of $w_c$.
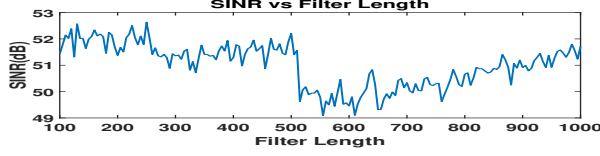
The overall output of the two microphone GJBF is

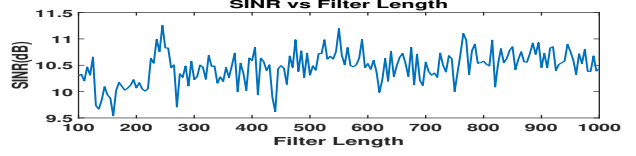$$z(n) = y_f(n) - y_b(n), \tag{8}$$

where $y_f$ has the target signal, noise and interference with response determined by $w_c^H$, and $y_b$ has only the noise and interference. The filter weights $w_a$ is found based on minimizing the power contained in $z(n)$. The power will be minimum when $y_b(n)$ closely models the interference and noise present in the $y_f(n)$. In particular, Frequency Domain Adaptive Filter (FDAF) approach [7] using overlap-save method has been utilized for computing $w_a$. The FDAF approach has logarithmic complexity when compared to polynomial complexity for the classical LMS as utilized in the original GJBF.

### 2.2.1. LMS Filter Length

The GJ beamformed output quality depends on the length of LMS filter. The filter length depends on the environmental conditions.

(a) Anechoic recording.



(b) A Meeting room recording.

**Fig. 3**: Dynamic LMS Filter Length Estimation (a) The SINR is max at filter length of 250 (b) The SINR is maximum for filter length 245.

The length of the filter can be dynamically decided based on the highest SINR in GJBF output. The SINR is computed in Short Time Fourier Transform (STFT) domain as $\frac{|Z(f,\nu)|^2 - \sigma^2(f,\nu)}{\sigma^2(f,\nu)}$, where the estimation of residual noise variance is detailed in Section 2.3. SINR at GJBF output is plotted in Figure 3 with the LMS filter length for anechoic chamber and reverberant room recording. It is to note that the optimum filter length for anechoic set-up is 250 while that of reverberant room is 245.

### 2.3. Adaptive Block Thresholding based Post-Filtering for Audio Zooming

In this Section, adaptive block thresholding based post-filtering is proposed for Audio Zooming effect creation. The estimate obtained by adaptive beamforming recovers the target with some amount of noise and interference. Additionally, the beamformed output might show transient or tonal behavior, or combination of the two. The residual interference along with transient and tonal behavior is now simply treated as single additive interference term for further processing. Mathematically, the beamformed output can now be written in STFT domain as

$$Z(f,\nu) = S(f,\nu) + I(f,\nu) \quad (9)$$

where $S(f,\nu)$ is STFT coefficient of desired signal and $I(f,\nu)$ is STFT coefficient of residual noise and interference. The total variance of such additive interference can be computed as

$$\hat{\sigma}^2(f,\nu) = \frac{1}{2}\{(Y_1(f,\nu) - Z(f,\nu))^2 + (Y_2(f,\nu) - Z(f,\nu))^2\} \quad (10)$$

where $Y_1(f,\nu)$, $Y_2(f,\nu)$ are the STFT coefficients of microphone 1 and 2 respectively. This is the best residual interference variance that can be estimated considering the beamformed output is close to the target signal. This estimate works well with practical scenarios. More accurate estimation can provide better result. It is to be noted that because of limited degrees of freedom (two, utilized both) in the spatial domain, post-filtering is now attempted in frequency domain. For this purpose, the output of the either beamformer is considered in short-time Fourier transform (STFT) domain, with suitably chosen block sizes in the time and frequency, as detailed below.

The time-frequency coefficients are modified by multiplying each of them by an attenuation factor to attenuate the interference dominated components as

$$\hat{S}(f,\nu) = a(f,\nu) * Z(f,\nu) \quad (11)$$

and creating the zooming effect. The attenuation factor $a(f,\nu)$ depends upon the values $Z(f',\nu')$ for all $[f',\nu']$ in the neighborhood of $[f,\nu]$. The signal estimate $S(\hat{f},\nu)$ is computed from the noisy data $Z(f,\nu)$ with a constant attenuation factor $a_i$ over the sub-block $B_i$ as

$$\hat{S}(f,\nu) = a_i \, Z(f,\nu) \, \forall (f,\nu) \in B_i \quad (12)$$
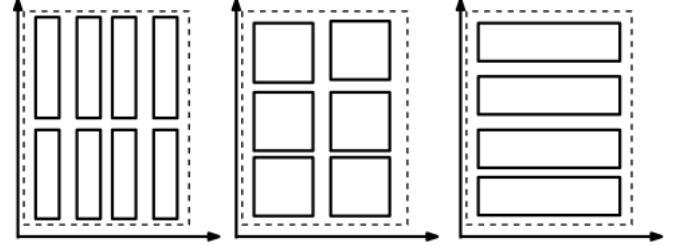


**Fig. 4**: Example of dividing a macro-block into sub-blocks

Selection of the block is done as follows. The entire STFT matrix $\mathbf{Z}$ is divided into macro-blocks of size $P \times Q$. Each macro-block is further divided into sub-blocks of size $2^{H-v} \times 2^v$ where $2^H = L$ and $2^V = W$ and $v \in \{0, 1.....H\}$. Various ways of dividing a macro-block into sub-blocks is shown in Figure 4. A Threshold block SNR is chosen to differentiate higher and lower block SNR values. Out of the various structure of sub-blocks, the one having highest SNR is chosen.

The mean square error can be written as

$$r = E(|\hat{S} - S|^2) \leq \frac{1}{A} \sum_{i=1}^{I} \sum_{f,\nu \in B_i} E\{(a_i \, Z(f,\nu) - S(f,\nu))^2\} \quad (13)$$

The error can be minimized by choosing [8]

$$a_i = \left(1 - \frac{1}{\zeta_i + 1}\right)_+^3 \quad (14)$$

where $\zeta_i = \frac{\bar{S}^2}{\bar{\sigma}^2}$ is the average apriori SNR in the sub-block $B_i$, that is computed from

$$\bar{S}^2 = \frac{1}{B_i^0} \sum_{f,\nu \in B_i} S^2(f,\nu) \text{and}, \quad \bar{\sigma}^2 = \frac{1}{B_i^0} \sum_{f,\nu \in B_i} \sigma^2(f,\nu) \quad (15)$$

Here $B_i^0$ is number of coefficients in the sub-block. As $S(f,\nu)$ is unknown, the apriori SNR $\zeta_i$ can be computed alternatively using

$$\hat{\zeta}_i = \frac{\bar{Z}_i^2}{\bar{\sigma}_i^2} - 1 \quad (16)$$

that can be derived from (9). $\bar{Z}_i^2$ can be computed as in (15) using the beamformed signal.

It is to be noted that the selection of $a_i$ ensures that the target signal is enhanced corresponding to the high SNR sub-blocks while suppressing the low-SNRs sub-blocks. This results in the required audio-zooming effect.

---

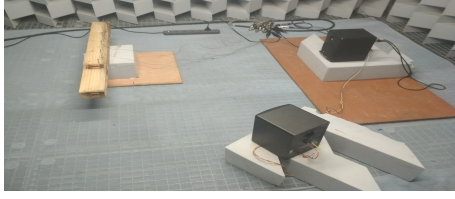[3] a = $(x)_+ \implies a = x$ if $x > 0$ else $a = 0$

**Fig. 5**: Anechoic Chamber Recording Setup

## 3. PERFORMANCE EVALUATION

The performance evaluation of the proposed audio zooming systems is presented herein using simulation and real data experiments. Four objective measures that include Mean Square Error(MSE) [9], Output SINR (OSINR) [9], Perceptual Evaluation of Speech Quality (PESQ) [10] and Short-time Objective Intelligibility (STOI) [11] are utilized for the simulation experiments. Experiments were additionally conducted on real data recorded in anechoic chamber and in the field. The Mean Opinion Score (MOS) [12] measure is utilized to evaluate the performance of experiments with real recorded data. Non-refference PESQ score is also given to evaluate the speech intelligibility. The proposed time and frequency domain audio zooming systems are compared with RMVDR based system [3] for two microphones.

### 3.1. Simulation Experiments

Two microphones separated by 10 cm were taken for the simulation experiments. The objective parameters were computed for fifty target speech files from TIMIT database [13]. The target source was taken at azimuth $90°$, while the interference was assumed to be at $60°$. Signal-to-Interference Ratio (SIR) was taken to be 0dB. The experiments were performed in reverberant condition. For reverberant condition, the reverberation time was taken to be 100ms. The results are presented in Table 1. It is to be noted that the proposed MPDR beamforming based system outperforms the other methods. The performance of GJBF based system is comparable to RMVDR based system.

| Zooming System | OSINR(dB) | PESQ | STOI | MSE(dB) |
|---|---|---|---|---|
| RMVDR | 10.56 | 3.093 | 0.8935 | -25.83 |
| Griffith Jim | 9.35 | 3.0572 | 0.8655 | -23.56 |
| MPDR | 13.85 | 3.203 | 0.9175 | -33.453 |

**Table 1**: Objective Evaluation of proposed audio zooming systems

### 3.2. Experiments on Real Data

The data recording was done was in anechoic chamber and open field. Subjective listening tests were conducted. Mean Opinion Score (MOS) measure is presented to evaluate the proposed audio zooming system. Fifteen subjects in the age group of $20 - 25$ years were invited for listening the zoomed audio. The subjects listened the mixed received signal and the output. The rating was given for the quality of the output audio on the scale of 0 to 5 [14].

#### 3.2.1. Microphone Array Recordings in Anechoic Chamber

A uniform linear microphone array was utilized for recording in anechoic chamber with inter-element spacing as 4.5cm. The array consists of Sennheiser HSP 2 microphones. The target speaker

was placed at $90°$ azimuth. The interference was kept at $40°$ azimuth.Various kind of interference were selected as shown in table 2. Data of two channel separated by 9cm was utilized for evaluating the audio zooming systems. The MOS measures being close to or more than 4 shows a good perception of the output.

| Zooming Technique | Interference | MOS | |
|---|---|---|---|
| | | *Interference suppression* | *Speech quality* |
| RMVDR | Train | 3.25 | 4 |
| | Vacuum | 3.5 | 4 |
| | Speech | 3.25 | 3.75 |
| | Sonic | 3.25 | 4 |
| Griffith-Jim | Train | 4.16 | 3.75 |
| | Vacuum | 4.25 | 4 |
| | Speech | 4 | 3.75 |
| | Sonic | 4.16 | 3.85 |
| MPDR | Train | 4 | 4.25 |
| | Vacuum | 4.25 | 4.16 |
| | Speech | 4.16 | 4 |
| | Sonic | 4.16 | 4.25 |

**Table 2**: Performance evaluation for anechoic chamber recording

#### 3.2.2. Smartphone Recordings in Open Ground

As the application target is smartphone, we used two channel recording from Samsung Galaxy A5(2017) smartphone. The recording was performed in a open ground scenarios. The Two orators were located at $90°$ and $45°$ to the microphone array. MOS and non-reference PESQ (NR-PESQ) [15] measures are given in Table 3. High MOS and NR-PESQ measures shows the practical applicability of the proposed systems. The more field recording results are made available online at http://web.iitd.ac.in/~lalank/msp/demo.html for reviewers to evaluate.

| Zooming Technique | MOS | | NR-PESQ |
|---|---|---|---|
| | *Interference suppression* | *Speech quality* | |
| RMVDR | 3.5 | 4 | 2.835 |
| Griffith-Jim | 4.08 | 3.66 | 2.756 |
| MPDR | 4.25 | 4.166 | 3.021 |

**Table 3**: Evaluation for open ground recording

## 4. CONCLUSIONS

In this paper, two channel audio zooming systems are proposed for smartphone for the first time. Two channel time and frequency domain beamforming based target extraction is explored. A novel block-thresholding based post-filtering is utilized for creating the audio zooming effect. The proposed MPDR and GJ beamformer based systems are compared with two channel RMVDR based system. The MPDR based audio zooming system outperforms while performance of GJBF based system is comparable with RMVDR. The proposed systems are tested on smartphones for the expected result. MPDR based system can be utilized for smartphones with more than two microphones with better output. A mobile app is being developed for the same. Subjective and objective measures suggest rich user experience.

# References

[1] Keansub Lee, Hosung Song, Yonghee Lee, SON Youngjoo, and KIM Joontae, "Mobile terminal and audio zooming method thereof," Jan. 26 2016, US Patent 9,247,192.

[2] Carlos Avendano and Ludger Solbach, "Audio zoom," Dec. 8 2015, US Patent 9,210,503.

[3] Ngoc QK Duong, Pierre Berthet, Sidkieta Zabre, Michel Kerdranvat, Alexey Ozerov, and Louis Chevallier, "Audio zoom for smartphones based on multiple adaptive beamformers," in *International Conference on Latent Variable Analysis and Signal Separation*. Springer, 2017, pp. 121–130.

[4] J. Capon, "High-resolution frequency-wavenumber spectrum analysis," *Proceedings of the IEEE*, vol. 57, no. 8, pp. 1408–1418, 1969.

[5] Lloyd Griffiths and CW Jim, "An alternative approach to linearly constrained adaptive beamforming," *IEEE Transactions on antennas and propagation*, vol. 30, no. 1, pp. 27–34, 1982.

[6] Guoshen Yu, Stéphane Mallat, and Emmanuel Bacry, "Audio denoising by time-frequency block thresholding," *IEEE Transactions on Signal processing*, vol. 56, no. 5, pp. 1830–1839, 2008.

[7] David Mansour and A Gray, "Unconstrained frequency-domain adaptive filter," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 30, no. 5, pp. 726–734, 1982.

[8] D. Donoho and I. Johnstone, "Idea spatial adaptation via wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455.

[9] Valentin Emiya, Emmanuel Vincent, Niklas Harlander, and Volker Hohmann, "Subjective and objective quality assessment of audio source separation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 7, pp. 2046–2057, 2011.

[10] Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra, "Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs," in *Acoustics, Speech, and Signal Processing, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on*. IEEE, 2001, vol. 2, pp. 749–752.

[11] Cees H Taal, Richard C Hendriks, Richard Heusdens, and Jesper Jensen, "A short-time objective intelligibility measure for time-frequency weighted noisy speech," in *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*. IEEE, 2010, pp. 4214–4217.

[12] Mahesh Viswanathan and Madhubalan Viswanathan, "Measuring speech quality for text-to-speech systems: development and assessment of a modified mean opinion score (mos) scale," *Computer Speech & Language*, vol. 19, no. 1, pp. 55–83, 2005.

[13] Victor Zue, Stephanie Seneff, and James Glass, "Speech database development at mit: Timit and beyond," *Speech Communication*, vol. 9, no. 4, pp. 351–356, 1990.

[14] ITUT Rec, "P. 85 (1994) a method for subjective performance assessment of the quality of speech voice output devices," *International Telecommunication Union, Geneva Google Scholar*.

[15] International Telecommunication Union, "Non refference pesq score calculation," https://www.itu.int/rec/T-REC-P.562-200511-I!Amd2/.