

Blind Separation of Speech Mixtures via Time-Frequency Masking

Özgür Yılmaz and Scott Rickard

Abstract—Binary time-frequency masks are powerful tools for the separation of sources from a single mixture. Perfect demixing via binary time-frequency masks is possible provided the time-frequency representations of the sources do not overlap, a condition we call *W-disjoint orthogonality*. We introduce here the concept of *approximate W-disjoint orthogonality* and present experimental results demonstrating the level of approximate W-disjoint orthogonality of speech in mixtures of various orders. The results demonstrate that ideal binary time-frequency masks exist which can separate several speech signals from one mixture. While determining these masks blindly from just one mixture is an open problem, we show that we can approximate the ideal masks in the case where two anechoic mixtures are provided. Motivated by the maximum likelihood mixing parameter estimators, we define a power weighted two-dimensional histogram constructed from the ratio of the time-frequency representations of the mixtures which is shown to have one peak for each source with peak location corresponding to the relative attenuation and delay mixing parameters. The histogram is used to create time-frequency masks which partition one of one mixture into the original sources. Experimental results on speech mixtures verify the technique. Example demixing results can be found online: <http://alum.mit.edu/www/rickard/bss.html>

I. INTRODUCTION

The goal in blind source separation (BSS) is to determine the original sources given mixtures of those sources. When the number of sources is greater than the number of mixtures, the problem is degenerate in that traditional matrix inversion demixing cannot be applied. However, when a representation of the sources exists such that the sources have disjoint support in that representation, it is possible to partition the support of the mixtures and obtain the original sources. One solution to the problem of degenerate demixing is thus to (1) determine an appropriate disjoint representation of the sources and (2) determine the partitions in this representation which demix. In this paper, we show that the Gabor expansion (i.e., the discrete short-time (or windowed) Fourier transform) is a good representation for demixing speech mixtures. Specifically, we show that partitions of the time-frequency lattice exist that can demix mixtures of several speech signals from one mixture. Determining the partition blindly from one mixture is an open problem, but, given a second mixture, we describe a method for partitioning the time-frequency lattice which separates the sources.

Özgür Yılmaz (oyilmaz@math.umd.edu) is with the Department of Mathematics, University of Maryland.

Scott Rickard (rickard@ieee.org) is with the Department of Electronic and Electrical Engineering, University College Dublin, Ireland.

This work was partially funded by Siemens Corporate Research, Princeton, NJ.

Submitted to IEEE Transactions on Signal Processing, November 4, 2002. Revised March 27, 2003. Second revision July 24, 2003.

Formally, let \mathcal{S} be the family of signals of interest. Typically \mathcal{S} will be some collection of square integrable bandlimited functions. Suppose there exists some linear transformation $T : s_j \in \mathcal{S} \mapsto S_j$ (where T maps the set \mathcal{S} to another family of functions) with the following properties:

- (i) T is invertible on \mathcal{S} (i.e., $T^{-1}(Ts) = s, \forall s \in \mathcal{S}$).
- (ii) $\Lambda_j \cap \Lambda_k = \emptyset$ for $j \neq k$, where Λ_j is the support of S_j , i.e., $\Lambda_j := \text{supp } S_j := \{\lambda : S_j(\lambda) \neq 0\}$.

For example, we can consider the case where \mathcal{S} is a collection of square integrable functions with mutually disjoint supports in the Fourier domain; any two functions s_1 and s_2 in \mathcal{S} satisfy $\hat{s}_1(\omega)\hat{s}_2(\omega) = 0$ for all ω , where \hat{s}_j denotes the Fourier transform of s_j . Then if we define T on \mathcal{S} as $Ts := \hat{s}$, it is clear that T satisfies (i) and (ii).

For any T with properties (i) and (ii), we can demix a mixture x_1 of signals in \mathcal{S} , $x_1(t) = \sum_{j=1}^N s_j(t)$, via

$$s_j = T^{-1}(1_{\Lambda_j} T x_1) \quad (1)$$

where 1_{Λ_j} is the indicator function of the set Λ_j , i.e.,

$$1_{\Lambda_j}(\lambda) := \begin{cases} 1 & \lambda \in \Lambda_j \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Going back to our example above, this corresponds to s_j being equal to the inverse Fourier transform of $1_{\Lambda_j} \hat{x}_1$ which is certainly true since the functions in \mathcal{S} satisfy (ii).

Suppose now that we have another mixture $x_2(t) = \sum_{j=1}^N a_j s_j(t - \delta_j)$, which is the case in anechoic environments when we have two microphones. In the mixing, a_j and δ_j are the relative attenuation and delay parameters respectively corresponding to the j^{th} source. Assume

- (iii) $\text{supp } Ts(\cdot - \delta) = \text{supp } Ts$ for any $s \in \mathcal{S}$, $\forall |\delta| < \Delta$, and
- (iv) there exist functions F and G such that $a_j = F(Tx_1(\lambda), Tx_2(\lambda))$ and $\delta_j = G(Tx_1(\lambda), Tx_2(\lambda))$ for $\lambda \in \Lambda_j$ for $j = 1, \dots, N$,

where Δ is the maximum possible delay between mixtures due to the distance separating the sensors. Using (iii) and (iv), we can label each $\lambda \in \text{supp } Tx_1$ with the pair $(F(Tx_1(\lambda), Tx_2(\lambda)), G(Tx_1(\lambda), Tx_2(\lambda)))$, and Λ_j is exactly the set of all points with the label (a_j, δ_j) . It follows that given the mixtures $x_1(t)$ and $x_2(t)$, we can demix via

$$s_j = T^{-1}(1_{\Lambda_j} T x_1). \quad (3)$$

Clearly, (iii) will be satisfied for the example above since the Fourier transform of $s(\cdot - \delta)$ will be just a modulated version of the Fourier transform of s and thus it will have the

same support as s . As to the existence of functions F and G , one can show that $F(\hat{x}_1(\omega), \hat{x}_2(\omega)) = |\hat{x}_2(\omega)/\hat{x}_1(\omega)|$ and $G(\hat{x}_1(\omega), \hat{x}_2(\omega)) = -\frac{1}{\omega} \angle(\hat{x}_2(\omega)/\hat{x}_1(\omega))$ where $\angle z$ denotes the phase of the complex number z taken between $-\pi$ and π , satisfies (iv).

The general algorithm explained above mainly depends on two major points: (a) the existence of an invertible transformation T that transforms the signals to a domain on which they have disjoint representations (properties (i), (ii), and (iii)), and (b) finding functions F and G that provide the means of labeling on the transform domain (property (iv)). Note that in the description above we required F and G to yield the exact mixing parameters. Although this is desired since the mixing parameters provide the perfect labels and can also be used for various other purposes (e.g., direction-of-arrival determination), it is not necessary for the demixing algorithm to work. Some function that provides a unique labeling on the transform domain is sufficient. Moreover, requirement (ii) that the transformation T is “disjoint” is very strong. In practice, one is usually more interested in transforms that satisfy (ii) in some approximate sense. Transforms that result in sparse representations of the signals of interest, representations where a small percentage of the signal coefficients capture a large percentage of the signal energy, can lead to (ii) being approximately satisfied.

There are many examples in the literature that use this type of approach with various choices of T for various mixing models and demixing methods [1–12]. The mixing model in [1–3, 5, 8, 9, 11] is “instantaneous” (sources have different amplifications in different mixtures) while [4, 6, 7, 10, 12] use an anechoic mixing model (sources have different amplifications and time delays in different mixtures). [1–3, 11] consider the time domain sampling operator as T . The general assumption in these is that at any given time at most one source is non-zero. [4–7, 9, 10, 12] use the short-time Fourier transform (STFT) operator as T . Condition (ii) is satisfied in this case, at least approximately, because of the sparsity of the time-frequency representations of speech signals. Empirical support for this can be found in [7, 13], and a more extensive discussion is given in Section II-A. [8] chooses T depending on the signal class of interest in such a way that it yields a sparse representation. In principle, [1–12] all use some clustering algorithm for estimating the mixing parameters, although there are several different approaches to demixing. [1, 3, 4, 6, 7, 9–11] use a labeling scheme based on the estimated mixing parameters and thus demix in the above described way by creating binary masks in the transform domain corresponding to each source. That is, given the mixtures x_1 and x_2 , demixing is done by grouping the clusters of points in (Tx_1, Tx_2) space, although different techniques are used to detect these clusters. For example, [4, 6, 7, 9, 10] demix essentially by constructing binary time-frequency masks that partition the time-frequency plane such that each partition corresponds to the time-frequency points that “belong” to a particular source. The fact that such a mask exists has been observed also in [14] in the context of BSS of speech signals from *one* mixture, and in [15] in the context of source localization. In [2, 8, 11, 12], the demixing is done by making additional assumptions on the statistical properties of the sources and using a maximum a posteriori (MAP) estimator.

[5, 11] demix by assuming that the number of sources active in the transform domain at any given point is equal to the number of mixtures. They then demix by inverting the now non-degenerate M -by- M mixing matrices and appropriately combining the outputs. The above comparison is summarized in Table I.

TABLE I
A COMPARISON OF DEGENERATE DEMIXING METHODS USING DISJOINT REPRESENTATIONS.

mixing model	T operator	demixing
instantaneous [1–3, 5, 8, 9, 11]	sampling [1–3, 11]	masking [1, 3, 4, 6, 7, 9–11]
anechoic [4, 6, 7, 10, 12]	STFT [4–7, 9, 10, 12]	MAP [2, 8, 11, 12]
	signal dependent [8]	matrix masking [5, 11]

In this paper, for the linear transform T , we use the short-time Fourier transform (STFT) and Gabor expansions (the discrete version of the STFT) of speech signals. We present extensive empirical evidence that speech signals indeed satisfy (ii) in an approximate sense when T is the STFT with an appropriate window function. Based on this, we extend the Degenerate Unmixing Estimation Technique (DUET), originally presented in [4] for sources with disjointly supported STFTs, to anechoic mixtures of speech signals. The algorithm we propose relies on estimating the mixing parameters via maximum likelihood motivated estimators and constructing binary time-frequency masks using these estimates. Thus the method presented here: (1) uses an anechoic mixing model, (2) uses the STFT as T , and (3) performs demixing via masking.

In Section II we introduce a way of measuring the degree of “approximate” W-disjoint orthogonality, WDO_M , of a signal in a given mixture for a given mask M . We construct a family of time-frequency masks, Φ^x , that correspond to the indicator functions of the time-frequency points in which one source dominates the others by x dB. We test the demixing performance of these masks experimentally and illustrate that WDO_{Φ^x} is indeed a good measure of the demixing performance of the masks Φ^x . The results show that binary time-frequency masks exist that are capable of demixing several speech signals from just a single mixture. At present, there is no known robust technique for determining these masks blindly from *one* mixture. However, in Section III we derive a technique that given a *second* anechoic mixture can approximate these demixing masks blindly. We first derive the maximum likelihood estimators for the delay and attenuation coefficients. We then compare the performance of these with other estimators motivated by the maximum likelihood estimators. The modified delay and attenuation estimators are weighted averages of the instantaneous time-frequency delay and attenuation estimates. We combine the delay and attenuation estimators and show that a weighted two-dimensional histogram can be used to enumerate the sources, determine the mixing parameters, and demix the sources. The number of peaks in the histogram is the number of sources, the peak locations reveal the mixing parameters, and the mixing parameters can be used to partition the time-frequency representation of one of the mixtures to obtain estimates of the original sources. In Section IV, we verify the

method presenting demixing results for speech signals mixed synthetically and in both anechoic and echoic rooms.

II. W-DISJOINT ORTHOGONALITY

In this section, we focus on showing that binary time-frequency masks exist which are capable of separating multiple speech signals from one mixture. Our goal is, given a mixture

$$x_1(t) := \sum_{j=1}^N s_j(t) \quad (4)$$

of sources $s_j(t)$, $j = 1, \dots, N$, to recover the original sources. In order to accomplish this, we exploit the fact that the sources are pairwise approximately W-disjoint orthogonal. In this section we will define a qualitative measure of W-disjoint orthogonality and relate this measure to demixing performance.

We call two functions s_1 and s_2 **W-disjoint orthogonal** (W-DO) if, for a given a window function W , the supports of the short-time Fourier transforms (STFTs) of s_1 and s_2 are disjoint [4]. The STFT of s_j is defined

$$F^W[s_j](\tau, \omega) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} W(t - \tau) s_j(t) e^{-i\omega t} dt \quad (5)$$

which we will refer to as $\hat{s}_j(\tau, \omega)$. For a detailed discussion of the properties of this transform consult [16]. The W-disjoint orthogonality assumption can be stated concisely

$$\hat{s}_1(\tau, \omega) \hat{s}_2(\tau, \omega) = 0, \forall \tau, \omega. \quad (6)$$

The two limiting cases for W , namely $W = 1$ and $W(t) = \delta(t)$, result in interesting sets of W-DO signals. In the $W = 1$ case, the τ argument in (6) is irrelevant because the windowed Fourier transform is simply the Fourier transform. In this case, the condition is satisfied by signals which are frequency disjoint, such as frequency division multiplexed signals. In the other extreme, when $W(t) = \delta(t)$, signals which are time disjoint such as time-division multiplexed signals satisfy the condition. For window functions which are well localized in time and frequency, the W-disjoint orthogonality condition leads to signals such as those used in frequency-hopped multiple access systems [17]. Indeed, the method presented here could be applied to time domain multiplexed, frequency domain multiplexed, or frequency-hopped multiple access signals; however, in this paper we exclusively consider speech signals.

Unfortunately, (6) will not be satisfied for simultaneous speech signals because the time-frequency representation of active speech is rarely zero. However, speech is sparse in that a small percentage of the time-frequency coefficients in the Gabor expansion of speech capture a large percentage of the overall energy. In other words, the magnitude of the Gabor coefficients of speech is often small. For different speech signals, it is unlikely that the large Gabor coefficients will coincide, which leads to the signals being W-disjoint orthogonality in an approximate sense. The goal of this section is to show that speech signals satisfy a weakened version of (6) and are thus approximately W-DO. The higher the degree of approximate

W-disjoint orthogonality, the better separation results are possible. Figure 1 illustrates that speech signals have sparse time-frequency representations and satisfy a weakened version of (6), in that the product of their time-frequency representations is almost always small. A condition similar to (6) is also considered in [18], the only difference being that the time-frequency transform used was the Wigner distribution. Signals satisfying (6) for the Wigner distribution were called “time-frequency disjoint.”

The approximate W-disjoint orthogonality of speech has been described as the “sparsity” and “disjointness” of the short-time Fourier transform of the sources [5], “when one source has large energy the other does not” and “harmonic components” which “hardly overlap” [7], “when a datapoint is large the most likely decomposition is to assume that it belongs to a single source” [12], “spectra [that] are non-overlapping” [14], and “useful” time-frequency points containing a “contribution of one speaker...significantly higher than the energy of the other speaker” [19]. A quantitative measure of approximate W-disjoint orthogonality is discussed later in this section.

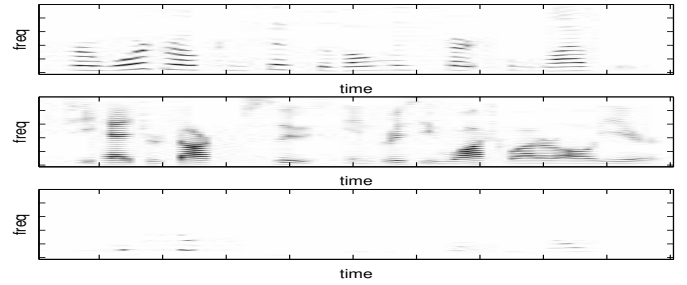


Fig. 1. A picture of W-disjoint orthogonality. The three figures are grayscale images of $|\hat{s}_1(\tau, \omega)|$ (top), $|\hat{s}_2(\tau, \omega)|$ (middle), and $|\hat{s}_1(\tau, \omega) \hat{s}_2(\tau, \omega)|$ (bottom) for two speech signals $s_1(t)$ and $s_2(t)$, sampled at 16 kHz, normalized to have unit energy. A Hamming window of length 64 ms was used as $W(t)$ and all signals had a length of three seconds. $|\hat{s}_1(\tau, \omega) \hat{s}_2(\tau, \omega)|$ contains fewer large components than $|\hat{s}_1(\tau, \omega)|$ or $|\hat{s}_2(\tau, \omega)|$. Further analysis of these signals reveals that the time-frequency points that contain 90% of the energy of s_1 contain only 1.1% of the energy of s_2 . Similarly, the time-frequency points that contain 90% of the energy of s_2 contain only 0.6% of the energy of s_1 . Thus, these speech signals approximately satisfy the W-disjoint orthogonality condition.

We can rewrite the model from (4) in the time-frequency domain

$$\hat{x}_1(\tau, \omega) = \hat{s}_1(\tau, \omega) + \dots + \hat{s}_N(\tau, \omega). \quad (7)$$

Assuming the sources are pairwise W-DO, at most one of the N sources will be non-zero for a given (τ, ω) , and thus

$$\hat{x}_1(\tau, \omega) = \hat{s}_{J(\tau, \omega)}(\tau, \omega) \quad (8)$$

where $J(\tau, \omega)$ is the index of the source active at (τ, ω) . To demix, one creates the time-frequency mask corresponding to each source and applies each mask to the mixture to produce the original source time-frequency representations. For example, defining

$$M_j(\tau, \omega) := \begin{cases} 1 & \hat{s}_j(\tau, \omega) \neq 0 \\ 0 & \text{otherwise,} \end{cases} \quad (9)$$

the indicator function for the support of s_j , one obtains the time-frequency representation of s_j from the mixture via

$$\hat{s}_j(\tau, \omega) = M_j(\tau, \omega) \hat{x}_1(\tau, \omega), \forall \tau, \omega. \quad (10)$$

A. Measuring the W-Disjoint Orthogonality of Speech

Clearly, the W-disjoint orthogonality assumption is not strictly satisfied for our signals of interest. We introduce here a measure of approximate W-disjoint orthogonality based on the demixing performance of time-frequency masks created using knowledge of the instantaneous source and interference time-frequency powers. In order to measure W-disjoint orthogonality for a given mask, we combine two important performance criteria: (1) how well the mask preserves the source of interest, and (2) how well the mask suppresses the interfering sources. These two criteria, the preserved-signal ratio (PSR) and the signal-to-interference ratio (SIR), are introduced below.

First, given a time-frequency mask M such that $0 \leq M(\tau, \omega) \leq 1$ for all (τ, ω) , we define PSR_M , the PSR of the mask M , as

$$\text{PSR}_M := \frac{\|M(\tau, \omega) \hat{s}_j(\tau, \omega)\|^2}{\|\hat{s}_j(\tau, \omega)\|^2} \quad (11)$$

which is the portion of energy of the j th source remaining after demixing using the mask. Note that $\text{PSR}_M \leq 1$ with $\text{PSR}_M = 1$ only if $\text{supp } M_j \subseteq \text{supp } M$. Here $\|f(x, y)\|^2 := \iint |f(x, y)|^2 dx dy$. Now, we define

$$y_j(t) := \sum_{\substack{k=1 \\ j \neq k}}^N s_k(t) \quad (12)$$

so that $y_j(t)$ is the summation of the sources interfering with the j th source. Then, we define the signal-to-interference ratio of time-frequency mask $M(\tau, \omega)$

$$\text{SIR}_M := \frac{\|M(\tau, \omega) \hat{s}_j(\tau, \omega)\|^2}{\|M(\tau, \omega) \hat{y}_j(\tau, \omega)\|^2} \quad (13)$$

which is the output signal-to-interference ratio after using the mask to demix.

We now combine the PSR_M and SIR_M into one measure of approximate W-disjoint orthogonality. We propose the normalized difference between the signal energy maintained in masking and the interference energy maintained in masking as a measure of the W-disjoint orthogonality associated with a particular mask:

$$\text{WDO}_M := \frac{\|M(\tau, \omega) \hat{s}_j(\tau, \omega)\|^2 - \|M(\tau, \omega) \hat{y}_j(\tau, \omega)\|^2}{\|\hat{s}_j(\tau, \omega)\|^2} \quad (14)$$

$$= \text{PSR}_M - \text{PSR}_M / \text{SIR}_M. \quad (15)$$

For signals which are W-DO, using the mask $M_j(\tau, \omega)$ defined in (9), we note that $\text{PSR}_{M_j} = 1$, $\text{SIR}_{M_j} = \infty$, and $\text{WDO}_{M_j} = 1$. This is the maximum obtainable WDO value because $\text{WDO}_M \leq 1$ for all M such that $0 \leq M(\tau, \omega) \leq 1$. Moreover, for any M , $\text{WDO}_M = 1$ implies that $\text{PSR}_M = 1$, $\text{SIR}_M = \infty$, and that (6) is satisfied. That is, $\text{WDO}_M = 1$ implies that the signals are W-DO and that mask M perfectly separates the j th source from the mixture. In order for a mask to have $\text{WDO}_M \approx 1$, i.e., good demixing performance, it must simultaneously preserve the energy of the signal of interest while suppressing the energy of the interference. The failure of a mask to accomplish either of the goals can result in a small,

even negative, WDO value. For example, $\text{WDO}_M = 0$ implies either that $\text{PSR}_M = 0$ (the mask kills all the energy of the source of interest) or that $\text{SIR}_M = 1$ (the mask results in equal energy for source and interference). Masks with $\text{SIR}_M < 1$ have associated $\text{WDO}_M < 0$.

Now we establish that binary time-frequency masks exist which are capable of demixing speech signals from one mixture and detail their performance in relation to the three presented measures. Consider the following family of time-frequency masks

$$\Phi_j^x(\tau, \omega) := \begin{cases} 1 & 20 \log(|\hat{s}_j(\tau, \omega)| / |\hat{y}_j(\tau, \omega)|) \geq x \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

which is the indicator function for the time-frequency points where s_j dominates the interference in the mixture by x dB. We will use $\text{PSR}_j(x)$ and $\text{SIR}_j(x)$ as shorthand for $\text{PSR}_{\Phi_j^x}$ and $\text{SIR}_{\Phi_j^x}$, respectively.

To determine the demixing ability of the above mask type, the masks for various x were applied to speech mixtures of various orders and the demixing performance measures, $\text{PSR}_j(x)$ and $\text{SIR}_j(x)$, were determined. We refer to a mixture of N sources as a *mixture of order N* and the mixtures used in these tests had orders $N = 2, 3, \dots, 10$. The demixed speech was then rated by the authors as falling into one of five subjective categories. The speech signals were selected from 16 male and 16 female continuous speech segments of three seconds taken from the TIMIT database and normalized to unit energy. The time-frequency representation of the 16kHz sampled data was created using a Hamming window of 1024 samples with 50% overlap. The results of the 333 listening tests are displayed in Figure 2. We note that there is a fairly accurate relationship between the WDO performance measure and the subjective ratings listed in the table under the figure.

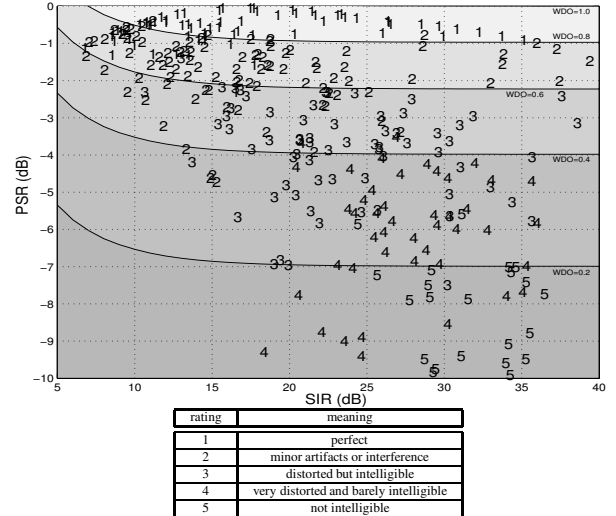


Fig. 2. Results of subjective listening test performed by the authors. For example, $0.8 > \text{WDO} \geq 0.6$ implies a “minor artifacts or interference” rating or better

Now that we have some idea how PSR, SIR, and WDO map to demixing performance, we analyze the demixing performance of the masks described in (16). Figure 3 shows plots of $\text{PSR}_j(x)$ versus $\text{SIR}_j(x)$ and a table of $(\text{PSR}_j(x), \text{SIR}_j(x))$

pairs averaged for groups of speech mixtures of different orders. For $N = 2$, each source was compared against each of the remaining 31 sources, resulting in $32 \times 31 = 992$ tests being averaged for each data point. For larger N , each source was compared against a random mixing of $N - 1$ of the remaining 31 sources. This was done 31 times per source in order to keep the number of tests per data point constant at 992. As we tested mixtures from $N = 2$ to $N = 10$, a total of $9 \times 992 = 8928$ mixtures were created to generate the data for Figure 3. Figure 3 demonstrates that time-frequency masks exist which exhibit excellent demixing performance; For example, considering the 0 dB mask Φ_j^0 , we see that on average this mask produces demixtures with WDO measure greater than 0.6 for mixtures of up to ten sources.

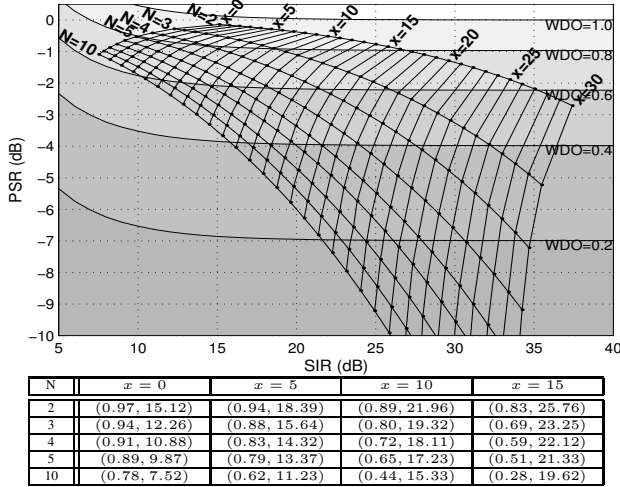


Fig. 3. Time-Frequency Mask Demixing Performance. Plot contains $\text{PSR}_j(x)$ (in dB) versus $\text{SIR}_j(x)$ (in dB) for $x = 0, 1, \dots, 30$ for $N = 1, 2, \dots, 10$. Table contains $(\text{PSR}_j(x), \text{SIR}_j(x))$ (in dB) for $N = 2, 3, 4, 5, 10$ for $x = 0, 5, 10, 15$ dB. The different gray regions correspond to different regions of approximate W-disjoint orthogonality as determined by the lines of constant WDO. For example, using the $x = 5$ dB mask in mixtures of four sources yields 14.32 dB output SIR while maintaining 83% of the desired source energy. This $(\text{PSR}, \text{SIR}) = (0.83, 14.32)$ dB pair results in $\text{WDO} = 0.80$, which from Figure 2 implies perfect demixing performance. In other words, if we can correctly map time-frequency points with 5 dB or more single source dominance to the correct corresponding output partition, we can recover 83% of the energy of each of the original sources and produce demixtures with 14.32 dB output SIR from a mixture of four sources.

Now that we know that good time-frequency masks exist, we wish to determine the dependence of these performance measures on the window function $W(t)$ and window size. For this task, we examine the performance of the 0 dB mask, Φ_j^0 . Figure 4 shows PSR, SIR, and WDO for pairwise mixing for various window sizes and types. Each data point in the figure represents the average of the results for 992 mixtures. In all measures, the Hamming window of size 1024 samples performed the best. Note, however, that the performance of the other masks (with the exception of the rectangle) was extremely similar and exhibited better than 90% W-disjoint orthogonality for pairwise mixing across a wide range of window sizes (from roughly 500 to 4000 samples). Other mixture orders and masks (i.e., Φ_j^x for $x > 0$) exhibited similar performance and in all cases the Hamming window of size 1024 had the best performance. A similar conclusion regarding the optimal time-

frequency resolution of a window for speech separation was arrived at in [7]. Note that even when the window size is 1 (i.e., T is sampling), the mixtures still exhibit a high level of PSR, SIR, and WDO. This fact was exploited by those methods that used the time-disjoint nature of speech [1–3, 11]. However, Figure 4 clearly shows the advantage of moving from the time domain to the time-frequency domain: the speech signals are more disjoint in the time-frequency domain provided the window size is sufficiently large. Choosing the window size too large, however, results in reduced W-disjoint orthogonality.

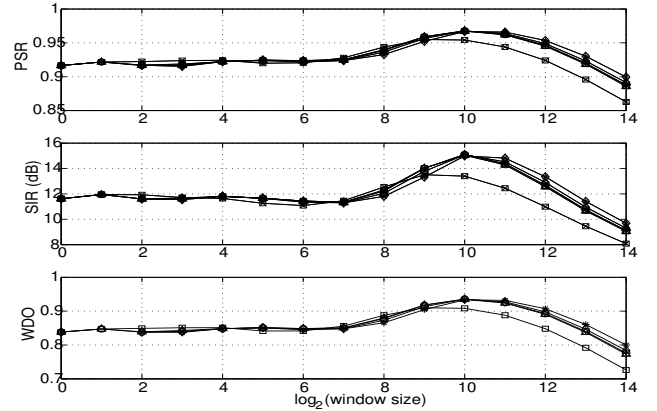


Fig. 4. Window size and type comparison. Hamming (\circ), Blackman ($*$), Hann (\diamond), Triangle (\triangle), and Rectangle (\square). PSR, SIR, and WDO for the 0 dB mask for window size = 1, 2, 4, \dots , 16384 samples for various window types for pairwise mixing of speech signals sampled at 16 kHz. The Hamming window of size 1024 has the best performance.

We close this section by proposing WDO_M with $M = \Phi_j^0$ as the general measure of W-disjoint orthogonality. Table II shows $\text{WDO}_{\Phi_j^0}$ values for mixtures of various orders. Again, each data point represents the average measurement over 992 mixtures. It can be shown using (14) that the 0 dB mask, Φ_j^0 , maximizes WDO, and thus the 0 dB mask line represents the upper bound of WDO for any mask. We thus say that, for example, *speech signals in pairwise mixtures are 93.6% W-disjoint orthogonal*.

TABLE II
WDO FOR THE 0 dB MASK FOR MIXTURES OF VARIOUS ORDERS.

N	2	3	4	5	6	7	8	9	10
% WDO	93.6	88.0	83.4	79.2	75.6	72.3	69.3	66.6	64.0

III. PARAMETER ESTIMATION AND DEMIXING

In this section, we will present a demixing algorithm that separates an arbitrary number of sources using two mixtures. We start by describing our anechoic mixing model. Suppose we have N sources $s_1(t), \dots, s_N(t)$. Let $x_1(t)$ and $x_2(t)$ be the mixtures such that

$$x_k(t) = \sum_{j=1}^N a_{kj} s_j(t - \delta_{kj}), \quad k = 1, 2 \quad (17)$$

where parameters a_{kj} and δ_{kj} are the attenuation coefficients and the time delays associated with the path from the j^{th} source to the k^{th} receiver. Without loss of generality we set $a_{1j} = 1$

and $\delta_{1j} = 0$ for $j = 1, \dots, N$; for simplicity we rename a_{2j} as a_j and δ_{2j} as δ_j . In addition we assume that the windowed Fourier transform of any source function, $F^W[s_j](\tau, \omega)$ satisfies the narrowband assumption for array processing, i.e.,

$$\begin{aligned} F^W[s_j(\cdot - \delta)](\tau, \omega) &= \exp(-i\omega\delta)F^W[s_j](\tau - \delta, \omega) \\ &\approx \exp(-i\omega\delta)F^W[s_j](\tau, \omega). \end{aligned} \quad (18)$$

This assumption is realistic as long as the window function W is chosen appropriately. A detailed discussion about this assumption can be found in [20].

Now we go back to discussing the mixing model, described in (17). We take the STFT of x_1 and x_2 with an appropriate choice of W . Using the assumption discussed above, the mixing model (17) reduces to

$$\begin{bmatrix} \hat{x}_1(\tau, \omega) \\ \hat{x}_2(\tau, \omega) \end{bmatrix} = \begin{bmatrix} 1 & \dots & 1 \\ a_1 e^{-i\omega\delta_1} & \dots & a_N e^{-i\omega\delta_N} \end{bmatrix} \begin{bmatrix} \hat{s}_1(\tau, \omega) \\ \vdots \\ \hat{s}_N(\tau, \omega) \end{bmatrix}. \quad (19)$$

A. Parameter Estimation and Demixing for W-DO sources

To motivate the Degenerate Unmixing Estimation Technique (DUET), which we will describe in the next section, we first consider the case where the sources are W-DO, i.e.,

$$\hat{s}_j(\tau, \omega)\hat{s}_k(\tau, \omega) = 0, \quad \forall(\tau, \omega), \quad \forall j \neq k. \quad (20)$$

This condition is the idealization of the properties of speech signals discussed in Section II. We now construct the parameter estimators and the demixing algorithm for W-DO signals. Clearly, when the sources are W-DO, at most one source will be active at any time-frequency point (τ, ω) ; in particular for any (τ, ω) at which $\hat{x}_1(\tau, \omega) \neq 0$, there exists a j such that $\hat{s}_j(\tau, \omega) \neq 0$ and $\hat{s}_k(\tau, \omega) = 0$ for $j \neq k$. Recalling the definition of the time-frequency mask M_j in (10), we note that $M_j(\tau, \omega)M_k(\tau, \omega) = 0$ for all (τ, ω) if $j \neq k$. From (19) we deduce that

$$\hat{s}_j = M_j \hat{x}_1. \quad (21)$$

This shows that we can demix an arbitrary number of sources from only one of the mixtures if we can construct the corresponding mask M_j for each source. Next we will describe how to construct the masks M_j using the mixtures x_1 and x_2 .

Let j be arbitrary, and define $\Omega_j := \{(\tau, \omega) : M_j(\tau, \omega) = 1\}$ so that $M_j = 1_{\Omega_j}$. Note that the Ω_j are pairwise disjoint. Now consider

$$R_{21}(\tau, \omega) := \frac{\hat{x}_2(\tau, \omega)}{\hat{x}_1(\tau, \omega)}. \quad (22)$$

Clearly, on Ω_j

$$R_{21}(\tau, \omega) = a_j e^{-i\delta_j \omega}. \quad (23)$$

In this case $|R_{21}(\tau, \omega)| = a_j$ and $-\frac{1}{\omega} \angle R_{21}(\tau, \omega) = \delta_j$, where $\angle z$ denotes the phase of the complex number z taken between $-\pi$ and π .

The observation above yields a way of constructing the sets Ω_j and thus a demixing algorithm: we simply label each time-frequency point (τ, ω) with the pair $(|R_{21}(\tau, \omega)|, -\frac{1}{\omega} \angle R_{21}(\tau, \omega))$. Since the sources are W-DO, there will be N distinct labels. By grouping the time-frequency

points (τ, ω) with the same label, we construct the sets Ω_j , thus the masks $M_j = 1_{\Omega_j}$.

The above described demixing algorithm is the motivation behind DUET. Note that the algorithm separates the sources without inverting the mixing matrix, which makes it possible to deal with mixtures of an arbitrary number of sources. Aside from demixing, it also yields the mixing parameters: the labels $(|R_{21}(\tau, \omega)|, -\frac{1}{\omega} \angle R_{21}(\tau, \omega))$ which we used to construct the masks are exactly the mixing parameters a_j and δ_j . Motivated by this fact we define the instantaneous DUET attenuation and delay parameter estimators as

$$\tilde{a}(\tau, \omega) := |R_{21}(\tau, \omega)| \quad (24)$$

$$\tilde{\delta}(\tau, \omega) := -\frac{1}{\omega} \angle R_{21}(\tau, \omega) \quad (25)$$

respectively. We will use these estimators in the next section.

In summary, the **DUET algorithm for demixing W-DO sources** is,

- 1) From mixtures $x_1(t)$ and $x_2(t)$ construct time-frequency representations $\hat{x}_1(\tau, \omega)$ and $\hat{x}_2(\tau, \omega)$.
- 2) For each non-zero time-frequency point, calculate $(\tilde{a}(\tau, \omega), \tilde{\delta}(\tau, \omega))$.
- 3) Take the union of the $(\tilde{a}(\tau, \omega), \tilde{\delta}(\tau, \omega))$ pairs, $S = \bigcup_{(\tau, \omega)} \{(\tilde{a}(\tau, \omega), \tilde{\delta}(\tau, \omega))\}$. Note S will be equal to $\{(a_j, \delta_j) : j = 1, \dots, N\}$.
- 4) For each (a_j, δ_j) in S , $j = 1, \dots, N$, note that $\hat{s}_j(\tau, \omega) = 1_{\{(\tilde{a}(\tau, \omega), \tilde{\delta}(\tau, \omega)) = (a_j, \delta_j)\}}(\tau, \omega) \hat{x}_1(\tau, \omega)$ for (τ, ω) with $\hat{x}_1(\tau, \omega) \neq 0$ and $\hat{s}_j(\tau, \omega) = 0$ otherwise. Clearly, $\hat{s}_j(\tau, \omega)$ will be the time-frequency representations of one of the original sources. The numbering of the sources is arbitrary.
- 5) Convert each $\hat{s}_j(\tau, \omega)$ back into the time domain.

Remark 1: Note that the instantaneous DUET delay estimator yields a meaningful estimate at a time-frequency point $(\tau, \omega) \in \Omega_j$ only if

$$|\omega\delta_j| < \pi. \quad (26)$$

This follows from the periodicity of the complex exponential. For $(\tau, \omega) \in \Omega_j$, we have $-\frac{1}{\omega} \angle R_{21}(\tau, \omega) = r(\omega\delta_j)\delta_j$, with $r(u) := \frac{\angle u}{u}$ where $\angle u := (u + \pi) \pmod{2\pi} - \pi$. When (26) is not satisfied, the delay estimate obtained using the instantaneous DUET estimator will be a fraction of its true value. Let ω_{\max} be the maximum element of $\{|\omega| : (\omega, \tau) \in \bigcup_j \Omega_j, \text{ for some } \tau\}$, which is the maximum frequency present in the sources, and denote by ω_s the sampling rate. Let $\delta_{j\max} := \max_j |\delta_j|$. Clearly, (26) is guaranteed for all j and for all $\omega \in \bigcup_j \Omega_j$ if

$$\omega_{\max} \delta_{j\max} < \pi. \quad (27)$$

Now define $\delta_{\omega\max} := \pi/\omega_{\max}$. Any delay parameter with modulus less than $\delta_{\omega\max}$ can be estimated correctly. Clearly, (27) is equivalent to the condition $\delta_{j\max} < \delta_{\omega\max}$. If $\omega_{\max} = \omega_s/2$, the Nyquist frequency, then this means that the maximum delay, $\delta_{\omega\max} = \frac{2\pi}{\omega_s}$, is exactly equal to the sampling period. In other words, as long as the delay between the two microphone readings is less than a sample, the estimated phase will be accurate. While the ω_{\max} is determined by the characteristics of speech signals, the maximum physically possible delay, which we will

denote by $\delta_{d\max}$, is determined by the microphone spacing. For two microphones separated by a distance d , $\delta_{d\max} = d/c$ where c is the speed of sound. Clearly, we have $\delta_{j\max} < \delta_{d\max}$, and therefore (27) will be satisfied if $\delta_{d\max} < \delta_{\omega\max}$. This suggests that one can guarantee (27) simply by choosing d , and thus $\delta_{d\max}$, sufficiently small. For example, for a sampling rate $\omega_s/(2\pi) = 16$ kHz, assuming $\omega_{\max} = \omega_s/2$ and $c = 344$ m/s, we obtain that $\delta_{d\max} \leq \delta_{\omega\max}$ as long as $d \leq 2.15$ cm. If we knew, however, that $\omega_{\max}/(2\pi) = 4$ kHz, then this distance would be increased by a factor of 4 to 8.60 cm. The smaller the largest frequency present in the signal, the larger the allowable microphone separation (or equivalently the larger we can choose $\delta_{j\max}$) that guarantees accurate phase parameter estimates. The demixing technique presented in this paper does not require knowledge of the value of d , but we assume that d is small enough such that (27) is satisfied.

B. Parameter Estimation and Demixing for Approximately W-DO sources

In Section II, we illustrated that the time-frequency representations of speech signals are nearly disjoint and demonstrated that we can indeed recover a speech signal from one mixture of an arbitrary number of sources if we can construct an appropriate time-frequency mask. This suggests that a weakened W-DO condition holds for speech signals: if at a time-frequency point one of the sources has considerable power, the contribution of all the other sources at that time-frequency point is likely to be small. This observation is the key to the demixing algorithm we propose in this section. First we shall discuss how to estimate the mixing parameters.

From this point on, instead of the continuous STFT, we use the equivalent discrete counterpart¹

$$\hat{s}_j[k, l] = \hat{s}_j(k\tau_0, l\omega_0) \quad (28)$$

where τ_0 and ω_0 are the time-frequency lattice spacing parameters. We define the instantaneous DUET delay estimate, the discrete version of (25),

$$\tilde{\delta}[k, l] := -\frac{1}{l\omega_0} \angle R_{21}[k, l]. \quad (29)$$

where $R_{21}[k, l] := \frac{\hat{x}_2[k, l]}{\hat{x}_1[k, l]}$. For convenience, we define $\tilde{\delta}[k, l] = 0$ if $\hat{x}_1[k, l] = 0$ or $\hat{x}_2[k, l] = 0$. Similarly, we define the instantaneous DUET attenuation estimate

$$\tilde{a}[k, l] := |R_{21}[k, l]| \quad (30)$$

which is the discrete version of (24). For convenience, we define $\tilde{a}[k, l] = 1$ if $\hat{x}_1[k, l] = 0$ or $\hat{x}_2[k, l] = 0$. For reasons discussed in Appendix I, we choose to estimate

$$\alpha_j := a_j - 1/a_j \quad (31)$$

instead of directly estimating a_j , and define the instantaneous DUET symmetric attenuation estimate

$$\tilde{\alpha}[k, l] := \tilde{a}[k, l] - 1/\tilde{a}[k, l]. \quad (32)$$

¹The equivalence is nontrivial and only true for appropriately chosen window functions W with sufficiently small τ_0 and ω_0 . An illustrative discussion can be found in [16].

We will say that s_j is *dominant* at $[k, l]$ if $|\hat{s}_j[k, l]| \geq |\hat{y}_j[k, l]|$, where \hat{y}_j is as in (12). Note, the 0 dB mask, Φ_j^0 , in (16) is the indicator function for the dominant time-frequency points of the j th source. In Appendix I, we derive that under certain assumptions the maximum likelihood (ML) estimates for the mixing parameters a_j and δ_j can be determined via certain weighted averages (see (57) and (59)) of the instantaneous DUET delay ($\tilde{\delta}[k, l]$) and instantaneous DUET attenuation ($\tilde{a}[k, l]$) estimates, at the time-frequency points at which s_j is dominant. In Appendix II we compare the performance of the weighted estimators suggested by the ML derivation with other empirically motivated weighted estimators, and we illustrate that more accurate estimates ($\alpha_j^{(p)}, \delta_j^{(p)}$) of the true parameters (α_j, δ_j) are determined by

$$\alpha_j^{(p)} = \frac{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k, l] \hat{x}_2[k, l]|^p \tilde{\alpha}[k, l]}{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k, l] \hat{x}_2[k, l]|^p} \quad (33)$$

$$\delta_j^{(p)} = \frac{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k, l] \hat{x}_2[k, l]|^p \tilde{\delta}[k, l]}{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k, l] \hat{x}_2[k, l]|^p} \quad (34)$$

when $p = 2$. To employ these estimates however, we need to first construct the sets $\Lambda_j = \{[k, l] : |\hat{s}_j(k, l)| > |\hat{y}_j(k, l)|\}$ for each j . Note that these sets would be the discrete version of Ω_j of Section III-A if the sources are W-DO. In the W-DO case we used the instantaneous DUET estimates as labels for each time-frequency point (τ, ω) , and each Ω_j consisted of the points with identical labels. In the approximately W-DO case the instantaneous DUET estimates for the time-frequency points in Λ_j will not be identical anymore. However, we claim that we can still use these estimates as a means of labeling, and thus construct the sets Λ_j , at least approximately. Once we know the Λ_j , we demix simply by partitioning the support of $\hat{x}_1[k, l]$ using Λ_j and converting the resulting time-frequency representations back into the time domain. In order to determine the Λ_j , we rely on three observations which lead us to create a smoothed two-dimensional power weighted histogram of the $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$ pairs. Enumerating the peaks in this histogram estimates the number of sources, the peak centers estimate the mixing parameters, and the set of time-frequency points which contribute to a given peak provide an estimate for the associated Λ_j .

Observation 1: *The time-frequency points with instantaneous DUET estimates $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$ inside a small rectangle centered on the true mixing parameter pair (α_j, δ_j) contain most of the source energy.*

We wish to show that the time-frequency points which yield instantaneous DUET estimates that are in close proximity to the true mixing parameters contain most of the energy of the source. Let,

$$M_{\alpha, A}[k, l] = \begin{cases} 1 & \text{if } |\tilde{\alpha}[k, l] - \alpha| < A \\ 0 & \text{otherwise} \end{cases} \quad (35)$$

be the indicator function for time-frequency points with instantaneous DUET symmetric attenuation estimate within A of α where A is a resolution parameter. We are interested in,

$$\text{PSR}_{M_{\alpha, A}} = \frac{\sum_{(k,l)} M_{\alpha, A}[k, l] |s_j[k, l]|^2}{\sum_{(k,l)} |s_j[k, l]|^2} \quad (36)$$

which will show the portion of the energy of s_j contained in time-frequency points with corresponding $\tilde{\alpha}[k, l]$ within A of

the true mixing parameter value α_j . Figure 5 shows $\text{PSR}_{M_{\alpha_j,A}}$ averaged over 100 randomly selected speech signals taken from the TIMIT database. The curves represent the expected energy contained in time-frequency points with instantaneous DUET symmetric attenuation estimates close to the true symmetric attenuation. For example, with $A = 0.1$ we expect more than 60% of the source energy to come from time-frequency points with corresponding $\tilde{\alpha}[k, l]$ located within 0.1 of the true value α_j in mixtures of five sources. In this section, the model described by (47) in Appendix I and discussed in Appendix II was used to simulate mixtures of $N = 2, 3, 5, 10$ sources.

Similarly, for the delay, we define

$$M_{\delta_j,D}[k, l] = \begin{cases} 1 & \text{if } |\tilde{\delta}[k, l] - \delta_j| < D \\ 0 & \text{otherwise} \end{cases} \quad (37)$$

where D is a resolution parameter. Then we are interested in

$$\text{PSR}_{M_{\delta_j,D}} = \frac{\sum_{(k,l)} M_{\delta_j,D}[k, l] |s_j[k, l]|^2}{\sum_{(k,l)} |s_j[k, l]|^2} \quad (38)$$

which will show the portion of energy of s_j with instantaneous DUET delay estimates within D of the true mixing parameter value δ_j . Figure 5 also shows $\text{PSR}_{M_{\delta_j,D}}$ as a function of D for various mixture orders. For example, 70% of the energy of the source is expected to be contained in time-frequency points with corresponding $\tilde{\delta}[k, l]$ within 0.1 samples of the true value of δ_j in pairwise mixing.

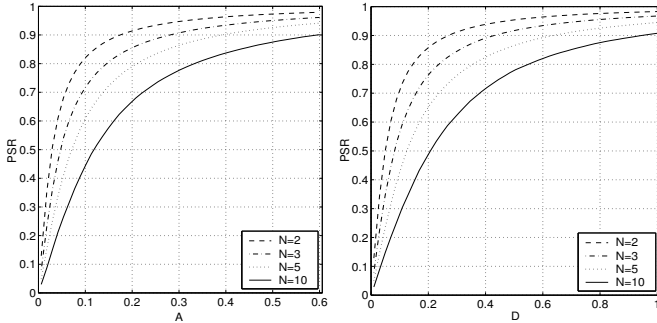


Fig. 5. Energy distribution $\text{PSR}_{M_{\alpha_j,A}}$ (left) and $\text{PSR}_{M_{\delta_j,D}}$ (right) of instantaneous DUET estimates around the true mixing parameters. Note, for W-DO signals, the corresponding source energy portion would be 1.0 for all distances from α_j (and δ_j).

We now show that the source energy is localized simultaneously around (α_j, δ) . To do so, we look at

$$\text{PSR}_{M_{\alpha_j,A} M_{\delta_j,D}} = \frac{\sum_{(k,l)} M_{\alpha_j,A}[k, l] M_{\delta_j,D}[k, l] |s_j[k, l]|^2}{\sum_{(k,l)} |s_j[k, l]|^2} \quad (39)$$

which measures the portion of source energy for time-frequency points with instantaneous DUET symmetric attenuation estimates within A of α_j and instantaneous DUET delay estimates within D of δ_j . Before we examine $\text{PSR}_{M_{\alpha_j,A} M_{\delta_j,D}}$, we need to determine the appropriate A to D ratio. Plotting (A, D) pairs for $\text{PSR}_{M_{\alpha_j,A}} = \text{PSR}_{M_{\delta_j,D}}$ for the same mixture order reveals that the (A, D) lie essentially

along a line. The least-mean-square fit of this line determines a ratio of $A/D = 1/1.65$ samples. This means that, for example, $\text{PSR}_{M_{\alpha_j,A}} \approx 0.6$ for $A = 0.1$ for $N = 5$ implies that $\text{PSR}_{M_{\delta_j,D}} \approx 0.6$ for $D = 1.65 \times 0.1 = 0.165$ samples for $N = 5$, a property which can be verified from the data displayed in Figure 5. Figure 6 shows $\text{PSR}_{M_{\alpha_j,A} M_{\delta_j,D}}$ versus A and D for $1.65A = D$. Note that the A axis is at the bottom and the D axis is at the top. For example, almost 80% of the energy of the source is contained by time-frequency points with corresponding $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$ falling in a rectangle with dimensions 0.2-by-0.33 centered on (α_j, δ_j) , i.e., $[\alpha_j - 0.1, \alpha_j + 0.1] \times [\delta_j - 0.165, \delta_j + 0.165]$, for mixtures of three sources. As the number of sources increases, the energy spreads over a wider area, but remains relatively well localized around the source's mixing parameters.

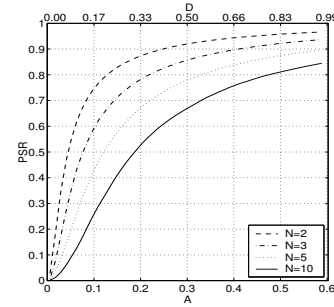


Fig. 6. Energy distribution $\text{PSR}_{M_{\alpha_j,A} M_{\delta_j,D}}$ of instantaneous DUET estimates in a rectangle centered on the true mixing parameter pair (α_j, δ_j) .

Observation 2: *Observation 1 is true for the individual sources in mixtures.*

This observation is based on the fact that, from the experiments with speech mixtures (see Figure 3), we know that the time-frequency points at which one source dominates maintain a significant percentage of the dominating source's energy. For $N = 2, 3, 4, 5$, and 10 the percentage source energy preserved when only considering dominant time-frequency points is 97%, 94%, 91%, 89%, and 78%, respectively. Considering the time-frequency points when one source dominates, Figure 6 shows that the instantaneous DUET estimates falling in a rectangle of dimension 0.2-by-0.33 centered on the true mixing parameter pair. Thus, in pairwise mixing we would expect the time-frequency points which yield estimates $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$ inside a 0.2-by-0.33 rectangle centered on (α_1, δ_1) to contain the product of 87% (from Figure 6) and 97% (from Figure 3) for a total of $87\% \times 97\% = 84\%$ of the energy of the first source. Similarly, we would expect 84% of the energy of the second source to come from time-frequency points which have instantaneous DUET estimate pairs within a 0.2-by-0.33 rectangle centered on (α_2, δ_2) . As N increases, the source energy percentage we expect to see in a fixed size rectangle centered on each source's mixing parameters decreases (it is 39% for $N=10$); nevertheless, Observation 1 will still hold.

Observation 3: *The peaks in a smoothed two-dimensional*

power weighted histogram of the instantaneous DUET estimates will be in one-to-one correspondence with the rectangle centers in Observation 2.

One way to determine the mixing parameters for multiple sources is to look at the two-dimensional weighted histogram $\sum_{k,l} M_{\alpha,A}[k,l] M_{\delta,D}[k,l] |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p$ as a function of α and δ , for some fixed p . If (A, D) is chosen large enough to capture a large portion of the source energy, as determined by Figure 6, yet small enough so the (A, D) rectangle does not contain significant energy contributions from multiple sources, we would expect the local maxima to occur around the true mixing parameter pairs (α_j, δ_j) , $j = 1, \dots, N$. Therefore, one way of determining the mixing parameters would be to calculate $\sum_{k,l} M_{\alpha,A}[k,l] M_{\delta,D}[k,l] |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p$ for the range of interest of (α, δ) pairs and select the local maxima. A computationally efficient way of doing this is to construct a two-dimensional weighted histogram at a high resolution, and then smooth that histogram with a kernel of the dimensions of the desired (A, D) rectangle. We perform smoothing to group the time-frequency points which are likely to correspond to one source. Recall that the estimators (33) and (34) averaged the instantaneous estimates over all time-frequency points where the source of interest was dominant. We know from the results shown in Figure 6 that a rectangle centered on the true mixing parameters will capture most of the corresponding source's energy. By smoothing, we locate the rectangle centers that capture locally the largest energy contribution, and thus estimate the mixing parameters. Histograms have been used previously for parameter estimation of voice mixtures; for example, [21] clusters onset arrival difference to determine the time delays of the various sources.

Now we construct a two dimensional weighted histogram for $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$, where $\tilde{\alpha}[k, l]$ and $\tilde{\delta}[k, l]$ are the instantaneous DUET estimates, with the weights $|\hat{x}_1[k, l] \hat{x}_2[k, l]|^p$ for some p . The weighted histogram with resolution widths β and Δ and weighting exponent p , is defined as

$$h(\alpha, \delta) := \sum_{k,l} M_{\alpha,\beta/2}[k,l] M_{\delta,\Delta/2}[k,l] |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p. \quad (40)$$

which we will smooth with a rectangular kernel $r(\alpha, \delta)$,

$$r(\alpha, \delta) := \begin{cases} 1/AD & (\alpha, \delta) \in [-A/2, A/2] \times [-D/2, D/2] \\ 0 & \text{otherwise,} \end{cases} \quad (41)$$

to produce the smoothed histogram

$$H(\alpha, \delta) := [h * r](\alpha, \delta) \quad (42)$$

where $*$ denotes two-dimensional convolution.

Figure 7 shows an example histogram before and after smoothing generated using the dominant time-frequency points of a speech signal generated using the model of five source mixing. For the mixing model, $(\alpha_j, \delta_j) = (0.2, 0.5)$ which match well with the peak location. The importance of the smoothing is clear in that it combines all the energy from the estimates in a local region and results in a clear single peak and thus mixing parameter estimate. In Appendix III, we compare the performance of histogram-based parameter estimators to the performance of the ML estimators and discuss the choice of p .

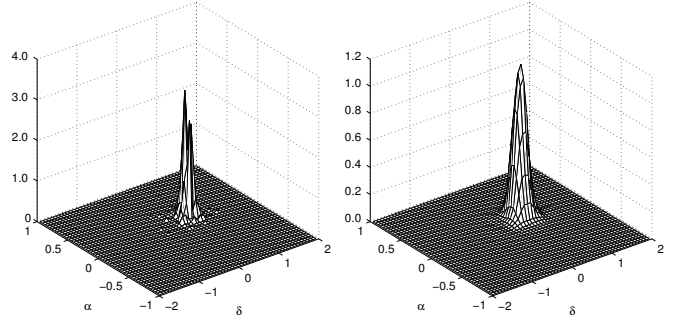


Fig. 7. Example raw (left) and smoothed (right) power weighted ($p = 2$) histograms for one speech signal in a mixture of five. The peak location of the smoothed histogram corresponds to the mixing parameters $(\alpha_j, \delta_j) = (0.2, 0.5)$.

C. Demixing Algorithm for Approximately W-DO Sources

Recall that in the W-DO case, sources were demixed using time-frequency masks that were constructed by grouping the time-frequency points that yield the same instantaneous parameter estimates. We demix in a similar way for approximately W-DO sources. First we estimate the mixing parameters, for example, using the histogram method described in the previous section. Then, we group time-frequency points that yield instantaneous parameter estimates that are “close” to these estimated mixing parameters. One natural definition of closeness is the instantaneous likelihood function for the j th source

$$L_j[k, l] := p(\hat{x}_1[k, l], \hat{x}_2[k, l] | a_j, \delta_j) = \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2} |a_j e^{-i\delta_j l \omega_0} \hat{x}_1[k, l] - \hat{x}_2[k, l]|^2 / (1+a_j^2)} \quad (43)$$

obtained by substituting the instantaneous ML source estimate (53) into the likelihood function in (48) modified to consider only time-frequency point $[k, l]$. $L_j[k, l]$ is, in a sense, the likelihood that the j th source is dominant at time-frequency point $[k, l]$. One way to demix the mixtures is to construct a time frequency mask for s_j by taking those time-frequency points for which $L_j[k, l] \geq L_i[k, l]$, $\forall i \neq j$. The time-frequency mask for demixing s_j is thus

$$\tilde{M}_j := 1_{\{[k,l]: j=\arg \max_m L_m[k,l]\}} \quad (44)$$

and defining

$$\tilde{\Lambda}_j = \{[k, l] : \arg \max_m L_m[k, l]\} \quad (45)$$

the estimate of the time-frequency points for which the j th source is dominant, we can relate this demixing mask to those that were used in the W-DO case. There are many other ways we can envision using these likelihoods, for example, some type of relative weighting resulting in fractional masks instead of the binary winner-take-all masks created by the scheme we have proposed. However, we have shown in Section II that the 0-dB binary masks exhibit excellent demixing performance and maximize the WDO performance measure so we consider exclusively binary time-frequency masks in this paper.

As before, we estimate the source by converting

$$\tilde{s}_j[k, l] := \tilde{M}_j[k, l] \hat{x}_1[k, l] \quad (46)$$

into the time domain. Note, we could apply the mask to x_2 as well, and, could combine the two demixtures using the ML estimate of the source as in (53). However, in order to compare with the results obtained in Section II, the experimental results presented in the next section will use (46).

In summary, the **DUET algorithm for demixing Approximately W-DO sources** is,

- 1) From mixtures $x_1(t)$ and $x_2(t)$ construct time-frequency representations $\hat{x}_1[k, l]$ and $\hat{x}_2[k, l]$.
- 2) For each time-frequency point, calculate $(\tilde{\alpha}[k, l], \tilde{\delta}[k, l])$ using (32) and (29).
- 3) Construct histogram and locate peaks:
 - a) Construct a high resolution histogram as in (40)
 - b) Smooth the histogram as in (42)
 - c) Locate peaks in histogram. There will be N peaks, one for each source, with peak locations approximately equal to the true mixing parameter pairs, $\{(\alpha_j, \delta_j) : j = 1, \dots, N\}$.
- 4) For the N pairs of (α_j, δ_j) estimates, construct the time-frequency masks corresponding to each pair using the ML partitioning as in (44) and apply these masks to one of the mixtures as in (46) to yield estimates of the time-frequency representations of the original sources.
- 5) Convert each estimate back into the time domain.

IV. EXPERIMENTS

In order to demonstrate the technique, we present results in this section for both synthetic and real mixtures. One issue that we have not addressed is how the histogram peaks are automatically enumerated and identified. For the following demonstration, we used an ad-hoc technique that iteratively selected the highest peak and removed a region surrounding the peak from the histogram. Peaks were removed as long as the histogram maintained a threshold percentage of its original weight. The threshold percentage and region dimensions had to be occasionally altered in the course of the tests to ensure the correct number of sources was found. Indeed, peak enumeration and identification remains a topic of future research. In all examples, we used histograms with $p = 1$, as suggested in Appendix III.

A. Synthetic mixtures

Figure 8 shows the smoothed histogram (42) for a six source synthetic mixing example with histogram resolution widths $(\beta, \Delta) = (0.05, 0.12)$ samples and smoothing kernel dimensions $(A, D) = (0.12, 0.2)$ samples. The six sources were taken from the TIMIT database and the stereo mixture was created using (symmetric attenuation, delay) mixing parameters pairs $\{(1, -2), (3/2, -1), (3/2, 1), (1, 2), (2/3, 1), (2/3, -1)\}$. It is clear given only the stereo mixture, one can determine how many sources were used to create the mixture by enumerating the peaks in the histogram. Using the ML partitioning, the first channel of the mixture was demixed and the SIR, PSR, and WDO were measured; the results are shown in Table III. For comparison, WDO_{Φ0}, the optimal WDO created using the 0 dB mask is shown in the last column. The demixtures average over 13 dB SIR gain and the WDO numbers indicate demixtures which would rate right on the border

between “minor artifacts or interference” and “distorted but intelligible.” Note that even though the blind method performs reasonably well, the performance of the 0 dB mask shows that there exist time-frequency masks which would further improve the performance. Figure 9 shows the original six sources, the two mixtures, and the six demixtures.

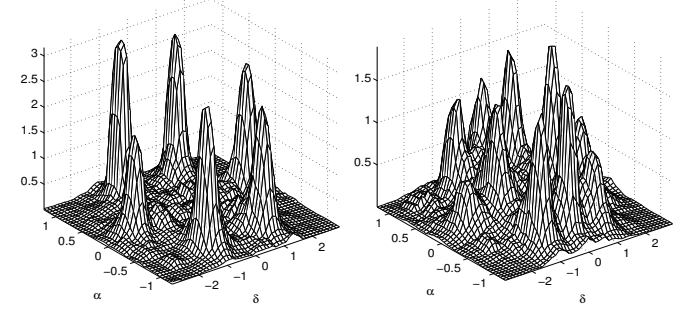


Fig. 8. Six and Ten Source Synthetic Mixing Smoothed Histograms ($p = 1$). Each peak corresponds to one source and the peak location corresponds to the associated source’s mixing parameters.

TABLE III
SIX AND TEN SOURCE DEMIXING PERFORMANCE. PERFORMANCE OF THE BLIND TECHNIQUE IS COMPARED AGAINST THE OPTIMAL TIME-FREQUENCY MASK, THE 0 dB MASK.

source	SIR in (dB)	SIR out (dB)	SIR gain (dB)	PSR	WDO DUET	WDO 0dB
s_1	-7.29	5.92	13.21	0.76	0.57	0.80
s_2	-7.29	5.24	12.53	0.78	0.55	0.78
s_3	-5.08	6.60	11.67	0.80	0.62	0.81
s_4	-9.29	5.35	14.63	0.79	0.56	0.69
s_5	-5.03	7.06	12.09	0.78	0.63	0.81
s_6	-9.28	5.47	14.75	0.77	0.55	0.66
s_1	-9.74	-0.32	9.42	0.58	-0.04	0.70
s_2	-7.73	3.14	10.87	0.66	0.34	0.77
s_3	-11.64	3.43	15.06	0.68	0.37	0.64
s_4	-9.72	-0.60	9.13	0.58	-0.09	0.67
s_5	-7.73	3.93	11.66	0.66	0.39	0.73
s_6	-11.61	3.14	14.75	0.56	0.29	0.51
s_7	-7.75	2.57	10.31	0.56	0.25	0.74
s_8	-11.62	1.36	12.98	0.61	0.16	0.62
s_9	-9.72	4.70	14.42	0.60	0.39	0.67
s_{10}	-9.74	3.33	13.07	0.60	0.32	0.64

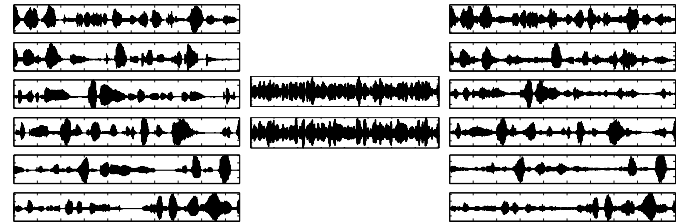


Fig. 9. Six Sources, Stereo Mixture, and Six Demixtures.

To show the limits of this technique, a ten source stereo mixture was synthetically mixed. The smoothed histogram for the mixture is shown in Figure 8 and Table III contains the demixing performance. The SIR gains are still high, the average gain above 12 dB, however, the WDO performance has dropped to “very distorted and barely intelligible.” However, as we are trying to demix ten sources from just two mixtures, these results are promising. More promising indeed is the fact that the 0 dB

mask's performance is significantly better showing that there is room for improvement.

B. Anechoic and Echoic Mixing Results

We also tested DUET on speech mixtures recorded in an anechoic room. For the tests, each speech signal was recorded separately at 16 kHz and then the signals were mixed additively to generate the mixtures for the tests. Knowledge of the actual signals present in each mixture allows us to calculate the performance measures exactly. For the recordings, the omnidirectional microphones were separated by 1.75 cm and the speech signals were played from various positions on a 1.5 meter radius semicircle around the microphones with the microphone axis along the line from the 0° position to the 180° position. Two female (F1 and F2) and one male (M1) TIMIT sound files were used for the tests. Pairwise mixing results for female-female and male-female mixtures are shown in Table IV. Again, for comparison purposes, the WDO obtained by the DUET algorithm is compared to the optimal WDO which is obtained using the 0 dB mask. The separation obtained by DUET is nearly perfect and in all but the 30° case: the DUET mask's performance is essentially the same as the performance of the optimal mask. The reason for the slight fall in performance in the 30° case is that the source peak regions begin to overlap and as a result some time-frequency points are misassigned.

TABLE IV
PAIRWISE ANECHOIC DEMIXING PERFORMANCE.

test	SIR in (dB)	SIR out (dB)	SIR gain (dB)	PSR	WDO DUET	WDO 0dB
F1 0°	-0.58	12.69	13.26	0.92	0.87	0.96
F2 30°	0.58	11.25	10.68	0.96	0.89	0.96
F1 0°	-0.54	15.97	16.51	0.98	0.95	0.96
F2 60°	0.54	17.21	16.68	0.98	0.96	0.96
F1 0°	-0.62	15.29	15.91	0.97	0.94	0.94
F2 90°	0.62	15.69	15.07	0.98	0.95	0.95
F1 0°	-0.49	17.50	17.99	0.98	0.96	0.96
F2 120°	0.49	17.36	16.87	0.98	0.97	0.97
F1 0°	-0.50	15.79	16.29	0.97	0.94	0.94
F2 150°	0.50	15.51	15.01	0.98	0.95	0.95
F1 0°	-0.44	16.29	16.73	0.96	0.94	0.94
F2 180°	0.44	14.49	14.05	0.98	0.94	0.95
F1 0°	3.54	13.99	10.46	0.96	0.92	0.97
M1 30°	-3.54	10.35	13.88	0.91	0.83	0.94
F1 0°	3.60	18.42	14.81	0.99	0.97	0.98
M1 60°	-3.60	15.41	19.01	0.97	0.94	0.95
F1 0°	3.63	18.92	15.29	0.99	0.98	0.98
M1 90°	-3.63	15.91	19.54	0.97	0.95	0.95
F1 0°	3.69	19.91	16.22	0.99	0.98	0.98
M1 120°	-3.69	15.79	19.48	0.98	0.95	0.95
F1 0°	3.75	19.57	15.82	0.99	0.98	0.98
M1 150°	-3.75	16.37	20.12	0.97	0.95	0.95
F1 0°	3.90	18.47	14.57	0.99	0.97	0.98
M1 180°	-3.90	15.51	19.41	0.97	0.94	0.94

Higher order mixing results (i.e., $N > 2$) are listed in Table V. In addition to the three, four, and five source anechoic mixtures tested, a three source echoic mixture was tested. All of the speech signals, three female (F1, F2, and F3) and two male (M1 and M2), were taken from the TIMIT database. The echoic recording was made in an echoic office environment with an approximate reverberation time of 500 ms. As the number of sources increases, the demixing performance decreases, although the performance is still acceptable in the five source mixture. As expected, the performance drops off significantly when switching from the anechoic to the echoic environment as the method is based on an anechoic mixing model. However, some separation is still achieved.

Figure 10 compares the one source histograms for anechoic and echoic recordings for sources at three different angles. The histograms corresponding to the summation of the three sources are also shown. The anechoic histograms are well localized and the peak regions are clearly distinct, even in the histogram corresponding to the summation of the sources. The peak regions in the echoic histograms are spread out and overlap with one another. This overlap results in reduced demixing performance. Note, however, that the 0 dB mask still performs well in the echoic case, so there remains a gap between what we can separate blindly and what we can separate with knowledge of the instantaneous time-frequency attenuations when using time-frequency masking to demix.

TABLE V
HIGHER ORDER DEMIXING PERFORMANCE. RESULTS FOR THREE SOURCE, FOUR SOURCE, AND FIVE SOURCE ANECHOIC MIXTURES, AS WELL AS THREE SOURCE ECHOIC MIXING.

Anechoic						
test	SIR in (dB)	SIR out (dB)	SIR gain (dB)	PSR	WDO DUET	WDO 0dB
M1 0°	-2.72	13.67	16.39	0.92	0.88	0.90
F1 90°	-2.05	7.96	10.00	0.96	0.80	0.93
M2 180°	-4.37	13.32	17.70	0.88	0.84	0.87
M1 0°	-6.93	9.89	16.83	0.78	0.70	0.80
F1 60°	-3.19	7.11	10.30	0.92	0.74	0.91
M2 120°	-4.37	6.98	11.35	0.85	0.68	0.89
F2 180°	-5.05	10.08	15.12	0.86	0.78	0.90
F1 0°	-9.77	7.97	17.74	0.73	0.62	0.76
M1 60°	-4.30	7.16	11.46	0.83	0.67	0.86
F2 90°	-3.77	5.99	9.76	0.91	0.68	0.91
M2 120°	-5.60	7.05	12.65	0.80	0.65	0.85
F3 180°	-8.59	8.53	17.11	0.76	0.65	0.82

Echoic						
test	SIR in (dB)	SIR out (dB)	SIR gain (dB)	PSR	WDO DUET	WDO 0dB
M1 0°	-5.20	5.38	10.58	0.56	0.40	0.81
M2 90°	0.07	4.33	4.26	0.89	0.56	0.91
F1 180°	-4.48	6.03	10.51	0.65	0.49	0.87

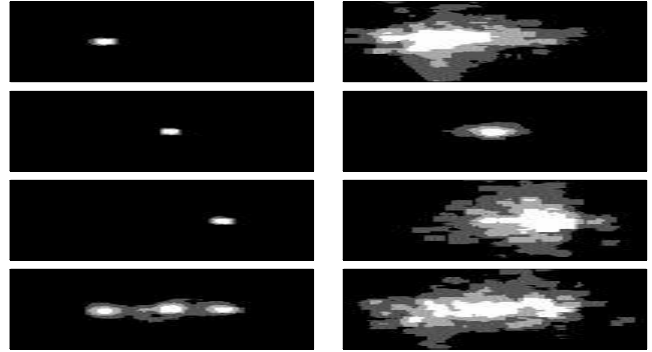


Fig. 10. Anechoic vs. Echoic Histogram Comparison. The left column images are of the histograms for three anechoic sources at 0° , 90° , 180° , and their mixture. The histogram of the mixture is essentially the summation of the individual histograms and the peak regions in the histogram are clearly separated. The right column images are of the histograms for three echoic sources 0° , 90° , 180° , and their mixture. While the individual histograms show some level of localization (left, center, right), peak regions in each histogram overlap and the peaks are difficult to identify in the summation image. Thus, the algorithm performs worse on echoic mixtures. In all images, the x-axis is delay ranging from -2 to 2 samples and the y-axis is symmetric attenuation ranging from -1 to 1. All histograms used $p = 1$.

Figure 11 shows the histogram for the Te-Won Lee real office room recording consisting of two speakers [22]. The histogram shows a number of peaks, the peaks with $\alpha > -0.5$ are all associated with the Spanish speaker, and those along the $\alpha = -1.0$ line correspond to the English speaker. Note that for this

recording, it is the attenuation direction in the histogram that allows for the separation and that a method that only relied on delays would not be able to separate the sources. Demixtures generated from this recording using the DUET algorithm are compared to several other BSS techniques in [23].

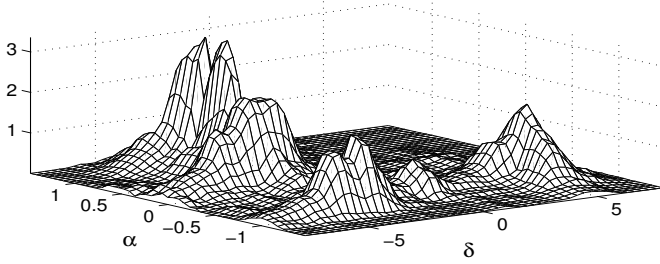


Fig. 11. Histogram ($p = 1$) for Te-Won Lee's "A real Cocktail Party Effect" Echoic Mixing Example.

V. CONCLUSIONS

In this paper, we presented a method to blindly separate mixtures of speech signals. We first illustrated experimentally that binary time-frequency masks exist that can separate as many as 10 different speech signals from one mixture. This relies upon a property of the Gabor expansions of speech signals, which we refer to as W-disjoint orthogonality. W-disjoint orthogonality in the strict sense is satisfied by signals which have disjoint time-frequency supports. Speech signals, as a result of the sparsity of their Gabor expansions, satisfy an approximate version of the W-disjoint orthogonality property. In Section II-A, we introduced a means of measuring the degree of W-disjoint orthogonality of a signal in a given mixture with respect to a windowing function W . Listening experiments showed that there is a fairly accurate relationship between the WDO value of a particular signal in a mixture for a given time-frequency mask and the subjective performance of the time-frequency mask to separate the signal from the mixture.

Next, we addressed the problem of blindly constructing binary time-frequency masks that demix. The solution we presented in this paper considered the two mixture case. For strictly W-disjoint orthogonal signals, we showed that the instantaneous DUET attenuation and delay estimators are the anechoic mixing parameters and, using this fact, described a simple algorithm to construct a binary time-frequency mask that demixes perfectly. Next we showed, by modeling the contributions of the interfering sources as independent Gaussian white noise, that the ML estimators for the mixing parameters are given by weighted averages of the instantaneous DUET estimates. Motivated by this, we constructed a weighted histogram which was used to enumerate the sources and partition the time-frequency representation of one of the mixtures to demix.

In Section IV we presented experimental results demonstrating that the algorithm works extremely well for synthetic mixtures of speech as well as for speech mixtures recorded in an anechoic room and produces near perfect demixtures. That is, the performance of the mask generated by the DUET algorithm was close to the performance of the ideal mask. In an echoic

room the anechoic model is violated and the quality of the demixing is reduced. In the echoic case, the demixtures contain some crosstalk and distortion, but are intelligible. There is, in the echoic case, a performance gap between the ideal binary time-frequency mask and the mask generated using the DUET algorithm. Closing this gap is one goal of our future work. As we mentioned in Section IV, the enumeration and identification of the histogram peaks are also topics for future research.

APPENDIX I

DERIVATION OF THE ML ESTIMATORS

Let us concentrate on one source, say s_j . Let Λ_j be the set of time-frequency points $[k, l]$ at which s_j is dominant as defined in Section III-B. On Λ_j , we model the mixtures as follows:

$$\begin{aligned}\hat{x}_1[k, l] &= \hat{s}_j[k, l] + n_1[k, l] \\ \hat{x}_2[k, l] &= a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l] + n_2[k, l]\end{aligned}\quad (47)$$

where n_1 and n_2 are i.i.d. white complex Gaussian noise signals with mean zero and variance σ^2 . Here n_1 and n_2 model the contributions of other sources at the time-frequency points where s_j is the dominant source. We model the interfering sources as independent Gaussian noise in order to obtain simple closed-form source and mixing parameter estimators. In reality, the interference in the different mixtures will be correlated and may not be Gaussian distributed. However, the estimates we obtain here are used simply to *motivate* the weighted histograms used in our algorithm discussed in Section III-B and the model is sufficient for that purpose.

For the model in (47), we want to employ an ML estimate to find the parameter pair $(a_j, \delta_j) \in \mathbb{R}^2$ as well as $\hat{s}_j[k, l]$ which has the maximum likelihood. To that goal, we define the likelihood, L_0 , of (s_j, a_j, δ_j) , where $s_j = (\hat{s}_j[k, l])_{(k, l) \in \Lambda_j}$ with each $\hat{s}_j[k, l] \in \mathbb{C}$, given the data $\hat{x}_1[k, l]$ and $\hat{x}_2[k, l]$, by

$$\begin{aligned}L_0(s_j, a_j, \delta_j) &:= p(\mathbf{x}_1, \mathbf{x}_2 | s_j, a_j, \delta_j) \\ &= \prod_{(k, l) \in \Lambda_j} f_{N_1, N_2}(\hat{x}_1[k, l] - \hat{s}_j[k, l], \hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l]) \\ &= C \exp \left(-\frac{1}{2\sigma^2} \sum_{(k, l) \in \Lambda_j} |\hat{x}_1[k, l] - \hat{s}_j[k, l]|^2 + \right. \\ &\quad \left. |\hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l]|^2 \right)\end{aligned}\quad (48)$$

where $\mathbf{x}_i = (\hat{x}_i[k, l])_{(k, l) \in \Lambda_j}$. The last equality holds because we assume n_1 and n_2 are i.i.d. complex Gaussian noise signals. Clearly, maximizing L_0 is equivalent to maximizing

$$L(s_j, a_j, \delta_j) := - \sum_{(k, l) \in \Lambda_j} |\hat{x}_1[k, l] - \hat{s}_j[k, l]|^2 + |\hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l]|^2. \quad (49)$$

We want to solve the equations $\frac{\partial L}{\partial \hat{s}_j^R[k, l]} = 0$, $\frac{\partial L}{\partial \hat{s}_j^I[k, l]} = 0$ for all $(k, l) \in \Lambda_j$, $\frac{\partial L}{\partial a_j} = 0$, and $\frac{\partial L}{\partial \delta_j} = 0$ simultaneously, where $\hat{s}_j^R[k, l]$ and $\hat{s}_j^I[k, l]$ denote the real and imaginary parts of $\hat{s}_j[k, l]$ respectively. We start with $\frac{\partial L}{\partial \hat{s}_j^R[k, l]}$. For any $(k, l) \in \Lambda_j$, we have

$$\begin{aligned}\frac{\partial L}{\partial \hat{s}_j^R[k, l]} &= \frac{\partial}{\partial \hat{s}_j^R[k, l]} \left(|\hat{x}_1[k, l] - \hat{s}_j^R[k, l] - i\hat{s}_j^I[k, l]|^2 + \right. \\ &\quad \left. |\hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j^R[k, l] + i\hat{s}_j^I[k, l]|^2 \right)\end{aligned}\quad (50)$$

We solve $\frac{\partial L}{\partial \hat{s}_j^R[k, l]} \Big|_{\hat{s}_j^R[k, l] = \hat{s}_j^{R*}[k, l]} = 0$ for $\hat{s}_j^{R*}[k, l]$, the ML estimate of $\hat{s}_j^R[k, l]$, and obtain

$$\hat{s}_j^{R*}[k, l] = \mathbf{Re} \left\{ \frac{\hat{x}_1[k, l] + a_j e^{i\delta_j l \omega_0} \hat{x}_2[k, l]}{1 + a_j^2} \right\}. \quad (51)$$

Similarly, solving $\frac{\partial L}{\partial \hat{s}_j^I[k, l]} \Big|_{\hat{s}_j^I[k, l] = \hat{s}_j^{I*}[k, l]} = 0$ for $\hat{s}_j^{I*}[k, l]$, the ML estimate of $\hat{s}_j^I[k, l]$, yields

$$\hat{s}_j^{I*}[k, l] = \mathbf{Im} \left\{ \frac{\hat{x}_1[k, l] + a_j e^{i\delta_j l \omega_0} \hat{x}_2[k, l]}{1 + a_j^2} \right\} \quad (52)$$

which we combine with (51) to get the ML estimate \mathbf{s}_j^* for \mathbf{s}_j :

$$\mathbf{s}_j^*[k, l] = \frac{\hat{x}_1[k, l] + a_j e^{i\delta_j l \omega_0} \hat{x}_2[k, l]}{1 + a_j^2}. \quad (53)$$

Next, we consider $\frac{\partial L}{\partial \delta_j}$. We have

$$\begin{aligned} \frac{\partial L}{\partial \delta_j} &= \frac{\partial}{\partial \delta_j} \left(\sum_{(k, l) \in \Lambda_j} |\hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l]|^2 \right) \\ &= 2a_j \sum_{(k, l) \in \Lambda_j} l \omega_0 \mathbf{Im} \left\{ \hat{x}_2[k, l] \overline{\hat{s}_j[k, l]} e^{i\delta_j l \omega_0} \right\}. \end{aligned} \quad (54)$$

where $\overline{\hat{s}_j[k, l]}$ is the complex conjugate of $\hat{s}_j[k, l]$. We now plug in $\hat{s}_j[k, l] = \mathbf{s}_j^*[k, l]$ in (54), which yields

$$\begin{aligned} \frac{\partial L}{\partial \delta_j} &= \frac{2a_j}{1 + a_j^2} \sum_{(k, l) \in \Lambda_j} l \omega_0 |\hat{x}_1[k, l]|^2 \mathbf{Im} \left\{ R_{21}[k, l] e^{i\delta_j l \omega_0} \right\} \\ &= \frac{2a_j}{1 + a_j^2} \sum_{(k, l) \in \Lambda_j} l \omega_0 |\hat{x}_1[k, l] \hat{x}_2[k, l]| \sin(\angle R_{21}[k, l] + \delta_j l \omega_0). \end{aligned} \quad (55)$$

We assume that $|\angle R_{21}[k, l] + \delta_j l \omega_0| = |l \omega_0 (\delta_j - \tilde{\delta}[k, l])|$ is small which is reasonable because we are considering only the $[k, l]$ where s_j is dominant and we make the approximation

$$\sin(\angle R_{21}[k, l] + \delta_j l \omega_0) \approx \angle R_{21}[k, l] + \delta_j l \omega_0. \quad (56)$$

After plugging (56) into (55), we solve the equation $\frac{\partial L}{\partial \delta_j} \Big|_{\delta_j = \delta_j^*} = 0$ for δ_j^* and obtain (using the definition in (29))

$$\delta_j^* = \frac{\sum_{(k, l) \in \Lambda_j} \tilde{\delta}[k, l] l^2 \omega_0^2 |\hat{x}_1[k, l] \hat{x}_2[k, l]|}{\sum_{(k, l) \in \Lambda_j} l^2 \omega_0^2 |\hat{x}_1[k, l] \hat{x}_2[k, l]|}. \quad (57)$$

Note that δ_j^* , the ML estimate for the parameter δ_j , is a weighted average of the instantaneous DUET delay estimates, with each estimate weighted by the product magnitude of the mixtures as well as $(l \omega_0)^2$. We also observe that the ML estimate δ_j^* does not depend on the attenuation parameter a_j .

Finally, we will solve $\frac{\partial L}{\partial a_j} \Big|_{a_j = a_j^*}$ for a_j^* . We have

$$\begin{aligned} \frac{\partial L}{\partial a_j} &= \frac{\partial}{\partial a_j} \left(\sum_{(k, l) \in \Lambda_j} |\hat{x}_2[k, l] - a_j e^{-i\delta_j l \omega_0} \hat{s}_j[k, l]|^2 \right) \\ &= \sum_{(k, l) \in \Lambda_j} 2a_j |\hat{s}_j[k, l]|^2 - 2 \mathbf{Re} \left\{ \hat{x}_2[k, l] \overline{\hat{s}_j[k, l]} e^{i\delta_j l \omega_0} \right\} \end{aligned} \quad (58)$$

After setting $\hat{s}_j[k, l] = \mathbf{s}_j^*[k, l]$ and some algebra we get (using the definition in (30))

$$\alpha_j^* := a_j^* - \frac{1}{a_j^*} = \frac{\sum_{(k, l) \in \Lambda_j} |\hat{x}_1[k, l] \hat{x}_2[k, l]| (\tilde{\alpha}[k, l] - 1/\tilde{\alpha}[k, l])}{\sum_{(k, l) \in \Lambda_j} \mathbf{Re} \left\{ \hat{x}_2[k, l] \overline{\hat{x}_1[k, l]} e^{i\delta_j^* l \omega_0} \right\}} \quad (59)$$

The estimate for $a_j^* - \frac{1}{a_j^*}$ is symmetric in \mathbf{x}_1 and \mathbf{x}_2 : swapping the mixture labels will only result in a sign change of this quantity (i.e., $(1/a_j) - (1/(1/a_j)) = -(a_j - 1/a_j)$). In the original presentation of DUET [4], the logarithm of the attenuation estimates was used solely because it has the same property (i.e., $\log(1/a_j) = -\log(a_j)$). However, motivated by its appearance in the ML estimator (59), we will replace the role of the logarithm with the DUET symmetric attenuation estimator defined in (32).

Remark 2: Although the estimate given in (59) is not a weighted average, it is interesting to note that if we replace δ_j^* in (59) with $\tilde{\delta}[k, l]$ we obtain that

$$\mathbf{Re} \left\{ \hat{x}_2[k, l] \overline{\hat{x}_1[k, l]} e^{i\tilde{\delta}[k, l] l \omega_0} \right\} = |\hat{x}_1[k, l] \hat{x}_2[k, l]| \quad (60)$$

and in this case, (59) becomes (33) with $p = 1$, a weighted average of $\tilde{\alpha}[k, l]$.

APPENDIX II

EXPERIMENTAL EVALUATION OF THE ML ESTIMATORS

In this section, we experimentally evaluate the ML estimators as well as other estimators motivated by the previous section. In order to simulate mixtures, we use the model in (47) and adjust the noise energy to model the different number of interfering sources. The model in (47) is valid for the dominant time-frequency points of one source. In order to determine the set of dominant time-frequency points Λ_j , a speech signal taken from the TIMIT database was compared to a random mixture of 1, 2, 4, or 9 TIMIT speech signals to model $N = 2, 3, 5, 10$, and in each case, the time-frequency points corresponding to the 0-dB mask were selected. The mixtures of interfering sources were only used to determine Λ_j and were discarded after the dominant time-frequency points were identified. In order to simulate the presence of interfering sources, i.i.d. Gaussian white noise was added to the dominant time-frequency points of source s_j on both channels. The added noise was amplified to produce a 15.12 dB, 12.26 dB, 9.87 dB, or 7.52 dB SNR so as to model mixing of order $N = 2, 3, 5$, or 10, respectively. These SNR's were selected to model different mixture orders because they match the average SIR's for the 0-dB mask from Figure 3. That is, 15.12 dB, 12.26 dB, 9.87 dB, and 7.52 dB are the expected SIR's after applying the 0-dB mask to mixtures of order $N = 2, 3, 5$, and 10, and thus in order to model these mixture orders for the dominant time-frequency points of one source, we add noise to one source to produce the corresponding SNR's. Note that the dominant time-frequency points are precisely the support of the source's 0 dB mask. Thus the performance of the estimator evaluated with this model at these SNR's should approximate the true performance of the estimator in speech mixtures of order $N = 2, 3, 5$, and 10. All results in the remainder of this section are obtained using this model.

We choose to experimentally evaluate the estimators using the model as described above as opposed to creating synthetic mixtures of multiple speech signals because we (1) wanted to prevent the results from depending on the specific choice of mixing parameters of the interfering sources and (2) wanted to evaluate the estimators using the model that motivated them. The disadvantage of modeling the presence of interfering sources in this way is that the interference should be correlated and this correlation is lost when the interference is modeled as independent noise. Our desire in this section is to explore the qualitative performance of a family of estimators to motivate the demixing algorithm, and modeling the interference as noise is sufficient for this purpose.

Figure 12 shows the ML estimate δ_j^* from (57) versus δ as δ ranges linearly from -5 to 5 samples with $a_j = 1$. We can see that the ML delay estimator exhibits bad performance outside the -1 to 1 sample range, and biased performance inside this range. The bad performance for larger delays is due to the phase wrap around problem discussed in Remark 1 in Section III-A. The squared frequency weighting factor in the ML delay estimator accentuates this problem. In addition, such a frequency weighting would make signals with higher frequency content have higher likelihood estimates of their delay parameters. In the next sections, we will be using these weightings to construct weighted histograms for source separation and it is undesirable to assign more likelihood to one set of parameters simply because their associated source contains higher frequencies. While methods for unwrapping the phase do exist, these methods are inappropriate for our purposes as different sources may be active from one frequency to the next. In order to see if we could reduce the bias, eliminate the wrap-around effect, and remove the high frequency weighting, we removed the squared frequency weighting factor in the ML delay estimator and considered estimators of the following form

$$\delta_j^{(p)} := \frac{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p \tilde{\delta}[k,l]}{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p}. \quad (61)$$

The free parameter p in (61) determines how strongly the estimates obtained from time-frequency points $[k,l]$ with large $|x_j[k,l]|$ are weighted relative to those obtained from time-frequency points with small $|x_j[k,l]|$. Note that $p = 0$ corresponds to (61) being a plain average and $p = \infty$ results in $\delta_j^{(p)}$ being the instantaneous DUET estimate from the time-frequency point $[k,l]$ with the largest coefficient magnitude. Moreover, with $p = 1$, this estimator is the ML delay estimator with the squared frequency weighting factor removed. The estimator in (61) for $p = 1/2, 1, 2$ is compared with the ML estimator in Figure 13. For this test, δ_j ranges linearly from -5 to 5 samples while α_j ranges linearly from 0.15 to -0.15 and the SNR and Λ_j were selected to model mixture orders of 2, 3, 5, and 10. The $p = 1/2$ estimator suffers similar deficiencies as the ML estimator. The $p = 1$ estimator is clearly biased outside the -1 to 1 delay range, but is monotonic with increasing δ_j and exhibits good (although biased) performance inside -1 to 1 sample delay. The $p = 2$ estimator exhibits near perfect estimates.

Figure 13 also shows α_j^* versus α_j for the same data used for the delay estimates. Similar to the delay case, in addition to the

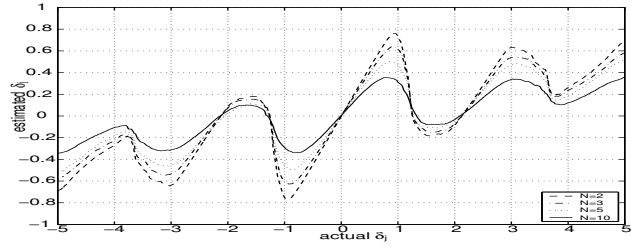


Fig. 12. Maximum Likelihood Delay Estimator. The plot compares estimated δ_j versus true δ_j for the ML estimator for δ_j ranging linearly from -5 to 5 samples when $\alpha_j = 0$ for mixture models of 2, 3, 5, and 10 sources.

ML estimator, we consider estimators of the following form,

$$\alpha_j^{(p)} := \frac{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p \tilde{\alpha}[k,l]}{\sum_{(k,l) \in \Lambda_j} |\hat{x}_1[k,l] \hat{x}_2[k,l]|^p}. \quad (62)$$

Note that with $p = 1$, this estimator is the ML symmetric estimator with the substitution described in Remark 2 previously. The ML and $p = 1/2$ symmetric attenuation estimators are clearly biased. The $p = 1$ symmetric attenuation estimator is also biased, although less so. The $p = 2$ symmetric attenuation estimator exhibits near perfect performance.

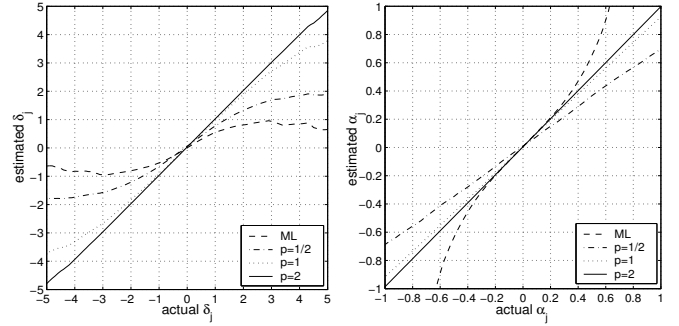


Fig. 13. Delay (left) and Attenuation (right) Estimator Comparison. The graph on the left compares estimated δ_j versus true δ_j for the ML delay estimator and the weighted average DUET delay estimators with $p = 1/2, 1, 2$ as δ_j ranges linearly from -5 to 5 samples while α_j ranges linearly from 0.15 to -0.15 for a mixture model of 5 sources. The graph on the right compares estimated α_j versus true α_j for the ML and weighted average DUET symmetric attenuation estimators for the same experimental data used to generate the delay graph.

APPENDIX III

HISTOGRAM-BASED PARAMETER ESTIMATION

In order to evaluate the usefulness of the histogram as a parameter estimator, a smoothed histogram was created for each of the tests used to generate Figure 13 and the peak location of the histogram was used as the symmetric attenuation and delay estimate. Figure 14 contains the results of these tests. Each estimator histogram consisted of 401-by-401 points with a delay range from -6 to 6 samples and symmetric attenuation from -1.2 to 1.2, and the smoothing kernel had parameters $(A, D) = (0.12, 0.2)$. Comparing Figure 13 and Figure 14, we conclude that the histogram based estimators are more accurate than the previously considered ML motivated estimators.

The similar estimator performance for different choices of p in Figure 14 suggests that the choice of p should be driven by

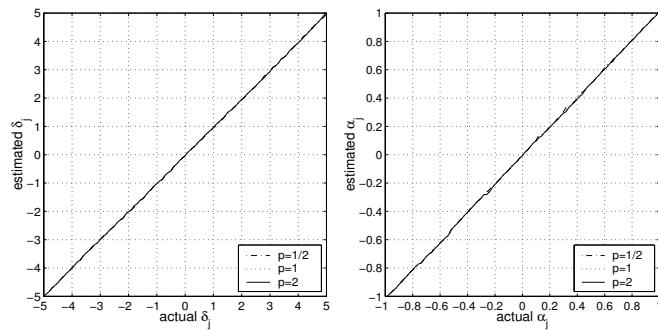


Fig. 14. Histogram Delay (left) and Attenuation (right) Estimator Comparison. The graph on the left compares estimated δ_j versus true δ_j for the smoothed histogram peak estimators with $p = 1/2, 1, 2$ for δ_j ranging linearly from -5 to 5 samples as α_j ranges linearly 0.15 to -0.15. The graph on the right compares estimated α_j to the true α_j . Both graphs were generated using a model of 5 source mixing.

other concerns. Identifying the peaks in the histogram is the crucial step in the separation process. Two important criteria for the weighting exponent p selection are (1) the shape around the peak (the “peak shape”) and (2) the relative peak heights. In order to aid in peak identification, we want the peak shape to be narrow and tall, and we want the peaks to be roughly of the same height. Figure 15 compares the histogram peak shapes for $p = 1/2, 1, 2$ for both the α and δ axes by taking the summation along the other axis. That is, Figure 15 contains 1-D weighted histograms for both α and δ . As p increases, the peak shape becomes narrower and taller. This would suggest that we should select p as large as possible. However, the larger we choose p , the more the peak heights depend only on the largest instantaneous product power time-frequency components of each source. If these components have different magnitude distributions for different sources, the resulting peaks heights can vary by several orders of magnitude making identification of the smaller peaks impossible. While $p = 2$ results in the best peak shape, smaller choices of p may result in easier peak identification. The choice of p is thus data dependent, however, motivated once again by the form of the ML estimators, we will suggest $p = 1$ as the default choice.

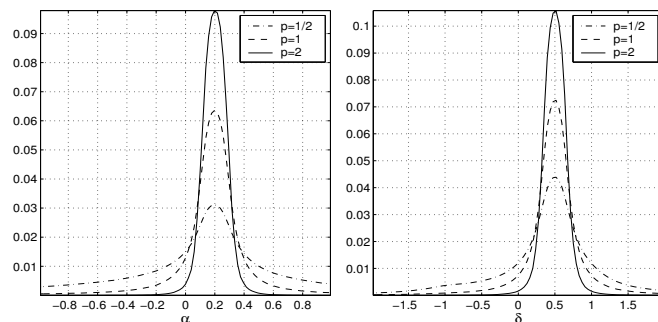


Fig. 15. Peakshape for $p=1/2, 1$, and 2 for mixing with $(\alpha_j, \delta_j) = (0.2, 0.5)$.

REFERENCES

[1] J.-K. Lin, D. G. Grier, and J. D. Cowan, “Feature extraction approach to blind source separation,” in *IEEE Workshop on Neural Networks for*

Signal Processing (NNSP), Amelia Island Plantation, Florida, September 24–26 1997, pp. 398–405.

[2] T.-W. Lee, M. Lewicki, M. Girolami, and T. Sejnowski, “Blind source separation of more sources than mixtures using overcomplete representations,” *IEEE Signal Proc. Letters*, vol. 6, no. 4, pp. 87–90, April 1999.

[3] M. V. Hulle, “Clustering approach to square and non-square blind source separation,” in *IEEE Workshop on Neural Networks for Signal Processing (NNSP)*, Madison, Wisconsin, August 23–25 1999, pp. 315–323.

[4] A. Jourjine, S. Rickard, and O. Yilmaz, “Blind separation of disjoint orthogonal signals: Demixing N sources from 2 mixtures,” in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 5, Istanbul, Turkey, June 5–9 2000, pp. 2985–2988.

[5] P. Bofill and M. Zibulevsky, “Blind separation of more sources than mixtures using sparsity of their short-time Fourier transform,” in *International Workshop on Independent Component Analysis and Blind Source Separation (ICA)*, Helsinki, Finland, June 19–22 2000, pp. 87–92.

[6] R. Balan and J. Rosca, “Statistical properties of STFT ratios for two channel systems and applications to blind source separation,” in *International Workshop on Independent Component Analysis and Blind Source Separation (ICA)*, Helsinki, Finland, June 19–22 2000, pp. 429–434.

[7] M. Aoki, M. Okamoto, S. Aoki, H. Matsui, T. Sakurai, and Y. Kaneda, “Sound source segregation based on estimating incident angle of each frequency component of input signals acquired by multiple microphones,” *Acoustical Science & Technology*, vol. 22, no. 2, pp. 149–157, Feb. 2001.

[8] M. Zibulevsky and B. A. Pearlmutter, “Blind source separation by sparse decomposition in a signal dictionary,” *Neural Computation*, vol. 13, no. 4, pp. 863–882, 2001.

[9] L.-T. Nguyen, A. Belouchrani, K. Abed-Meraim, and B. Boashash, “Separating more sources than sensors using time-frequency distributions,” in *International Symposium on Signal Processing and its Applications (ISSPA)*, Kuala Lumpur, Malaysia, August 13–16 2001, pp. 583–586.

[10] S. Rickard, R. Balan, and J. Rosca, “Real-time time-frequency based blind source separation,” in *International Workshop on Independent Component Analysis and Blind Source Separation (ICA)*, San Diego, CA, December 9–12 2001, pp. 651–656.

[11] L. Vielva, D. Erdogmus, C. Pantaleon, I. Santamaria, J. Pereda, and J. C. Principe, “Underdetermined blind source separation in a time-varying environment,” in *IEEE Int. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 3, Orlando, Florida, May 13–17 2002, pp. 3049–3052.

[12] P. Bofill, “Underdetermined blind separation of delayed sound sources in the frequency domain,” *preprint*, 2002.

[13] S. Rickard and O. Yilmaz, “On the approximate W-disjoint orthogonality of speech,” in *IEEE Int. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 1, Orlando, Florida, May 13–17 2002, pp. 529–532.

[14] S. T. Roweis, “One microphone source separation,” in *Neural Information Processing Systems 13 (NIPS)*, 2000, pp. 793–799.

[15] K. Suyama, K. Takahashi, and R. Hirabayashi, “A robust technique for sound source localization in consideration of room capacity,” in *IEEE Workshop on the Applications of Signal Processing to Audio and Acoustics*, New Paltz, New York, October 21–24 2001, pp. 63–66.

[16] I. Daubechies, *Ten Lectures on Wavelets*. Philadelphia, PA: SIAM, 1992.

[17] R. Mersereau and T. Seay, “Multiple access frequency hopping patterns with low ambiguity,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. AES-17, no. 4, July 1981.

[18] W. Kozek, “Time-frequency signal processing based on the Wigner-Weyl framework,” *Signal Processing*, vol. 29, pp. 77–92, 1992.

[19] B. Berdugo, J. Rosenhouse, and H. Azhari, “Speakers’ direction finding using estimated time delays in the frequency domain,” *Signal Processing*, vol. 82, pp. 19–30, 2002.

[20] R. Balan, J. Rosca, S. Rickard, and J. O’Ruanaidh, “The influence of windowing on time delay estimates,” in *Conf. on Info. Sciences and Systems (CISS)*, vol. 1, Princeton, NJ, March 2000, pp. WP1–(15–17).

[21] J. Huang, N. Ohnishi, and N. Sugie, “A biomimetic system for localization and separation of multiple sound sources,” *IEEE Transactions on Instrumentation and Measurement*, vol. 44, no. 3, pp. 733–738, June 1995.

[22] [Online]. Available: http://www.cnl.salk.edu/~tewon/Blind/blind_audio.html

[23] [Online]. Available: <http://alum.mit.edu/www/rickard/bss.html>