

Multichannel Speech Enhancement Based on Time-frequency Masking Using Subband Long Short-Term Memory

Xiaofei Li and Radu Horaud
Inria Grenoble, France

1 Introduction

2 Subband LSTM network

3 Experiments

Introduction

Scenario

- Single speaker
- Ambient noise
- Microphone array

Task

- Speech denoising

Conventional unsupervised subband methods

- Single-channel: noise estimation + spectral subtraction ¹
- Multichannel: beamforming (spatial filtering) ²

¹Y. Ephraim, et al, "Speech enhancement using a minimum-meansquare error short-time spectral amplitude estimator," *TASP*, 1984.

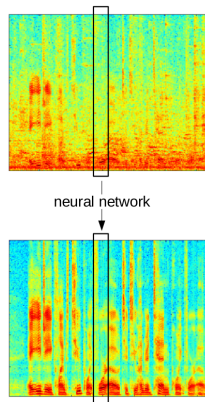
²S. Gannot, et al, "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *TASLP*, 2017.

Introduction

Deep neural network methods ³

Single-channel: spectral regression

- Learn using full-band spectra and dynamics (past, present and future frames)
- Sensitive to speaker identities and noise types
- Input/output (noisy speech / clean speech) : high dimensional vectors, therefore a very large network is necessary



³D. Wang, et al, "Supervised speech separation based on deep learning: An overview," *TASLP*, 2018.

Deep neural network methods

Multichannel

- Single-channel spectral regression + beamforming ⁴
- Full-band spatial information learning ⁵

⁴C. Boeddeker, et al, “Exploring practical aspects of neural mask-based beamforming for far-field speech recognition,” *ICASSP*, 2018.

⁵S. Chakrabarty, et al, “Time-frequency masking based online speech enhancement with multi-channel data using convolutional neural networks,” *IWAENC*, 2018.

Proposed: subband LSTM (Long Short-Term Memory)

Motivation: unsupervised methods (spectral subtraction and beamforming) show that subband-based methods are powerful to discriminate between speech and noise:

- Speech is temporally nonstationary / noise exhibits more stationarity
- Speech is spatially coherent / noise is diffuse
- Spectral subtraction and spatial filtering are performed in subbands

Objective: train an LSTM network that exploits subband information.

1 Introduction

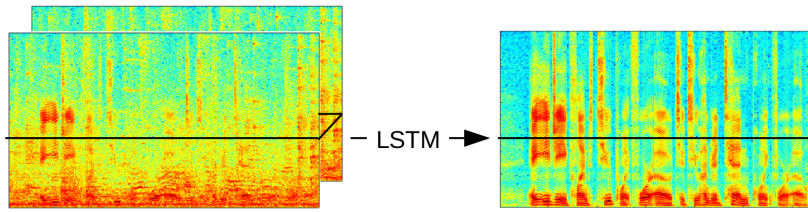
2 Subband LSTM network

3 Experiments

Subband LSTM network

Train one LSTM common to all frequencies

- Learn spatial and temporal information common to all frequencies
- Exclude cross-band spectral information
- Map *multichannel-sequence-to-speech-sequence* for each subband



Subband LSTM network

Signal formulation in the STFT domain: $x_{t,k}^i = s_{t,k}^i + u_{t,k}^i$

- Input: STFT Fourier coefficients

$$\mathbf{x}_{t,k} = [\mathcal{R}(x_{t,k}^1), \mathcal{I}(x_{t,k}^1) \dots \mathcal{R}(x_{t,k}^l), \mathcal{I}(x_{t,k}^l)]^T \in \mathbb{R}^{2l}$$

- Target: rectified magnitude ratio mask (MRM)

$$M_{t,k} = \min \left(\frac{|s_{t,k}^1|}{|x_{t,k}^1|}, 1 \right) \in [0, 1]$$

- Loss:

$$(\hat{M}(\mathbf{x}_{t,k}, \theta) - M_{t,k})^2 \quad \forall k$$

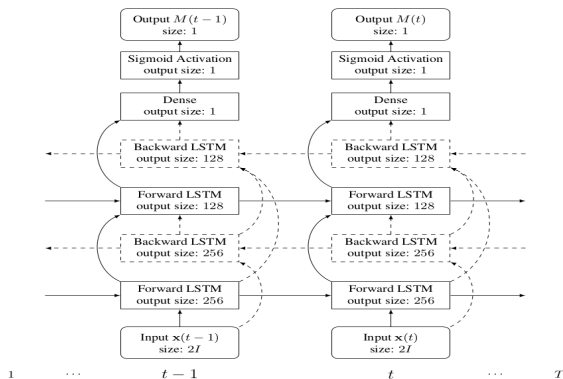
- Test: denoising using the predicted soft mask

$$\begin{aligned} |\hat{s}_{t,k}| &= \hat{M}_{t,k} |x_{t,k}^1|, \\ \arg(\hat{s}_{t,k}) &= \arg(x_{t,k}^1) \end{aligned}$$

Subband LSTM network

Bidirectional LSTM network

- Sequence-to-sequence
- Small-sized network: two BLSTM layers, 1.19 M learnable parameters



1 Introduction

2 Subband LSTM network

3 Experiments

CHiME3 dataset ⁶, six microphones, four noise types

- BUS: bus
- CAF: coffee shop
- PED: pedestrian area
- STR: street junction

Data generation: multichannel *booth speech* with added background noise

- Four speakers for training, one for validation, three for test
- Training data: about 11 hours generated data, cut into sequences with length of 192 frames (LSTM learns 192 time steps of memory)
- Test data: original CHiME utterances with varying lengths, SNR is 0 dB

⁶J. Barker, et al, "The third 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," *ASRU*, 2015.

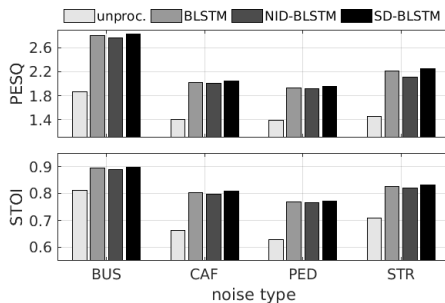
Performance metrics

- PESQ (perceptual evaluation of speech quality): speech quality
- STOI (short-time objective intelligibility): speech intelligibility

Experiments

Network generalization performance in terms of speaker and noise type

- BLSTM: speaker-independent, noise-type-dependent
- NID-BLSTM: speaker-independent, noise-type-independent
- SD-BLSTM: speaker-dependent, noise-type-dependent
- Results: two microphones



Baseline methods

- BeamformIt: unsupervised filter-and-sum beamforming ⁷
- Multichannel full-band CNN ⁸
 - Input: multichannel, single-frame, full-band STFT coefficients
 - Output: full-band time-frequency mask
 - CNN network

⁷X. Anguera, et al, “Acoustic beamforming for speaker diarization of meetings,” *TASLP*, 2007.

⁸S. Chakrabarty, et al, “Time-frequency masking based online speech enhancement with multi-channel data using convolutional neural networks,” *IWAENC*, 2018.

Audio examples ⁹ (four microphones)

	unproc.	BeamformIt	CNN	BLSTM
CAF	audio	audio	audio	audio
STR	audio	audio	audio	audio

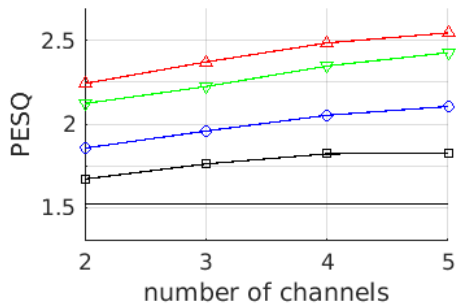
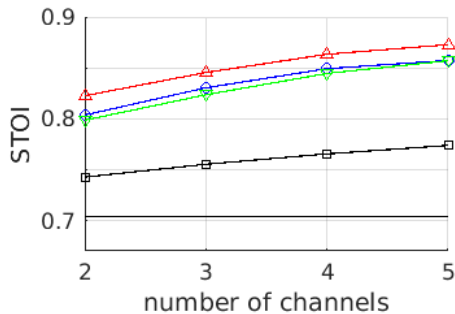
⁹<https://team.inria.fr/perception/research/mse-lstm/>

Experiments

Quantitative results (averaged over all noise types)

- The proposed achieves better speech quality

— unproc. —□— BeamformIt —◇— CNN —▽— LSTM —△— BLSTM



Multichannel subband LSTM

- Generalizes well to unseen speakers and noise types
- Complementary to full-band spectral regression methods

In the future we plan to experiment with other targets, such as

- *complex ideal ratio mask* or *Fourier coefficients* to predict both the module and the phase,
- *spatial filtering* to mimic beamforming.

`https://team.inria.fr/perception/research/mse-lstm/`