

Continuous Time-Frequency Masking Method for Blind Speech Separation with Adaptive Choice of Threshold Parameter Using ICA

Zbyněk Koldovský, Jan Nouza, and Jan Kolorenc

Department of Electronics and Signal Processing
Faculty of Mechatronics and Interdisciplinary Engineering Studies
Technical University of Liberec, Liberec, Czech Republic
{zbynek.koldovsky, jan.nouza, jan.kolorenc}@tul.cz

Abstract

We propose a novel method for blind speech separation using continuous time-frequency masking. The method is equipped with an adaptive choice of a threshold parameter that is based on utilization of ICA methods. We present a direct application that consists in the speech segregation for automatic transcription of spoken broadcasts disturbed by background music. Experimental results show improved performance in comparison with traditionally used binary masking methods.

Index Terms: Independent Component Analysis, time-frequency masking, speech recognition, automatic transcription.

1. Introduction

Blind Source Separation (BSS), which consists in recovering original signals from their mixtures when the mixing process is unknown, has been widely studied problem in last two decades. Independent Component Analysis (ICA) is one of most popular methods for BSS based on the assumption of independence of the unknown original signals [1]. The underlying task is an instantaneous linear mixing model

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) \quad t = 1, \dots, N, \quad (1)$$

where $\mathbf{x}(t) = [x_1(t), \dots, x_m(t)]^T$ and $\mathbf{s}(t) = [s_1(t), \dots, s_n(t)]^T$ denote the realizations of the mixed signals and of the original signals, respectively, \mathbf{A} is an $m \times n$ full-rank mixing matrix, and N is the number of samples.

Numerous methods have been proposed in the literature that are able to estimate the original signals \mathbf{s} when their number does not exceed the number of the mixtures \mathbf{x} , i.e. $m \geq n$ [2, 3]. A typical feature of the ICA methods is that the original signals can be retrieved up to their original order, scales, and signs [4] unless any prior information is available.

The separation of the *underdetermined* model, i.e. when $m < n$, is often based on the assumption of sparseness of the original signals either in the time, frequency, or time-frequency domain [5]. Many methods utilize the property together with the assumption of W-disjoint orthogonality [6, 7] of the original sources, and the separation proceeds by 1) identification of system parameters via ICA and 2) sources segregation by means of binary time-frequency masking [8, 9].

This paper focuses on the masking stage only while the mixing system and its parameters are known. Specifically, our direct

application consists in separation of a speech signal from its stereo mixture with background noise (music). The model has a form

$$\begin{aligned} x_1(t) &= s(t) + y_1(t) \\ x_2(t) &= s(t) + y_2(t) \quad t = 1, \dots, N, \end{aligned} \quad (2)$$

where s denotes the speech signal, and y_1 and y_2 denotes left and right interfering noise, respectively. We further assume that $s(t)$ is independent of $y_1(t)$ and $y_2(t)$ for all $t = 1, \dots, N$, and all signals have zero mean, i.e., $E[s(t)] = E[y_1(t)] = E[y_2(t)] = 0$. Here, E stands for the expectation operator.

We will further process signals

$$u(t) = \frac{1}{2}(x_1(t) + x_2(t)) = s(t) + \frac{1}{2}y_1(t) + \frac{1}{2}y_2(t) \quad (3)$$

$$v(t) = x_1(t) - x_2(t) = y_1(t) - y_2(t) \quad (4)$$

that are clearly defined thanks to known mixing parameters. However, if they were not known, it would be highly expectable that most ICA methods applied on signals x_1 and x_2 result in scaled, reordered and resigned copies of u and v . This is justified by the fact that u and v are not correlated when $E[y_1(t)]^2 = E[y_2(t)]^2$. Next, v does not contain the signal s , thus, u and v “usually” form the most independent pair of signals [13]. We claim this since it is the starting point for further generalization of our work. In addition, we utilize residual interference estimator, that is a by-product of the hypothetically applied ICA algorithm with known performance [11, 12], for masking threshold parameter choice.

It was already suggested in [13, 14] that a continuous masking can be used instead of the binary one [7], especially, when the sparsity or the W-disjoint orthogonality of the signals are violated. Even if they are not so, the binary masking segregates a source via zeroing rejected frequencies of the processed signal (its finite-sample short-time Fourier Transformation (STFT)). This results in unwanted distortion of the segregated signal.

The masking is often driven by a nuisance threshold parameter whose optimum value is unknown. By means of suitable criteria, we will show that a bad choice of the parameter may seriously degrade the quality of the segregated signal. Based on this, we propose a novel continuous masking method being less sensitive to the parameter choice, and further we introduce a heuristically derived trick using ICA for an adaptive choice of the parameter.

The paper is organized as follows. Next section introduces performance measures that are very important for objective evaluation of performances of the masking methods and for their comparisons. In Section 3, our continuous mask is proposed, and the adap-

⁰This work was partly supported by the Czech Science Foundation (GA ČR) through the project 102/05/0278.

tive choice of the threshold parameter is derived. In Section 4, better performance of our masking method is validated by computer simulations, and results of application in automatic transcription of spoken broadcasts disturbed by background music are presented.

2. Performance Measurement

2.1. Short-time Fourier transformation

We define the L -point short-time Fourier Transformation of a signal $x(t)$, $t = 1, \dots, N$, by

$$\begin{aligned} \text{STFT}[x(t)](\omega_k, \ell) &= x(\omega_k, \ell) = \\ &= \sum_{i=\ell M+1}^{\ell M+L} x(i)w(i)e^{-\frac{2\pi j}{L}(i-\ell M-1)(\omega_k-1)}, \end{aligned} \quad (5)$$

where L is length of the time-window, M is length of the non-overlap segment, $\ell = 0, \dots, (p-1)r$, $\omega_k = 1, \dots, L$, and $r = L/M$ and $p = N/L$ (we assume that r and p are integers). An inverse transformation to the STFT will be denoted by ISTFT. We have used $L = 1024$, $M = 128$, and the rectangular window function $w \equiv 1$ in (5) since the windowing effect was shown to be negligible [6].

2.2. Performance measures

Let $\mathcal{M}(\omega_k, \ell)$ be a positive real function representing a mask, and let $s(\omega_k, \ell)$, $y_1(\omega_k, \ell)$, $y_2(\omega_k, \ell)$, $u(\omega_k, \ell)$, and $v(\omega_k, \ell)$ be, respectively, the short-time Fourier transformations of the signals s , y_1 , y_2 , u , and v . Masked versions of s , y_1 , and y_2 are defined, respectively, by

$$\tilde{s}^{\mathcal{M}}(t) = \text{ISTFT}[\mathcal{M}(\omega_k, \ell)s(\omega_k, \ell)](t) \quad (6)$$

$$\tilde{y}_1^{\mathcal{M}}(t) = \text{ISTFT}[\mathcal{M}(\omega_k, \ell)y_1(\omega_k, \ell)](t) \quad (7)$$

$$\tilde{y}_2^{\mathcal{M}}(t) = \text{ISTFT}[\mathcal{M}(\omega_k, \ell)y_2(\omega_k, \ell)](t) \quad (8)$$

Finally, $\hat{s}^{\mathcal{M}}$ denotes the resulting estimated signal s arisen from masking of the known signal u , i.e.,

$$\hat{s}^{\mathcal{M}}(t) = \text{ISTFT}[\mathcal{M}(\omega_k, \ell)u(\omega_k, \ell)](t) \quad (9)$$

In order to measure the quality of the signal $\hat{s}^{\mathcal{M}}$ we use *interference-plus-distortion-to-signal ratio* (IDSR)

$$\text{IDSR}^{\mathcal{M}} = \frac{\min_{\alpha} \mathbb{E}[s(t) - \alpha \hat{s}^{\mathcal{M}}(t)]^2}{\mathbb{E}[s(t)]^2}. \quad (10)$$

We present also partial criteria: *distortion-to-signal ratio* (DSR) and *interference-to-signal ratio* (ISR), respectively, defined by

$$\text{DSR}^{\mathcal{M}} = \frac{\min_{\alpha} \mathbb{E}[s(t) - \alpha \tilde{s}^{\mathcal{M}}(t)]^2}{\mathbb{E}[s(t)]^2} \quad (11)$$

$$\text{ISR}^{\mathcal{M}} = \frac{\min_{\alpha} \mathbb{E}[\tilde{s}^{\mathcal{M}}(t) - \alpha \tilde{s}^{\mathcal{M}}(t)]^2}{\mathbb{E}[\tilde{s}^{\mathcal{M}}(t)]^2}. \quad (12)$$

Thanks to linearity of the (inverse) Fourier transformation it holds

$$\hat{s}^{\mathcal{M}}(t) = \tilde{s}^{\mathcal{M}}(t) + \frac{1}{2}\tilde{y}_1^{\mathcal{M}}(t) + \frac{1}{2}\tilde{y}_2^{\mathcal{M}}(t)$$

Hence, using the fact that s is independent of y_1 and y_2 , and the expectation values can be estimated via sample mean, denoted by

$\hat{\mathbb{E}}[\cdot]$, the criteria (10), (11) and (12) can be estimated as

$$\text{IDSR}^{\mathcal{M}} = 1 - \frac{\hat{\mathbb{E}}^2[s(t)\tilde{s}^{\mathcal{M}}(t)]}{\hat{\mathbb{E}}[s(t)]^2\hat{\mathbb{E}}[\tilde{s}^{\mathcal{M}}(t)]^2} \quad (13)$$

$$\text{DSR}^{\mathcal{M}} = 1 - \frac{\hat{\mathbb{E}}^2[s(t)\tilde{s}^{\mathcal{M}}(t)]}{\hat{\mathbb{E}}[s(t)]^2\hat{\mathbb{E}}[\tilde{s}^{\mathcal{M}}(t)]^2} \quad (14)$$

$$\text{ISR}^{\mathcal{M}} = \frac{\hat{\mathbb{E}}[\tilde{y}_1^{\mathcal{M}}(t) + \tilde{y}_2^{\mathcal{M}}(t)]^2}{4 \cdot \hat{\mathbb{E}}[\tilde{s}^{\mathcal{M}}(t)]^2}. \quad (15)$$

The drawback of the criterion (13) is that it evaluates bad estimates $\hat{s}^{\mathcal{M}}$ coarsely [15] ($\text{IDSR} \leq 1$), however, we use it since it provides clear and fair criterion for comparisons.

3. Proposal of the masking method

The idea of the binary masking comes from the assumption of W-disjoint orthogonality (W-DO) of the original signals [6, 7]. In the time-frequency domain, it means that for each pair (ω_k, ℓ) only one of the original signals has nonzero STFT representation. In case of our model, this assumption can be relaxed so that either only $s(\omega_k, \ell) \neq 0$ or only $y_1(\omega_k, \ell) \neq 0 \vee y_2(\omega_k, \ell) \neq 0$.

The binary mask, that separates signal s using known signals u and v only, can be defined as

$$\mathcal{M}_{\tau}^b(\omega_k, \ell) = \begin{cases} 1 & |u(\omega_k, \ell)| > \tau|v(\omega_k, \ell)| \\ 0 & |u(\omega_k, \ell)| \leq \tau|v(\omega_k, \ell)|, \end{cases} \quad (16)$$

where τ is a threshold parameter. This utilizes the W-DO assumption through the fact that frequencies of the signal s are included in u but not included in v .

The drawback of this approach is that the W-DO assumption does not hold exactly in real applications, especially, when processing non-speech signals (music). Moreover, $s(\omega_k, \ell) = 0$ with zero probability even if the signal is sparse. This suggest using continuous mask instead [13, 14].

Our proposal of the continuous mask is

$$\mathcal{M}_{\tau}^c(\omega_k, \ell) = \frac{|u(\omega_k, \ell)|^2}{|u(\omega_k, \ell)|^2 + \tau|v(\omega_k, \ell)|^2}. \quad (17)$$

The choice was inspired by an ideal mask \mathcal{M}^i that minimizes $|s(\omega_k, \ell) - \mathcal{M}^i(\omega_k, \ell)u(\omega_k, \ell)|^2$. Simple calculus gives

$$\mathcal{M}^i(\omega_k, \ell) = \frac{|s(\omega_k, \ell)|^2 + \Re(\overline{s(\omega_k, \ell)}y(\omega_k, \ell))}{|u(\omega_k, \ell)|^2} \quad (18)$$

where $y(\omega_k, \ell) = (y_1(\omega_k, \ell) + y_2(\omega_k, \ell))/2$, and $\Re(z)$ stands for the real part of a complex number z .

Both masks (16) and (17) involve the threshold parameter τ that significantly affects the quality of the separation. Typical behaviors of the criteria (10)-(12) are demonstrated in Figure 1. For $\tau = 0$ the signal u remain unchanged, and DSR is zero. For τ increasing the DSR grows, and ISR usually decays (depends on the validity of the W-DO assumption). As can be seen, only the IDSR criterion can identify the optimum value of τ . Next, an important observation is that the DSR of the binary mask rapidly grows with τ (and for that reason the IDSR also), which brings on the sensitivity to the threshold parameter choice.

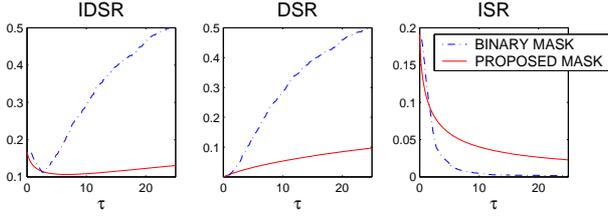


Fig. 1 A speech signal s ($\hat{E}[s(t)]^2 = 1$) of length $N = 2^{15}$ was separated from its stereo mixture with a stereo music signal ($\hat{E}[y_1(t)]^2 = \hat{E}[y_2(t)]^2 = 0.25$). Graphs show IDSR, DSR, and ISR achieved by the binary mask (16) and by the proposed continuous mask (17) for $\tau \in [0, 25]$.

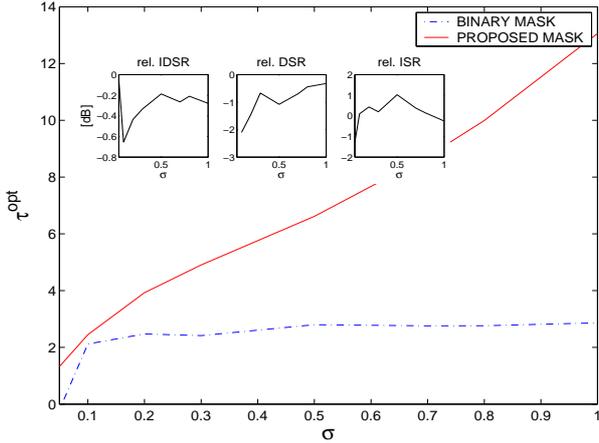


Fig. 2 Optimum value of τ versus σ , where $\sigma^2 = \hat{E}[y_1(t)]^2 = \hat{E}[y_2(t)]^2$; $\hat{E}[s(t)]^2 = 1$. Small graphs show relative IDSR, DSR, and ISR defined as $\text{IDSR}^{\mathcal{M}_\tau^c} / \text{IDSR}^{\mathcal{M}_\tau^b}$, etc.

3.1. Adaptive choice of the threshold parameter τ

It is difficult to consistently estimate the optimum value of τ (denoted by τ^{opt}), moreover, this may not be possible without any prior information about the input sources s , y_1 , and y_2 . Here we exploit a quite obvious fact that the value of τ^{opt} significantly depends on input interference-to-signal ratio, in our case, defined by

$$\text{ISR}^{\text{in}} = (\text{E}[y_1(t)]^2 + \text{E}[y_2(t)]^2) / \text{E}[s(t)]^2, \quad (19)$$

because the masking must be stronger for higher ISR^{in} . This is achieved by taking higher values of τ .

In simulations, we compute the optimum value τ^{opt} through minimization of the IDSR (13) by means of the function `fminsearch` in MatlabTM. The dependence of τ^{opt} on ISR^{in} (through the parameter $\sigma^2 = \hat{E}[y_1(t)]^2 = \hat{E}[y_2(t)]^2$) is demonstrated in Figure 2, where the same signals were used like in Figure 1. Figure 3(a) demonstrates the dependence when random signals from a database were used in the same experiment.

The leading idea of our proposal for estimation of τ^{opt} is that the signals u and v are dependent in spite of being uncorrelated (or almost uncorrelated), and that the dependency should be higher for higher ISR^{in} . We believe that the dependency can be reflected by the residual interference-to-signal ratio (residual ISR) of ICA components separated from signals $\mathbf{x} = [u, v]^T$ (as it was suggested in the introduction section the separated signals can be expected to be

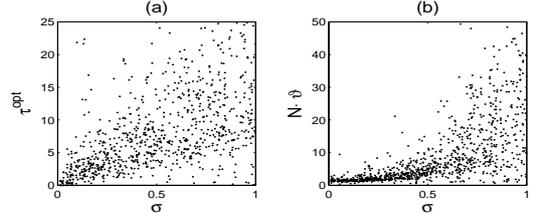


Fig. 3 Optimum τ of the continuous mask (17) and the dependence measure (21) versus σ ($\sigma^2 = \hat{E}[y_1(t)]^2 = \hat{E}[y_2(t)]^2$; $\hat{E}[s(t)]^2 = 1$) from 1000 independent experiments. In each trial, σ was randomly chosen from $[0, 1]$, signal s of length $N = 2^{15}$ was randomly taken from training database #1 described in Section 4. y_1 and y_2 were, respectively, left and right channel of a random piece of music. Results with $\tau^{opt} > 25$ and $N \cdot \vartheta > 50$ were discarded.

estimates of u and v). Another approach would be to use a direct estimation of mutual information of u and v [16], but here we may profit from previous utilization of ICA in the general case when the mixing parameters are not known. Here, u and v are explicitly defined, thus, the running of the ICA method is not necessary.

The residual ISR can be measured in case of ICA algorithms with known performance [11, 12]. For this purpose, we utilize algorithm EFICA [11] whose residual interference-to-signal ratio between k -th and ℓ -th estimated components (denoted by \hat{s}_k, \hat{s}_ℓ) is estimated via

$$\text{ISR}_{k\ell}^{EF} = \frac{1}{N} \frac{\gamma_k(\gamma_\ell + \tau_\ell^2)}{\tau_\ell^2 \gamma_k + \tau_k^2(\gamma_\ell + \tau_\ell^2)}, \quad (20)$$

where

$$\begin{aligned} \gamma_k &= \beta_k - \mu_k^2 & \mu_k &= \hat{E}[\hat{s}_k g_k(\hat{s}_k)] \\ \tau_k &= |\mu_k - \rho_k| & \rho_k &= \hat{E}[g'_k(\hat{s}_k)] \\ & & \beta_k &= \hat{E}[g_k^2(\hat{s}_k)], \end{aligned}$$

and $g_k(\cdot)$ is a nonlinear function chosen for k -th signal in the algorithm. Hence, we define a heuristic measure of the dependence of u and v as the sum of residual interference of signals $\hat{s}_1 = u$ and $\hat{s}_2 = v$, i.e.

$$\vartheta = \text{ISR}_{12}^{EF} + \text{ISR}_{21}^{EF}. \quad (21)$$

The results in Fig. 3(b) corroborate legitimacy of our approach.

Finally, we propose an ad-hoc choice of the parameter τ for the proposed masking (17) using the measure ϑ by

$$\tau = \min\{0.5 + a \cdot N \cdot \vartheta, 25\}, \quad (22)$$

where $a = 0.5689$ was determined experimentally as a regression coefficient from results presented in Figure 3.

4. Simulations

In our simulations, we utilize two databases #1 and #2 of, respectively, 215 and 503 utterances of various length recorded from Czech spoken broadcasts. The stereo music signals are taken as a random piece of Mike Oldfield's Ommadawn (part 1), which is very multifarious and dynamical instrumental composition suitable for our simulations. We demonstrate the masking method in two experiments.

In the first experiment, signal s of length $N = 2^{15}$ was randomly taken from database #2 and mixed via (2) with a randomly

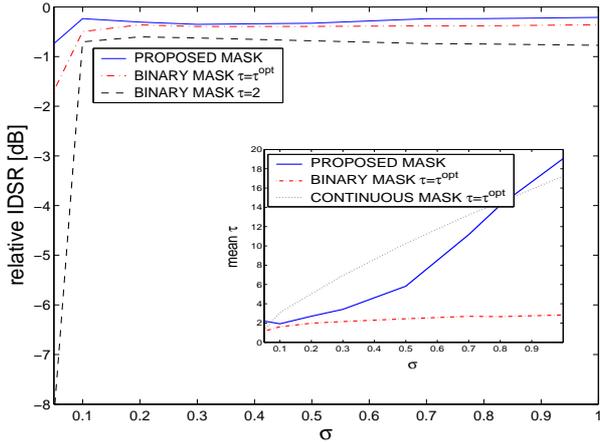


Fig. 4 Relative IDSR and average τ versus σ .

Table I. Accuracy achieved by the speech recognizer [17].

	$\sigma = 0.3$	$\sigma = 0.7$	$\sigma = 0.9$
clear signals	81.35	81.35	81.35
mixed signals	75.26	60.76	53.90
binary masking, $\tau = 2$	75.91	66.25	60.05
proposed masking	76.63	67.05	61.64

chosen music signal (left channel y_1 and right channel y_2) in each of 500 independent trials. The signals were normalized so that $\hat{E}[s(t)]^2 = 1$ and $\hat{E}[y_1(t)]^2 = \hat{E}[y_2(t)]^2 = \sigma^2$. In each trial, the speech signal s was separated using the continuous mask (17) with optimum τ , the same mask with the choice (22), the binary mask (16) with optimum τ , and the binary mask with $\tau = 2$. Relative IDSRs related to that of the optimum continuous masking $\text{IDSR}_{\tau^{\text{opt}}}^{\text{C}}$ are shown in Figure 4 together with averaged values of τ s used in forenamed masks. The proposed masking method outperforms the binary mask even if its optimum τ is used.

The second experiment was done with the whole utterances from database #2 that were mixed with a random piece of stereo music of the same length. Here, the data were not normalized before mixing; ISR^{in} was only roughly controlled by multiplying y_1 and y_2 with σ . Each utterance was separated by the proposed masking method and by the binary mask with $\tau = 2$. In the proposed method, the adaptive choice of τ (22) was done separately for each segment of length $N = 2^{15}$.

The original utterances s , the mixed signals u , and the separated signals were passed through the automatic continuous Czech speech recognizer [17]. Performance was measured in terms of accuracy defined as $100 \cdot (C - D - I - S)/C$, which is computed via comparison of a reference text with the recognized one. Here, C is the number of words in the reference text, D is the number of deletions (words omitted by speech recognizer), I is the number of insertions (new words added on recognizer's output), and S is the number of substitutions (words exchanged with another). The whole database #2 contain 10322 words. Results in Table I demonstrate the achieved improvement via the proposed masking method in terms of the accuracy.

5. References

[1] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*, Wiley-Interscience, New York, 2001.

[2] J.-F. Cardoso, "High-order Contrasts for Independent Component Analysis", *Neural Computat.*, vol. 11, no. 1, pp. 157-192, 1999.

[3] A. Hyvärinen and E. Oja, "A Fast Fixed-Point Algorithm for Independent Component Analysis", *Neural Computation*, 9(7):1483-1492, 1997.

[4] P. Tichavský and Z. Koldovský, "Optimal Pairing of Signal Components Separated by Blind Techniques", *IEEE Signal Processing Letters*, vol. 11, no. 2, pp. 119-122, 2004.

[5] P. Bolfill and M. Zibulevsky, "Blind Separation of More Sources Than Mixtures using Sparsity of Their Short-time Fourier Transform", in *Proc. ICA 2000*, Helsinki, Finland, pp. 87-92, June 19-22, 2000.

[6] R. Balan and J. Rosca, "Statistical Properties of STFT Ratios for Two Channel Systems and Applications to Blind Source Separation", in *Proc. ICA 2000*, Helsinki, Finland, pp. 429-434, June 19-22, 2000.

[7] Ö. Yilmaz and S. Rickard, "Blind Separation of Speech Mixtures via Time-Frequency Masking", *IEEE Trans. on Signal Processing*, vol. 52, no. 7, July 2004.

[8] M. S. Pedersen, D. Wang, J. Larsen, and U. Kjems, "Over-complete Blind Source Separation by Combining ICA and Binary Time-Frequency Masking", in *Proc. of 2005 IEEE Workshop on Machine Learning and Signal Processing*, pp. 15-20, 28-30 Sept. 2005.

[9] S. Araki, S. Makino, H. Sawada, and R. Mukai, "Underdetermined blind separation of convolutive mixtures of speech with directivity pattern based mask and ICA", in *Proc. ICA2004*, pp. 898-905, Sept. 2004.

[10] Z. Koldovský and P. Tichavský, "Methods of Fair Comparison of Performance of Linear ICA Techniques in Presence of Additive Noise", to be presented at ICASSP-2006.

[11] Z. Koldovský, P. Tichavský and E. Oja, "Efficient Variant of Algorithm FastICA for Independent Component Analysis Attaining the Cramér-Rao Lower Bound", to be published in *IEEE Trans. on Neural Networks*, Sept. 2006.

[12] A. Yeredor, "Blind separation of Gaussian sources via second-order statistics with asymptotically optimal weighting," *IEEE Signal Processing Letters*, vol. 7, pp. 197-200, July 2000.

[13] H. Sawada, S. Araki, R. Mukai, and S. Makino, "Blind extraction of dominant target sources using ICA and time-frequency masking", in *Proc. of ISCAS 2005*, vol. 6, pp. 5882-5885, 23-26 May 2005.

[14] S. Araki, S. Makino, H. Sawada, and R. Mukai, "Underdetermined blind speech separation with directivity pattern based continuous mask and ICA", in *Proc. EUSIPCO2004*, pp. 1991-1994, Sept. 2004.

[15] E. Vincent, R. Gribonval, and C. Févotte, "Performance Measurement in Blind Audio Source Separation", to be published in *IEEE Trans. on Speech and Audio Processing*, 2006.

[16] G. A. Darbellay and P. Tichavský, "Independent Component Analysis through Direct Estimation of the Mutual Information", *Proc. of ICA'2000*, Helsinki, Finland, pp. 69-75, June 2000.

[17] J. Nouza, J. Žďánský, P. David, P. Červa, J. Kolorenč, and D. Nejedlová, "Fully Automated System for Czech Spoken Broadcast Transcription with Very Large (300K+) Lexicon", in *Proc. Interspeech2005*, Lisboa, Portugal, pp. 1681-1684, Sept. 2005.