

Acoustic Source Separation with Microphone Arrays



Lucas C. Parra
Biomedical Engineering Department
City College of New York

CCNY

Craig Fancourt
Clay Spence
Chris Alvino



Montreal Workshop, Nov 6, 2004

Blind Source Separation and ICA

The goal of Independent Component Analysis (ICA) is to factorize an observation matrix \mathbf{X} into a basis vectors \mathbf{S} and coefficients \mathbf{A} such that the columns of \mathbf{S} are statistically independent:

$$\begin{array}{c} \text{sensors} \downarrow \\ \text{samples} \rightarrow \\ \mathbf{X} \end{array} = \begin{array}{c} \text{sensors} \downarrow \\ \text{Components} \rightarrow \\ \mathbf{A} \end{array} * \begin{array}{c} \text{components} \downarrow \\ \text{samples} \rightarrow \\ \mathbf{S} \end{array}$$

The notion is that the observations $\mathbf{x}(t)$ have been generated by linear mixing \mathbf{A} of independent sources $\mathbf{s}(t)$, and the recovered component are the originating sources.

Separation Based on Independence

Statistical independence implies for all $i \neq j, t, l, n, m$:

$$E[s_i^n(t) s_j^m(t+l)] = E[s_i^n(t)]E[s_j^m(t+l)]$$

For M sources and N sensors each tuple $\{t, l, n, m\}$ this equation gives $M(M-1)/2$ conditions for the NM unknowns in \mathbf{A} .

Sufficient conditions if we use multiple:

<u>use</u>	<u>sources assumed*</u>	<u>condition</u>	<u>statistic</u>	<u>algorithm</u>
t	non-stationary	$\mathbf{R}_x(t) = \mathbf{A} \mathbf{R}_s(t) \mathbf{A}^T$	covariance	decorrelation
l	non-white	$\mathbf{R}_x(l) = \mathbf{A} \mathbf{R}_s(l) \mathbf{A}^T$	cross-correlation	SOBI
n, m	non-Gaussian	$\mathbf{C}_x(i,j) = \mathbf{A} \mathbf{C}_s(i,j) \mathbf{A}^T$	4th cumulants	JADE (ICA)

* zero mean

'Quickie BSS' – GEV of two Cross-Statistics

The independence assumption establishes that there are some cross-statistics of the observations say \mathbf{R}_1 , \mathbf{Q}_2 that are diagonalized by $\mathbf{W}=\mathbf{A}^{-1}$:

$$\mathbf{R}_S = \mathbf{W} \mathbf{R}_X \mathbf{W}^T = \text{diag}$$

$$\mathbf{C}_S = \mathbf{W} \mathbf{C}_X \mathbf{W}^T = \text{diag}$$

This can be combined to a standard generalized Eigenvalue equation

$$\mathbf{R}_X^{-1} \mathbf{C}_X \mathbf{W} = \mathbf{W} \mathbf{D}$$

```
>> [W,D] = eig(X*X',C);
```

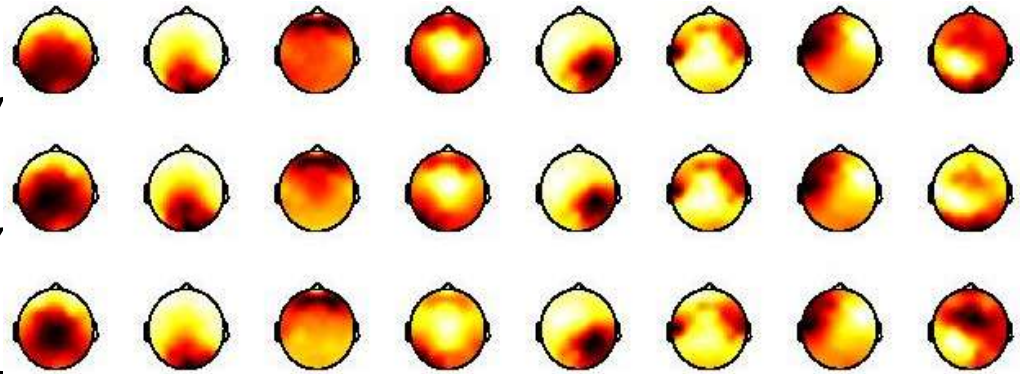
```
>> S = W' * X;
```

Example: ICA in Electro-Encephalography

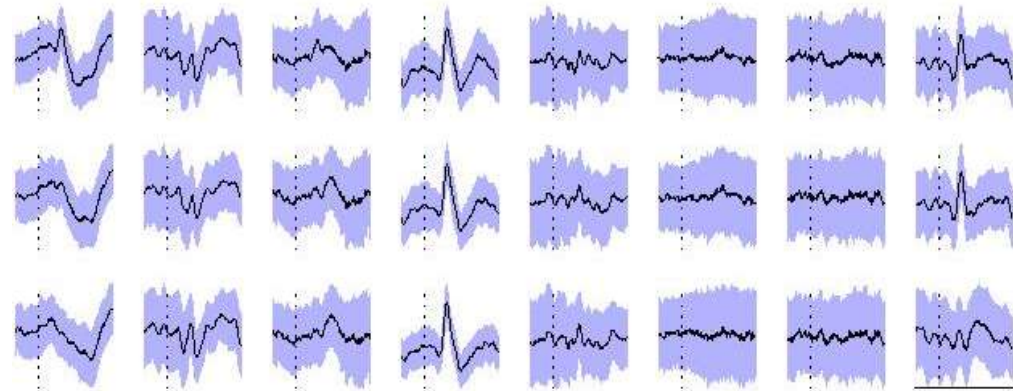
Interestingly, diagonalization algorithms based on these three different properties of the signals give often similar results in EEG:

- non-white
- Non-Gaussian
- non-stationary

EEG sensor projections $\mathbf{A} = \mathbf{W}^{-1}$



Trial averaged components $\mathbf{s}(t) = \mathbf{W} \mathbf{x}(t)$



0.8 s

This supports the interpretation of EEG as linear superposition of independent signals.

Acoustic Mixing and Separation

Acoustic environments are characterized by multi-path propagation which can be represented as a convolutive mixture:

$$\mathbf{x}(t) = \sum_l \mathbf{A}(l) \mathbf{s}(t-l)$$

Source $s_i(t)$ is coupled to microphone $x_j(t)$ by the room response $A_{ji}(l)$.

The goal of acoustic source separation is to find filters $W_{ij}(l)$ that generate model sources $y_i(t)$, corresponding to the sources $s_i(t)$

$$\mathbf{y}(t) = \sum_l \mathbf{W}(l) \mathbf{x}(t-l)$$

Frequency Domain Separation

For efficiency reasons this problem is often transformed into the frequency domain:

$$\mathbf{x}(t) = \mathbf{A}(t) * \mathbf{s}(t) \quad \Leftrightarrow \quad \mathbf{x}(\omega) = \mathbf{A}(\omega) \mathbf{s}(\omega)$$

Source separation aims to find **filters** $\mathbf{W}(\omega)$ such that model source $\mathbf{y}(\omega)$ correspond to the original sources $\mathbf{s}(\omega)$.

$$\mathbf{y}(\omega) = \mathbf{W}(\omega) \mathbf{A}(\omega) \mathbf{s}(\omega)$$

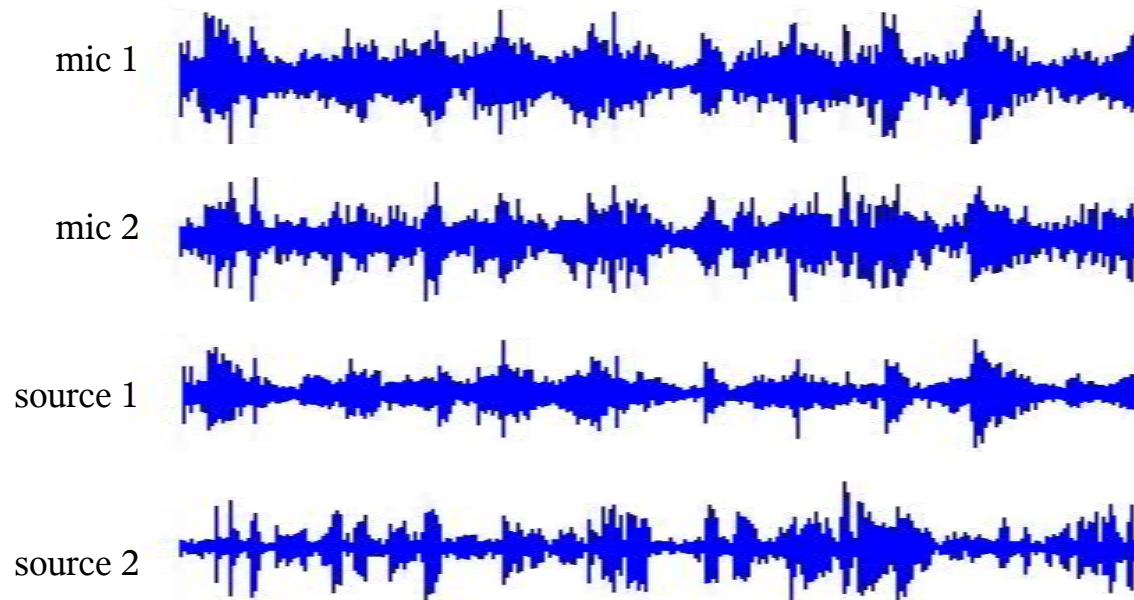
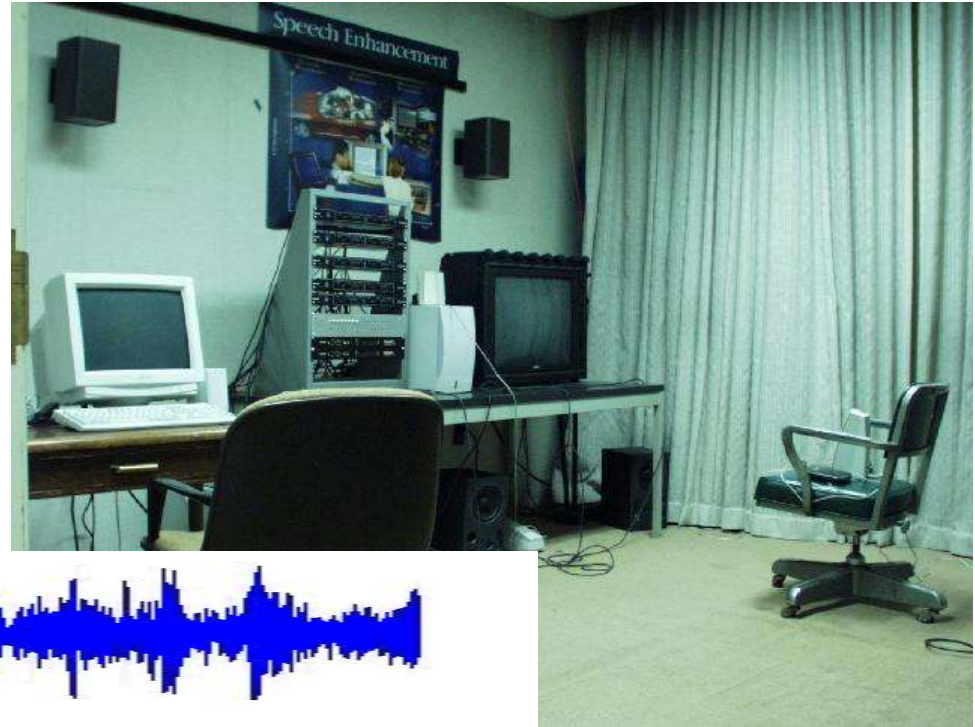
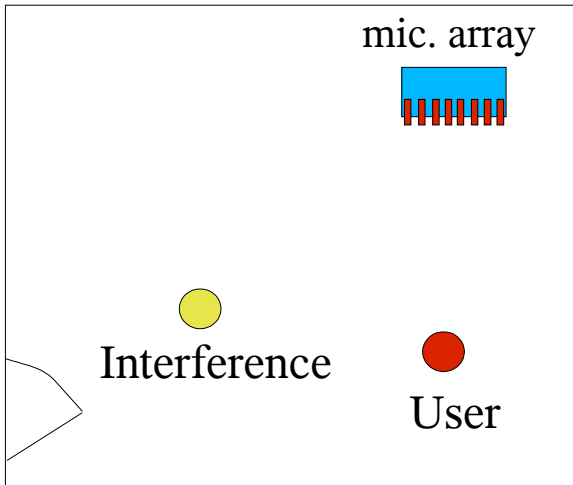
Separation based on Non-Stationarity

For non-stationary signals this can be achieved by minimizing cross-powers $\mathbf{R}_{yy}(t, \omega) = \mathbb{E}[\mathbf{y}(\omega)\mathbf{y}^H(\omega) | t]$ estimated over multiple estimation periods:

$$\underset{\mathbf{W}(\omega)}{\operatorname{argmin}} \sum_t \left\| \operatorname{offdiag}(\mathbf{R}_{yy}(t, \omega)) \right\|^2$$

$$\mathbf{R}_{yy}(t, \omega) = \mathbf{W}(\omega) \mathbf{R}_{xx}(t, \omega) \mathbf{W}^H(\omega)$$

Separation for Speech Controlled TV



Ambiguities of the Separation Criteria

Independence criteria specifies result up to permutation $\mathbf{P}(\omega)$ and scaling $\mathbf{D}(\omega)$

$$\mathbf{W}(\omega) \mathbf{A}(\omega) = \mathbf{P}(\omega) \mathbf{D}(\omega)$$

- Convolution: Scaling $\mathbf{D}(\omega)$ corresponds to convolution of each source.
- Permutation: $\mathbf{P}(\omega)$ may be different for every frequency bin!
- Subspace: $N \leq M$ dimensional space of equivalent solutions $\mathbf{W}(\omega)$.

Convolutional BSS Summary

State-of-the-art (?): **Frequency domain SOS** achieving **15dB-20dB** interference reduction for mixtures of up to **3 sources**.

Second order statistics (SOS) sufficient for non-stationary signals such as speech.

SOS algorithms easier to implement.

Frequency domain formulations give fast algorithms.

However SOS in the frequency domain leads to frequency permutation ambiguity.

Criteria that have been proposed to fix this problem are:

- Limit support in time domain (smooth spectrum)
- Exploit frequency co-modulation
- Impose geometric constraints

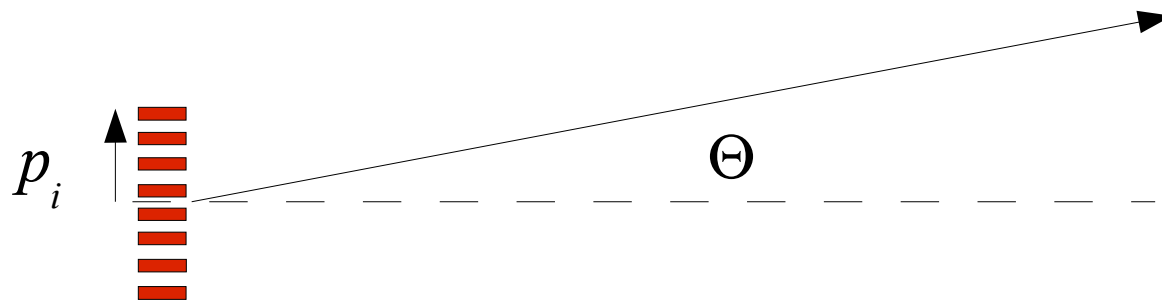
Farfield Beam Response

In order to interpret the resulting filter geometrically consider the farfield array response $\mathbf{r}(\Theta, \omega)$ for filters $\mathbf{W}(\omega)$:

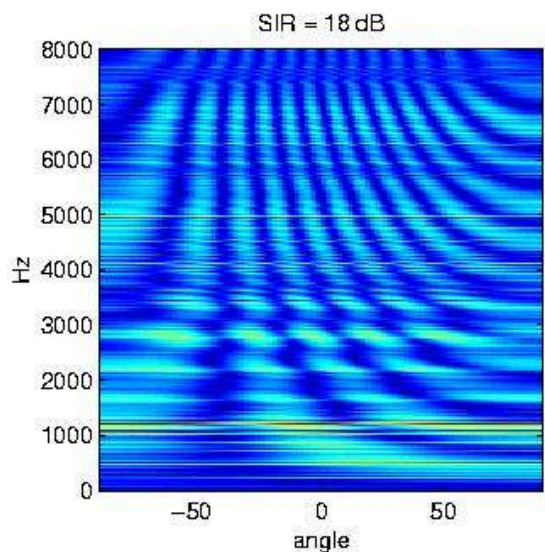
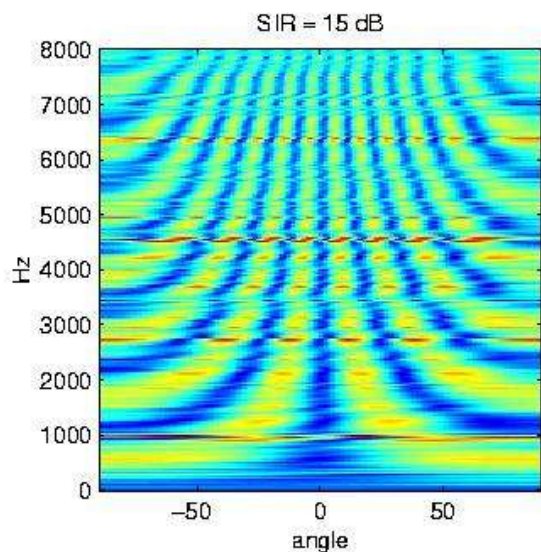
$$\mathbf{r}(\Theta, \omega) = \mathbf{W}(\omega) \mathbf{d}(\Theta, \omega)$$

Array vector for microphones at positions p_i :

$$d_i(\Theta, \omega) = \exp\left(j \frac{2\pi\omega}{f} \frac{p_i}{c} \sin(\Theta)\right)$$

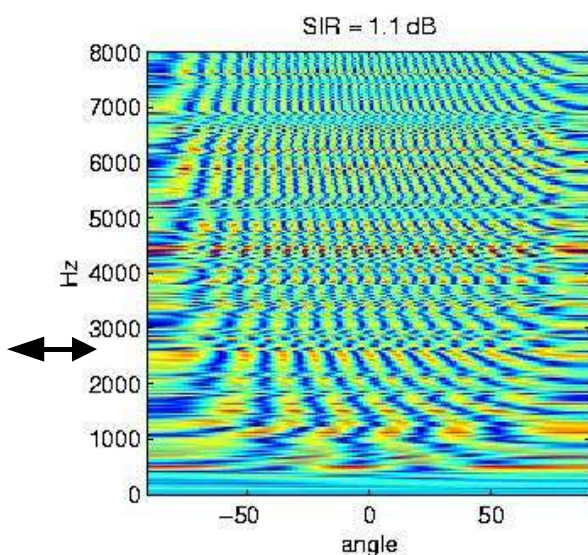
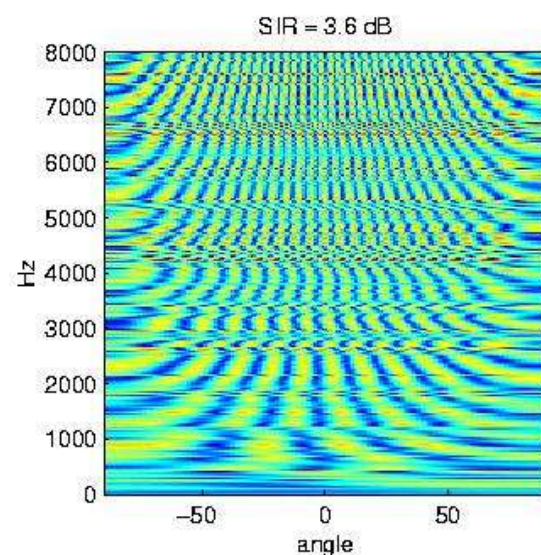


Example Beam Patterns - 2 Microphones



Two users at 2.5 m distance from array and 2 mics with 30 cm aperture.

Notice nulls consistent across frequency bands at user positions 0° and -45° .



Same signal and user positions with 70cm aperture.

Notice "permutation" in lower frequency band.

Numerical kurtosis measured on 15s of data is higher than above.

Beamforming: Geometry as Linear Constraints

To avoid trivial solution, $W(\omega)=0$, normalization is required, e.g.

In MSC and often in BSS we keep gain of one microphone constant:

$$W(\omega) \mathbf{e}_i = 1$$

In linearly constraint minimum variance (LCMV) or generalized sidelobe canceling (GSC) keep one orientation constant:

$$W(\omega) \mathbf{d}(\omega, \Theta) = 1$$

Or in some robust adaptive beam-forming the angle response is constraint to change slowly with direction:

$$W(\omega) \partial_{\Theta} \mathbf{d}(\omega, \Theta) = 0$$

Power vs. Cross-Power minimization

Second order separation and adaptive beamforming have similar criteria:

Beam power and cross-power is given by

$$R_{yy}(\omega) = W(\omega) R_{xx}(\omega) W^H(\omega)$$

Adaptive beam forming minimizes power - *does not allow crosstalk*

$$\underset{W(\omega)}{\operatorname{argmin}} \operatorname{diag}(R_{yy}(\omega))$$

Decorrelation algorithms minimizes cross-power - *allows crosstalk*

$$\underset{W(\omega)}{\operatorname{argmin}} \left\| \operatorname{offdiag}(R_{yy}^t(\omega)) \right\|^2$$

Combination of BSS with Beamforming

Source Separation based on second order non-stationarity

Minimum Cross-Power over time:
$$\underset{W(\omega)}{\operatorname{argmin}} \sum_t \left\| \operatorname{offdiag} [R_{yy}^t(\omega)] \right\|^2$$

"Geometrically *initialized* Source Separation":

Delay-sum initialization
$$W_{init}(\omega) = \mathbf{d}^H(\omega, \Theta)$$

"Geometrically *constraint* Source Separation":

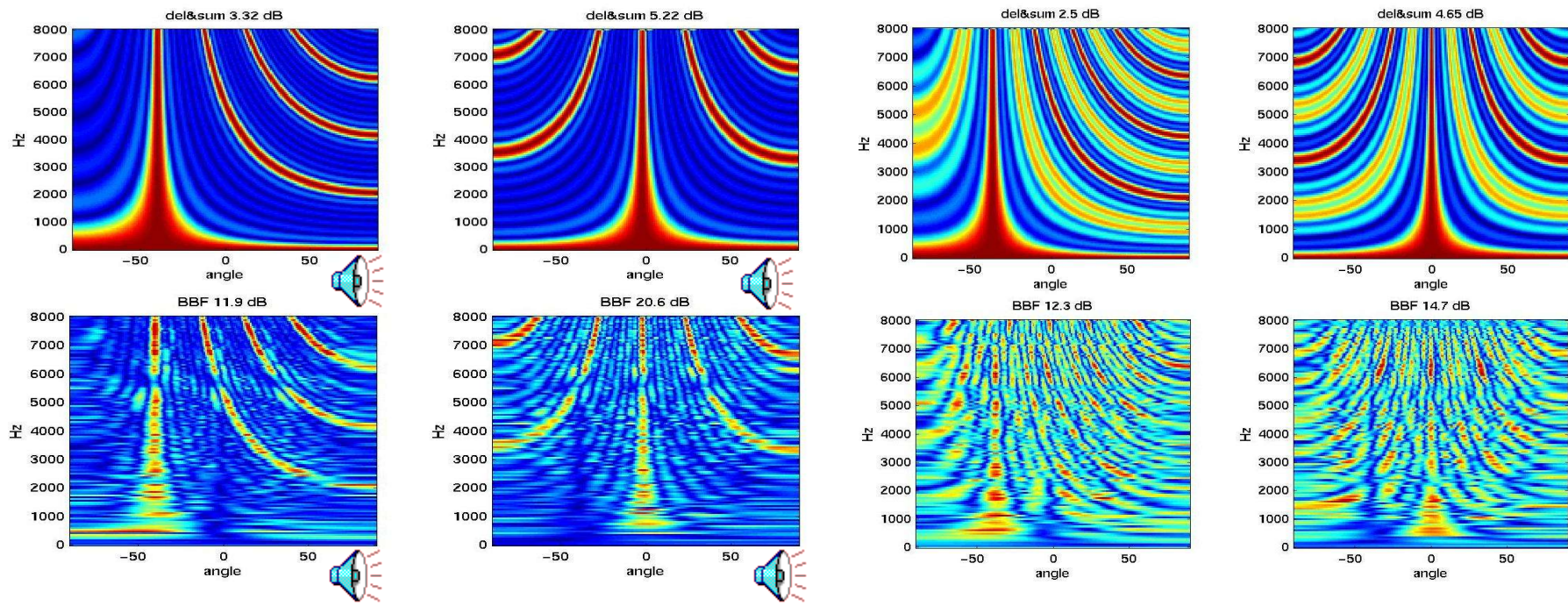
Unit-gain constrained in Θ
$$W(\omega) \mathbf{d}(\omega, \Theta) = 1$$

BSS with Geometric Initialization

Cross-power minimization across time can be made robust by initializing the filter structure with delay-sum beams, i.e. for each channel a maxima in the orientation of one of the sources.

8 microphones 

4 microphones

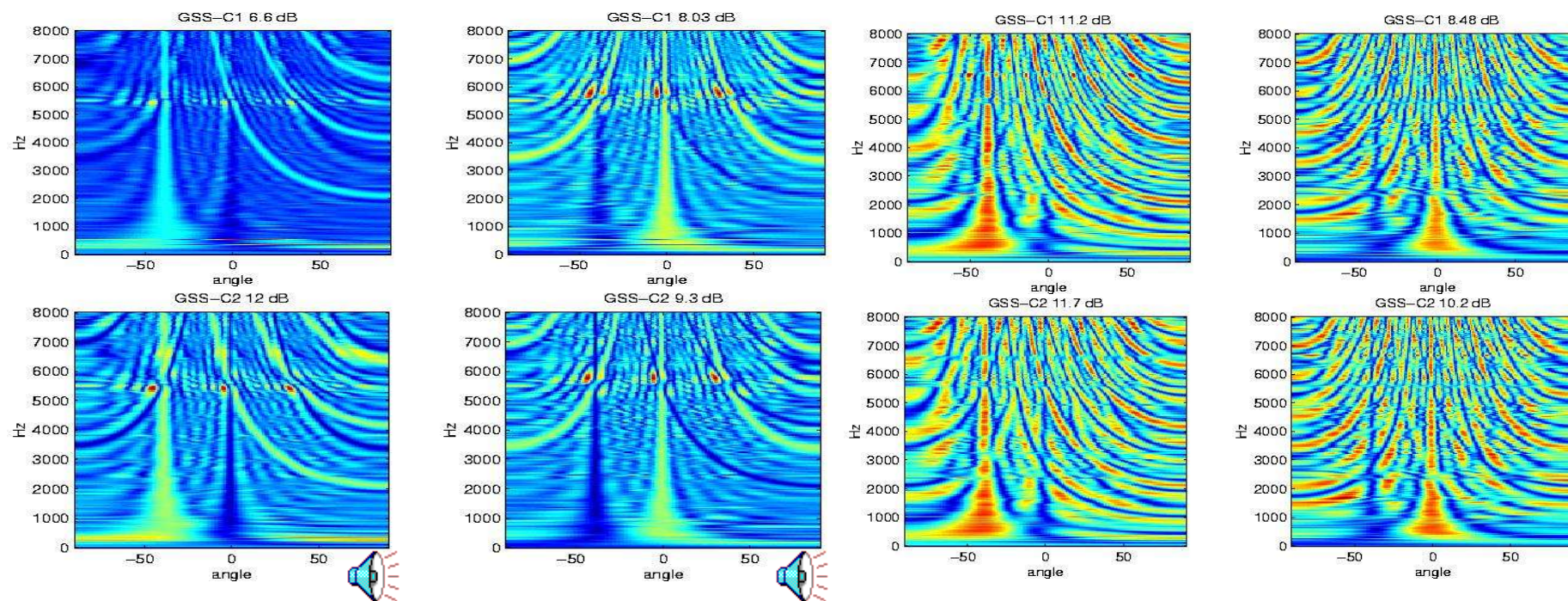


BSS with 'soft' Geometric Constraints

Cross-power minimization across time with 'soft' constraints on the response for given angles. Here we used constant unit-gain for user locations, and zero-gain for interference locations.

8 microphones 

4 microphones



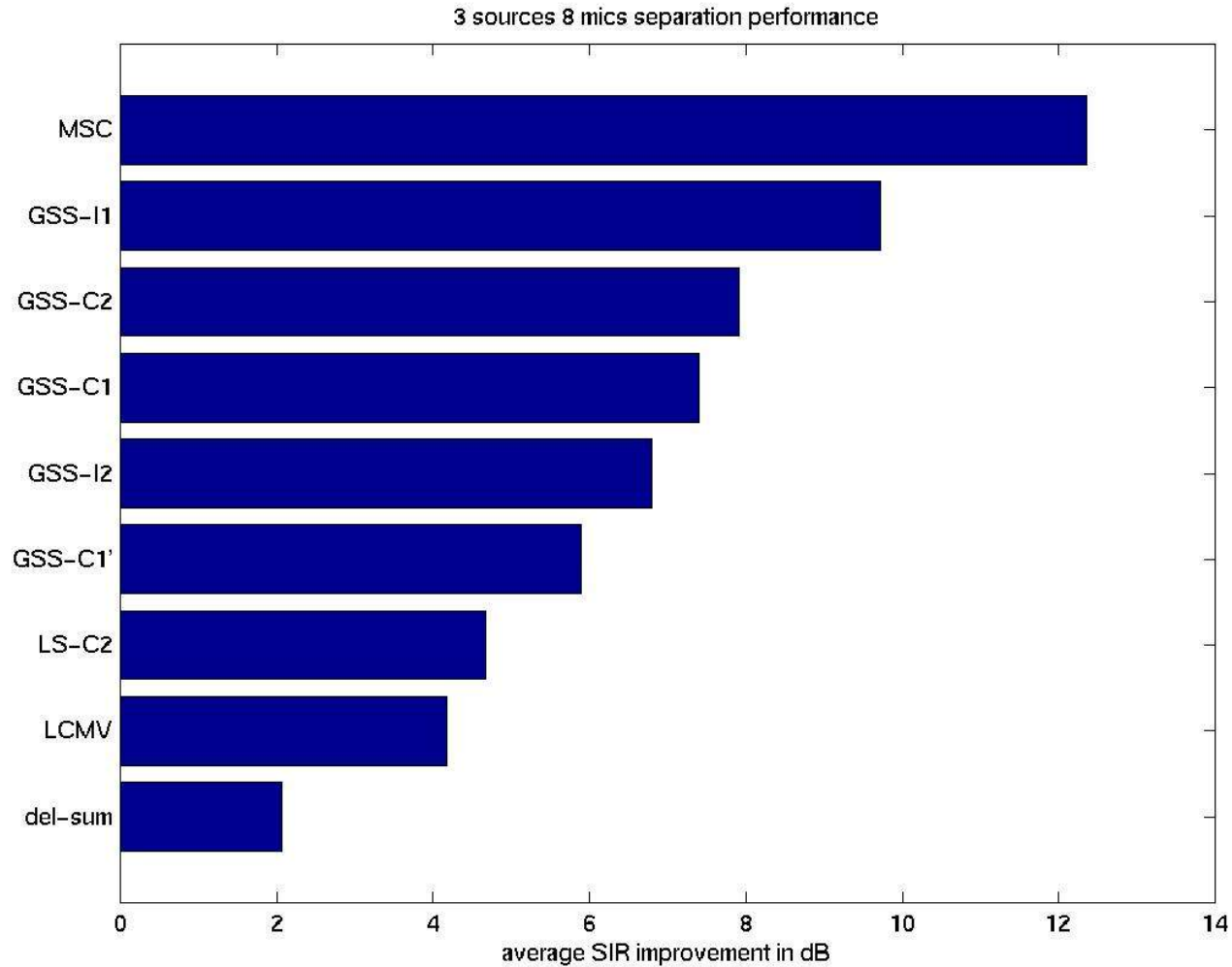
Acknowledgement

This work was performed 2000-2002 while working at



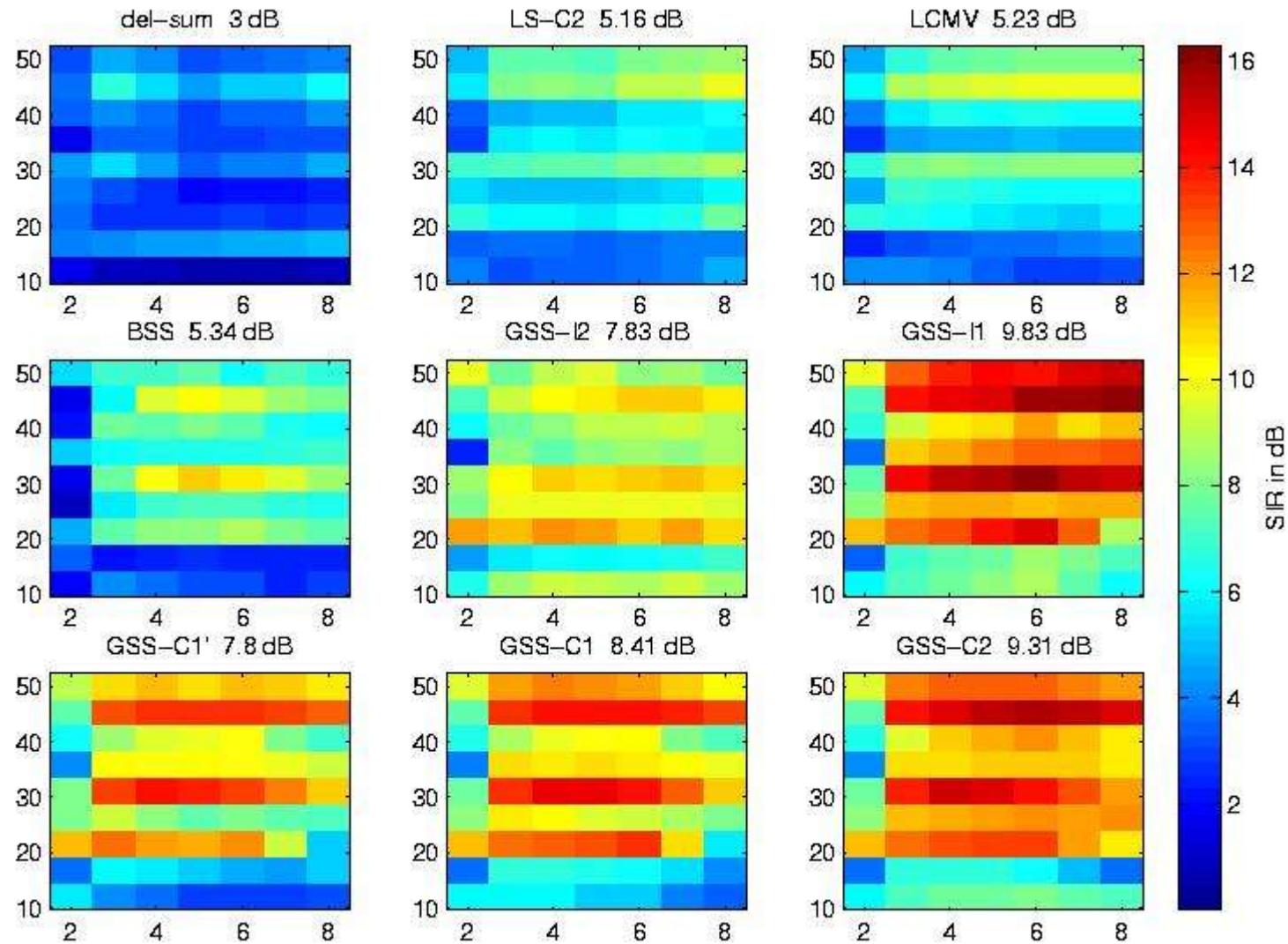
Clay Spence, Sarnoff Corporation
Chris Alvino, Rutgers
Craig Fancourt, Sarnoff Corporation

Performance Comparison 3 Sources



Real room performance comparison

... 2 sources for variable locations and number of microphones ...



Separation Based on Independence

Second order non-stationary

$$\mathbf{R}_x(t) = E[\mathbf{x}(t) \mathbf{x}^T(t)] = \mathbf{A} E[\mathbf{s}(t) \mathbf{s}^T(t)] \mathbf{A}^T = \mathbf{A} \mathbf{R}_s(t) \mathbf{A}^T$$

Second order non-white

$$\mathbf{R}_x(l) = E[\mathbf{x}(t+l) \mathbf{x}^T(t)] = \mathbf{A} E[\mathbf{s}(t+l) \mathbf{s}^T(t)] \mathbf{A}^T = \mathbf{A} \mathbf{R}_s(l) \mathbf{A}^T$$

4th order non-Gaussian

Assuming independence one can show that for any $k, l, i \neq j$

$$\begin{aligned} \text{Cum}(s_i, s_j, s_k, s_l) &= E[s_i s_j s_k s_l] - E[s_i s_j] E[s_k s_l] \\ &\quad - E[s_i s_k] E[s_j s_l] - E[s_i s_j] E[s_j s_k] = 0 \end{aligned}$$

Hence for any $i \neq j$ $c_{i,j}(\mathbf{M}) = \sum_{kl} \text{Cum}(s_i, s_j, s_k, s_l) m_{kl} = 0$

or in matrix notation $\mathbf{C}_s(\mathbf{M}) = \text{diag}$

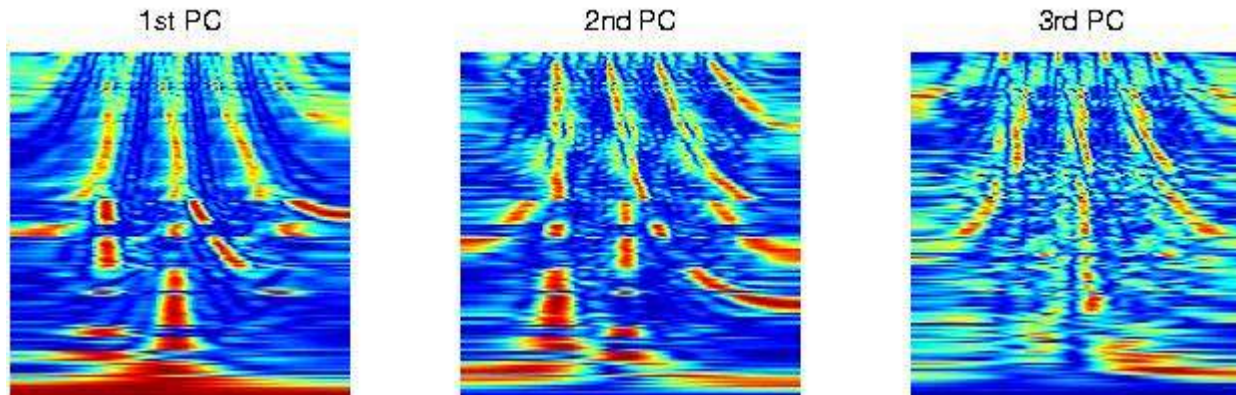
after some algebra $\mathbf{C}_x(\mathbf{I}) = \mathbf{A} \mathbf{C}_s(\mathbf{A}^T \mathbf{A}) \mathbf{A}^T$

Source Subspace

... extract signal from noise from additional microphones ...

$N \times M$ dimensional space of equivalent solutions for $W(\square)$ because there are more observations than signals. Signal exists in a M dimensional subspace.

Recover source subspace with PCA for each $\square \circ$



Problem:

- Orthogonality within source subspace may be meaningless in terms of beam geometry.
- Sources of interest may have small powers.