

Unscented Transform-Based Dual-Channel Noise Estimation: Application to Speech Enhancement on Smartphones

Iván López-Espejo¹, Juan M. Martín-Doñas², Angel M. Gomez², and Antonio M. Peinado²

¹VeriDas | das-Nano, Navarre, Spain

²Dept. of Signal Theory, Telematics and Communications, University of Granada, Granada, Spain
ilopez@das-nano.com, {mdjuamart, amgg, amp}@ugr.es

41st International Conference on Telecommunications and Signal Processing
7-2018, Athens (Greece)

Outline

- 1 Introduction
 - Motivation
 - Objectives
- 2 UT-Based Dual-Channel Noise Estimation
 - Dual-Channel MMSE Estimation
 - Dual-Channel Distortion Model
 - Unscented Transform Application
 - Implementation Issues
- 3 Experimental Evaluation
 - The AURORA2-2C-CT Database
 - Estimation Error Results
 - Speech Quality Results
 - Speech Intelligibility Results
- 4 Conclusions

Introduction

Motivation

New speech processing boom

The use of speech processing applications has notably increased due to the latest smartphones:

- Great amount of apps (search-by-voice, dictation, voice biometrics, etc.).

Noise-robust speech processing in smartphones

- It is crucial to deal with noisy environments.
- We can take benefit of the dual-mic feature.



Introduction

Objectives

- 1 **To estimate the noise power spectrum** at the primary channel by exploiting the information at both channels.
 - Dual-channel MMSE estimation.
 - Unscented transform (UT) vs. vector Taylor series (VTS).
- 2 To explore a particular application of interest: **speech enhancement on a dual-mic smartphone** in close-talk conditions.



UT-Based Dual-Channel Noise Estimation

Dual-Channel MMSE Estimation

- $k = 1$ is the primary mic while $k = 2$ is the extra mic:

$$\mathbf{y}_k(t) = (|Y_k(0, t)|^2, \dots, |Y_k(\mathcal{M} - 1, t)|^2)^\top$$

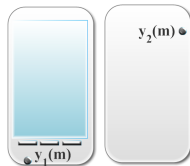
$$\mathbf{n}_k(t) = (|N_k(0, t)|^2, \dots, |N_k(\mathcal{M} - 1, t)|^2)^\top$$

- The MMSE estimate of $\mathbf{n}_1(t)$ is

$$\hat{\mathbf{n}}_1(t) = \mu_{n_1}(t) + \underbrace{\Sigma_{n_1 y_1}(t) \Sigma_{y_1}^{-1}(t)}_{\text{Through the unscented transform!}} (\underbrace{\mathbf{y}_1(t) - \mu_{y_1}(t)}_{\text{Through the unscented transform!}})$$

Through the unscented transform!

- $\mathbf{y}_2(t)$ will play here a supporting role.



UT-Based Dual-Channel Noise Estimation

Dual-Channel Distortion Model

- Dual-channel speech distortion model:

$$\mathbf{y}_1(t) = \mathbf{x}_1(t) + \mathbf{n}_1(t)$$

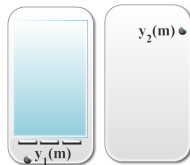
$$\mathbf{y}_2(t) = \mathbf{x}_2(t) + \mathbf{n}_2(t)$$

- We relate the speech power at both channels through speech gains:

$$\mathbf{x}_2(t) = \mathbf{a}_{21}(t) \odot \mathbf{x}_1(t)$$

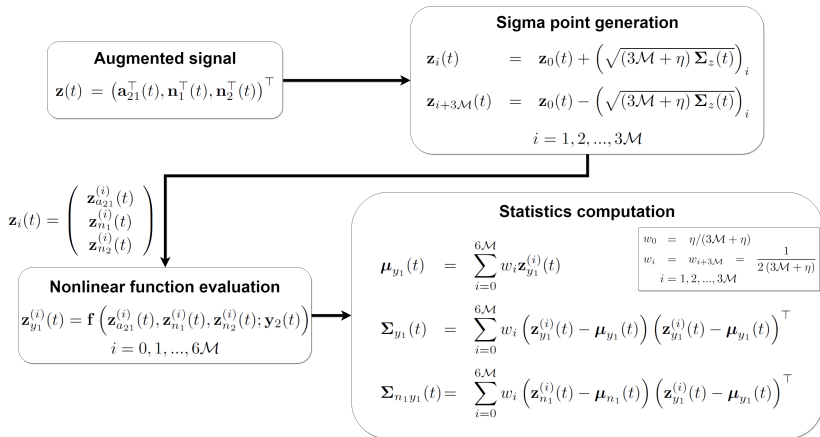
- We sample the next $\mathbf{y}_1(t)$ model through **UT** to get $\boldsymbol{\mu}_{y_1}(t)$, $\boldsymbol{\Sigma}_{y_1}(t)$ and $\boldsymbol{\Sigma}_{n_1 y_1}(t)$:

$$\begin{aligned} \mathbf{y}_1(t) &= \mathbf{f}(\mathbf{a}_{21}(t), \mathbf{n}_1(t), \mathbf{n}_2(t); \mathbf{y}_2(t)) \\ &= \mathbf{a}_{21}^{-1}(t) \odot (\mathbf{y}_2(t) - \mathbf{n}_2(t)) + \mathbf{n}_1(t) \end{aligned}$$



UT-Based Dual-Channel Noise Estimation

Unscented Transform Application



UT-Based Dual-Channel Noise Estimation

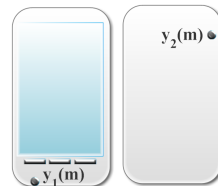
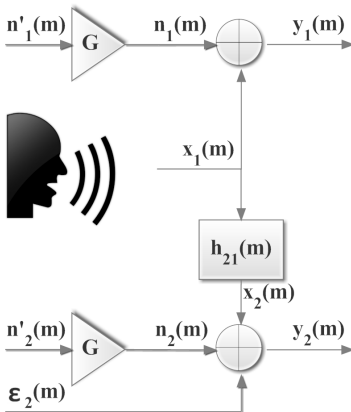
Implementation Issues

Relevant practical issues

- We assume that both $\mathbf{a}_{21}(t)$ and $\mathbf{n}_k(t)$ ($k = 1, 2$) are wide-sense stationary random processes:
 - $\boldsymbol{\mu}_{a_{21}}$ and $\boldsymbol{\Sigma}_{a_{21}}$ are obtained from a development dataset.
 - Noise statistics are obtained from the first and last frames of each utterance.
- η is set to 0.
- Negative estimated bins are replaced by those from a noise estimation based on linear interpolation (around 2% of cases).

Experimental Evaluation

The AURORA2-2C-CT Database

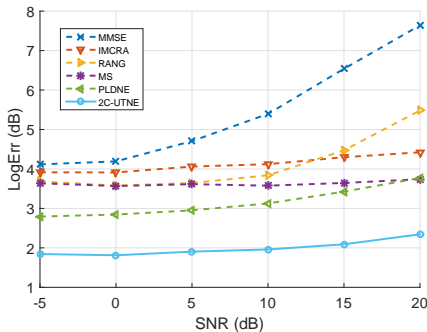


- **Test A:** Bus, babble, car and pedestrian street
- **Test B:** Café, street, bus and train stations
- **SNRs:** $\{-5, 0, 5, 10, 15, 20\}$ dB and clean

López-Espejo I., et al.: "Feature Enhancement for Robust Speech Recognition on Smartphones with Dual-Microphone". In: *EUSIPCO*, Lisbon (2014)

Experimental Evaluation

Estimation Error Results



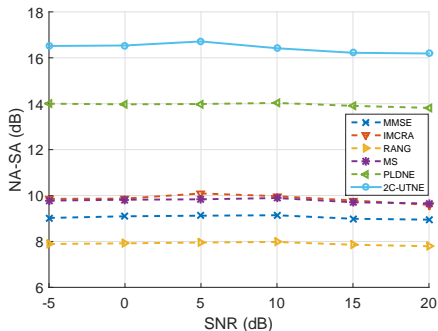
- **2C-UTNE - Our proposal**
- **PLDNE** - Power Level Difference Noise Estimator
- **MS** - Minimum Statistics
- **IMCRA** - Improved Minima Controlled Recursive Averaging
- **RANG** - Rangachari's algorithm
- **MMSE** - MMSE noise estimator

$$\text{LogErr} = \frac{1}{\mathcal{M}T} \sum_{f=0}^{\mathcal{M}-1} \sum_{t=0}^{T-1} \left| 10 \log_{10} \left[\frac{\Phi_{n_1}(f, t)}{\hat{\Phi}_{n_1}(f, t)} \right] \right|$$

$$\Phi_{n_1}(f, t) = 0.9\Phi_{n_1}(f, t-1) + 0.1|N_1(f, t)|^2$$

Experimental Evaluation

Speech Quality Results



- Noise estimation methods are combined with **dual-channel Wiener filtering**.
- **NA-SA**: Noise Attenuation minus Speech Attenuation.

- **2C-UTNE** - Our proposal
- **PLDNE** - Power Level Difference Noise Estimator
- **MS** - Minimum Statistics
- **IMCRA** - Improved Minima Controlled Recursive Averaging
- **RANG** - Rangachari's algorithm
- **MMSE** - MMSE noise estimator

Experimental Evaluation

Speech Intelligibility Results

	SNR (dB)						Avg.
	-5	0	5	10	15	20	
Noisy	0.280	0.415	0.556	0.693	0.819	0.909	0.612
MMSE	0.321	0.478	0.629	0.761	0.867	0.936	0.665
IMCRA	0.313	0.470	0.621	0.753	0.858	0.928	0.657
RANG	0.323	0.477	0.627	0.760	0.866	0.935	0.665
MS	0.319	0.476	0.627	0.760	0.867	0.936	0.664
PLDNE	0.346	0.496	0.640	0.767	0.868	0.933	0.675
2C-UTNE	0.350	0.506	0.653	0.780	0.879	0.941	0.685

- Noise estimation methods are combined with **dual-channel Wiener filtering**.
- **CSII**: Coherence Speech Intelligibility Index.

- **2C-UTNE** - Our proposal
- **PLDNE** - Power Level Difference Noise Estimator
- **MS** - Minimum Statistics
- **IMCRA** - Improved Minima Controlled Recursive Averaging
- **RANG** - Rangachari's algorithm
- **MMSE** - MMSE noise estimator

Conclusions

Some conclusions and future work

- We have achieved accurate noise estimates by taking advantage of...
 - 1 the unscented transform.
 - 2 the dual-channel observation (avoiding the use of a clean speech model while keeping a simple formulation).
- Results have shown the higher performance of our proposal with respect to state-of-the-art noise estimation methods.
- We will investigate on temporal dynamics modeling to further improve the performance of our method.

Thanks for your attention!



Iván López-Espejo
VeriDas | das-Nano, Navarre, Spain
ilopez@das-nano.com