

# Blind Source Separation of Audio Signals

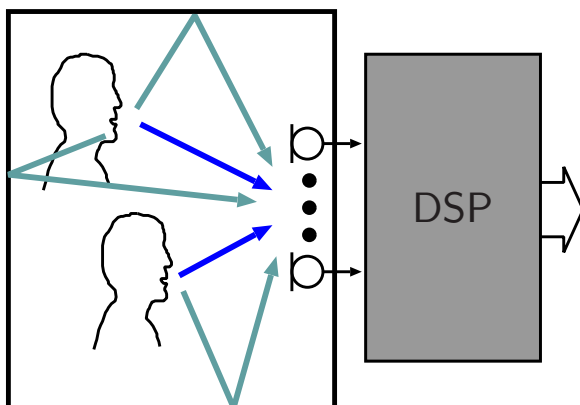
Herbert Buchner  
Technische Universität Berlin  
Lecture Machine Learning II

May 3, 2011

## Introduction - Broadband Adaptive MIMO Filtering

**Motivation:** Signal acquisition by sensor arrays in convolutive environments

**Example:** Speech capture by microphone arrays in reverberant rooms

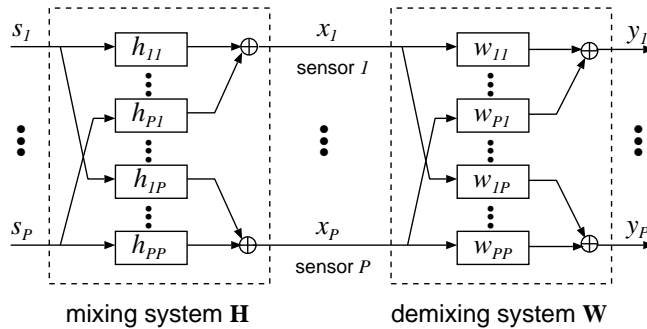


**Applications:** Teleconferencing, hands-free speech recognition, modern hearing aids,...

# Introduction - Broadband Adaptive MIMO Filtering

**Motivation:** Signal aquisition by sensor arrays in convolutive environments

## MIMO model



$$x_p(n) = \sum_{q=1}^P \sum_{\kappa=0}^{M-1} h_{qp}(\kappa) s_q(n - \kappa),$$

$$y_q(n) = \sum_{p=1}^P \sum_{\kappa=0}^{L-1} w_{pq}(\kappa) x_p(n - \kappa)$$

## Adaptive Signal Processing Tasks:

- **Signal Separation**

$$\check{\mathbf{y}} = \check{\mathbf{W}} * \check{\mathbf{H}} * \check{\mathbf{s}} \stackrel{!}{=} \text{diag}\{\check{\mathbf{W}} * \check{\mathbf{H}}\} * \check{\mathbf{s}}$$

- **Deconvolution:**

$$\check{\mathbf{y}} = \check{\mathbf{W}} * \check{\mathbf{H}} * \check{\mathbf{s}} \stackrel{!}{=} \alpha \cdot \delta(n - n_0) \mathbf{I} * \check{\mathbf{s}}$$

- **System Identification:**

estimate system  $\check{\mathbf{H}}$  using signals  $x_p$

- **Blind Estimation:** Propagation paths **and** original source signals are unknown
- **Supervised Estimation:** (Some) source signals and/or side info on paths are known

# Introduction - Broadband Adaptive MIMO Filtering

## Classification of the linear adaptive filtering problems

	supervised adaptive filtering problems	blind adaptive filtering problems
"direct adaptive filtering problems"	system identification  interference cancellation	blind system identification  blind source separation/ blind interference cancellation
"inverse adaptive filtering problems"	inverse modeling/equalization  linear prediction	blind (partial) deconvolution  linear prediction

## Semi-blind adaptive filtering:

**prominent example: adaptive beamforming** (=spatial filtering).

Can typically be implemented using (supervised) interference cancellation, but requires prior information on source locations for beamsteering.

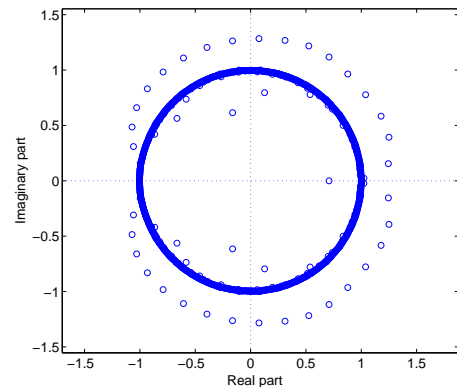
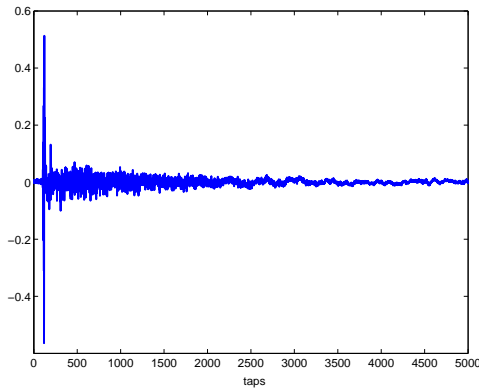
# Acoustic Environment

- **Reverberation time  $T_{60}$**   
(sound energy decayed by 60dB)
  - ▷ car  $\approx 50\text{ms}$
  - ▷ concert halls  $\approx 1 \dots 2\text{s}$

- **FIR models**

- ▷ typically  $L_H = T_{60} \cdot f_s / 3$  coefficients  
 $\Rightarrow$  hundreds...thousands of coeffs
- ▷ nonminimum-phase
- ▷ many zeros close to unit circle

Example: Office  $5.5\text{m} \times 3\text{m} \times 2.8\text{m}$ ,  $T_{60} \approx 300\text{msec}$ , sampling frequency  $f_s = 12\text{kHz}$ .



- **Excitation by speech and audio signals:**  
 $\Rightarrow$  nonwhiteness, nonstationarity, nongaussianity

## Overview

- **TRINICON - A Generic Concept for Adaptive MIMO Filtering**
  - ▷ Optimization Criterion and Generic Adaptation Algorithms
  - ▷ Incorporation of Stochastic Source Models
- **Applications to Signal Processing Problems for Speech Capture**
  - ▷ Blind Source Separation (BSS) / Interference Suppression
  - ▷ Supervised Separation and System Identification / Echo Cancellation
  - ▷ Blind System Identification (BSI) / Localization of Multiple Sources
  - ▷ \* Multich. Blind Deconvolution (MCBD) and  
Multich. Blind Partial Deconvolution (MCBPD) / Dereverberation
- **Concluding Remarks**

# TRINICON - Optimization Criterion

As a generic framework, **TRINICON** [Buchner et al., 2003-] exploits

**TRI**ple **N**, i.e., **N**onwhiteness, **N**onstationarity, and **N**ongaussianity of  $s_q(n)$ ,  
for **I**ndependent component analysis of **CON**vulsive mixtures

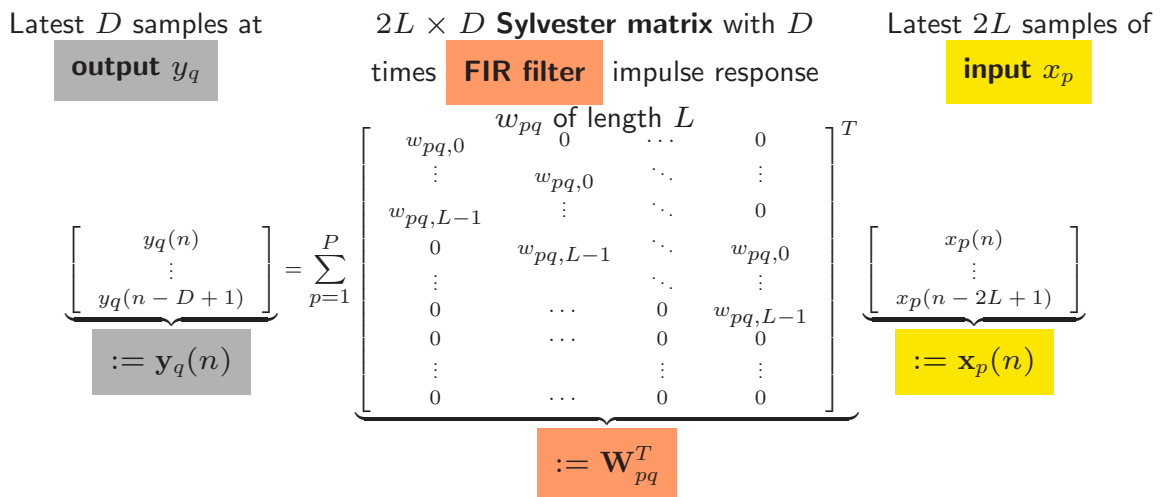
## Optimization Criterion

$$\mathcal{J}(m, \mathbf{W}) = - \sum_{i=0}^{\infty} \beta(i, m) \frac{1}{N} \sum_{j=iN_L}^{iN_L+N-1} \{ \log(\hat{p}_{s,PD}(\mathbf{y}(j))) - \log(\hat{p}_{y,PD}(\mathbf{y}(j))) \}$$

- **Nongaussianity** by minimizing Kullback-Leibler divergence (KLD) between the  $PD$ -variate probability density functions with data-dependent parameterizations
  - ▷  $\hat{p}_{s,PD}(\mathbf{y})$  for sources (*desired*)
  - ▷  $\hat{p}_{y,PD}(\mathbf{y})$  for outputs (*actual*)
- **Nonwhiteness** by simultaneous minimization for  $D$  time-lags (blocks of  $D$  output values per channel in vector  $\mathbf{y}$ )
- **Nonstationarity** by simultaneous minimization of  $N$  blocks (of  $D$  output samples per channel)

**Windowing**  $\beta(i, m)$  defines online, block-online, or offline adaptation

# TRINICON - MIMO Filtering of Broadband Signals



**Input/output of complete MIMO system:**

$$\underbrace{\begin{bmatrix} \mathbf{y}_1(n) \\ \vdots \\ \mathbf{y}_P(n) \end{bmatrix}}_{\mathbf{y}(n)} = \underbrace{\begin{bmatrix} \mathbf{W}_{11} & \cdots & \mathbf{W}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{W}_{P1} & \cdots & \mathbf{W}_{PP} \end{bmatrix}^T}_{\mathbf{W}^T} \cdot \underbrace{\begin{bmatrix} \mathbf{x}_1(n) \\ \vdots \\ \mathbf{x}_P(n) \end{bmatrix}}_{\mathbf{x}(n)}$$

## TRINICON - Euclidean Gradient-Based Coefficient Update

- Block-online update rule:  $\check{\mathbf{W}}^0(m) := \check{\mathbf{W}}(m-1)$ ,  
 $\check{\mathbf{W}}^\ell(m) = \check{\mathbf{W}}^{\ell-1}(m) - \mu \Delta \check{\mathbf{W}}^\ell(m), \quad \ell = 1, \dots, \ell_{\max},$   
 $\check{\mathbf{W}}(m) := \check{\mathbf{W}}^{\ell_{\max}}(m)$
- Coefficient updates: Euclidean gradient of  $\mathcal{J}$  w.r.t.  $\check{\mathbf{W}}$  yields

$$\Delta \check{\mathbf{W}}^\ell(m) = \frac{1}{N} \sum_{i=0}^{\infty} \beta(i, m) \mathcal{SC} \left\{ \sum_{j=iL}^{iL+N-1} \left[ \mathbf{x}(j) \Phi_{s,PD}^T(\mathbf{y}(j)) - \left( (\mathbf{W}^{\ell-1}(m))^T \right)^+ \right] \right\}$$

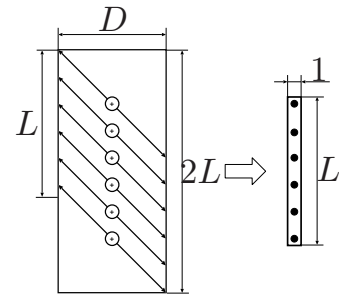
Selection of practical algorithms for specific applications by

- Sylvester constraint**  $\mathcal{SC}\{\bullet\}$  linking  $\mathbf{W}$  (in the cost function) and  $\check{\mathbf{W}}$  (in the optimization procedure).

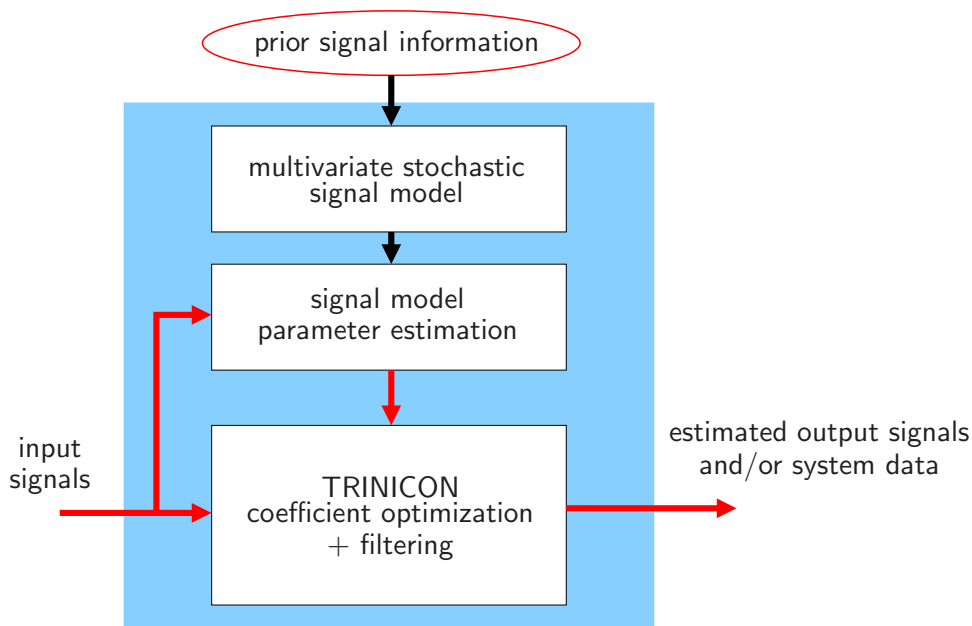
**General realization:** sums within Sylvester diagonals.

- multivariate score function**

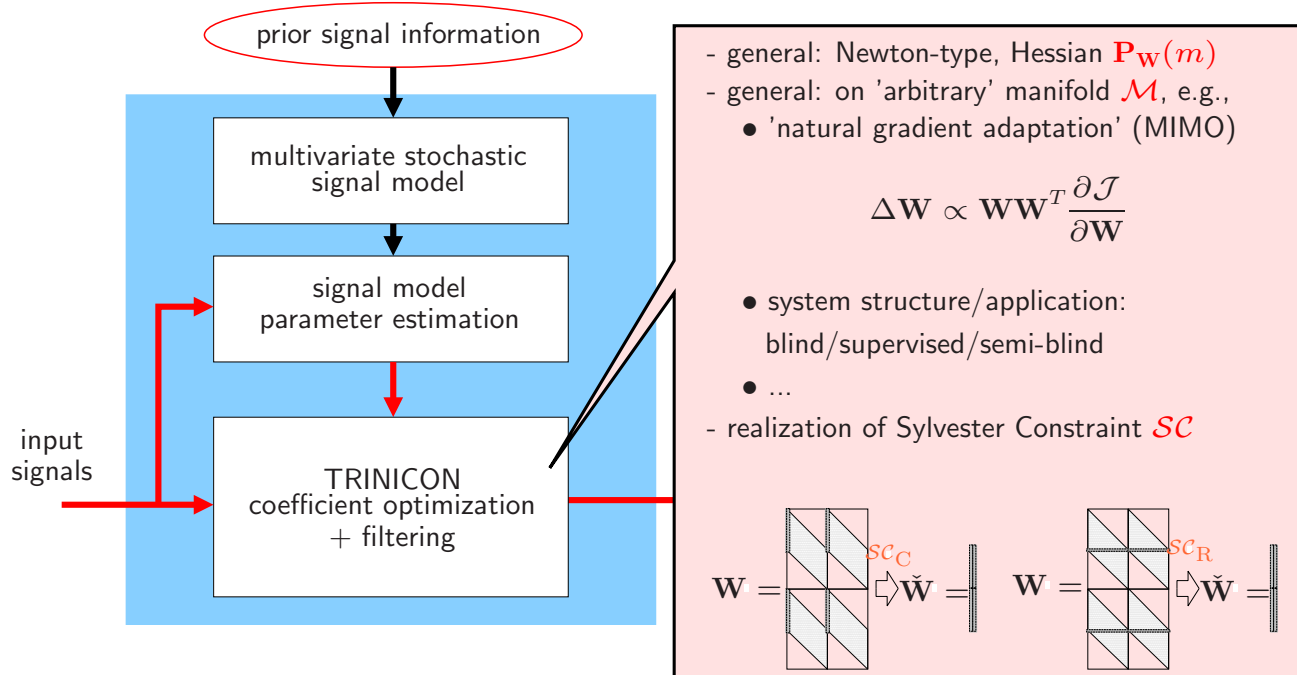
$$\Phi_{s,PD}(\mathbf{y}) = -\frac{\partial \log \hat{p}_{s,PD}(\mathbf{y})}{\partial \mathbf{y}}$$



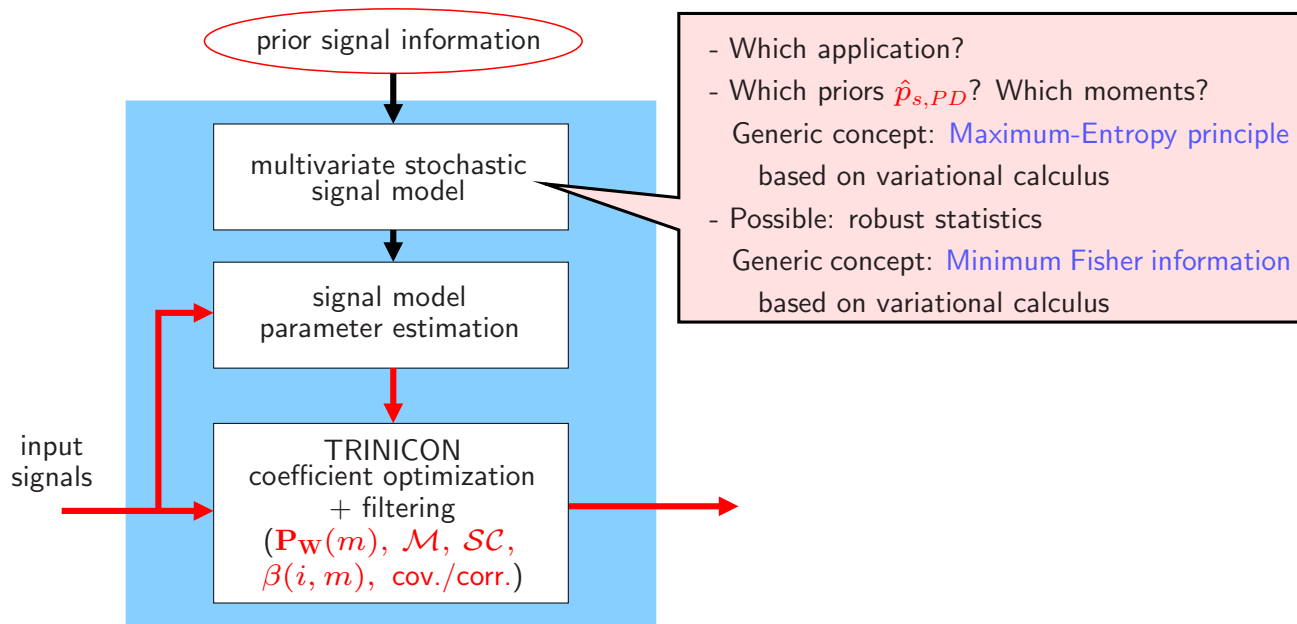
## TRINICON-Based Algorithm Design



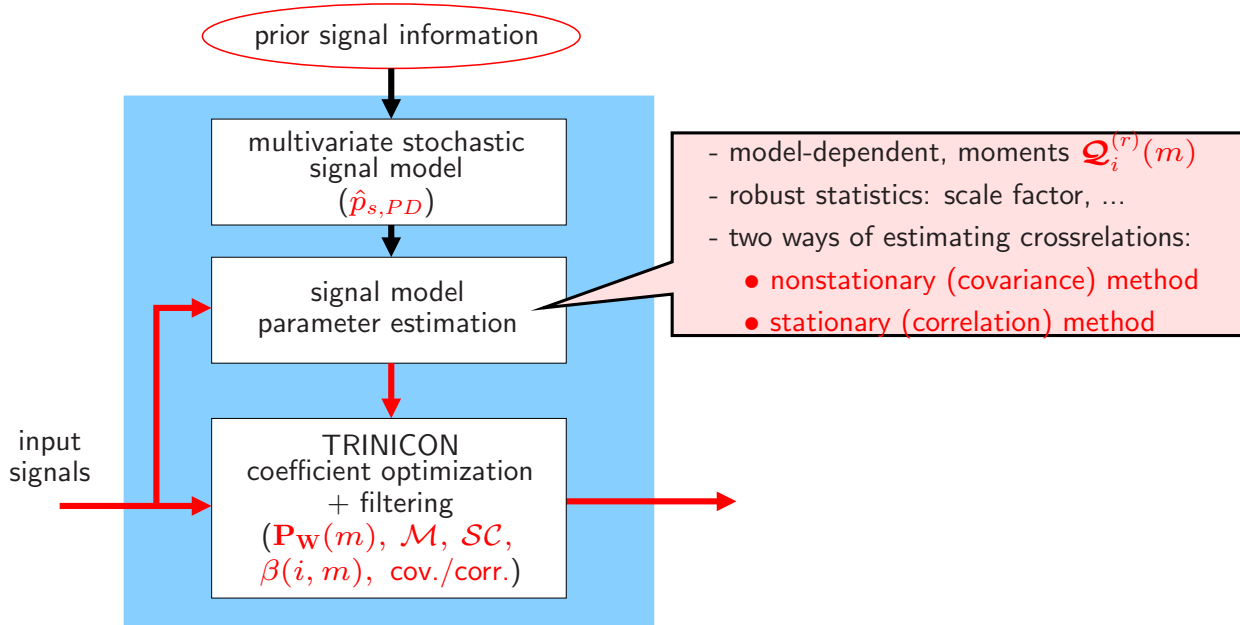
# TRINICON-Based Algorithm Design



# TRINICON-Based Algorithm Design



# TRINICON-Based Algorithm Design



## TRINICON - Exploiting HOS: SIRPs as Source Model

Spherically Invariant Random Processes (SIRPs) are described by multivariate pdfs of the form (with suitable function  $f_D$ ):

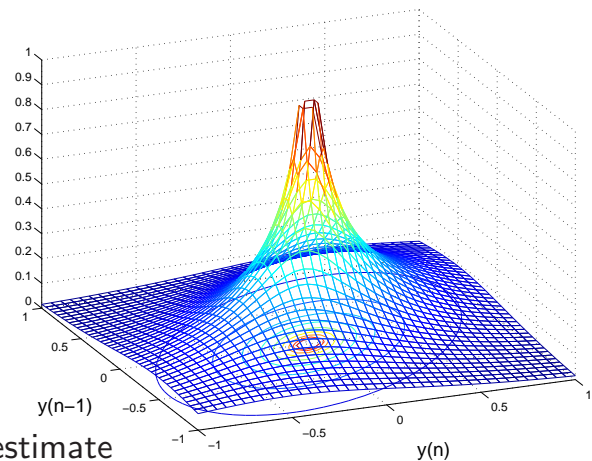
$$\hat{p}_D(\mathbf{y}_p) = \frac{1}{\sqrt{\pi^D \det(\hat{\mathbf{R}}_{pp})}} f_D \left( \mathbf{y}_p^T \hat{\mathbf{R}}_{pp}^{-1} \mathbf{y}_p \right)$$

Several attractive properties:

- Good model for speech signals
- Reduced number of model parameters to estimate
- Multivariate pdfs can be derived analytically from corresponding univariate pdfs
- Incorporation into TRINICON leads to inherent stepsize normalization of the update

equation:  $\Phi_{p,D}(\mathbf{y}_p(j)) = 2 \phi_{y_p,D} \left( \mathbf{y}_p^T(j) \mathbf{R}_{y_p y_p}^{-1}(i) \mathbf{y}_p(j) \right) \cdot \mathbf{R}_{y_p y_p}^{-1}(i) \mathbf{y}_p(j),$

*SIRP score:*  $\phi_{y_p,D}(u_p) = -\partial \log f_{p,D}(u_p) / \partial u_p \rightarrow$  *Gauss:*  $\phi_{y_p,D} = 1/2$



# Overview

- **TRINICON - A Generic Concept for Adaptive MIMO Filtering**
  - ▷ Optimization Criterion and Generic Adaptation Algorithms
  - ▷ Incorporation of Stochastic Source Models
- **Applications to Signal Processing Problems for Speech Capture**
  - ▷ Blind Source Separation (BSS) / Interference Suppression
  - ▷ Supervised Separation and System Identification / Echo Cancellation
  - ▷ Blind System Identification (BSI) / Localization of Multiple Sources
  - ▷ Multich. Blind Deconvolution (MCBD) and Multich. Blind Partial Deconvolution (MCBPD) / Dereverberation
- **Concluding Remarks**

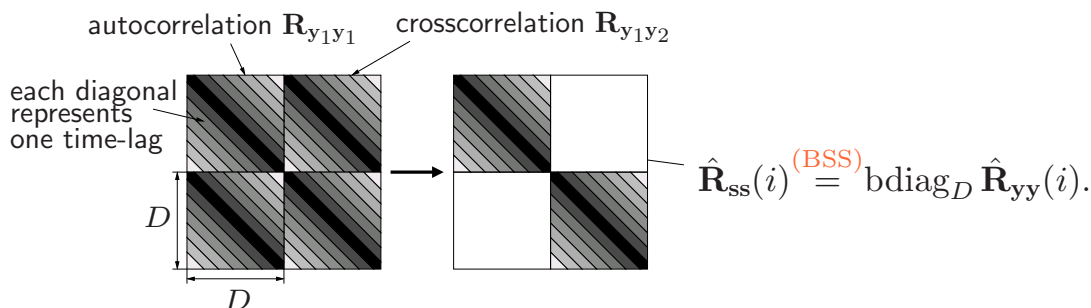
## TRINICON for Blind Source Separation

- **Desired score function for BSS:** Independence between channels

$$\hat{p}_{s,PD}(\mathbf{y}(j)) \stackrel{(\text{BSS})}{=} \prod_{q=1}^P \hat{p}_{y_q,D}(\mathbf{y}_q(j))$$

- **Illustration for SOS** (Gaussian source model  $\Rightarrow \Phi_{p,D}(\mathbf{y}_p(j)) = \mathbf{R}_{y_p y_p}^{-1}(i) \mathbf{y}_p(j)$ ): Natural gradient-based update

$$\Delta \check{\mathbf{W}}(m) = 2 \sum_{i=0}^{\infty} \beta(i, m) \mathcal{SC} \left\{ \mathbf{W}(i) \left\{ \hat{\mathbf{R}}_{yy} - \hat{\mathbf{R}}_{ss} \right\} \hat{\mathbf{R}}_{ss}^{-1} \right\}$$

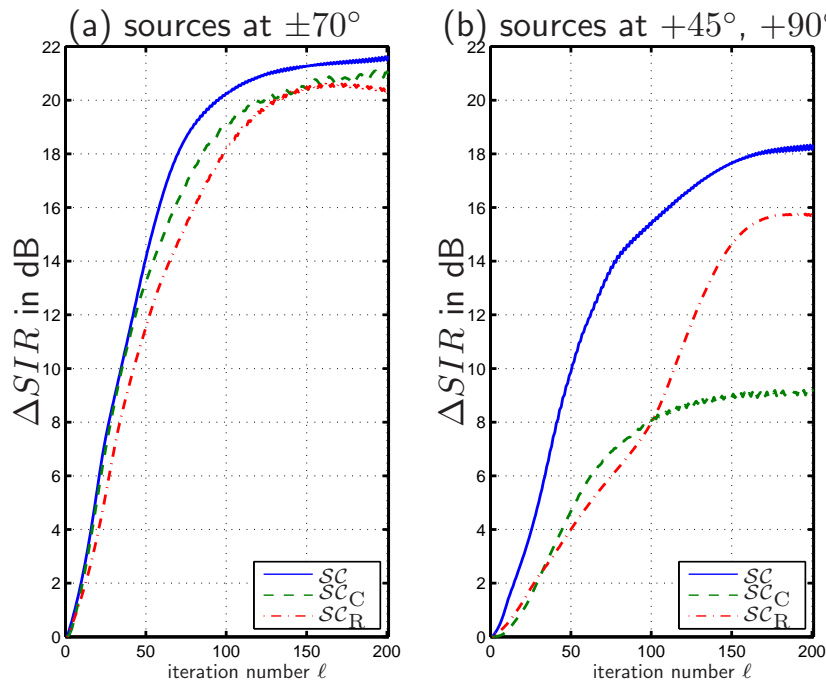


- **Inherent normalization ( $\rightarrow$  stepsize control)** in each channel by  $\text{bdiag}_D \hat{\mathbf{R}}_{yy}^{-1}$



## BSS - Results: Generic SOS + Sylvester Constraint

Offline BSS generic SOS natural gradient adaptation, cov. method,  $L = 256$

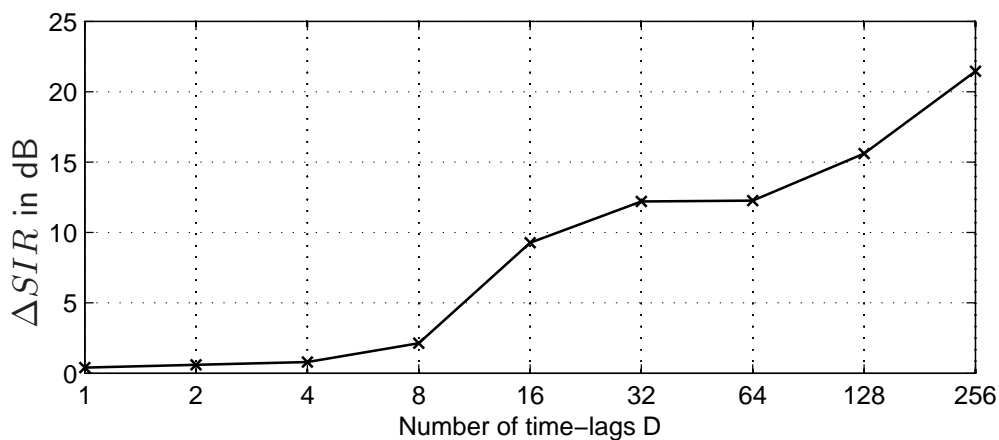


⇒ **Exact  $\mathcal{SC}$  combines versatility ( $\mathcal{SC}_R$ ) and robust convergence ( $\mathcal{SC}_C$ )**

## TRINICON for BSS - Exploitation of Nonwhiteness

**Influence of the number  $D$  of simultaneously optimized time lags for exploiting nonwhiteness**

Offline BSS generic SOS natural gradient adaptation,  $\mathcal{SC}_R$ , cov. method,  $L = 256$

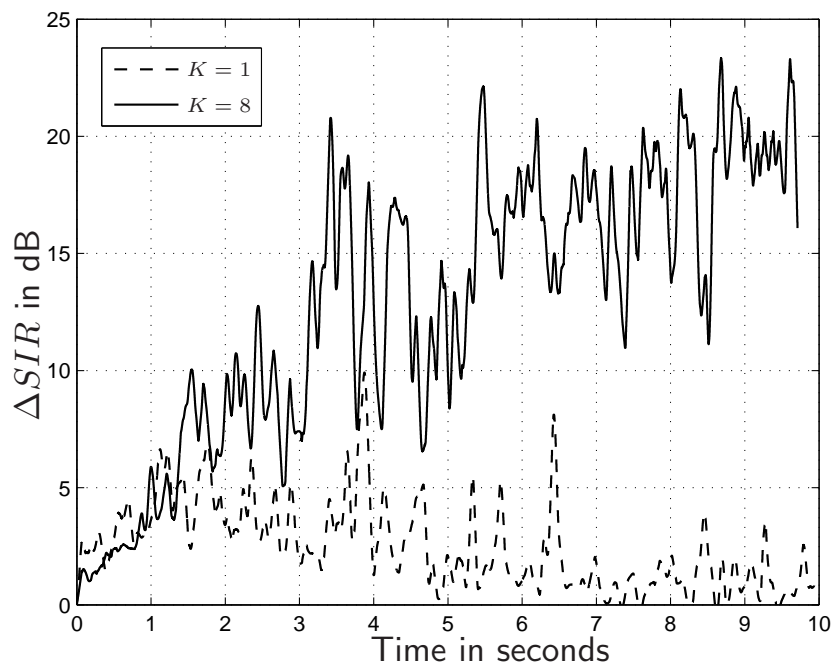


(no further improvement can be achieved for  $D > L = 256$ )

# TRINICON for BSS - Exploitation of Nonstationarity

## Influence of the number $K$ of simultaneously optimized blocks

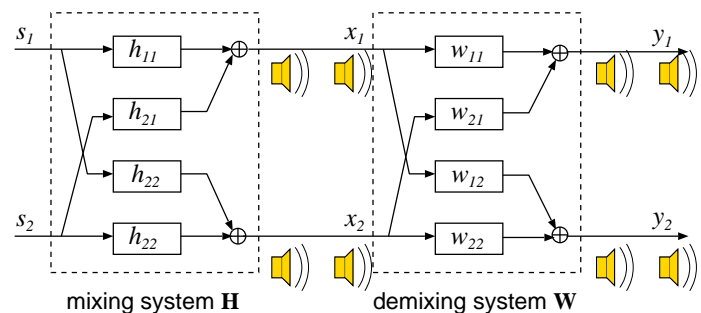
Block-offline BSS generic SOS natural gradient adaptation,  $\mathcal{SC}_R$ , cov. method,  $L = 256$



## TRINICON for BSS - Results

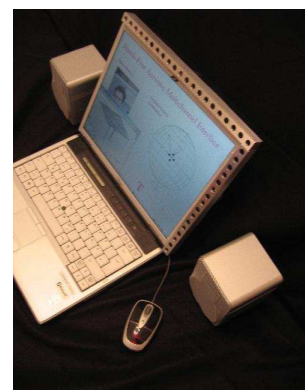
- Example: Convergence behaviour

- ▷ SOS-BSS with diagonal power normalization
- ▷  $f_S = 16\text{kHz}$
- ▷  $T_{60} = 50, 200\text{msec}$
- ▷  $D = L = 1024$
- ▷ block-online adaptation

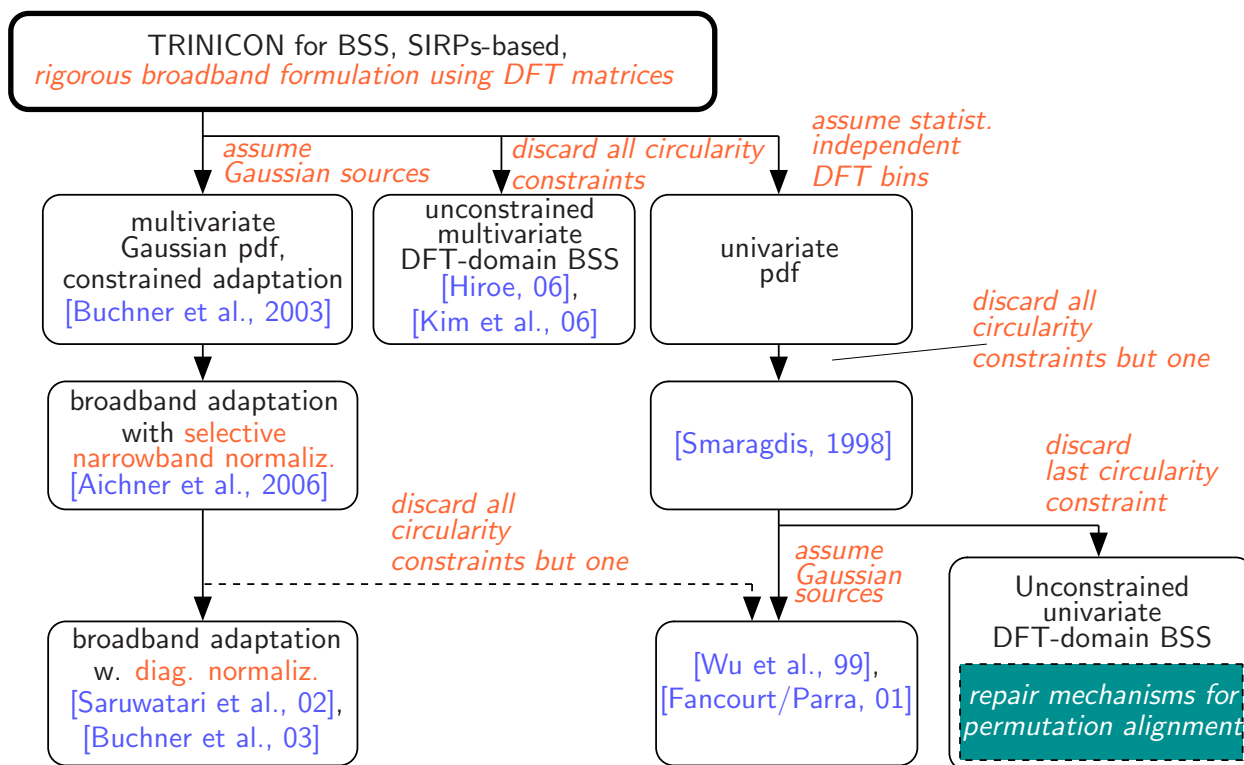


- Real-time implementation

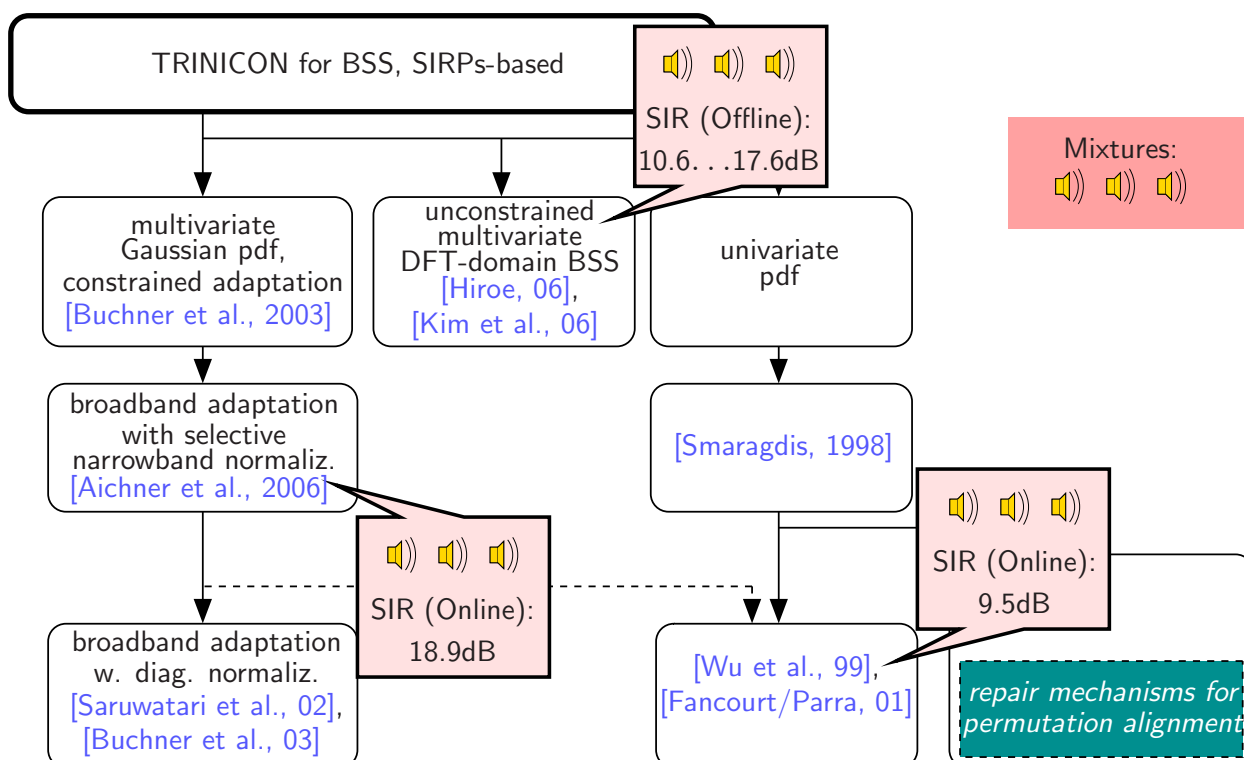
Implementation in the DFT domain equivalent to time-domain  
↔ Broadband algorithm, no internal permutation or circular convolution



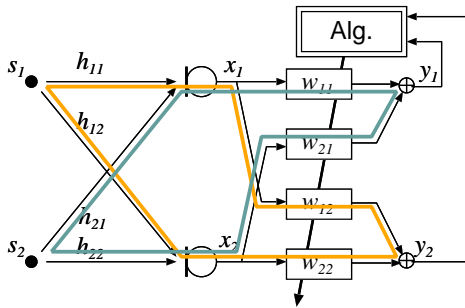
## TRINICON for BSS - Efficient Frequency-Domain Realizations



## TRINICON for BSS - Efficient Frequency-Domain Realizations



# TRINICON for Blind System Identification



Equilibria of overall system  $\mathbf{C} = \mathbf{H}\mathbf{W}$  in Sylvester structure (Buchner et al., 2005):

$$\text{boff}\{\mathbf{C}\} = \mathbf{0} \quad \text{iff } D = L$$

$$\Rightarrow \begin{aligned} h_{11} * w_{12} &= -h_{12} * w_{22} \\ h_{21} * w_{11} &= -h_{22} * w_{21} \end{aligned}$$

Ideal separation filters: 
$$\begin{bmatrix} \mathbf{W}_{11} & \mathbf{W}_{12} \\ \mathbf{W}_{21} & \mathbf{W}_{22} \end{bmatrix} = \begin{bmatrix} \alpha_1 \mathbf{H}_{22} & -\alpha_2 \mathbf{H}_{12} \\ -\alpha_1 \mathbf{H}_{21} & \alpha_2 \mathbf{H}_{11} \end{bmatrix}$$

Unique solution up to scaling  $(\alpha_1, \alpha_2)$  iff

- No common zeros in  $H_{11}(z), H_{12}(z)$ /no common zeros in  $H_{21}(z), H_{22}(z)$
- Demixing filter length  $L \leq \text{length } M$  of mixing system

$\Rightarrow$  Application to localization of simultaneously active sources

## TRINICON for Localization $\rightarrow$ TDOA estimation

Simultaneous Localization of **Multiple** Sound Sources in **Reverberant** Environments

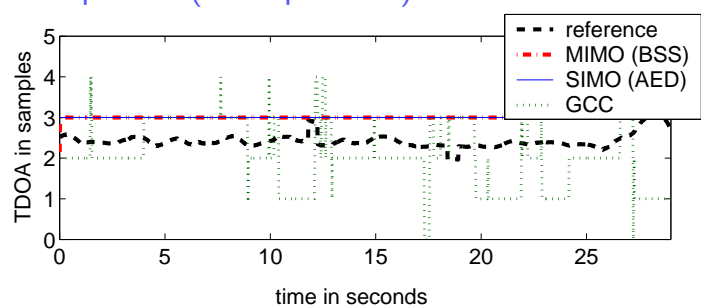
### Setup:

Two speakers recorded in a TV studio,  $T_{60} \approx 700\text{ms}$ ,  $f_s = 48\text{kHz}$

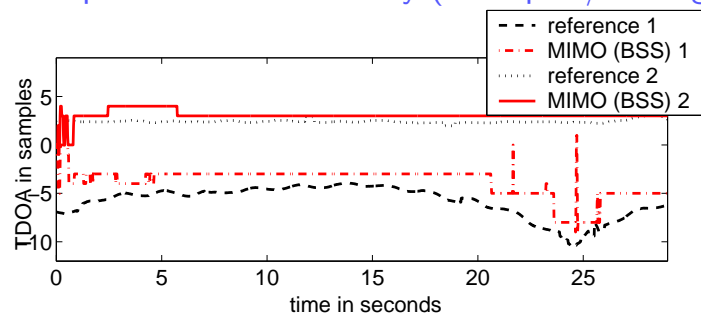
### TDOA estimation:

- Generalized cross-correlation (GCC) with phase-transform (PHAT) weighting + VAD
- **SIMO-based BSI** + VAD (AED, Benesty et al., 1999)
- **MIMO-based BSI** (may be seen as a generalization of AED)

one speaker (fixed position)

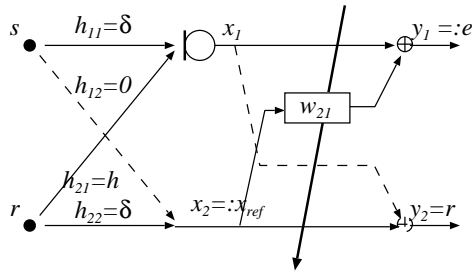


two speakers simultaneously (fixed pos./moving)



# TRINICON for Supervised Adaptive Filtering Problems

Example: acoustic echo cancellation (AEC) problem from a specialized mixing model



Ideal filters:

$$\begin{bmatrix} \mathbf{w}_{11} & \mathbf{w}_{12} \\ \mathbf{w}_{21} & \mathbf{w}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{h}_{22} & -\mathbf{h}_{12} \\ -\mathbf{h}_{21} & \mathbf{h}_{11} \end{bmatrix}$$

Here:

$$\begin{bmatrix} \mathbf{w}_{11} & \mathbf{w}_{12} \\ \mathbf{w}_{21} & \mathbf{w}_{22} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_1 & \mathbf{0} \\ -\mathbf{h} & \mathbf{1}_1 \end{bmatrix}$$

Lower left sub-matrix of simple gradient-based SIRPs TRINICON update:

$$\hat{\mathbf{h}}^\ell(m) = \hat{\mathbf{h}}^{\ell-1}(m) + \frac{\mu}{N} \sum_{i=0}^{\infty} \beta(i, m) \mathcal{SC} \left\{ \sum_{j=iL}^{iL+N-1} \mathbf{x}_{\text{ref}}(j) \mathbf{e}^T(j) \mathbf{R}_{\text{ee}}^{-1}(i) \underbrace{\phi_{e,D}(\mathbf{e}^T(j) \mathbf{R}_{\text{ee}}^{-1}(i) \mathbf{e}(j))}_{\text{'SIRP score'}} \right\}$$

⇒ Generalization of Least-Mean-Squares (LMS): **inherent 'stepsize control'**

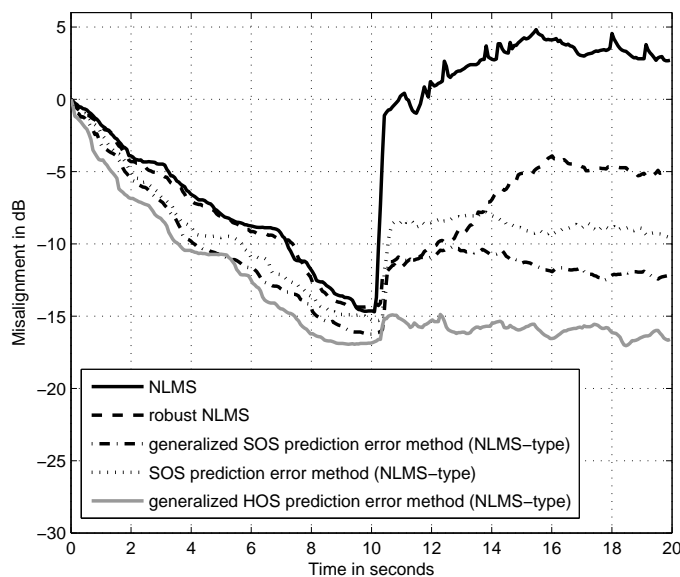
- **HOS case:** nongaussianity of local speech, **incorporation of 'robust statistics'**
- **Analogously:** generalization of other supervised algorithms (NLMS, Newton, RLS,...)

# TRINICON for Supervised Adaptive Filtering Problems

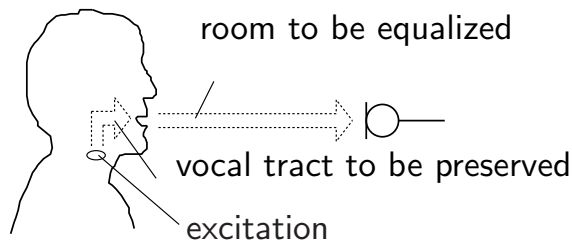
**Misalignment convergence ( $\|\mathbf{h} - \hat{\mathbf{h}}\|^2 / \|\mathbf{h}\|^2$ ) for gradient-based updates with NLMS-type normalization**

$L = 1024$ ,  $f_s = 16\text{kHz}$ , without double-talk detector

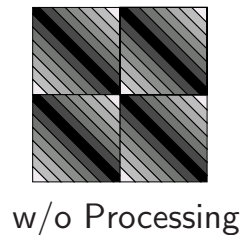
Single talk during the first 10sec, double talk starts after 10sec



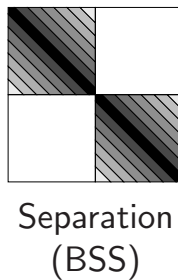
# TRINICON for Dereverberation = Partial Deconvolution



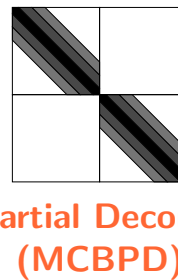
## Desired output statistics - illustration for SOS:



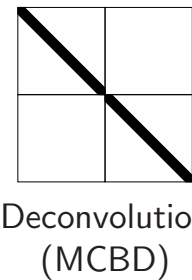
w/o Processing



Separation  
(BSS)



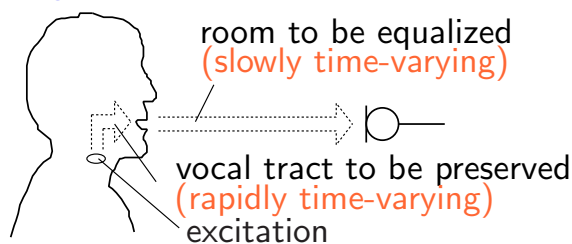
Partial Deconv.  
(MCBPD)



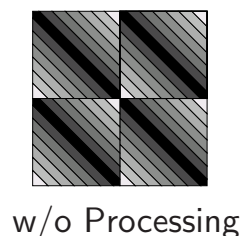
Deconvolution  
(MCBD)

# TRINICON for Dereverberation = Partial Deconvolution

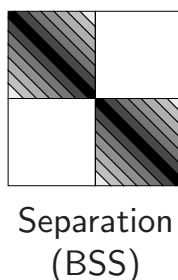
## Nonstationary:



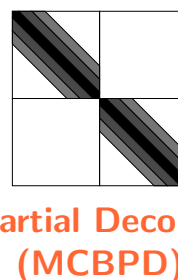
## Nonwhiteness:



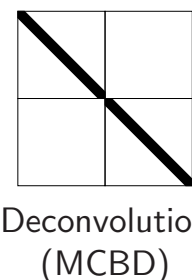
w/o Processing



Separation  
(BSS)



Partial Deconv.  
(MCBPD)



Deconvolution  
(MCBD)

## Nongaussianity:

- Probability density of speech signal: supergaussian
- Room acoustics described by convolutional sum → mic signals closer to Gaussian
- Aim: maximize nongaussianity of demixing filter outputs → e.g., max. **kurtosis**

# TRINICON - Nearly Gaussian Densities as Source Model

Expansions based on Chebyshev-Hermite polynomials  $P_{H,n}(\cdot)$

**Here: Gram-Charlier expansion**

**Univariate Example:** fourth-order approximation for a zero-mean process

$$\hat{p}_{y,1}(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-y^2/2\sigma^2} \left( 1 + \underbrace{\frac{\kappa_3}{3!\sigma^3} P_{H,3}\left(\frac{y}{\sigma}\right)}_{\kappa_3=\text{skewness, negligible}} + \underbrace{\frac{\kappa_4}{4!\sigma^4} P_{H,4}\left(\frac{y}{\sigma}\right)}_{\kappa_4=\text{kurtosis, nongaussianity}} \right)$$

# TRINICON - Nearly Gaussian Densities as Source Model

Expansions based on Chebyshev-Hermite polynomials  $P_{H,n}(\cdot)$

**Here: Gram-Charlier expansion**

**Multivariate Case:**

$$\begin{aligned} \hat{p}_{y_p,D}(\mathbf{y}_p(j)) &= \frac{1}{\sqrt{(2\pi)^D \det \mathbf{R}_{y_p y_p}(i)}} e^{-\frac{1}{2} \mathbf{y}_p^T(j) \mathbf{R}_{y_p y_p}^{-1}(i) \mathbf{y}_p(j)} \\ &\cdot \sum_{n_1=0}^{\infty} \cdots \sum_{n_D=0}^{\infty} a_{n_1 \dots n_D, p} P_{H,n_1} \left( [\mathbf{L}_p^{-1}(i) \mathbf{y}_p(j)]_1 \right) \cdots P_{H,n_D} \left( [\mathbf{L}_p^{-1}(i) \mathbf{y}_p(j)]_D \right) \end{aligned}$$

**For speech signals:** we introduce the factorization

$$\mathbf{R}_{y_p y_p}^{-1}(i) = \mathbf{A}_p(i) \mathbf{\Sigma}_{\tilde{\mathbf{y}}_p \tilde{\mathbf{y}}_p}^{-1}(i) \mathbf{A}_p^T(i)$$

Unit lower triangular matrix  $\mathbf{A}_p(i) \rightarrow$  interpreted as a whitening convolution matrix

- model  $y_p(n)$  as an AR process of order  $n_A = D - 1$
- shift prefilter matrix into the data terms:

$$\tilde{\mathbf{y}}_p := \mathbf{A}_p^T \mathbf{y}_p = [\tilde{y}_p(n), \tilde{y}_p(n-1), \dots, \tilde{y}_p(n-D+1)]^T \quad (0)$$

# TRINICON - Nearly Gaussian Densities as Source Model

⇒ **Multivariate speech model based on fourth-order approximation:**

$$\hat{p}_{y_p, D}(\mathbf{y}_p(j)) = \prod_{d=1}^D \frac{1}{\sqrt{2\pi} \hat{\sigma}_{\tilde{y}_p}^2(j-d+1)} e^{-\frac{\tilde{y}_p^2(j-d+1)}{2\hat{\sigma}_{\tilde{y}_p}^2(j-d+1)}} \left( 1 + \frac{\hat{\kappa}_{4, \tilde{y}_p}}{4! \hat{\sigma}_{\tilde{y}_p}^4(j-d+1)} P_{H, n_d} \left( \frac{\tilde{y}_p(j-d+1)}{\hat{\sigma}_{\tilde{y}_p}(j-d+1)} \right) \right).$$

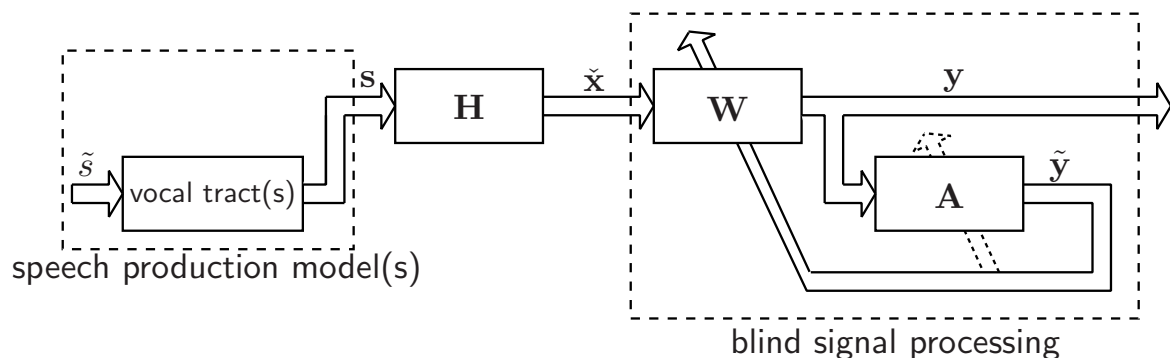
**Resulting generalized multivariate score function:**

$$\Phi_{y, PD}(\mathbf{y}(j)) = \mathbf{A}(i) \left[ \frac{\tilde{y}_p(j-d+1)}{\hat{\sigma}_{\tilde{y}_p}^2(j-d+1)} - \left( \frac{\sum_{j=iN_L}^{iN_L+N-1} \tilde{y}_p^4(j-d+1)}{3 \left( \sum_{j=iN_L}^{iN_L+N-1} \tilde{y}_p^2(j-d+1) \right)^2} - 1 \right) \right. \\ \left. \cdot \left( \frac{\tilde{y}_p^3(j-d+1)}{\hat{\sigma}_{\tilde{y}_p}^4(j-d+1)} - \frac{\tilde{y}_p(j-d+1) \sum_{j=iN_L}^{iN_L+N-1} \tilde{y}_p^4(j-d+1)}{\hat{\sigma}_{\tilde{y}_p}^6(j-d+1)} \right) \right] \quad (0)$$

# TRINICON - Nearly Gaussian Densities as Source Model

**Filtered-x-type interpretation** (filtered versions of microphone signals and output signals in the coeff. update):

Inversion of the speech production models within the blind signal processing

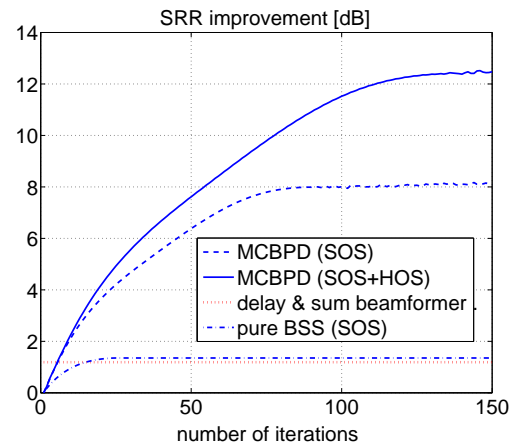
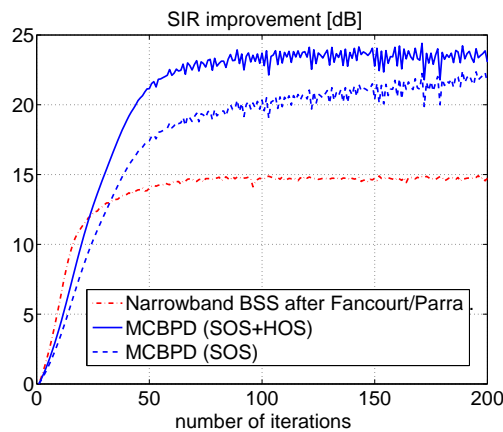


- The coefficients in  $\mathbf{W}$  and  $\mathbf{W}$  are estimated in an alternating way.
- Note: The filtered input vector is calculated using the LP of the output signals.



# Experimental Results – TRINICON for Dereverberation

- MCBPD, 2 sources, 4 mics.,  $L = 3000$ , offline adaptation;  $T_{60} \approx 700\text{ms}$ ,  $f_s = 16\text{kHz}$ ,  $n_A = 32$



Sources:



Mixture (mic 1):



Outputs:



## Summary

- In acoustic preprocessing for hands-free speech communication, **Blind MIMO signal processing** seems appropriate for practical separation and dereverberation tasks.
- **From a generic framework** for adaptive signal processing, we can establish **links between various known algorithms**, and, so far, it has led to various **novel algorithms** for
  - ▷ robust BSS real-time system
  - ▷ robust TDOA estimation for localization of multiple sources
  - ▷ dereverberation without artifacts
  - ▷ improved supervised adaptive filtering, e.g., for acoustic echo cancellation