2020　2022

A Unified Computational Model for Audio-Visual Recognition of Human Social Interaction

NUGRAHA, Aditya Arie

3,300,000

Normalizing Flow　　　　　IVA　　　　FastMNMF　　　　　　　　　　　BSS

BSS

One key achievement is the prototype of adaptive speech enhancement for real-time speech transcription with head-worn smart glasses. It involves challenging egocentric information processing with non-stationary sensors. This technology may benefit older adults and people with hearing impairment.

We aimed for a probabilistic computational model of audio-visual information processing for understanding human verbal communication. We proposed a model for generating speech signals from speaker labels controlling the voice characteristics and phone labels controlling the speech content. For speech enhancement, it potentially improves not only the signal quality but also the speech intelligibility. We also introduced principled time-varying extensions, based on a novel deep generative model called normalizing flow, of time-invariant blind source separation (BSS) methods, including the classical independent vector analysis and the state-of-the-art FastMNMF. Finally, we developed adaptive audio-visual speech enhancement with augmented reality smart glasses. Camera images allow speakers of interest to be identified to control direction-aware enhancement. We achieve robust low-latency enhancement via a fast environment-sensitive beamforming governed by a slow environment-agnostic BSS.

Audio-visual speech enhancement for smart glasses

Audio-visual processing　Smart glasses　Adaptive system　Blind source separation　Speech enhancement　Speech recognition　Neural spatial model　Generative model

## 1.　研究開始当初の背景

**Human communication relies on information obtained by the different senses in a complementary manner.** Although human senses include sight, hearing, taste, smell, and touch, visual and auditory systems are typically predominant. A human can interact with a group of people, even in an unfamiliar crowd (environment), by visually localizing and identifying the sound sources, especially the active speakers, to improve the listening focus. **Could we formulate a computational model of audio-visual information processing so that a computer can robustly do the same in different environments?**

While conventional studies on audio-visual information typically address different tasks separately, such as speaker identification, speech separation, and lip reading, **jointly considering multiple tasks based on audio-visual information in a complementary manner should be beneficial**. For example, visual-based speaker localization is more reliable than audio-based one when a speaker is seen. In contrast, audio-based localization can still work when a person is not seen but actively speaking. Associating speech characteristics obtained from speaker identification to facial features from face detection could improve the robustness of audio-visual speaker localization in unideal cases that trouble localization based on audio or visual data only, such as when speakers are occluded or silent at times. Unified processing of both modalities promises a superior performance that benefits applications, e.g., adaptive speech enhancement.

**Many studies on audio-visual information processing employ deep neural networks (DNNs) trained in a fully supervised manner, so the robustness of the DNNs to unseen inputs would be limited**, depending on the parallel training data diversity. We, therefore, want to formulate a probabilistic model to cope with the variability of audio-visual data and its interrelationship. In the DNN-based generative modeling framework called the variational autoencoder (VAE), **a generation DNN and a recognition DNN are jointly trained on the observations in an unsupervised manner to discover the latent variables**. This unsupervised training does not need parallel data, so it is practical for real-world applications.

## 2.　研究の目的

This research aimed **to formulate and realize a probabilistic computational model of audio-visual information processing**. We wanted to learn the underlying representation (latent variables), such as *message (idea)* $\mathbf{z}_m$, *voice feature* $\mathbf{z}_s$, and *facial feature* $\mathbf{z}_f$, from the audio-visual inputs, such as *speech signal* $\mathbf{s}$ and *facial expression* $\mathbf{f}$, in an unsupervised manner. It translates to modeling the joint probability distribution given by:

$$p(\mathbf{s}, \mathbf{f}, \mathbf{z}_s, \mathbf{z}_f, \mathbf{z}_m) \;=\; p(\mathbf{s}|\mathbf{z}_s, \mathbf{z}_m)\, p(\mathbf{f}|\mathbf{z}_f, \mathbf{z}_m)\, p(\mathbf{z}_s)\, p(\mathbf{z}_f)\, p(\mathbf{z}_m).$$

**This unified model may simultaneously take into account different tasks** that are often separately considered in existing studies. The tasks may include *recognition* ones, such as speaker identification (speech signal $\mathbf{s} \rightarrow$ speaker identity $\mathbf{z}_s$), face identification (face image $\mathbf{f} \rightarrow$ human identity $\mathbf{z}_f$), and speech recognition (speech signal $\mathbf{s} \rightarrow$ text $\mathbf{z}_m$), and also *generation* ones, such as speech synthesis (text $\mathbf{z}_m$ + speaker identity $\mathbf{z}_s \rightarrow$ speech $\mathbf{s}$) and voice conversion (text $\mathbf{z}_m$ + another speaker identity $\mathbf{z}_s \rightarrow$ speech $\mathbf{s}$). **We then aim to use the model for downstream tasks, such as audio-visual speech enhancement,** and expect the model competes well with models trained in a supervised manner.

## 3. 研究の方法

This research wanted **to elucidate whether our idea of human-inspired computational audio-visual information processing is achievable**. The simple model above may not be able to capture the complex relationship between the latent and observed variables because it is merely based on intuition about human capability in exploiting audio-visual inputs.

By having audio-visual speech enhancement and speaker diarization as the target downstream tasks, we planned to do research in the following order. We first wanted to develop *speech generative models*. Building upon the existing model, in which time-varying speech is generated by time-varying latent variables, we wanted to have time-invariant voice representation. We then wanted to develop *face- or lip-movement generative models*. The main challenge for all of these models is the disentanglement of the time-varying latent variables from the time-invariant ones. Finally, we wanted to have *joint speech and face- or lip-movement generative models*.

In FY2020, we developed a generative model of time-varying speech from disentangled representations of time-invariant speaker labels and time-varying phone labels. In addition to this model based on the variational autoencoder (VAE), we also developed another one based on a combination of VAE and normalizing flow (NF). See **Section 4(1)** for further description.

In FY2021, we improved blind source separation (BSS) techniques, which are useful for speech enhancement. We integrated normalizing flow (NF) into the state-of-the-art joint diagonalization techniques for spatial covariance matrices. See **Section 4(2)** for further description.

In FY2021 and FY2022, we worked on the visual aspect of audio-visual information processing. Learning from recently published works on lip-movement-informed speech enhancement, we decided to shift our focus to tackling issues in real-world scenarios where lip movement detection is often not reliable, e.g., because the target speaker is too far away from the camera or the image resolution is too low due to the camera hardware limitations. In this case, instead of lip movement, human faces or bodies can still be detected from the camera images to inform speech enhancement about the locations of the target speakers. See **Section 4(3)** for further description.

## 4. 研究成果

Due to space limitations, we can only briefly describe a few representative research results. We grouped the results into speech generative models, time-varying blind source separation, and adaptive audio-visual speech enhancement with smart glasses.

## (1) Speech Generative Models

*Du et al.* (2020) proposed **a phone- and speaker-aware speech generative model based on variational autoencoder (VAE)** given three latent variables: (1) time-invariant speaker label, encoding the voice characteristic, (2) time-varying phone label, conveying the message, and (3) additional time-varying latent variables, holding other unspecified aspects in speech. The generated speech sounded natural with a good set of latent variables, and voice conversion could be done by modifying the speaker label. The model also worked well for speech separation outperforming the classical model without phone and speaker information. Additionally, *Nugraha et al.* (2020b) presented **another speech generative model based on a combination of VAE and the normalizing flow (NF)**. We showed that this novel model represents and produces better speech harmonics and improves a speech enhancement system utilizing it.

## (2)  Time-Varying Blind Source Separation Based on Normalizing Flow

*Nugraha et al.* (2020a) introduced **NF-IVA, a time-varying extension of independent vector analysis (IVA) based on the normalizing flow (NF) for determined BSS of multi-channel audio signals**. It estimates demixing matrices that transform mixture spectra to source spectra in the complex-valued spatial domain such that the likelihood of those matrices for the mixture spectra is maximized under some non-Gaussian source model by gradient descent in an unsupervised manner. *Nugraha et al.* (2022) then proposed **NF-FastMNMF that integrates NF into the multichannel nonnegative matrix factorization with jointly-diagonalizable spatial covariance matrices, a.k.a. FastMNMF, for determined and non-determined BSS**. NF-FastMNMF successfully performed separations of multiple speech utterances by stationary or non-stationary speakers in noisy-reverberant environments.

## (3)  Adaptive Audio-Visual Speech Enhancement with Smart Glasses

The *dual process theory* in human cognition study postulates that human thought can arise from two different processes: a fast, intuitive process and a slow, deliberate process. **We adopt the concept to develop dual-process systems for audio-visual speech enhancement with augmented reality (AR) head-mounted display, a.k.a. smart glasses.** Using smart glasses equipped with cameras and microphones, sources (speakers) of interest can be identified from the camera images. We use this directional information to drive speech enhancement. The downstream task is a streaming automatic speech recognition (ASR) for cocktail-party conversation assistance that would benefit older adults and people with hearing impairment.

*Nugraha et al.* (2022) proposed **a practical dual-process visually-informed speech enhancement system that adapts environment-sensitive frame-online beamforming (front end) with help from environment-agnostic block-online BSS (back end)**. To use minimum variance distortionless response (MVDR) beamforming, one may train a deep neural network (DNN) that estimates time-frequency masks used for computing the covariance matrices of sources (speech and noise). Instead, one may try directly estimating the source covariance matrices with a BSS method, such as the state-of-the-art FastMNMF. In practice, however, neither the DNN nor the FastMNMF can be updated in a frame-online manner due to its computationally-expensive nature. Our DNN-free system, as shown in **Figure 1**, leverages the posteriors of the latest source spectrograms given by block-online FastMNMF to derive the current source covariance matrices for frame-online beamforming. This method is applied with a blind dereverberation method called weighted prediction error (WPE). The evaluation shows that our system can respond to scene changes due to interfering speaker movements and outperformed a system with DNN-based beamforming in terms of an ASR performance metric called word error rate (WER).
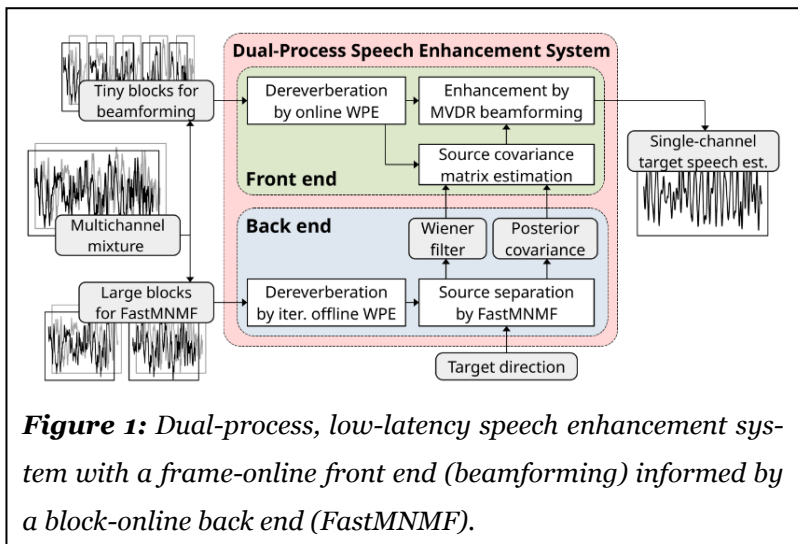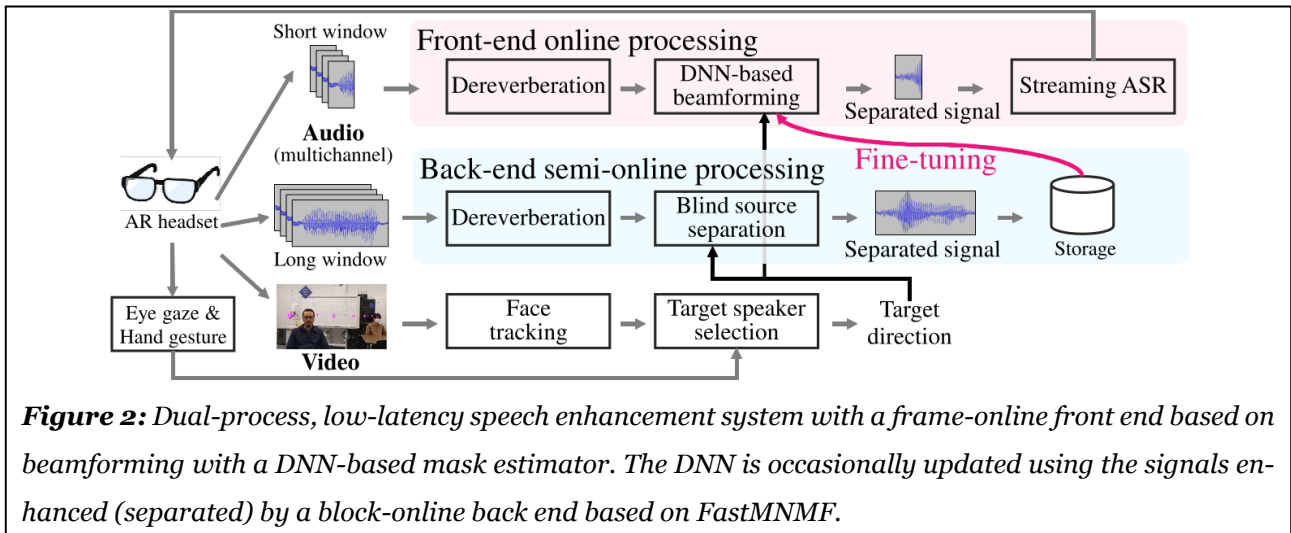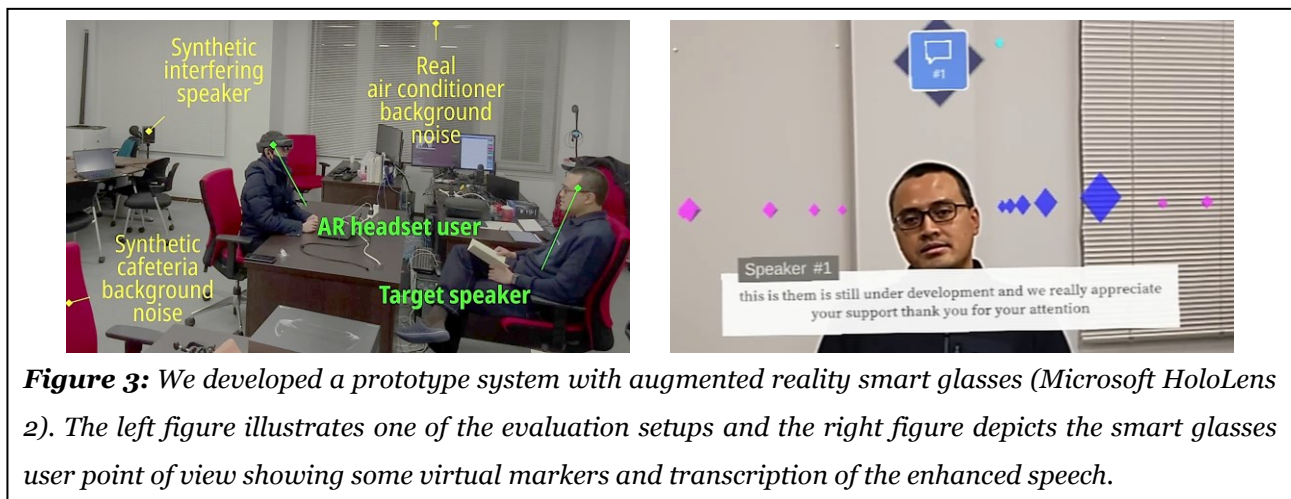


**Figure 1:** *Dual-process, low-latency speech enhancement system with a frame-online front end (beamforming) informed by a block-online back end (FastMNMF).*

**Figure 2:** *Dual-process, low-latency speech enhancement system with a frame-online front end based on beamforming with a DNN-based mask estimator. The DNN is occasionally updated using the signals enhanced (separated) by a block-online back end based on FastMNMF.*

*Sekiguchi et al.* (2022) introduced another **dual-process visually-informed speech enhancement system that utilizes block-online FastMNMF (back end) for run-time adaptation of frame-online beamforming with DNN-based mask estimator (front end)**. FastMNMF works well in various environments thanks to its unsupervised nature. In contrast, beamforming that uses a pre-trained DNN for estimating spatial information of sources (speech and noise) suffers from drastic performance degradation in mismatched conditions. We thus propose speech enhancement based on DNN-based beamforming with FastMNMF-guided adaptation, as shown in **Figure 2**. Pairs of a noisy speech and its enhanced speech obtained by FastMNMF (back end) are used together with the original parallel training data for updating the DNN (front end) at a computationally-allowable interval. Our experiments showed that a run-time adaptation using only 12 minutes of observation could improve the WER by more than 10 points. By taking the idea of model adaptation further, *Du et al.* (2022) jointly fine-tuned the DNN-based mask estimator and ASR models to maximize the ASR objective function.

We developed **a prototype system using Microsoft HoloLens 2**, which was the only developer-friendly holographic smart glasses available in the market. Our system uses visual information from the camera images to provide possible target speakers from which the user can select using eye gaze or hand gestures. The speech enhancement front end extracts possibly multiple speech signals coming from the target directions corresponding to the user selection and input the extracted speech signals into an ASR system. The generated transcriptions are then displayed as virtual content on the AR headset, as shown in **Figure 3**.



**Figure 3:** *We developed a prototype system with augmented reality smart glasses (Microsoft HoloLens 2). The left figure illustrates one of the evaluation setups and the right figure depicts the smart glasses user point of view showing some virtual markers and transcription of the enhanced speech.*

| | 6 | 6 | 6 | 4 |
|---|---|---|---|---|

| | |
|---|---|
| Fontaine Mathieu  Sekiguchi Kouhei  Nugraha Aditya Arie  Bando Yoshiaki  Yoshii Kazuyoshi | 30 |
| Generalized Fast Multichannel Nonnegative Matrix Factorization Based on Gaussian Scale Mixtures for Blind Source Separation | 2022 |
| IEEE/ACM Transactions on Audio, Speech, and Language Processing | 1734  1748 |
| DOI<br>10.1109/TASLP.2022.3172631 | |
| | |

| | |
|---|---|
| Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Fontaine Mathieu  Yoshii Kazuyoshi  Kawahara Tatsuya | 30 |
| Autoregressive Moving Average Jointly-Diagonalizable Spatial Covariance Analysis for Joint Source Separation and Dereverberation | 2022 |
| IEEE/ACM Transactions on Audio, Speech, and Language Processing | 2368  2382 |
| DOI<br>10.1109/TASLP.2022.3190734 | |
| | |

| | |
|---|---|
| Bando Yoshiaki  Sekiguchi Kouhei  Masuyama Yoshiki  Nugraha Aditya Arie  Fontaine Mathieu  Yoshii Kazuyoshi | 28 |
| Neural Full-Rank Spatial Covariance Analysis for Blind Source Separation | 2021 |
| IEEE Signal Processing Letters | 1670  1674 |
| DOI<br>10.1109/LSP.2021.3101699 | |
| | |

| | |
|---|---|
| Nugraha Aditya Arie  Sekiguchi Kouhei  Yoshii Kazuyoshi | 28 |
| A Flow-Based Deep Latent Variable Model for Speech Spectrogram Modeling and Enhancement | 2020 |
| IEEE/ACM Transactions on Audio, Speech, and Language Processing | 1104  1117 |
| DOI<br>10.1109/TASLP.2020.2979603 | |
| | |

| | |
|---|---|
| Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Yoshii Kazuyoshi  Kawahara Tatsuya | 28 |
| Fast Multichannel Nonnegative Matrix Factorization With Directivity-Aware Jointly-Diagonalizable Spatial Covariance Matrices for Blind Source Separation | 2020 |
| IEEE/ACM Transactions on Audio, Speech, and Language Processing | 2610  2625 |
| DOI<br>10.1109/TASLP.2020.3019181 | |
| | |

| | |
|---|---|
| Nugraha Aditya Arie  Sekiguchi Kouhei  Fontaine Mathieu  Bando Yoshiaki  Yoshii Kazuyoshi | 27 |
| Flow-Based Independent Vector Analysis for Blind Source Separation | 2020 |
| IEEE Signal Processing Letters | 2173  2177 |
| DOI<br>10.1109/LSP.2020.3039944 | |
| | |

14          0          11

| |
|---|
| Nugraha Aditya Arie  Sekiguchi Kouhei  Fontaine Mathieu  Bando Yoshiaki  Yoshii Kazuyoshi |
| Flow-Based Fast Multichannel Nonnegative Matrix Factorization for Blind Source Separation |
| IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) |
| 2022 |

| |
|---|
| Fontaine Mathieu  Di Carlo Diego  Sekiguchi Kouhei  Nugraha Aditya Arie  Bando Yoshiaki  Yoshii Kazuyoshi |
| Elliptically Contoured Alpha-Stable Representation for MUSIC-Based Sound Source Localization |
| European Signal Processing Conference (EUSIPCO) |
| 2022 |

Sumura Yoshiaki  Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Yoshii Kazuyoshi

Joint Localization and Synchronization of Distributed Camera-Attached Microphone Arrays for Indoor Scene Analysis

International Workshop on Acoustic Signal Enhancement (IWAENC)

2022

Nugraha Aditya Arie  Sekiguchi Kouhei  Fontaine Mathieu  Bando Yoshiaki  Yoshii Kazuyoshi

DNN-Free Low-Latency Adaptive Speech Enhancement Based on Frame-Online Beamforming Powered by Block-Online FastMNMF

International Workshop on Acoustic Signal Enhancement (IWAENC)

2022

Du Yicheng  Nugraha Aditya Arie  Sekiguchi Kouhei  Bando Yoshiaki  Fontaine Mathieu  Yoshii Kazuyoshi

Direction-Aware Joint Adaptation of Neural Speech Enhancement and Recognition in Real Multiparty Conversational Environments

Annual Conference of the International Speech Communication Association (Interspeech)

2022

Sekiguchi Kouhei  Nugraha Aditya Arie  Du Yicheng  Bando Yoshiaki  Fontaine Mathieu  Yoshii Kazuyoshi

Direction-Aware Adaptive Online Neural Speech Enhancement with an Augmented Reality Headset in Real Noisy Conversational Environments

IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)

2022

| |
|---|
| Fontaine Mathieu  Sekiguchi Kouhei  Nugraha Aditya Arie  Bando Yoshiaki  Yoshii Kazuyoshi |
| Alpha-Stable Autoregressive Fast Multichannel Nonnegative Matrix Factorization for Joint Speech Enhancement and Dereverberation |
| Annual Conference of the International Speech Communication Association (Interspeech) |
| 2021 |

| |
|---|
| Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Fontaine Mathieu  Yoshii Kazuyoshi |
| Autoregressive Fast Multichannel Nonnegative Matrix Factorization For Joint Blind Source Separation And Dereverberation |
| IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) |
| 2021 |

| |
|---|
| Nugraha Aditya Arie  Sekiguchi Kouhei  Fontaine Mathieu  Bando Yoshiaki  Yoshii Kazuyoshi |
| Determined Blind Source Separation Based on NF-IVA with Time-Varying Linear Transformations |
| The Spring Meeting of the Acoustical Society of Japan (ASJ) |
| 2021 |

| |
|---|
| Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Fontaine Mathieu  Yoshii Kazuyoshi |
| Joint Blind Source Separation and Dereverberation Based on ARMA-FastMNMF |
| The Spring Meeting of the Acoustical Society of Japan (ASJ) |
| 2021 |

| Nugraha Aditya Arie  Sekiguchi Kouhei  Fontaine Mathieu  Bando Yoshiaki  Yoshii Kazuyoshi |
| --- |
| Unsupervised Source Separation with Deep Spatial Models |
| RIKEN AIP Open Seminar |
| 2021 |

| Fontaine Mathieu  Sekiguchi Kouhei  Nugraha Aditya Arie  Yoshii Kazuyoshi |
| --- |
| Unsupervised Robust Speech Enhancement Based on Alpha-Stable Fast Multichannel Nonnegative Matrix Factorization |
| Annual Conference of the International Speech Communication Association (Interspeech) |
| 2020 |

| Yoshii Kazuyoshi  Sekiguchi Kouhei  Bando Yoshiaki  Fontaine Mathieu  Nugraha Aditya Arie |
| --- |
| Fast Multichannel Correlated Tensor Factorization for Blind Source Separation |
| European Signal Processing Conference (EUSIPCO) |
| 2020 |

| Du Yicheng  Sekiguchi Kouhei  Bando Yoshiaki  Nugraha Aditya Arie  Fontaine Mathieu  Yoshii Kazuyoshi  Kawahara Tatsuya |
| --- |
| Semi-supervised Multichannel Speech Separation Based on a Phone- and Speaker-Aware Deep Generative Model of Speech Spectrograms |
| European Signal Processing Conference (EUSIPCO) |
| 2020 |

o

Demo web page for NF-FastMNMF: https://aanugraha.github.io/demo/nffastmnmf/
Demo web page for Neural FCA: https://ybando.jp/projects/spl2021/
Demo web page for NF-IVA: https://aanugraha.github.io/demo/nfiva/
Demo web page for GF-VAE: https://aanugraha.github.io/demo/gfvae/

| | | | |
|---|---|---|---|
| | | | |

o

| | |
|---|---|
| | |