

A New Perspective of Auxiliary-Function-Based Independent Component Analysis in Acoustic Echo Cancellation

Yueyue Na, Ziteng Wang, Zhang Liu, Yun Li, Biao Tian, Qiang Fu

Machine Intelligence Technology, Alibaba Group

{yueyue.nyy, ziteng.wzt, yinan.lz, yl.yy, tianbiao.tb, fq153277}@alibaba-inc.com

Abstract

For better human-machine or human-human voice communication, acoustic echo cancellation (AEC) is required to suppress echo from observed signals. Traditional AEC's performance is restricted when double-talk occurs. Blind source separation (BSS) based AEC has good performance in double-talking since BSS is inherently based on the full duplex signal mode, which both far-end and near-end signals coexist. In this paper, the auxiliary-function-based independent component analysis (Aux-ICA) algorithm is used to solve the AEC problem. After mathematical simplification, it can be seen that the Aux-ICA AEC is equivalent to weighted recursive least square (RLS) algorithm, with the ICA nonlinearity as the weighting function. The performance of the proposed approach is verified by both simulated and real-world experiments. In addition, a speech enhancement front-end is unified in the BSS framework for the application of smart TV keyword spotting.

Index Terms: acoustic echo cancellation, blind source separation, independent component analysis, double-talk

1. Introduction

Microphone array and loudspeakers are often installed on modern smart devices for better human-machine or human-human voice communication. In such cases, sound played by device itself will be collected back by its own microphone, i.e., the acoustic echo, which will contaminate the useful signal from real user. Thus, acoustic echo cancellation (AEC) technique is required to suppress echo from original signals.

Adaptive filtering [1] is a common way to perform AEC. Normalized least mean square (NLMS) approach [1-3] is a typical algorithm of this type, which minimizes the mean square error between echo estimation and microphone signal by gradient descent method. In order to prevent filter divergence, double-talk detector (DTD) [4], or adaptive step size policy [5] is used to freeze or slow down the filter adaptation when double-talk occurs. Recursive least square (RLS) [6, 7] is another class of algorithm which can be used to solve the AEC problem. Compared with NLMS, RLS has faster convergence speed but higher computational complexity (section 2.4 of [7]).

In addition to the linear echo, due to the sound effect and the dynamic range compression of the playback system, as well as the device vibration, nonlinear echo also presents in real applications, which cannot be cancelled directly by linear adaptive filtering. Post filtering and acoustic echo suppression (AES) techniques, such as [8-10], are usually applied after adaptive filtering to reduce nonlinear echo. Deep learning based approaches, like [11, 12], are also developed for nonlinear echo suppression with the help of deep neural network's powerful nonlinear modeling ability.

To achieve better voice communication quality, not only acoustic echo, but also environmental noise must be handled. Blind source separation (BSS) is a technique which performs speech enhancement by separating sound components from observed noisy signals. Independent component analysis (ICA) [13, 14] and independent vector analysis (IVA) [15-19] are typical BSS techniques. AEC can also be considered as a semi-BSS problem which separates echo and near-end components from microphone signals, with references (far-end) as the supervised information. Researches about BSS based AEC can be found in [20-23], etc.

Traditional AEC adaptive filtering plus post filtering or AES techniques have some drawbacks. First, for gradient descent based approaches, there always exists the balance between convergence speed and stability [14]. Second, although DTD and adaptive step size policy can work well in single-talking and sparse double-talking scenes, their performance may degrade in continuous double-talking scene where near-end signals always exist [21, 23]. This scene is often encountered in the application of smart TV, as introduced in section 4. Third, most post filtering and AES approaches are based on Wiener filtering or masking techniques, which suppress nonlinear echo at the price of near-end speech distortion.

Compared with traditional AEC approaches like [2, 3, 5-7], BSS based AEC algorithms have the advantage of good echo cancellation ability in continuous double-talking scene [23], since BSS signal model is a full duplex model, which both far-end and near-end signals coexist. In this paper, the auxiliary function based ICA (Aux-ICA) algorithm [14] is used to solve the AEC problem. In addition to the full duplex property, with the help of the auxiliary function technique, the explicit step size parameter is avoided, and the convergence speed is increased [14]. After mathematical simplification, it can be seen that the Aux-ICA AEC can be considered as a kind of weighted RLS algorithm, with the ICA nonlinearity as the weighting function. The following content is organized as follows: section 2 introduces the signal model; section 3 depicts the algorithm derivation; experiments and comparisons are given in section 4; at last, section 5 gives the conclusion of this paper.

2. Problem formulation

In this paper, lowercase italic letters denote scalars; lowercase bold italic letters denote vectors; uppercase bold italic letters denote matrices; the operator $*$ for convolution; the superscript T for transpose, H for Hermitian transpose or complex conjugate.

2.1. Signal model

The signal model used in this paper is depicted in Figure 1. The microphone signal x is consist of two parts: the linear echo e and the near-end signal s , as shown in (1), where k is the sample index.

$$x(k) = e(k) + s(k) \quad (1)$$

Supposing there are R references in the system, the linear echo e can be modeled as the sum of the convolution between individual reference r_1, \dots, r_R and the corresponding unknown echo path a_1, \dots, a_R , as shown in (2).

$$e(k) = \sum_{m=1}^R a_m(k) * r_m(k) \quad (2)$$

The near-end signal s can be considered as the combination of N source images s_1, \dots, s_N . Since nonlinear echo \tilde{e} cannot be cancelled by linear AEC, it is also classified to the near-end component in Figure 1.

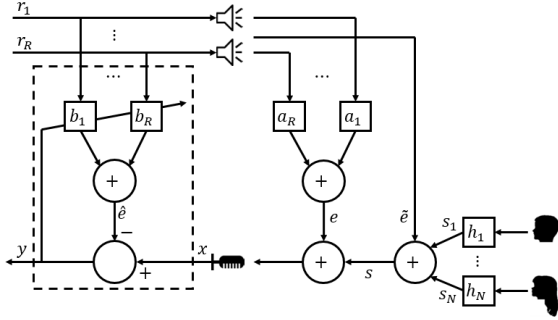


Figure 1: Signal model.

AEC performs echo cancellation by calculating the echo estimation \hat{e} from r_1, \dots, r_R and the estimated echo paths b_1, \dots, b_R , then, subtracting \hat{e} from x to get the output signal y , which is the estimation of the near-end signal s , as shown in (3) and (4), as well as the dotted box in Figure 1.

$$\hat{e}(k) = \sum_{m=1}^R b_m(k) * r_m(k) \quad (3)$$

$$y(k) = x(k) - \hat{e}(k) \quad (4)$$

2.2. The NLMS and the RLS algorithm

For the signal model in (1) to (4), equation (5) to (8) depict the common NLMS iteration policy, where μ is the step size, δ is a small positive number to prevent zero denominator, $\mathbf{w} = \mathbf{b}^H$ is the adaptive filter, \mathbf{r} is the reference vector with L samples lag [3], Table 2.2 of [7].

$$\mathbf{w}(k) = \mathbf{w}(k-1) + \frac{\mu \mathbf{y}^H(k) \mathbf{r}(k)}{\delta + \mathbf{r}^H(k) \mathbf{r}(k)} \quad (5)$$

$$\mathbf{y}(k) = x(k) - \mathbf{w}^H(k) \mathbf{r}(k) \quad (6)$$

$$\mathbf{w}(k) = \mathbf{b}^H(k) = [b_1^H(k), \dots, b_R^H(k), \dots, b_1^H(k-L+1), \dots, b_R^H(k-L+1)]^T \quad (7)$$

$$\mathbf{r}(k) = [r_1(k), \dots, r_R(k), \dots, r_1(k-L+1), \dots, r_R(k-L+1)]^T \quad (8)$$

The RLS approach is based on the so called normal equation in (9), where \mathbf{R} is the reference auto-correlation matrix, \mathbf{p} is the reference-mic cross correlation vector, $0 \ll \alpha < 1$ is the forgetting factor [6, 7].

$$\mathbf{w}(k) = \mathbf{R}^{-1}(k) \mathbf{p}(k) \quad (9)$$

$$\mathbf{R}(k) = \alpha \mathbf{R}(k-1) + (1-\alpha) \mathbf{r}(k) \mathbf{r}^H(k) \quad (10)$$

$$\mathbf{p}(k) = \alpha \mathbf{p}(k-1) + (1-\alpha) \mathbf{r}(k) x^H(k) \quad (11)$$

Solving (9) directly is computationally intensive because of the matrix inversion operation. Several policies can be used to speed up the iteration process, e.g., the Woodbury matrix identity [24], and the QR decomposition (chapter 3 of [7]), etc.

In NLMS and RLS, the near-end signal s is not explicitly modeled, and it is often treated as noise in the system identification procedure [1, 7]. In equation (6), y is considered as the error term between desired and estimated signal, but not the estimation of the near-end component. Such modeling policy may not so appropriate when s is speech.

2.3. The BSS model

After subband decomposition techniques, such as short-time Fourier transform (STFT) is applied, time domain microphone signal $x(t)$ can be transformed to time-frequency domain $x(f, k)$, where f is the frequency bin index, k is the STFT frame index. Other signals have similar notations, too. Since all frequency bins have the same operation in the following algorithms, the frequency bin index f is omitted to simplify the notation. Equations in this paper are not distinguished between time and frequency domain since the theory works in both domains.

When STFT window is sufficiently longer than room impulse response, the time domain signal mode in equation (1) to (4) can approximately be transformed to time-frequency domain, as shown in equation (12) and (13), where $\mathbf{s} = [s, r_1, \dots, r_R]^T$, $\mathbf{x} = [x, r_1, \dots, r_R]^T$, and $\mathbf{y} = [y, r_1, \dots, r_R]^T$ are $1+R$ dimension vectors [20-23].

$$\mathbf{x}(k) = \mathbf{A}(k) \mathbf{s}(k) \quad (12)$$

$$\mathbf{y}(k) = \mathbf{B}(k) \mathbf{x}(k) \quad (13)$$

The mixing matrix \mathbf{A} and the demixing matrix \mathbf{B} have their special structure, as shown in (14) and (15), where \mathbf{I} denotes the identity matrix, a and b are the frequency domain version of real and estimated echo paths [20, 21].

$$\mathbf{A} = \begin{bmatrix} 1 & a_1 & \dots & a_R \\ 0 & \mathbf{I} & & \end{bmatrix} \quad (14)$$

$$\mathbf{B} = \begin{bmatrix} 1 & b_1 & \dots & b_R \\ 0 & \mathbf{I} & & \end{bmatrix} \quad (15)$$

Unlike NLMS and RLS, s is explicitly modeled as an independent component in the BSS based approaches. Thus, the BSS signal model can be considered as a full duplex model which both far-end and near-end signals coexist.

3. The proposed approach

In this paper, the Aux-ICA algorithm [14, 17] is used to perform the separation. For the detail theory and algorithm derivation of Aux-ICA (IVA), please refer to [14, 16-19].

3.1. Algorithm derivation

According to Aux-ICA, the weighted correlation matrix \mathbf{C} is updated first, as shown in (16) to (18), where $\varphi[\cdot]$ is a nonlinear

transform. The nonlinearity in (18) [19] is used in this paper, where γ is the sparse parameter.

$$\mathbf{C}(k) = \alpha \mathbf{C}(k-1) + \beta(k) \mathbf{x}(k) \mathbf{x}^H(k) \quad (16)$$

$$\beta(k) = (1 - \alpha) \varphi[y(k)] \quad (17)$$

$$\varphi[y] = (|y|^2 + \delta)^{(\gamma-2)/2} \quad (18)$$

Then, the demixing filter $\mathbf{v} = \mathbf{B}^H \mathbf{i}_1 = [1, b_1^H, \dots, b_R^H]^T$ can be calculated according to (19) and (20), where $\mathbf{i}_1 = [1, 0, \dots, 0]^T$ is a $1 + R$ dimension vector. Unlike the source separation problem, there are no permutation and scaling ambiguities in the BSS based AEC. As a result, \mathbf{v} is normalized according to (20) to make sure that the demixing matrix \mathbf{B} has the structure in (15).

$$\mathbf{v}(k) = [\mathbf{B}(k-1) \mathbf{C}(k)]^{-1} \mathbf{i}_1 \quad (19)$$

$$\mathbf{v} \leftarrow \mathbf{v}/v_1 \quad (20)$$

Equation (19) can be reformulated to (21), then, by noticing that \mathbf{B} has the special structure in (15), \mathbf{B}^{-1} must have the same structure. Thus, the inversion of \mathbf{B} is avoided, as shown in (22).

$$\mathbf{v}(k) = \mathbf{C}^{-1}(k) \mathbf{B}^{-1}(k-1) \mathbf{i}_1 \quad (21)$$

$$\mathbf{v}(k) = \mathbf{C}^{-1}(k) \mathbf{i}_1 \quad (22)$$

At this point, the demixing filter can already be solved by the Aux-ICA algorithm by first calculating (22), then, normalizing according to (20).

The correlation matrix \mathbf{C} in (16) and (22) can be partitioned according to (23) to (26).

$$\mathbf{C} = \begin{bmatrix} c_{11} & \mathbf{p}^H \\ \mathbf{p} & \mathbf{R} \end{bmatrix} \quad (23)$$

$$c_{11}(k) = \alpha c_{11}(k-1) + \beta(k) x(k) x^H(k) \quad (24)$$

$$\mathbf{p}(k) = \alpha \mathbf{p}(k-1) + \beta(k) \mathbf{r}(k) x^H(k) \quad (25)$$

$$\mathbf{R}(k) = \alpha \mathbf{R}(k-1) + \beta(k) \mathbf{r}(k) \mathbf{r}^H(k) \quad (26)$$

According to block matrix inversion, \mathbf{C}^{-1} can be denoted in (27).

$$\mathbf{C}^{-1} = \begin{bmatrix} \frac{1}{c_{11} - \mathbf{p}^H \mathbf{R}^{-1} \mathbf{p}} & \frac{-\mathbf{p}^H}{c_{11}} \left(\mathbf{R} - \frac{\mathbf{p} \mathbf{p}^H}{c_{11}} \right)^{-1} \\ \frac{-\mathbf{R}^{-1} \mathbf{p}}{c_{11} - \mathbf{p}^H \mathbf{R}^{-1} \mathbf{p}} & \left(\mathbf{R} - \frac{\mathbf{p} \mathbf{p}^H}{c_{11}} \right)^{-1} \end{bmatrix} \quad (27)$$

Comparing (27) with (22) and (20), it is easy to derive the simplified solution in (28). The negative sign is discarded from (27) to (28) since the BSS demixing model in (15) can be interpreted as adding the negative echo estimation to microphone signals, rather than subtracting it as in traditional AEC model.

$$\mathbf{w}(k) = \mathbf{R}^{-1}(k) \mathbf{p}(k) \quad (28)$$

3.2. Discussion

The objective function of the RLS algorithm is minimizing the mean square error between echo estimation and microphone signals, and the corresponding optimization policy is the least

square approach [6, 7]. Near-end signal is not explicitly modeled in the RLS signal model, and often be considered as noise in the system identification procedure. On the other hand, the objective function of Aux-ICA is minimizing the mutual information, which is measured by the KL divergence, and the auxiliary function technique is used for the optimization [14, 16-19]. Near-end signal is explicitly modeled as an independent component in the ICA model.

Although the basic theory of the two algorithms are different, from preceding subsection, it can be seen that the Aux-ICA AEC solution in (28) has the same form as the normal equation of the RLS algorithm in (9). However, from (10), (11) to (25), (26), the ICA nonlinearity β is used in the adaptation. Comparing β with the adaptive step size approach, like [5], although β and μ are updating in every step, β cannot be interpreted as a kind of adaptive step size. Since the normal equation is derived from the “gradient equals zero” assumption [7], RLS get to the minimum in one step, but not several steps as in the gradient descent approach. Thus, the Aux-ICA AEC can be considered as a kind of weighted RLS algorithm, with the ICA nonlinearity as the weighting function. From next section, it can be seen that the echo cancellation performance is improved by the nonlinear weighting.

As the conclusion to this section, the Aux-ICA AEC algorithm is summarized in Table 1. Fast implementation is not discussed here, readers can refer to [7] for further information.

Table 1: The Aux-ICA AEC algorithm.

Initialize: $\mathbf{w}(0) = \mathbf{0}$, $\mathbf{R}(0) = \mathbf{0}$, $\mathbf{p}(0) = \mathbf{0}$, $\alpha = 0.9999$, $\gamma = 0.2$, $\delta = 1 \times 10^{-10}$
Input: $x(k)$, $\mathbf{r}(k)$, output: $y(k)$
1. $y(k) = x(k) - \mathbf{w}^H(k-1) \mathbf{r}(k)$;
2. Calculate $\beta(k)$ according to (17) and (18);
3. Update $\mathbf{p}(k)$, $\mathbf{R}(k)$ according to (25) and (26);
4. Update $\mathbf{w}(k)$ according to (28);

4. Experiment

4.1. Simulated experiment

A simulated environment is established: one speech and one music signal of 16k Hz sampling rate are used as near-end and reference. Simulated data are generated according to (1) and (2), signal-to-echo ratio (SER) is adjusted to -20 dB. Echo paths of length 16 are randomly generated, and an echo path sudden change at 10 second is deliberately added. Three algorithms are compared: NLMS, RLS, and Aux-ICA AEC (Aux and AuxP). Filter length is set to 16 for all algorithms, $\mu = 0.05$ in NLMS. Misalignment in (29), and Perceptual Evaluation of Speech Quality (PESQ) [25] are used as performance indices.

$$\text{Misalignment} = 10 \log_{10} \left(\frac{\|\mathbf{b} - \mathbf{a}\|^2}{\|\mathbf{a}\|^2} \right) \quad (29)$$

Misalignment variation and corresponding near-end speech are shown in Figure 2. Several facts can be observed from Figure 2. First, the misalignment of NLMS and RLS is increased when double-talk occurs, however, the performance of Aux can keep relatively stable, which verifies the full duplex property of the Aux model. Second, original Aux is vulnerable to echo path sudden change, since more data is required to “forget” old states. This drawback can be avoided by adding some protection

mechanism (AuxP), e.g., reset algorithm context when $|y|^2 > |x|^2$ for several consecutive frames. Third, AuxP has the lowest misalignment and fastest convergence speed among compared approaches.

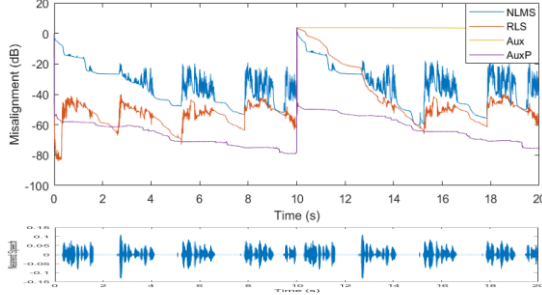


Figure 2: Misalignment variation and near-end speech (colored).

The average misalignment and PESQ are given in Table 2, which also supports the preceding observations.

Table 2: Average misalignment and PESQ.

	Average Misalignment (dB)	PESQ
NLMS	-33.87	1.34
RLS	-48.81	1.29
Aux	-32.27	1.10
AuxP	-64.36	4.27

4.2. Real-world experiment

Real recordings are captured in a testing room configured in Figure 3. The smart TV with 4 mics and 4 loudspeakers is the sound capturing device, the captured signals are 16 kHz, 8 channels data, with 4 mic and 4 reference channels. The 4 mics are arranged into a linear array with 3.5 cm mic spacing, which is located at the bottom of the TV. Target and noise loudspeakers also exist in Figure 3. The loudspeakers marked as “Noise 1”, “Noise 2”, and “Subwoofer” are connected to the same audio hub for noise playback, the “Noise 1” and “Noise 2” loudspeakers are facing to the wall to increase the diffusion effect.

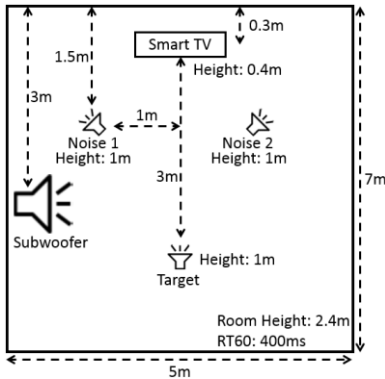


Figure 3: Real-world experimental environment.

Target speech, echo, and noise are captured individually in the environment of Figure 3, then, added together as the algo-

rithm input. The target source is 120 utterances of a certain keyword, pronounced by different people; the echo source is a music clip; the noise source is recorded in a TV store, which contains both diffuse noise and a news broadcast. The sound pressure level is tuned to 70, 90, and 75 dB(A) for target speech, echo, and noise, individually, according to a sound level meter placed close to the microphone array. Hence, the SER is -20 dB, and the signal-to-noise ratio (SNR) is -5 dB for the directly mixed signals.

The front-end used in this experiments utilizes the “AEC-BSS-AGC (Automatic Gain Control)” architecture. The three AEC algorithms in frequency domain are implemented for comparison. The FFT size and the frame shift is 640, 320 (not powers of 2 to compatible with the back-end system), AEC filter length is set to 10 in all frequency bins. After AEC, the Aux-IVA algorithm [17] is used to separate target speech, nonlinear echo, noise and interferences. Thus, the speech distortion is decreased for nonlinear echo cancellation. In order to perform objective comparison, the keyword spotting system introduced in [26] is concatenated to the front-end output. A keyword detected in any output channel is counted as one spot. Recall is used as the performance index, which is defined as: “number of spots / number of keywords”.

Table 3 gives the keyword spotting results in the continuous double-talking scene with speech, echo, and noise coexist. In addition to direct mixing, target speech power is adjusted to generate different SER and SNR. First, it can be seen from Table 3 that the keyword can hardly be detected when the AEC module is bypassed in the front-end pipeline, which verifies the importance of the AEC module. It also can be seen that the performance of Aux is better than NLMS and RLS in this experiment, which is accordance with the conclusion of the last subsection.

Table 3: Recall in continuous double-talking scene.

	Recall (%)				
SER (dB)	-25	-20	-19	-18	-17
SNR (dB)	-10	-5	-4	-3	-2
AEC Bypass	1.7	0.8	1.7	1.7	0.8
NLMS	16.7	57.5	69.2	75.0	88.3
RLS	36.7	73.3	85.0	87.5	88.3
Aux	71.7	95.0	94.2	94.2	95.0

5. Conclusion

Unlike traditional AEC, e.g., NLMS and RLS, BSS AEC’s signal model can be considered full duplex which both near-end and far-end signals coexists. Near-end speech is explicitly modeled as an independent component in the ICA model. Thus, the BSS AEC is expected to have better performance in continuous double-talking scene. In this paper, the Aux-ICA algorithm is used to solve the AEC problem. After mathematical simplification, it can be seen that the Aux-ICA AEC has strong connection with the RLS algorithm. Specifically, it can be considered as weighted RLS, with the ICA nonlinearity as the weighting function. The performance of the proposed approach is verified by both simulated and real-world experiments. In addition, a speech enhancement front-end is unified in the Aux-ICA/IVA framework for the real-world application of smart TV keyword spotting.

6. References

- [1] Haykin, Simon S. *Adaptive filter theory*. Pearson Education India, 2005.
- [2] Morgan, Dennis R., and Steven G. Kratzer. "On a class of computationally efficient, rapidly converging, generalized NLMS algorithms." *IEEE Signal Processing Letters* 3.8 (1996): 245-247.
- [3] Albu, Iuliana, Cristian Anghel, and Constantin Paleologu. "Adaptive filtering in acoustic echo cancellation systems—A practical overview." *2017 9th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)*. IEEE, 2017.
- [4] Yang, Jun. "Multilayer Adaptation Based Complex Echo Cancellation and Voice Enhancement." *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018.
- [5] Valin, Jean-Marc. "On adjusting the learning rate in frequency domain echo cancellation with double-talk." *IEEE Transactions on Audio, Speech, and Language Processing* 15.3 (2007): 1030-1034.
- [6] Buchner, Herbert, Jacob Benesty, and Walter Kellermann. "Generalized multichannel frequency-domain adaptive filtering: efficient realization and application to hands-free speech communication." *Signal Processing* 85.3 (2005): 549-570.
- [7] Apolinário, José Antonio, and R. Rautmann. *QRD-RLS adaptive filtering*. Ed. José Antonio Apolinário. New York: Springer, 2009.
- [8] Franzen, Jan, and Tim Fingscheidt. "An efficient residual echo suppression for multi-channel acoustic echo cancellation based on the frequency-domain adaptive Kalman filter." *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018.
- [9] Hong, Jungpyo. "Stereophonic Acoustic Echo Suppression for Speech Interfaces for Intelligent TV Applications." *IEEE Transactions on Consumer Electronics* 64.2 (2018): 153-161.
- [10] Huang, Hai, et al. "A multiframe parametric Wiener filter for acoustic echo suppression." *2016 IEEE International Workshop on Acoustic Signal Enhancement (IWAENC)*. IEEE, 2016.
- [11] Lee, Chul Min, Jong Won Shin, and Nam Soo Kim. "DNN-based residual echo suppression." *Sixteenth Annual Conference of the International Speech Communication Association*. 2015.
- [12] Carbajal, Guillaume, et al. "Multiple-input neural network-based residual echo suppression." *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018.
- [13] Hyvärinen, Aapo, and Erkki Oja. "Independent component analysis: algorithms and applications." *Neural networks* 13.4-5 (2000): 411-430.
- [14] Ono, Nobutaka, and Shigeki Miyabe. "Auxiliary-function-based independent component analysis for super-Gaussian sources." *International Conference on Latent Variable Analysis and Signal Separation*. Springer, Berlin, Heidelberg, 2010.
- [15] Kim, Taesu, et al. "Blind source separation exploiting higher-order frequency dependencies." *IEEE transactions on audio, speech, and language processing* 15.1 (2006): 70-79.
- [16] Ono, Nobutaka. "Stable and fast update rules for independent vector analysis based on auxiliary function technique." *2011 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2011.
- [17] Taniguchi, Toru, et al. "An auxiliary-function approach to online independent vector analysis for real-time blind source separation." *2014 4th Joint Workshop on Hands-free Speech Communication and Microphone Arrays (HSCMA)*. IEEE, 2014.
- [18] Scheibler, Robin, and Nobutaka Ono. "Independent vector analysis with more microphones than sources." *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*. IEEE, 2019.
- [19] Ono, Nobutaka. "Auxiliary-function-based independent vector analysis with power of vector-norm type weighting functions." *Proceedings of The 2012 Asia Pacific Signal and Information Processing Association Annual Summit and Conference*. IEEE, 2012.
- [20] Nesta, Francesco, Ted S. Wada, and Biing-Hwang Juang. "Batch-online semi-blind source separation applied to multi-channel acoustic echo cancellation." *IEEE transactions on audio, speech, and language processing* 19.3 (2010): 583-599.
- [21] Ikram, Muhammad Z. "Blind source separation and acoustic echo cancellation: A unified framework." *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012.
- [22] Buchner, Herbert, and Walter Kellermann. "A fundamental relation between blind and supervised adaptive filtering illustrated for blind source separation and acoustic echo cancellation." *2008 Hands-Free Speech Communication and Microphone Arrays*. IEEE, 2008.
- [23] Gunther, Jake. "Learning echo paths during continuous double-talk using semi-blind source separation." *IEEE transactions on audio, speech, and language processing* 20.2 (2011): 646-660.
- [24] Woodbury, Max A. "Inverting modified matrices." *Memorandum report* 42.106 (1950): 336.
- [25] Recommendation, ITU-T. "Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs." *Rec. ITU-T P. 862* (2001).
- [26] Chen, Mengzhe, et al. "Compact Feedforward Sequential Memory Networks for Small-footprint Keyword Spotting." *Inter-speech*. 2018.