# UC San Diego
## UC San Diego Electronic Theses and Dissertations

**Title**

Auxiliary Function Independent Vector Analysis Using a Harmonic Clique Dependence Model

**Permalink**

https://escholarship.org/uc/item/59c4p4pk

**Author**

Bezanson, Derrick English

**Publication Date**

2013

Peer reviewed|Thesis/dissertation

UNIVERSITY OF CALIFORNIA, SAN DIEGO

**Auxiliary Function Independent Vector Analysis Using a Harmonic Clique Dependence Model**

A Thesis submitted in partial satisfaction of the requirements for the degree Master of Science

in

Electrical Engineering (Signal and Image Processing)

by

Derrick English Bezanson

Committee in charge:

Professor William Hodgkiss, Chair
Professor Truong Nguyen
Professor Bhaskar D. Rao

2013

This thesis of Derrick English Bezanson is approved, and it is acceptable in quality and form for publication on microfilm:

_____

_____

_____
Chair

University of California, San Diego

2013

# Dedication

I dedicate this thesis to my wife, Kourtney Bezanson

and my daughter, Blayke Ann Bezanson.

# Table of Contents

# List of Figures

# List of Tables

# Acknowledgments

First and foremost, I would like to acknowledge the help and technical contributions I received from my colleague Amir Sarajedini. Amir helped me throughout the entire project to understand the concepts and background of BSS and IVA.

I would like to thank my chair, Professor William Hodgkiss. He took the time to help me define the project and through the end. Hodgkiss made sure I had a fundamental understanding of the theory behind BSS which prepared me for the defense. I would also like to thank my other committee members, Professor Rao and Professor Nguyen, for taking time to attend my defense.

My brother, Leverett Bezanson, also deserves big thanks. He is a UCSD student and not only helped me edit this Thesis, but has helped me with many of the classes we took together. And to everyone else who supported me and helped me edit this Thesis, Thank You to Tony Mauro, Chase Decker, Igor Fedorov, and my family.

ABSTRACT OF THE THESIS


Auxiliary Function Independent Vector Analysis Using a
Harmonic Clique Dependence Model


by


Derrick English Bezanson


Master of Science in Electrical Engineering (Signal and Image Processing)


University of California, San Diego, 2013


Professor William Hodgkiss, Chair


The problem of separating mixed signals using multiple sensors with little to no information about the source signals is known as Blind Source Separation (BSS). Many embedded systems, such as cell phones, have an audio environment which routinely suffers from multiple simultaneous audio or noise sources interfering with the desired user. One approach for improving audio quality is to use BSS techniques such as Independent Vector Analysis (IVA), an extension to the more common Independent

Component Analysis (ICA) approach, to remove interference or mitigate noise. However, these algorithms suffer from slow convergence rates, thus making it impractical to effectively separate sources in real time. This thesis explores Auxiliary Function Independent Vector Analysis (AuxIVA) with constraints on the frequency distribution of the audio components to improve convergence rates. AuxIVA has been shown to yield a faster convergence time and better results when separating audio sources compared to traditional IVA. By constraining input sources to be human speech, a harmonic frequency dependence model can be used to further improve convergence. We propose combining AuxIVA with a harmonic clique dependence model to achieve a more efficient algorithm, thus making a real time solution more viable. This thesis will demonstrate improved convergence rates and Signal to Interference Ratio (SIR) performance of the proposed technique relative to traditional IVA, AuxIVA, and AuxIVA with non-harmonic clique dependence.

# Chapter I  Introduction

Blind Source Separation (BSS) is the separation of a set of source signals from a set of mixed signals where only the statistical structure of the signal is assumed. BSS relies on the assumption that the source signals are independent of one another. Learning algorithms take advantage of this assumption by maximizing the statistical independence of each source from the multivariate input signal.

The classical example of BSS is the "cocktail party" problem where a number of people are talking simultaneously in a room. The earliest and most basic form of BSS problems started with a model of linear and instantaneous mixing of the sources. Independent Component Analysis (ICA) among one of the more successful approaches to this problem and has become widely adopted as a popular area of research [4]. BSS by ICA has received a lot of attention since the mid-1990s because of its potential applications in signal processing, such as speech recognition systems, telecommunications and medical signal processing. In contrast to correlation-based transformations, such as Principal Component Analysis (PCA), ICA decorrelates the signals but also reduces higher-order statistical dependencies, attempting to make the signals as independent as possible. Superficially, ICA is related to principal component analysis and factor analysis, but is a much more powerful technique capable of finding the underlying factors when these classic methods fail completely [9]. More recently a technique called Independent Vector Analysis (IVA), an extension of ICA, has been shown to improve the results of BSS [6]. IVA uses the same underlying assumptions of

the sources as ICA. However, IVA uses the entire spectrum as input to the objective function, creating a more efficient model.

An implementation of real-time solution to BSS is desirable in both the academic community and commercial industry because it could help reduce interfering audio sources. This is a challenging problem because of the complexity of IVA and its slow time to reach convergence. The goal of this thesis is to achieve a faster rate of convergence (separation of acoustic signals in a minimum number of iterations) with a high signal-to-interference ratio (SIR). This improvement will make real-time embedded solution, more viable.

A new set of update rules based on the auxiliary function technique (AuxIVA) is a recent (2011) improvement to IVA [13]. Since this new algorithm has already improved upon the traditional IVA algorithm, we will focus on optimizing AuxIVA for separating speech signals. Constraining the input sources to be speech allows us to use a harmonic dependence model to help the learning function converge faster. Since our goal in this thesis is to create a better model for real time separation of signals, we will constrain our test cases to simulate that of a small 2 microphone array (similar to modern cell phone architectures). A simulation test suite will be used to evaluate various types of speech and environments.

# Chapter II  Background

## II.1  Independent Component Analysis

### II.1.1  Overview

Independent Component Analysis (ICA) is a statistical and computational technique for revealing hidden factors that underlie sets of random variables, measurements, or any multivariate statistical data. This technique is commonly used to solve the BSS problem, otherwise known as the "cocktail party" problem. The cocktail party represents $N$ speakers and $M$ microphones in a room. For simplicity we will define $M = N$ however this is not a requirement for BSS. Acoustic sources (or speakers) are represented as $s_n, n = 1, \dots, N$ and microphone inputs as $x_n, n = 1, \dots, N$. We assume a spatial difference between each source and each microphone, thus there exists a linear instantaneous mixing of the sources [17]. This mixture is represented in the frequency domain as:

$$X(k) = A(k)S(k) \tag{1}$$

such that,

$$X(k) = [x_1(k) \dots x_N(k)]^T \tag{2}$$

$$S(k) = [s_1(k) \dots s_N(k)]^T \tag{3}$$

$$A(k) = \begin{bmatrix} a_{11}(k)a_{12}(k) \dots a_{1N}(k) \\ a_{21}(k)a_{22}(k) \dots a_{2N}(k) \\ \dots \\ a_{N1}(k)a_{N2}(k) \dots a_{NN}(k) \end{bmatrix} \tag{4}$$

where k is the frequency bin index, $^T$ is the transpose, and $A$ is the mixing matrix.

The goal of ICA is to estimate $S$ as best as possible when given only the observed inputs $X$. To do this, it is assumed that the sources are statistically independent of each other. ICA finds the independent components of a signal by maximizing the statistical independence of the estimated components. We may choose one of many ways to define independence, and this choice governs the form of the ICA algorithms. Traditionally (and for this thesis), ICA defines independence as the minimization of mutual information. This minimization will output a set of weight vectors called the unmixing matrix.

By multiplying the unmixing matrix (output of the ICA learning algorithm) with the mixed input we get the estimated sources represented as $y_n, n = 1, \dots, N$. The unmixing process is defined as:

$$Y(k) = W(k)X(k) \tag{5}$$

$$\begin{bmatrix} y_1(k) \\ \vdots \\ y_N(k) \end{bmatrix} = \begin{bmatrix} w_{11}(k)w_{12}(k) \dots w_{1N}(k) \\ w_{21}(k)w_{22}(k) \dots w_{2N}(k) \\ \vdots \\ w_{N1}(k)w_{N2}(k) \dots w_{NN}(k) \end{bmatrix} \begin{bmatrix} x_1(k) \\ \vdots \\ x_N(k) \end{bmatrix} \tag{6}$$

where $W$ is the unmixing matrix and ideally $W(k)A(k) = I \; (the \; identity \; matrix)$. The traditional approach to finding the unmixing matrix is to apply the natural gradient algorithm [16] (this will be explored in greater detail in the following section). Within the ICA framework, the learning is done over individual bins, thus no relationship between sequential bins is taken into account in this algorithm. It is for this reason that conventional ICA suffers from the "permutation problem" [3, 4].

*II.1.2 Permutation Problem*

ICA is effective at separating sources for each frequency bin. The permutation problem occurs when each unmixed source contains frequency bins which do not belong to it. Post processing was often used to solve the inconsistent permutation of the discrete Fourier transform bins after ICA has completed. Finding a solution to this problem has been an active area of research in recent years. A common approach to this problem is to solve it in the frequency domain and use all frequency components of each source signal together as a multivariate source. Rather than using a contrast function that measures the component-wise independence of each frequency bin, a contrast function that measures the whole independence among the multivariate sources can be applied (further explained in the next section). In general, ICA can be effective for BSS but suffers from the permutation problem and will not be studied in this thesis. For a more detailed discussion of ICA please refer to [4].

## II.2  Independent Vector Analysis

*II.2.1  Overview*

Independent vector analysis (IVA) for blind source separation is a more recent solution to the "cocktail party" problem. This technique is an extension of ICA but uses the entire frequency spectrum as an input to solve the permutation problem [3, 6]. In IVA a measure of independence is calculated from the entire spectrogram of each input. IVA fundamentally solves the permutation problem by separating whole spectrograms instead of binwise separation, as done by ICA. The calculations use multivariate probability density functions (PDFs) which take signal spectra across all frequency bins

as arguments. IVA models each individual source as a dependent multivariate symmetric super-Gaussian distribution (specifically a spherical Laplacian distribution as shown in Figure 1). This assumption holds while maintaining the fundamental assumption of BSS that each source is independent from the other [17].

In natural gradient IVA the signal is modeled by a bin-wise instantaneous mixture in the Short-Time Fourier Transform (STFT) domain. For the following equations we will omit the time frame index but note that the expectations are taken over time frames. Combining Eq. 1 with Eq. 5 we have:.

$$Y(k) = W(k)X(k) = W(k)A(k)S(k) \tag{7}$$

where $k$ is the frequency bin index.

In IVA unmixing is done over all frequency bins (bold is used to represent matrices which contain all sources and all frequency bins):

$$\boldsymbol{Y} = \boldsymbol{W}\boldsymbol{X} \tag{8}$$

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_N \end{bmatrix} = \begin{bmatrix} W_{11} \dots W_{1N} \\ \vdots \\ W_{N1} \dots W_{NN} \end{bmatrix} \begin{bmatrix} X_1 \\ \vdots \\ X_N \end{bmatrix} \tag{9}$$

$$\begin{bmatrix} \begin{pmatrix} y_1(1) \\ \vdots \\ y_1(K) \end{pmatrix} \\ \vdots \\ \begin{pmatrix} y_N(1) \\ \vdots \\ y(K) \end{pmatrix} \end{bmatrix} = \begin{bmatrix} \begin{pmatrix} w_{11}(1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_{11}(K) \end{pmatrix} & \cdots & \begin{pmatrix} w_{1N}(1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_{1N}(K) \end{pmatrix} \\ \vdots & \ddots & \vdots \\ \begin{pmatrix} w_{N1}(1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_{N1}(K) \end{pmatrix} & \cdots & \begin{pmatrix} w_{NN}(1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_{NN}(K) \end{pmatrix} \end{bmatrix} \begin{bmatrix} \begin{pmatrix} x_1(1) \\ \vdots \\ x_1(K) \end{pmatrix} \\ \vdots \\ \begin{pmatrix} x_N(1) \\ \vdots \\ x_N(K) \end{pmatrix} \end{bmatrix} \tag{10}$$

where N is the number of microphones (and sources) and $K$ is the number of frequency bins. Eq. (10) shows the coupling of the entire spectrum to each source.

*II.2.2 Objective function of IVA*

Similar to ICA a gradient-based algorithm is commonly used (i.e. natural gradient). The objective function uses multivariate activation functions derived from the PDFs to obtain an unmixing matrix which make the spectrograms independent. To find an optimal solution, the Kullback-Leibler Divergence (KLD) between $p(Y)$ and $\Pi_n p(Y_n)$ is used as a measure of independence in the whole spectrogram. The objective function is used to find the unmixing matrix $W$ which makes the output vectors $Y_1, \dots, Y_n$ independent and thus minimizing the KLD [6, 13].

$KLD(Y)$ is explained as the distance between the PDFs of $Y$ and the joint PDFs of $Y_n$ assuming independence. The unmixing matrices are estimated by minimizing the following objective function:

$$KLD(Y) = \sum_{n=1}^{N} E[G(Y_n)] - \sum_{k=1}^{K} \log |\det W(k)| \tag{11}$$

$$W(k) = [w_1(k) \dots w_N(k)]^H \tag{12}$$

where $w_n(k)$ is the column vector of weights for source n and $^H$ is the Hermitian transpose.

$$Y_n = [y_n(1) \dots y_n(K)]^T \tag{13}$$

$E[.]$ is the expectation over time frames and $G(Y_n)$ is the contrast function. The contrast function has the relationship $G(Y_n) = -\log p(Y_n)$ where $p(Y_n)$ represents a multivariate PDF for each source. This implementation uses a multivariate Laplacian PDF which has a spherical distribution property as shown in figure 1. The spherical contrast function used for IVA is represented as:

$$r_n = \|Y_n\|_2 = \sqrt{\sum_{k=1}^{K} |y_n(k)|^2} \tag{14}$$

where $\|.\|_2$ denotes the $L_2$-norm of a vector.

To minimize the objective function (Eq. 11) a learning algorithm is derived using a gradient decent method. The update equations for IVA result in a $\Delta W$ that ultimately steer a null towards the interfering source(s) and a beam towards the source. This is a convex optimization problem that is solved by applying the update rules based on the natural gradient [3,6,13]:

$$W(k) \leftarrow W(k) + \mu(I - E[\phi_k(Y)Y^H(k)])W(k) \tag{15}$$

$$\phi_k(Y) = [\phi_{1k}(Y_1) \ldots \phi_{Nk}(Y_N)] \tag{16}$$

$$\phi_{nk}(Y_n) = \frac{\partial G(Y_n)}{\partial y_n^*(k)} \tag{17}$$

where $^*$ is the complex conjugate, $\mu$ is the step size parameter and the expectation is over time frames. The step size is a tuning parameter that imposes a tradeoff between convergence speed and stability. Researching an optimum step size for a particular data set can be somewhat ambiguous, so for this thesis we will not use a step size value greater than 0.3 to avoid the algorithm diverging.

The phi function $\phi_k$ is the derivative of the log of the assumed distribution. This function is used in the update equation to match the assumed sources to our model:

$$\phi_k(Y) = -\frac{y_n(k)}{r_n} \tag{18}$$

where $r_n$ is defined in Eq. 14.

Figure 1 shows the comparison of an independent Laplacian join distribution (a) and a dependent spherical Laplacian distribution (b). ICA assumes (a) but suffers from the permutation problem. IVA assumes (b) which is a Laplacian distribution with a spherical property. By using this new dependent model for IVA each frequency bin depends on the entire spectrum, thus solving the permutation problem.



**Figure 1. Comparison between (a) An independent Laplacian distribution and (b) a dependent spherical Laplacian distribution [6]**

For a more complete discussion of IVA please refer to the original IVA paper by Kim, Lee, Attias [6] and the extension paper by Hiroe [3].

## II.3  Auxiliary Function IVA

### II.3.1  Overview

The natural gradient update for IVA described above has tradeoffs between convergence speed and stability based on the step size value used. If a large step size is used the algorithm could potentially diverge. Recently, an approach was developed to improve convergence time and eliminate the step size variable by using the auxiliary function technique to form a new objective function and new update equations. This

technique was first applied to ICA and later adapted for IVA. Auxiliary Function IVA (AuxIVA) is a framework to find an efficient iterative solution for nonlinear optimization problems and is an extension of the expectation-maximization (EM) algorithm (a common algorithm used for statistical inference problems in signal processing) [13, 14].

To introduce the auxiliary function technique we look at a general optimization problem where we want a vector $\theta = \theta^\dagger$ such that $\theta^\dagger = argmin_\theta J(\theta)$ where $J(\theta)$ is an objective function. The Auxiliary function technique defines a new equation:

$$J(\theta) = min_{\theta'} Q(\theta, \theta') \tag{19}$$

where $Q(\theta, \theta')$ is called the auxiliary function for function $J(\theta)$ and $\theta'$ is called an auxiliary variable. Instead of minimizing the objective function directly we minimize the auxiliary function $Q(\theta, \theta')$ in terms of $\theta$ and $\theta'$ alternatively. These variables are updated iteratively as:

$$\theta'^{(i+1)} = argmin_{\theta'} Q(\theta^{(i)}, \theta') \tag{20}$$

$$\theta^{(i+1)} = argmin_\theta Q(\theta, \theta'^{(i+1)}) \tag{21}$$

where $i$ is the iteration index. This process guarantees a monotonic decrease of $J(\theta)$. To find an appropriate $Q(\theta, \theta')$ is problem dependent. To extend IVA to use this technique we will use the objective function derived by Ono in [13].

*II.3.2  Objective Function*

The new objective function for AuxIVA is defined as:

$$Q(W,V) = \sum_{k=1}^{K} Q_k(W(k),V(k))$$

$$= \sum_{k=1}^{K} \left( \frac{1}{2} \sum_{n=1}^{N} w_n(k)V_n(k)w_n^H(k) - \log|\det W(k)| \right) + R \tag{22}$$

where $V_n(k)$ is the weighted covariance matrix at the $k$th frequency bin later defined in Eq. 24, $w_n(k)$ is the $n$th row of the unmixing matrix, and $R$ is a scalar constant term. AuxIVA can be thought of as taking the Taylor series expansion of the original IVA objective function but only keeping the $2^{nd}$ order term for the first update step we then use the output from this first step to minimize the weight vectors. It is because of this alternation of minimizing out new objective function in terms of $V$ and $W$ respectively that we can eliminate the step size parameter and still obtain a monotonic decrease.

The algorithm is summarized as the following alternative updates for all $k$ sources, which are applied in order until convergence. The axillary variable update step assumes the same spherical dependence $r_n$ as in IVA. The weighted covariance matrices $V_n(k)$ are updated for all $k$ as follows:

$$r_n = \sqrt{\sum_{k=1}^{K} |w_n^H(k)X(k)|^2} \tag{23}$$

$$V_n(k) = E\left[ \frac{G'(r_n)}{r_n} X(k)X^H(k) \right] \tag{24}$$

where the expectation is over time frames and $G'(r_n)$ is derived from the spherical Laplacian source PDF we are assuming.

Now we minimize $Q(W, V)$ in terms of $W$. Instead of updating all of $w_n(k)$ simultaneously, we update the weights from one source at a time keeping $w_l(k)$ fixed where $l \neq n$. From this we have the following equations:

$$w_n^H(k)V_n(k)w_n(k) = 1 \tag{25}$$

$$w_l^H(k)V_n(k)w_n(k) = 0 \ (l \neq n) \tag{26}$$

where Eq. 25 determines the scale of $w_n(k)$ and Eq. 26 determines the direction of $w_n(k)$. The weight vector is first updated and then normalized for all $n$. By combining Eqs. 25 and 26 we have:

$$w_n = (W(k)V_n(k))^{-1}e_n \tag{27}$$

where $e_n$ denotes the unit vector with the $n^{\text{th}}$ element unity. The last update step for a single iteration is to normalize by applying the following equation:

$$w_n(k) = \frac{w_n(k)}{\sqrt{w_n^H(k)V_n(k)w_n(k)}} \tag{28}$$

AuxIVA avoids the step size tuning problem in conventional IVA and gives effective iterative update rules which can guarantee the monotonic decrease of the objective function at each update. Similar to IVA, this method assumes a dependent spherical Laplacian PDF of the sources and thus does not suffer from the permutation problem. Recent studies have shown that this method can improve both SIR and convergence time when separating audio signals at a slight cost of added complexity. With this improvement, real-time IVA implementation is more feasible. For a more complete discussion of AuxIVA, please refer to the paper by Ono [13].

## II.4  Dependence Models Using Cliques

IVA and AuxIVA assume a spherical dependency model over all frequency bins (also known as radial symmetry). By using this model each frequency bin assumes dependence to every other bin equally. This dependency model is equivalent to a single clique in an undirected graph as shown in Figure 2 (top). For this thesis we define a clique as a subset of frequency bins of which each bin in the set assumes dependence to every other bin in the set.

The overlapping clique dependency model enables a more accurate model of statistical dependencies in accordance to the correlation coefficients observed in acoustic signals [11]. This model defines a fixed number of cliques which constitute a dependency graph such that neighboring frequency bins are assigned to the same clique while distant bins are assigned to different cliques. The permutation ambiguity is resolved by overlapped frequency bins between neighboring cliques. The clique sizes are set to be fixed with a fifty percent overlap. A recent study (2012) showed improved performance when for separating audio signals when compared to spherical dependency models in both IVA and AuxIVA [11]. Figure 2 shows the single clique distribution (top) as compared to the overlapping clique distribution (bottom).

**Figure 2. Single and Chain-like Overlapping Clique Models**

AuxIVA can be modified to use this model by changing the covariance matrix update step to sum over each clique:

$$\frac{G'_R(r_n)}{r_n} = \frac{1}{\sqrt{\sum_{k=k_{b1}}^{k_{c1}}|y_n(k)|^2}} + \frac{1}{\sqrt{\sum_{k=k_{b2}}^{k_{c2}}|y_n(k)|^2}} + \cdots$$

$$+ \frac{1}{\sqrt{\sum_{k=k_{bl}}^{k_{cl}}|y_n(k)|^2}}$$

(29)

where the set $[k_{bl}, \ldots, k_{cl}]$ represents the frequency bins indices for each of $l$ cliques. Subscript $b$ indicates the first bin index in the clique and subscript $c$ indicates the last. Although the overlapped clique dependence exhibits improved performance for some acoustic sources it may not match the characteristics of speech or other real-world signals with a strong harmonic structure. For a more complete discussion of this method refer to the paper by Lee [8].

## II.5  Pre and Post Processing

Pre-processing is not strictly required for IVA, however for this thesis we choose to pre-whiten our data to assist the learning steps and ultimately converge as fast as possible. Whitening removes the second order correlations to make the data uncorrelated while possibly leaving higher order statistics nonzero (thus not making the signals independent).

After convergence we have the optimum weights for each frequency bin. However, each bin contains an arbitrary scaling component which needs to be adjusted before source reconstruction. To overcome this we will perform rescaling of the weights based on the minimal distortion principle [12]. For IVA this is simply stated as:

$$W(k) \leftarrow diag\big(W^{-1}(k)\big)W(k) \tag{30}$$

Finally, we are ready to apply the unmixing matrix to our input data and produce an optimal estimation of each source. Lastly, the Inverse Short Time Fourier Transform (ISTFT) is performed on the frequency domain data. At this point we can listen to the output; if the separation worked well we hear each source individually instead of mixed. For comparison purposes each implementation of IVA throughout this thesis will perform the same pre and post-processing steps.

# Chapter III  The Problem

The goal of this thesis is to optimize AuxIVA with speech sources as input. Neither a single clique dependence model nor an overlapping clique model, as described in Chapter II, accurately match the structure of acoustic signals with strong harmonics. This thesis solves this problem by changing the objective function to use a harmonic dependence model as opposed to a spherical or overlapped clique model. We generate this model based on the fundamentals of human speech.

## III.1  Harmonics of Human Speech

A signal has harmonic structure when a 'fundamental' frequency component is accompanied by signal components at multiples of the fundamental frequency. Human speech, especially during voiced periods, exhibits a harmonic structure that we will exploit for improved convergence. This harmonic sound is produced from the vibration of vocal chords combined with air flowing out of the lungs. We discuss the assumed harmonic structure in the following section [16].

## III.2  Harmonic Frequency Dependence

The overlapped clique dependence can show improvement for some acoustic signals but may not match the characteristics of speech or other real-world signals. Specifically, voiced signals and other acoustic signals with strong harmonics require a dependence model more complex than just neighboring frequency bins. We can create a new dependence model by introducing a harmonic structure. This is an advanced frequency dependence model for IVA and should be more effective in separating sound

sources that have strong harmonic structures, such as speech and music signals. This method has been implemented for conventional IVA (as opposed to auxiliary IVA) and has been shown to yield better performance than overlapped clique dependence IVA when the input signals primarily are speech or music [2].

The clique structure we will analyze for this project is similar to the one described in [2], the only difference being that we do not have a clique that contains the entire spectrum. By omitting this clique we assume a stronger dependence to the harmonic model. The fundamental frequencies of each harmonic clique represent frequencies from 55 Hz to 880 Hz. The frequency is denoted by $F_h$ and defined as

$$F_h = F_1 \times 2^{(h-1)/12} \tag{31}$$

where $F_1 = 55Hz$ and harmonic clique $C_h$ varies over the range $h = \{1, ..., 49\}$ as described in [2].

$$C_h = \left\{ k \in \{1,..,K\} \middle| \frac{|f_k - mF_h|}{mF_h} < \delta \text{ for } \forall m \in \{1, ..., M\} \right\} \tag{32}$$

where $f_k$ is the frequency of the $k^{\text{th}}$ bin and $C_h$ is the $h^{\text{th}}$ local clique. Clique $C_h$ includes the frequency bins of the first eight multiples of $F_h$, i.e. $M = 8$. The bandwidth of the $m^{\text{th}}$ multiple of $F_h$, i.e. $mF_h$, is $2\delta$. Every two $mF_h$ consecutive harmonic cliques overlap by approximately 50%. We again use an overlapping model for sequential cliques to avoid the permutation problem as described in II.1.1. Figure 3 shows a graphical representation of the harmonic dependency model for FFT size of 1024 and a sampling frequency $F_s = 8000Hz$. This model should work well for most types of speech because the energy from any source should align with 1 or 2 of the cliques.

**Figure 3. Harmonic Clique Model**

## III.3  Solution

In this thesis we extend AuxIVA to use a harmonic dependency model as the input to the covariance matrix update step. This should produce improved results when separating two human speech signals. The auxiliary function IVA algorithm can be adapted to use a harmonic clique dependence model (AuxIVA_harmonic). The input sources will be constrained to have strong harmonics such as human speech. This approach is expected to yield a faster convergence time and improved SIR when compared to traditional AuxIVA or the overlapping clique dependency model.

The update rules proposed by [13] will be used, paired with the harmonic distribution proposed by [2]. We can modify AuxIVA to use the equation below for the

covariance matrix update step. We form our new harmonic dependent update for $r_n$ by combining Eqs. 29 and 32:

$$p([y_n(1), \dots, y_n(K)]) \propto \mathrm{E}\left[ -\sum_{h=1}^{H} \sqrt{\sum_{k \in C_h} \frac{|y_n(k)|^2}{(\sigma_h(k))^2}} \right] \Rightarrow r_n \tag{33}$$

where $H$ is total number of cliques, $\sigma_h(k)$ adjusts the variance of the variables (set to 1 as defined in [2]) and $C_h$ is the set of frequency bins that belongs to harmonic clique $h$.

By using a harmonic dependence model and the adjusted objective function, we expect a more accurate scaling factor when the inputs are constrained to be speech (or any acoustic source with strong harmonics). We will compare the performance of each of the following 4 algorithms:

1) Natural Gradient Independent Vector Analysis (IVA)

2) Auxiliary Function IVA (AuxIVA)

3) AuxIVA using the overlapped clique dependence (AuxIVA_overlap)

4) AuxIVA using the harmonic clique dependence (AuxIVA_harmonic)

Each algorithm will be implemented and compared to each other. Algorithms 1-3 have already been shown to produce increasingly better results in terms of performance and convergence. We expect the proposed algorithm (AuxIVA_harmonic) to produce better results respectively when given a constrained harmonic input. Each algorithm will be compared to each other in the following two ways:

1) Signal-to-Interference Ratio (SIR) after convergence

2) Rate of convergence – SIR as a function of number of iterations

## Chapter IV  Experiment

A custom simulation environment will evaluate the performance of the proposed algorithm. Speech and other audio samples were gathered from numerous online sources (primarily the open speech repository [15]). Other speech samples were recorded in house to be used as real world input. We collected 15 different speech samples (mixed male and female, Spanish and English) for use as input to the simulator. The samples generally are a single person talking or reading for about 7 to 10 seconds with no more than a 1-2 second pause between sentences. Because the speech samples have been pulled from different sources and recorded using unknown microphones we first normalize them by dividing the entire sample by its maximum amplitude before combining them into a single 2-channel mixed signal.

The simulator used for this experiment was derived from the Fast Image Source Method (Fast ISM) toolbox implemented in MATLAB [1, 10]. This toolbox uses the well-known image method for the purpose of simulating reverberant audio data in small-room acoustics. The image method simulates the impulse response between two points in a small rectangular room. This resulting impulse response is then convolved with any desired input signal (in our case a speech signal) which simulates room reverberation of the input.

Some modifications were made to this toolbox to allow multiple simultaneous speakers. The room is modeled in 3D with dimensions 5 x 5 x 4 meters. All simulations performed in this thesis use this room model with the microphone array and sources at the same height of 1.5 meters. Many variables can be adjusted, such as position of the

sources and microphone array spacing, which result in different room impulse responses (RIRs). For better comparison results, the random seed used in the reverberant calculations is kept constant when re-running the test for each algorithm. A constant reverberation time is set to T60 = 200ms (time required for reflections to decay 60 dB) and can be set to zero to test the anechoic case. Reflection coefficients for each wall can be adjusted to test different kinds of flooring and room configurations. We will use a constant set of reflection coefficients which simulate a carpeted floor and standard material on the ceiling and walls. The coefficient values are described as the following: $[0.9, 0.9, 0.9, 0.9, 1, 0.9]$ which correspond to the y-dimension walls, x-dimension walls, floor and ceiling, respectively.

A single random test will pull two arbitrary speech samples from the set and place them in a pseudo random location in the room. We limited the tests so that both sources would be the same distance away from the microphone-array (between 1 and 2 meters) with no less than 30 degrees of separation. Figure 4 shows an example of the simulated room as a top down view. Source 1 (labeled S1) was randomly selected to be *speech4.wav* and placed at -30$^o$ degrees relative to the microphone array. Source 2 (S2) was selected to be *speech5.wav* and set at 55$^o$. In the Chapter 5 we will look at the results from 100 random tests where the sources are pulled from a database of 15 different speech samples and are modeled in a similar random configuration as shown in Figure 4.

S1: speech4, S2: speech5, Distance: 1.40(m), lamda: 0.09, nmics: 2

**Figure 4. Simulation Environment Example Top View**

# Chapter V  Results

## V.1  Performance Measurement

Performance is measured primarily in two ways: Signal-to-Interference ratio improvement (SIR) and number of iterations until SIR convergence. SIR is similar to the signal-to-noise ratio, but in this case the interference is specific to co-channel interference from the other source. SIR improvement implies that we are computing the difference of the SIR from the mixed (input as recorded by microphone) signal to the original source signal (signal before mixing). To get an accurate SIR for comparison purposes we will use the BSS toolbox [19]. This tool box requires input (ground truth) sources before mixing. It takes into account the difference in signal power of each source and outputs measurements based on the decomposition of each estimated source signal into a number of contributions corresponding to the target source and interfering sources. For this thesis we will average the SIR for the 2 outputs together to give a single metric and more easily compare between algorithms.

When the algorithm is run for a large number of iterations, the weights will converge to their optimum values which correspond to a maximum SIR for any given test. We will consider the number of iterations to reach this convergence as the 'convergence rate.' When comparing performance of different algorithms, we will look at the converged SIR as well as the convergence rate to determine the performance.

**V.2  Experiment 1 (Anechoic)**

We conducted an experiment using the specifications listed in Table 1. This experiment consists of 100 tests using random input (pulled from a set of 15 speech samples). The two sources were placed at angles randomly chosen from the set of angles specified in Table 1. For each individual test all parameters are held constant for algorithm comparison. The output from this test shows the averaged SIR for all signals over all tests for each of the four algorithms as shown in Figure 5. A step size of 0.25 was used for the IVA tuning parameter. It has been shown [13] that a step size value of 0.3 or higher potentially can cause the IVA algorithm to diverge. IVA (red) is included in this experiment to show a clear improvement when using the auxiliary function update rules. For the overlapping clique dependence model (green) we used four 50% overlapped cliques as described in [8]. We see that the overlapped clique model we chose gave almost identical results on average when compared to traditional (one single spherical distribution) AuxIVA (blue) model.

The comparison we are most interested in is the one between our proposed AuxIVA_harmonic model (black) and AuxIVA (blue). We see an improvement in both convergence rate (speed) and SIR. AuxIVA_harmonic takes 16 iterations to reach 10dB while AuxIVA takes 23. We will consider the number of iterations for convergence to be about 45 for both AuxIVA and AuxIVA_harmonic. Our results show approximately 1 dB improvement after convergence over traditional AuxIVA. The SIR is relatively high in this case because we used the raw input without any added reflections. The

anechoic experiment is similar to a real world environment where there is little to no reverberation of the source signals, such as outdoors or in a very large room.

**Table 1. Parameters for 100 random tests**

| Parameter | Value(s) |
|---|---|
| Short-time Fourier Transform | Size: 1024 at 75% overlap |
| Random Angles | {-85, -60, -30, 0, 30, 60, 85} |
| Speech Samples | {*speech1.wav,…,speech15.wav*} |
| Distance from array | 1m – 2m |
| Number of microphones | 2 |
| Microphone spacing | 9cm |
| Duration of samples | Approximately 7 seconds |
| Reflections | Turned off for Experiment 1 Turned on for Experiment 2 |



**Figure 5. Averaged results from 100 random anechoic tests**

## V.3  Experiment 2 (Echoic)

The anechoic case showed the fundamental approach of using a harmonic dependence model and worked as expected with better average performance. Next we show the results from the same random test as described in Table 1 except we turn on reflections to simulate a more realistic reverberant room. For each individual test reverberations are calculated using the same random seed number to ensure accurate comparisons between algorithms. The results from this test are shown in Figure 6 and immediately we notice all algorithms have a much lower overall SIR than in the anechoic case. This is partially due to the fact that our original source did not contain any distortion or interference. The mixed source now has reverberations and the resulting unmixed source still contains some of the reverberations as well. Thus, the estimated unmixed signal will not match up to its ground truth signal as clearly as the anechoic case. Therefore, we will only compare the algorithms to each other for a specific test (reflections on or off) and not compare the tests themselves in terms of SIR values.

Figure 6 shows that natural gradient IVA (red) again has a very slow time of convergence, such that 100 iterations were not sufficient to see it converge when using a step size of 0.1. This illustrates the problem of having a slow convergence when using an arbitrary step size. We also note that AuxIVA_overlap has a slight performance increase over Aux_IVA which meets our expectations. As in the anechoic case, we are most interested in the performance of AuxIVA compared to AuxIVA_harmonic. Experiment 2 results in an improvement in both speed and SIR, but the improvement is

greater in terms of percentage of the anechoic test. AuxIVA_harmonic takes 30 iterations to reach 6dB whereas AuxIVA takes 50. We see about a 1-1.5 dB increase for SIR overall which is a 15% improvement, where the anechoic test showed approximately a 10% improvement.



**Figure 6. Averaged results from 100 random echoic tests**

## V.4  Case Study 1 - Harmonic Input

### V.4.1  Echoic Experiment

This experiment is a single test extracted from one of the averaged test cases above to further analyze the input and output signals. The test setup parameters are shown in Table 2. The two speech samples used are *speech1.wav* for Source #1 (S1), a

man reading from a book recorded by us, and *speech4.wav* for Source #2 (S2), a man

counting in Spanish found online at [15].

**Table 2. Parameters for case study 1**

| Parameter | Value(s) |
|---|---|
| Short-time Fourier Transform | Size: 1024 at 75% overlap |
| Speech Samples | S1: *speech1.wav*;   S2: *speech4.wav* |
| Angles | S1: -60$^o$;   S2: 30$^o$ |
| Distance from array | 1.5 meters |
| Number of microphones | 2 |
| Microphone spacing | 0.09 meters |
| Duration of samples | 7.3 seconds |
| Reflections | Turned On (echoic) |

S1: speech1, S2: speech4, Distance: 1.50(m), lamda: 0.09, nmics: 2

**Figure 7. Case Study 1 Top View of Simulation**

To verify that the source inputs are harmonic we look at the spectrogram shown in Figures 8 and 9. Harmonic content of speech predominantly shows up in vowels and is recognized by the ladder like shape in the frequency domain during a short time sample. Figure 9 shows a very strong harmonic structure for *speech4.wav*. Figure 8 shows harmonic content as well but is noisier; this is likely due to the recording environment.

**Figure 8. Spectrogram of Source 1: speech1.wav**



**Figure 9. Spectrogram of Source 2: speech4.wav**

To visually analyze the output, we look at the unmixing matrix returned by AuxIVA_harmonic. We expect the unmixing matrix to steer a null towards the interfering signal and a beam towards the desired signal. Beam patterns can be calculated from the unmixing matrix for each source and are used to verify that the algorithm is working as expected. Beam patterns are calculated by taking the inner product of each weight vector with each direction vector from $-90^o$ to $90^o$. We can do this for each source and expect to see a prominent null at the corresponding interfering source(s) [5].

To extract Source 1 the output weight vector will steer a null towards the interfering source [7]. In this case Source 2 is considered the interferer for Source 1 and vice-versa. Figure 10 shows the beam pattern of the weight vector for Source 1. We see a dominant null at $30^o$ over all of the frequency bins. The algorithm nulls out the interfering source and has a maximum at the angle it is located at. We are only using a 2 microphone array so the maximum is not as clear as the null. Source 2 has a very dominant harmonic structure and we can see that the null pattern in Figure 10 shows a similar looking harmonic pattern among nulls over frequency bins. Figure 11 shows a similar beam pattern for Source 2 with the null at $-60^o$ as expected.

**Figure 10. Beam Pattern of weights for Source 1 – Beam at -60 and null at 30$^{o}$**

**(echoic)**



**Figure 11. Beam Pattern of weights for Source 1 – Beam at 30 and null at -60$^{o}$**

**(echoic)**

The use of two strongly harmonic inputs in our test environment is expected to yield favorable results because of the harmonic dependence model. Figure 12 shows the results from this test up to 150 iterations for all four algorithms. As expected, AuxIVA has a better performance than IVA with a step size of 0.25. We also note that AuxIVA_overlap and AuxIVA have almost identical performance. This shows that a simple overlapping clique dependence model is not sufficient for improving the performance of separating two strongly harmonic speech signals. This is because an overlapped dependence assumes neighboring frequency bins are assigned to the same clique, while distant bins are assigned to different cliques, where our harmonic dependence is more flexible. As expected, AuxIVA_harmonic showed an improved SIR as well as a faster rate of convergence.



**Figure 12. Case Study 1 (echoic) SIR Results**

*V.4.2  Anechoic Experiment*

The following experiment used the same harmonic input as described above, however this time we examine the anechoic case. We look at this case to demonstrate that without reflections the signal separation performs as expected with ideal SIR and convergence. In a real world scenario this case can be thought of as being outside or in a soundproof room. Input spectrograms are the same as in Figures 8 and 9. We reiterate that the inputs have a strong harmonic structure and thus we expect to see improved performance from all algorithms yet AuxIVA_harmonic to stand out.

The follows figures show beam patterns obtained from the output weight vectors from the AuxIVA_harmonic algorithm. These beam patterns noticeably differ from the beam patterns in the echoic case. Figure 13 shows the beam pattern for source 1. There is a very strong null at 30 degrees where the interfering source is located and a beam at $-60^o$. The beam is harder to see because we only are using a 2 microphone array. We also notice a grating lobe on the bottom left part of the plot. This grating lobe appears here because of the microphone spacing. The spacing for this experiment is set to 9 cm which is larger than the half wavelength of our input, thus causing the lobe. Figure 14 shows a similar but opposite pattern for Source 2. There is a prominent null at $-60^o$ for the interfering source and a beam at $30^o$. This figure has a little more distortion around the null then in Figure 13 because Source 1 is a more broadband input in with weaker harmonics. This causes the separation SIR to be weaker for this particular source.

**Figure 13. Beam Pattern of weights for Source 1 – Null at 30$^\text{o}$ (anechoic)**



**Figure 14. Beam Pattern of weights for Source 2 – Null at -60$^\text{o}$ (anechoic)**

The SIR performance for AuxIVA_harmonic in Figure 15 shows a 4dB improvement over AuxIVA. Time of convergence improves as well. AuxIVA does not converge until around 20 iterations where AuxIVA_harmonic converges more quickly at about 13 iterations. For this specific test, AuxIVA_overlap has acceptable performance but has a lower SIR (comparatively) then in the echoic case. This lack in performance is due to the lack of reverberation and strong harmonic structure. The overlapped clique dependence model assumes neighboring bins are dependent and does not perform as well for this type of input. We also note that after enough time, natural gradient IVA performs the same as AuxIVA_harmonic. However, the time for IVA to converge is over 100 iterations.



**Figure 15. Case Study 1 (anechoic) SIR Results**

## V.5  Case Study 2 – Non-Harmonic Input

Case Study highlights why AuxIVA_harmonic performs supirior when the inputs are constrained to be harmonic in nature. Now we will look at a case study when the inputs do not have strong harmonics. With non-harmonic acoustic inputs, we would not expect AuxIVA_harmonic to perform as well as it did in Case Study 1. The two input sources for this experiment are both music sources. Some music can have a strong harmonic structure such as opera. However, we chose samples that are much more broadband.

**Table 3. Parameters for case study 2**

| Parameter | Value(s) |
|---|---|
| Short-time Fourier Transform | Size: 1024 at 75% overlap |
| Speech Samples | S1: *jazzTrio.wav*; S2: *harpsichord.wav* |
| Angles | S1: -10$^\text{o}$;   S2: 60$^\text{o}$ |
| Distance from array | 1.4 meters |
| Number of microphones | 2 |
| Microphone spacing | 0.09 meters |
| Duration of samples | 7.0 seconds |
| Reflections | Turned Off (anechoic) |

S1: JazzTrio, S2: Harpsichord, Distance: 1.40(m), lamda: 0.09, nmics: 2

S1: -10°

S2: 60°

**Figure 16. Case Study 2 Top View of Simulation**

Table 3 shows the test parameters for Case Study 2 and Figure 16 shows the simulated room setup. We choose Source 1 to be a jazz trio sample and Source 2 to be a harpsichord. The angles and distance away from the array were chosen at random. We decided to run this experiment with reflections turned off to better analyze each algorithm output based on the structure of their inputs without interference from reverberation or additive noise.

To verify that the source inputs do not have a harmonic input we look at their spectrograms. Figure 17 shows the spectrogram for the jazz trio sample (Source 1). Each of the three instruments in this sample have some harmonic content associated

with them. However because they are mixed together and have a somewhat random pitch, the spectrogram is more broadband and does not conform to a harmonic pattern as a whole. Figure 18 shows the spectrogram for the harpsichord sample (Source 2). This figure shows an even broader spectrum with almost no distinguishable harmonic structure as a whole.



**Figure 17. Spectrogram of Source 1 jazzTrio.wav**

**Figure 18. Spectrogram of Source 2 harpsichord.wav**

When compared to the other algorithms, AuxIVA_harmonic did not perform as well in terms of convergence rate. Figure 19 shows the separated signals for AuxIVA achieved an SIR of just under 16dB. We still achieve a separation of signals and eventually it reaches the same SIR as AuxIVA, but it takes close to 90 iterations to converge. The SIR does eventually converge to the same SIR as AuxIVA. This illustrates the importance of matching the clique distribution to the excepted input. While this is an acceptable level of separation to the human ear, when compared with AuxIVA_overlap it is 2dB lower. We can conclude that an overlapping clique model is a superior model for this specific input. Additionally, we can compare this performance to Case Study 1. In Case Study 1 with reflections turned off and voice as input we were able to obtain an SIR of over 26dB which is a noticeable difference of 10dB. It is

intuitive that a harmonic dependence model performs better when the input sources have a strong harmonic structure, and this test case shows that for non-harmonic input we still achieve a separation but with a decrease in performance.



**Figure 19. Case Study 2 SIR Results**

# Chapter VI Conclusions

## VI.1   Conclusion

The goal of this thesis was to identify an ideal solution to the Blind Source Separation problem that would make a real-time implementation feasible. We achieved this by applying a harmonic dependency model to Auxiliary Function IVA and constraining the input sources to be speech. We have shown that a harmonic dependency model increases performance when the source signals have a strong harmonic structure. We have also shown that Auxiliary Function IVA eliminates the step size parameter and converges much faster than IVA. To combine these two improvements we modified the objective function of AuxIVA to use a harmonic clique structure. This optimized algorithm converged faster and yielded a higher SIR when compared to traditional IVA, AuxIVA, and AuxIVA with an overlapping clique dependence model. The proposed method (AuxIVA_harmonic) has shown an improvement in convergence time (measured in number of iterations) and an improvement in SIR on average. AuxIVA_harmonic works well when separating sources that have a strong harmonic structure. However, it can result in worse performance when the acoustical input has a more broad (non-harmonic) structure.

A simulation environment was developed to simulate multiple speakers in a small room. We confirmed that the algorithm achieved the desired source separation by listening to the outputs and observing the beam patterns. The beam patterns verify our experiment is working as expected by pointing a null at the interfering source and a beam towards the desired source.

This thesis has shown that AuxIVA has superior performance to traditional IVA algorithms in general by using a more complex but effective update rules. BSS using IVA is a popular area of research and is constantly being improved and adapted for specific purposes. In many cases, the performance improvement of AuxIVA can be applied to other IVA related areas of research.

## VI.2  Future Work

A harmonic clique dependence model yields better performance when the input is constrained to speech signals. However, more research is needed to find a dependence model that works better for non-harmonic content or a mix of different types of acoustic signals. Through use of our custom test suite, various frequency dependence models could be developed and easily tested for different environments and situations. Additionally, more research and experimentation could be done to adapt to a changing environment such as non-stationary sources, more array elements, or more sources.

# References

[1]    J. Allen and D. Berkley, "Image method for efficiently simulating small room acoustics," J. Acoust. Soc. Amer., vol. 65, 1979.

[2]    C.H. Choi, W. Chang and S.Y. Lee, "Blind source separation of speech and music signals using harmonic frequency dependent independent vector analysis," Electronic Letters, 2012 Vol. 48 No. 2.

[3]    A. Hiroe, "Solution of permutation problem in frequency domain ICA, using multivariate probability density functions," ICA'06, 2006, pp. 601–608.

[4]    A. Hyvarinen, J. Karhunen, and E. Oja, "Indepedent Component Analysis," New York, Wiley Interscience, 2001.

[5]    M.Z. Ikram and D.R. Morgan, "A beamforming approach to permutation alignment for multichannel frequency-domain blind speech separation," Proc. ICASSP 2002, May 2002, pp. 881-884.

[6]    T. Kim, H.T. Attias, S.Y. Lee, and T.W. Lee, "Blind source separation exploiting higher-order frequency dependencies," IEEE Trans. Audio Speech Lang. Process., 2007, 15, (1), pp. 70–79.

[7]    S. Kurita, H. Samwatari, S. Kajita, K. Takeda, and F. Itahra, "Evaluation of blind signal separation method using directivity pattern under reverberant conditions," in Pmc. ICASSP 2000, June 2000, pp. 3140-3143.

[8]    I. Lee, G.J. Jang, and T.W. Lee, "Independent vector analysis using densities represented by chain-like overlapped cliques in graphical models for separation of convolutedly mixed signals," Electronic Letters, 2009, 45, (13), pp. 710–711.

[9]    T.W. Lee, "Introduction to Independent Component Analysis" http://cnl.salk.edu/~tewon/ICA/bkup/intro.htm

[10]   E. Lehmann and A. Johansson, "Diffuse Reverberation Model for Efficient Image-Source Simulation of Room Impulse Responses," IEEE Transactions on Audio, Speech and Language Processing, vol. 18, no. 6, pp. 1429-1439, August 2010.

[11]   Y. Liang, S.M. Naqvi and J. Chambers, "Overcoming block permutation problem in frequency domain blind source separation when using AuxIVA algorithm," Electronic Letters 12th April 2012 Vol. 48 No. 8.

[12]   K. Matsuoka and S. Nakashima, "Minimal distortion principle for blind source separation," Proc. ICA 2001, pp.722–727, Dec. 2001.

[13]    N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," 2011 IEEE Workshop on Applications of Signal Processing to Audio and Acoustic, New Paltz, USA, October 2011.

[14]    N. Ono and S. Miyabe, "Auxiliary-function-based Independent Component Analysis for Super-Gaussian Sources," Proc. LVA/ICA, pp.165-172, 2010.

[15]     "Open Speech Repository," http://www.voiptroubleshooter.com/open_speech

[16]    L.R. Rabiner and R.W. Schafer, "Digital Processing of Speech Signals," pp. 38-82 Prentice Hall, 1978.

[17]    A.M. Shirazi, "Blind Separation and Tracking of Sources with Spatial, Temporal, and Spectral Dynamics," Ph.D. University of California, San Diego, 2012.

[18]    P. Smaragdis, "Blind separation of convolved mixtures in the frequency domain," Neurocomputating, pp. 251–276, 1998.

[19]    E. Vincent, C. Fevotte, and R. Gribonval, "Performance Measurement in Blind Audio Source Separation," IEEE Trans. ASLP, vol. 14, no. 4, pp. 1462–1469, 2006.