

DNN-Based Speech Presence Probability Estimation for Multi-Frame Single-Microphone Speech Enhancement

ICASSP 2020

AUD-L7: Signal Enhancement and Restoration I

Marvin Tammen, Dörte Fischer, Bernd Meyer, Simon Doclo

Department of Medical Physics and Acoustics
Carl von Ossietzky University Oldenburg

May 2020

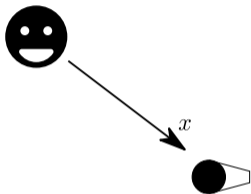
- 1 Problem Statement
- 2 Multi-Frame Filtering
- 3 Parameter Estimation
- 4 Experiments
- 5 Conclusions & Outlook

- 1 Problem Statement
- 2 Multi-Frame Filtering
- 3 Parameter Estimation
- 4 Experiments
- 5 Conclusions & Outlook

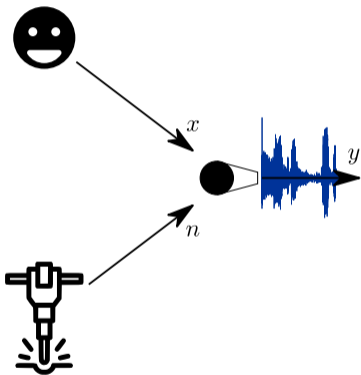
Problem Statement



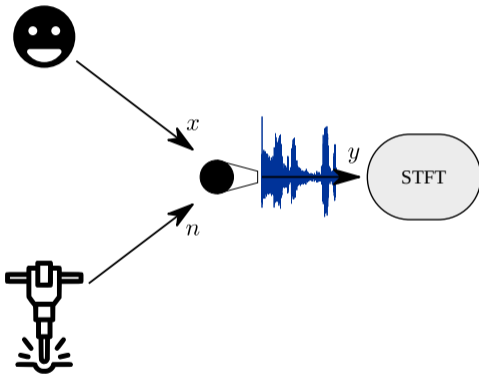
Problem Statement



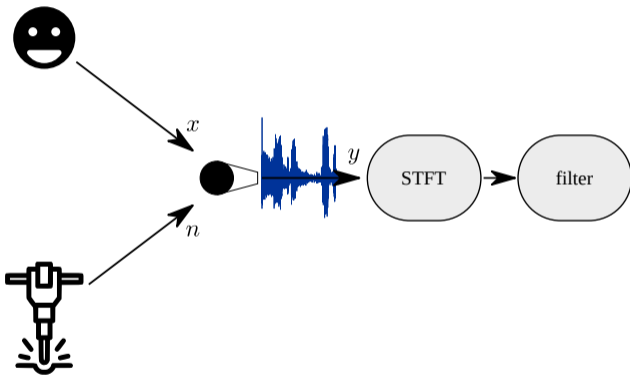
Problem Statement



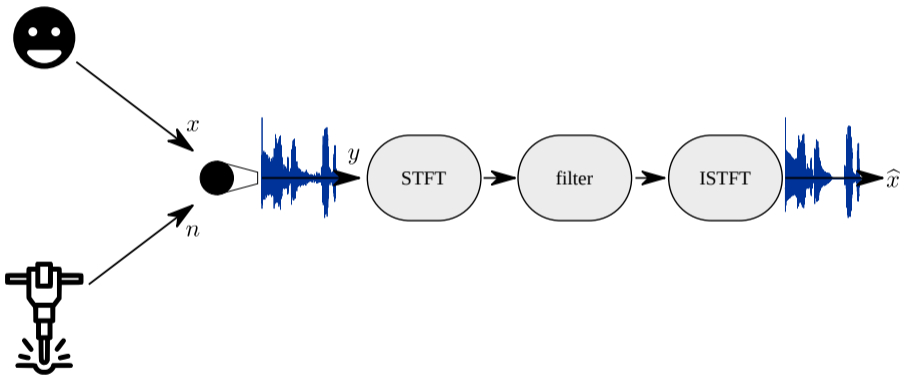
Problem Statement



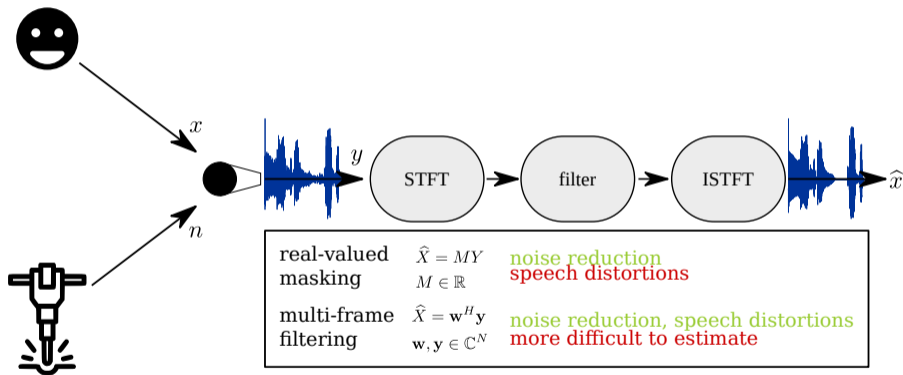
Problem Statement



Problem Statement



Problem Statement



- 1 Problem Statement
- 2 Multi-Frame Filtering**
 - Signal Model & Assumptions
 - Problem & Solution
- 3 Parameter Estimation
- 4 Experiments
 - Datasets & Settings
 - Results
- 5 Conclusions & Outlook

Signal Model & Assumptions

- ① additive noise:

$$\mathbf{y}(k, l) = \mathbf{x}(k, l) + \mathbf{n}(k, l)$$

¹Huang and Benesty 2012.

Signal Model & Assumptions

- 1 additive noise:

$$\mathbf{y}(k, l) = \mathbf{x}(k, l) + \mathbf{n}(k, l)$$

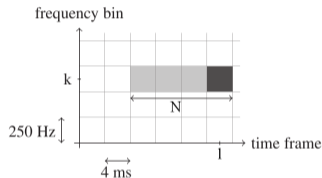


Figure: multi-frame filtering,
image adapted from Fischer
et al. 2016

¹Huang and Benesty 2012.

Signal Model & Assumptions

- 1 additive noise:

$$\mathbf{y}(k, l) = \mathbf{x}(k, l) + \mathbf{n}(k, l)$$

- 2 independent speech and noise:

$$\Phi_{\mathbf{y}}(l) = \mathcal{E}\{\mathbf{y}(l)\mathbf{y}^H(l)\} = \Phi_{\mathbf{x}}(l) + \Phi_{\mathbf{n}}(l)$$

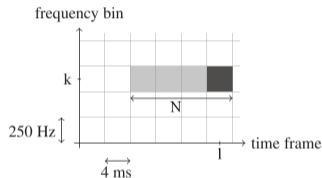


Figure: multi-frame filtering,
image adapted from Fischer
et al. 2016

¹Huang and Benesty 2012.

Signal Model & Assumptions

- ① additive noise:

$$\mathbf{y}(k, l) = \mathbf{x}(k, l) + \mathbf{n}(k, l)$$

- ② independent speech and noise:

$$\Phi_{\mathbf{y}}(l) = \mathcal{E}\{\mathbf{y}(l)\mathbf{y}^H(l)\} = \Phi_{\mathbf{x}}(l) + \Phi_{\mathbf{n}}(l)$$

- ③ speech has correlated and uncorrelated component¹:

$$\mathbf{x}(l) = \underbrace{\gamma_{\mathbf{x}}(l)\mathbf{X}(l)}_{\text{correlated}} + \underbrace{\mathbf{x}'(l)}_{\text{uncorrelated}}$$

$$\gamma_{\mathbf{x}}(l) = \frac{\mathcal{E}\{\mathbf{x}(l)\mathbf{X}^*(l)\}}{\mathcal{E}\{|\mathbf{X}(l)|^2\}} = \frac{\Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T\Phi_{\mathbf{x}}(l)\mathbf{e}}$$

$\gamma_{\mathbf{x}}(l)$: normalized speech interframe correlation vector

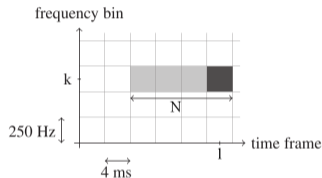


Figure: multi-frame filtering,
 image adapted from Fischer
 et al. 2016

¹Huang and Benesty 2012.

Optimization Problem and Solution

- optimization goals²:

²Huang and Benesty 2012.

Optimization Problem and Solution

- optimization goals²:
 - 1 minimize output power

$$\operatorname{argmin}_{\mathbf{w}(l) \in \mathbb{C}^N} \mathbf{w}^H(l) \Phi_y(l) \mathbf{w}(l)$$

²Huang and Benesty 2012.

Optimization Problem and Solution

- optimization goals²:
 - 1 minimize output power
 - 2 while preserving speech interframe correlations (IFCs)

$$\underset{\mathbf{w}(l) \in \mathbb{C}^N}{\operatorname{argmin}} \quad \mathbf{w}^H(l) \Phi_{\mathbf{y}}(l) \mathbf{w}(l), \quad \text{s.t.} \quad \mathbf{w}^H(l) \gamma_{\mathbf{x}}(l) = 1$$

²Huang and Benesty 2012.

Optimization Problem and Solution

- optimization goals²:
 - 1 minimize output power
 - 2 while preserving speech interframe correlations (IFCs)

$$\operatorname{argmin}_{\mathbf{w}(l) \in \mathbb{C}^N} \mathbf{w}^H(l) \Phi_{\mathbf{y}}(l) \mathbf{w}(l), \quad \text{s.t.} \quad \mathbf{w}^H(l) \gamma_{\mathbf{x}}(l) = 1$$

- solved by $\mathbf{w}_{\text{MFMPDR}}(l) = \frac{\Phi_{\mathbf{y}}^{-1}(l) \gamma_{\mathbf{x}}(l)}{\gamma_{\mathbf{x}}^H(l) \Phi_{\mathbf{y}}^{-1}(l) \gamma_{\mathbf{x}}(l)}$

²Huang and Benesty 2012.

Optimization Problem and Solution

- optimization goals²:
 - 1 minimize output power
 - 2 while preserving speech interframe correlations (IFCs)

$$\underset{\mathbf{w}(l) \in \mathbb{C}^N}{\operatorname{argmin}} \quad \mathbf{w}^H(l) \Phi_{\mathbf{y}}(l) \mathbf{w}(l), \quad \text{s.t.} \quad \mathbf{w}^H(l) \gamma_{\mathbf{x}}(l) = 1$$

- solved by $\mathbf{w}_{\text{MFMPDR}}(l) = \frac{\Phi_{\mathbf{y}}^{-1}(l) \gamma_{\mathbf{x}}(l)}{\gamma_{\mathbf{x}}^H(l) \Phi_{\mathbf{y}}^{-1}(l) \gamma_{\mathbf{x}}(l)}$
 → estimate $\Phi_{\mathbf{y}}(l), \gamma_{\mathbf{x}}(l)$

²Huang and Benesty 2012.

- 1 Problem Statement
- 2 Multi-Frame Filtering
 - Signal Model & Assumptions
 - Problem & Solution
- 3 Parameter Estimation**
- 4 Experiments
 - Datasets & Settings
 - Results
- 5 Conclusions & Outlook

Speech IFC Vector

- $$\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l) \mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l) \mathbf{e}}$$

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Speech IFC Vector

- $\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l)\mathbf{e}}$

- ① maximum likelihood estimator³

$$\hat{\gamma}_{\mathbf{x},\boldsymbol{\mu}}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)} \boldsymbol{\mu}_{\mathbf{n}} \quad (1)$$

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Speech IFC Vector

- $\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l)\mathbf{e}}$

- 1 maximum likelihood estimator³

$$\hat{\gamma}_{\mathbf{x},\mu}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)} \mu_{\mathbf{n}} \quad (1)$$

- 2 noise IFC vector-based estimator

$$\hat{\gamma}_{\mathbf{x},\gamma}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{n}}(l) \quad (2)$$

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Speech IFC Vector

- $\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l)\mathbf{e}}$

- 1 maximum likelihood estimator³

$$\hat{\gamma}_{\mathbf{x},\mu}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)}\mu_{\mathbf{n}} \quad (1)$$

- 2 noise IFC vector-based estimator

$$\hat{\gamma}_{\mathbf{x},\gamma}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{n}}(l) \quad (2)$$

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Speech IFC Vector

- $\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l)\mathbf{e}}$

- 1 maximum likelihood estimator³

$$\hat{\gamma}_{\mathbf{x},\mu}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)}\mu_{\mathbf{n}} \quad (1)$$

- 2 noise IFC vector-based estimator

$$\hat{\gamma}_{\mathbf{x},\gamma}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)}\hat{\gamma}_{\mathbf{n}}(l) \quad (2)$$

- a-priori SNR $\xi(l)$: decision-directed approach (DDA)⁴

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Speech IFC Vector

- $\gamma_{\mathbf{x}}(l) = \frac{1+\xi(l)}{\xi(l)}\gamma_{\mathbf{y}}(l) - \frac{1}{\xi(l)}\gamma_{\mathbf{n}}(l), \quad \xi(l) = \frac{\mathbf{e}^T \Phi_{\mathbf{x}}(l)\mathbf{e}}{\mathbf{e}^T \Phi_{\mathbf{n}}(l)\mathbf{e}}$

- 1 maximum likelihood estimator³

$$\hat{\gamma}_{\mathbf{x},\mu}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)} \mu_{\mathbf{n}} \quad (1)$$

- 2 noise IFC vector-based estimator

$$\hat{\gamma}_{\mathbf{x},\gamma}(l) = \frac{1 + \hat{\xi}(l)}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{y}}(l) - \frac{1}{\hat{\xi}(l)} \hat{\gamma}_{\mathbf{n}}(l) \quad (2)$$

- a-priori SNR $\xi(l)$: decision-directed approach (DDA)⁴

- requires $\mathbf{e}^T \hat{\Phi}_{\mathbf{n}}(l)\mathbf{e}$

³Schasse and Martin 2014.

⁴Ephraim and Malah 1984.

Correlation Matrices

- noisy: fixed recursive smoothing:

$$\hat{\Phi}_{\mathbf{y}}(l) = \lambda_{\mathbf{y}} \hat{\Phi}_{\mathbf{y}}(l-1) + (1 - \lambda_{\mathbf{y}}) \mathbf{y}(l) \mathbf{y}^H(l) \quad (3)$$

⁵Cohen 2003.

Correlation Matrices

- noisy: fixed recursive smoothing:

$$\hat{\Phi}_{\mathbf{y}}(l) = \lambda_y \hat{\Phi}_{\mathbf{y}}(l-1) + (1 - \lambda_y) \mathbf{y}(l) \mathbf{y}^H(l) \quad (3)$$

- noise: SPP-based adaptive recursive smoothing⁵:

$$\hat{\Phi}_{\mathbf{n}}(l) = \lambda_n(l) \hat{\Phi}_{\mathbf{n}}(l-1) + (1 - \lambda_n(l)) \mathbf{y}(l) \mathbf{y}^H(l), \quad \text{where} \quad (4)$$

$$\lambda_n(l) = \alpha_n + (1 - \alpha_n) \widehat{\text{SPP}}(l) \quad (5)$$

⁵Cohen 2003.

Correlation Matrices

- noisy: fixed recursive smoothing:

$$\hat{\Phi}_{\mathbf{y}}(l) = \lambda_y \hat{\Phi}_{\mathbf{y}}(l-1) + (1 - \lambda_y) \mathbf{y}(l) \mathbf{y}^H(l) \quad (3)$$

- noise: SPP-based adaptive recursive smoothing⁵:

$$\hat{\Phi}_{\mathbf{n}}(l) = \lambda_n(l) \hat{\Phi}_{\mathbf{n}}(l-1) + (1 - \lambda_n(l)) \mathbf{y}(l) \mathbf{y}^H(l), \quad \text{where} \quad (4)$$

$$\lambda_n(l) = \alpha_n + (1 - \alpha_n) \widehat{\text{SPP}}(l) \quad (5)$$

$$\Rightarrow \begin{array}{ll} \text{(i): } \widehat{\text{SPP}}(l) = 0 \rightarrow \lambda_n(l) = \alpha_n, & \text{update of } \hat{\Phi}_{\mathbf{n}}(l) \\ \text{(ii): } \widehat{\text{SPP}}(l) = 1 \rightarrow \lambda_n(l) = 1, & \text{no update of } \hat{\Phi}_{\mathbf{n}}(l) \end{array}$$

⁵Cohen 2003.

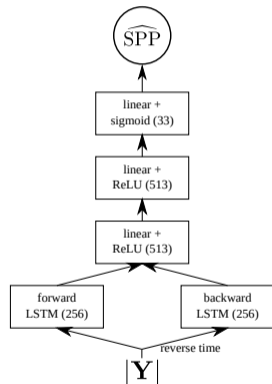
Speech Presence Probability

- 1 maximum likelihood-based⁶

⁶Gerkmann and Hendriks 2012.

Speech Presence Probability

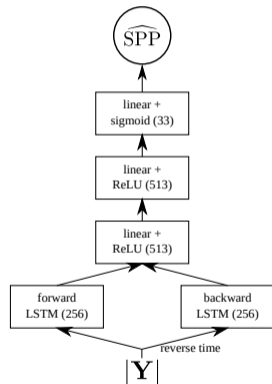
- 1 maximum likelihood-based⁶
- 2 deep recurrent neural network-based:



⁶Gerkmann and Hendriks 2012.

Speech Presence Probability

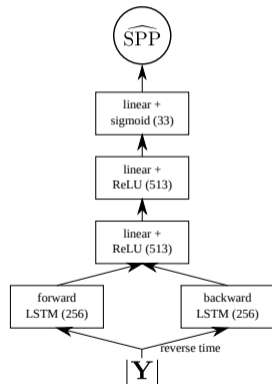
- 1 maximum likelihood-based⁶
- 2 deep recurrent neural network-based:
 - input feature: $|\mathbf{Y}| \in \mathbb{R}^{K \times L}$



⁶Gerkmann and Hendriks 2012.

Speech Presence Probability

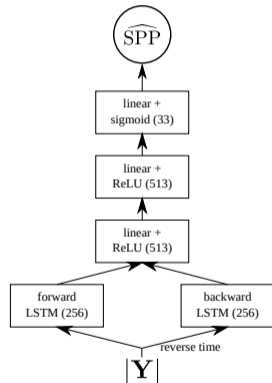
- 1 maximum likelihood-based⁶
- 2 deep recurrent neural network-based:
 - input feature: $|\mathbf{Y}| \in \mathbb{R}^{K \times L}$
 - sigmoid output to ensure $\widehat{SPP}(l) \in]0, 1[$



⁶Gerkmann and Hendriks 2012.

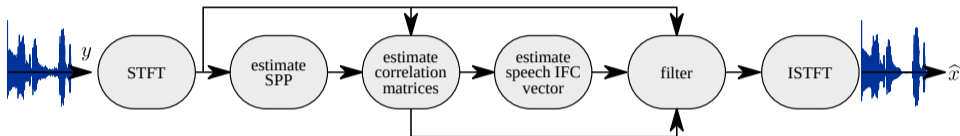
Speech Presence Probability

- 1 maximum likelihood-based⁶
- 2 deep recurrent neural network-based:
 - input feature: $|\mathbf{Y}| \in \mathbb{R}^{K \times L}$
 - sigmoid output to ensure $\widehat{SPP}(l) \in]0, 1[$
 - trained to minimize MSE loss between estimated and target SPP⁶

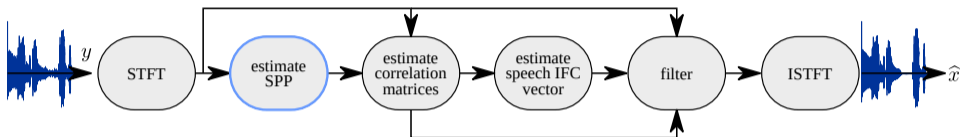


⁶Gerkmann and Hendriks 2012.

Method Overview

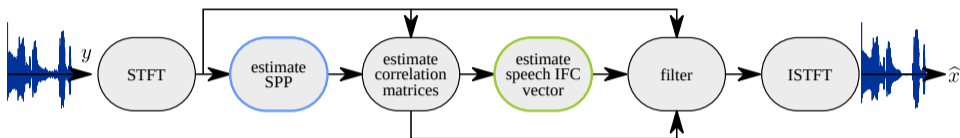


Method Overview



estimate SPP:
ML-based
vs. DNN-based

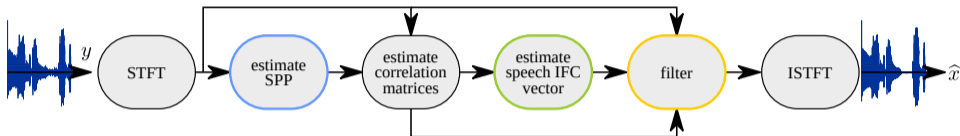
Method Overview



estimate SPP:
ML-based
vs. DNN-based

estimate speech IFC vector:
using ML constant noise IFC
vector (μ_n)
vs. using estimated noise
IFC vector ($\hat{\gamma}_n(l)$)

Method Overview



estimate SPP:
ML-based
vs. DNN-based

estimate speech IFC vector:
using ML constant noise IFC
vector (μ_n)
vs. using estimated noise
IFC vector ($\hat{\gamma}_n(l)$)

filter:
multi-frame (MPDR)
vs. single-frame Wiener gain
(WG)

- 1 Problem Statement
- 2 Multi-Frame Filtering
 - Signal Model & Assumptions
 - Problem & Solution
- 3 Parameter Estimation
- 4 Experiments**
 - Datasets & Settings
 - Results
- 5 Conclusions & Outlook

Datasets

	training	validation	testing
speech (WSJ0 ⁷)	101 speakers	20 speakers	4 speakers
noise	NOISEX92 ⁸		Aurora ⁹
SNRs / dB	[0, 20]		{-5, 0, ..., 20}

⁷Paul and Baker 1992.

⁸Varga and Steeneken 1993.

⁹Hirsch and Pearce 2000.

Datasets

	training	validation	testing
speech (WSJ0 ⁷)	101 speakers	20 speakers	4 speakers
noise	NOISEX92 ⁸		Aurora ⁹
SNRs / dB	[0, 20]		{-5, 0, ..., 20}

- disjoint training and testing datasets

⁷Paul and Baker 1992.

⁸Varga and Steeneken 1993.

⁹Hirsch and Pearce 2000.

Datasets

	training	validation	testing
speech (WSJ0 ⁷)	101 speakers	20 speakers	4 speakers
noise	NOISEX92 ⁸		Aurora ⁹
SNRs / dB	[0, 20]		{-5, 0, ..., 20}

- disjoint training and testing datasets
- one DNN is trained for all SNRs

⁷Paul and Baker 1992.

⁸Varga and Steeneken 1993.

⁹Hirsch and Pearce 2000.

Settings

- STFT and algorithm settings:
 - frame length 4 ms, shift length 1 ms, Hann windows
 - recursive smoothing constants $\alpha_n \hat{=} 50$ ms, $\lambda_y \hat{=} 12$ ms, $\lambda_{\text{DDA}} \hat{=} 33$ ms
 - Tikhonov regularization for matrix inversion
 - filter size $N = 18$ (correlations within 21 ms can be exploited)
 - minimum gain -17 dB
- DNN training
 - dropout, batch normalization, early stopping
 - Adam optimizer
- evaluation in terms of ΔPESQ

Results

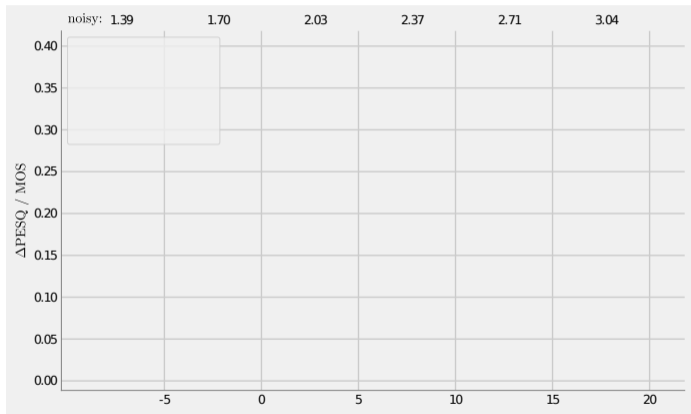


Figure: average PESQ improvement vs. input SNR / dB

Results

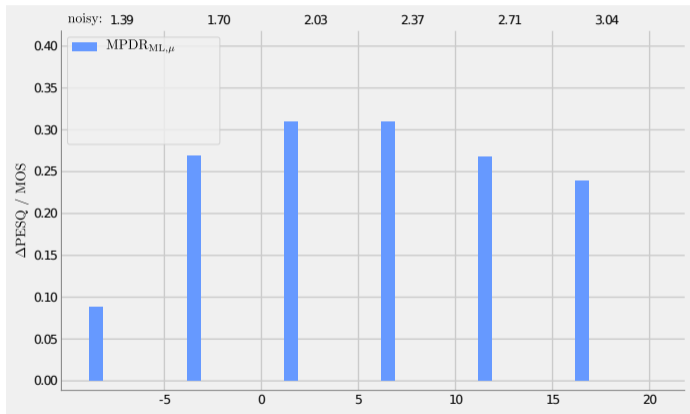


Figure: average PESQ improvement vs. input SNR / dB

Results

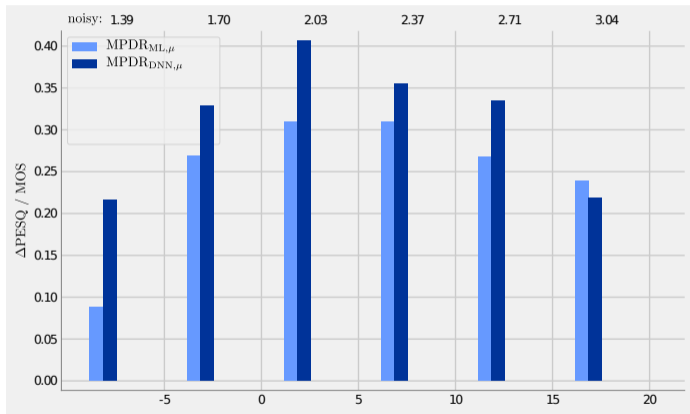


Figure: average PESQ improvement vs. input SNR / dB

Results

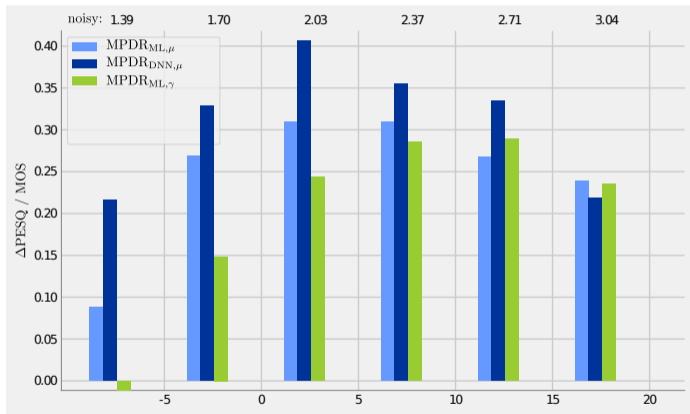


Figure: average PESQ improvement vs. input SNR / dB

Results

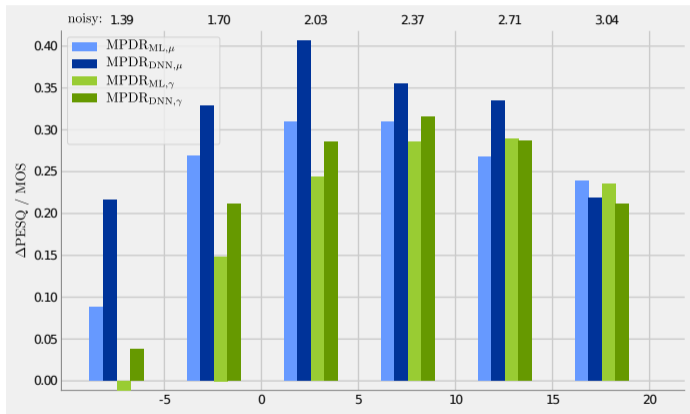


Figure: average PESQ improvement vs. input SNR / dB

Results

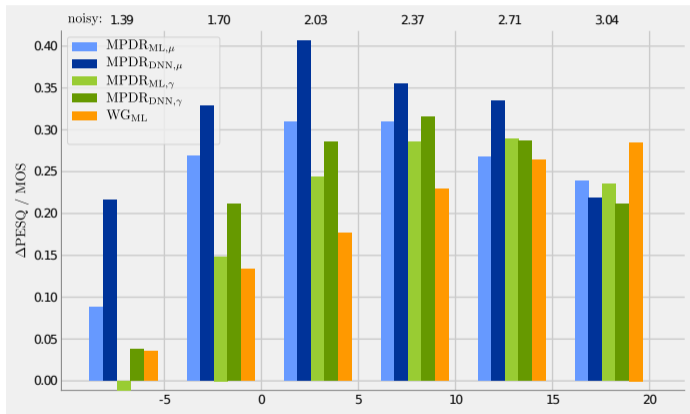


Figure: average PESQ improvement vs. input SNR / dB

Results

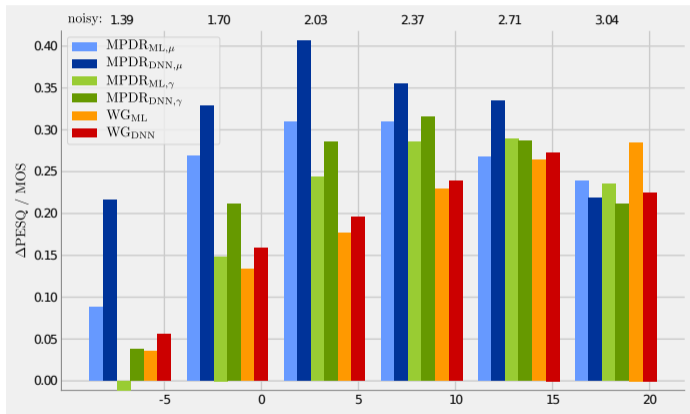


Figure: average PESQ improvement vs. input SNR / dB

Results

- $MPDR_{DNN,\mu}$ with highest performance

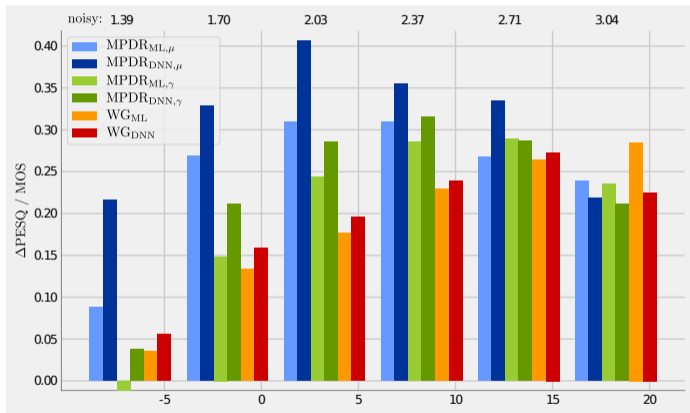


Figure: average PESQ improvement vs. input SNR / dB

Results

- $MPDR_{DNN,\mu}$ with highest performance
- SPP estimation: $DNN > ML$

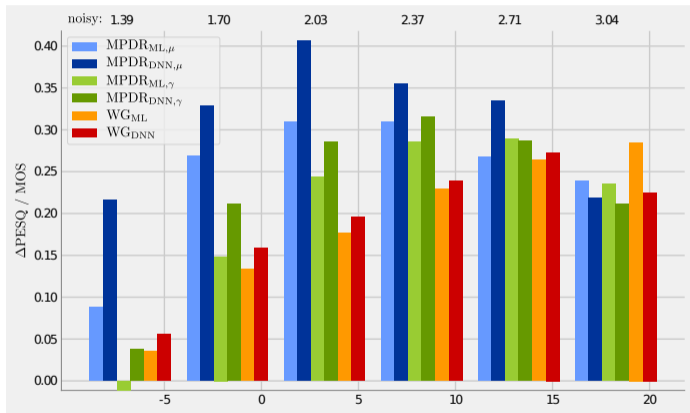


Figure: average PESQ improvement vs. input SNR / dB

Results

- $MPDR_{DNN,\mu}$ with highest performance
- SPP estimation: $DNN > ML$
- speech IFC estimation: using $\hat{\gamma}_{x,\mu}(l) > \hat{\gamma}_{x,\gamma}(l)$

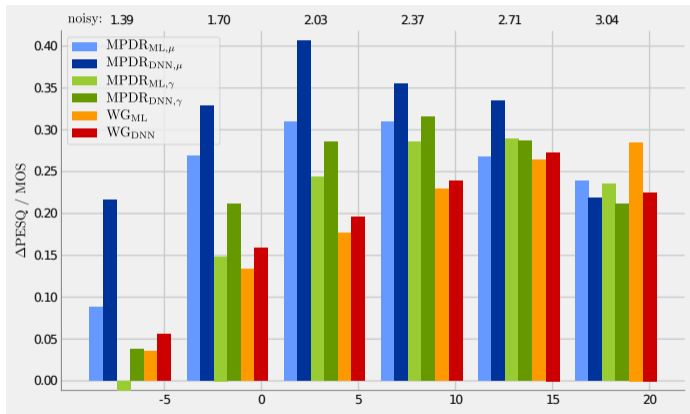


Figure: average PESQ improvement vs. input SNR / dB

Results

- $MPDR_{DNN,\mu}$ with highest performance
- SPP estimation:
DNN > ML
- speech IFC estimation:
using $\hat{\gamma}_{x,\mu}(l) > \hat{\gamma}_{x,\gamma}(l)$
- filtering:
multi-frame > single-frame

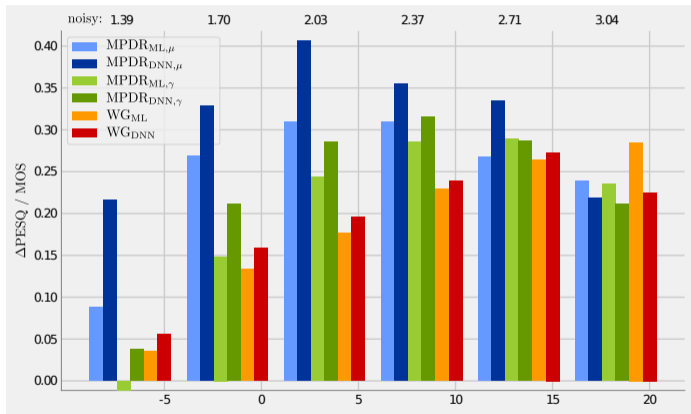


Figure: average PESQ improvement vs. input SNR / dB

- 1 Problem Statement
- 2 Multi-Frame Filtering
 - Signal Model & Assumptions
 - Problem & Solution
- 3 Parameter Estimation
- 4 Experiments
 - Datasets & Settings
 - Results
- 5 Conclusions & Outlook

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach
- **improved SE performance** by replacing conventional ML-based SPP estimator **with DNN-based SPP estimator**

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach
 - **improved SE performance** by replacing conventional ML-based SPP estimator **with DNN-based SPP estimator**
 - confirmed that **multi-frame filtering** can be **advantageous to single-frame filtering**
-

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach
 - **improved SE performance** by replacing conventional ML-based SPP estimator **with DNN-based SPP estimator**
 - confirmed that **multi-frame filtering** can be **advantageous to single-frame filtering**
-
- tight integration of deep learning with filtering

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach
 - **improved SE performance** by replacing conventional ML-based SPP estimator **with DNN-based SPP estimator**
 - confirmed that **multi-frame filtering** can be **advantageous to single-frame filtering**
-
- tight integration of deep learning with filtering
 - include filter in training process

Conclusions & Outlook

- considered **multi-frame single-microphone speech enhancement (SE)** approach
 - **improved SE performance** by replacing conventional ML-based SPP estimator **with DNN-based SPP estimator**
 - confirmed that **multi-frame filtering** can be **advantageous to single-frame filtering**
-
- tight integration of deep learning with filtering
 - include filter in training process
 - more strongly rely on DNN for parameter estimation

Thank you for your attention!

Any questions?

⇒ marvin.tammen@uol.de

Thank you for your attention!

Any questions?

⇒ marvin.tammen@uol.de

Maximum Likelihood-Based SPP Computation

- assumptions¹⁰:
 - complex Gaussian distributions for microphone, speech and noise STFT coefficients
 - equal a-priori probability of speech and noise
 - derive likelihood under speech presence (\mathcal{H}_1) and absence (\mathcal{H}_0)
 - typical SNR under speech presence: $\xi_{\mathcal{H}_1} = 15$ dB
- resulting a-posteriori SPP estimate:

$$\widehat{\text{SPP}}(l) = \left(1 + \frac{P(\mathcal{H}_0)}{P(\mathcal{H}_1)} (1 + \xi_{\mathcal{H}_1}) e^{-\frac{|Y(l)|^2}{\hat{\phi}_N(l-1)} \frac{\xi_{\mathcal{H}_1}}{1 + \xi_{\mathcal{H}_1}}} \right)^{-1} \quad (6)$$

¹⁰Gerkmann and Hendriks 2012.

Training and Testing

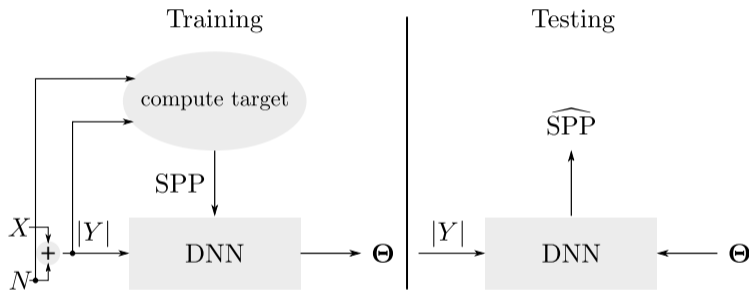


Figure: training and testing phase; Θ : DNN parameters

Wiener Gain Implementation

Wiener Gain

Input: noisy STFT coefficients $Y(k, l)$, noise PSD estimate $\hat{\phi}_n(k, l)$, smoothing constant α , minimum gain G_{\min}

Output: clean speech estimate $\hat{X}(k, l)$

foreach $k \in [0, K - 1]$ **do**

init. $\hat{X}(k, 0) = 0;$

foreach $l \in [1, L - 1]$ **do**

$$\hat{\xi}(k, l) = \alpha \frac{|\hat{X}(k, l-1)|^2}{\hat{\phi}_n(k, l-1)} + (1 - \alpha) \frac{|Y(k, l)|^2}{\hat{\phi}_n(k, l)}; \quad /* \text{ a-priori SNR estimation } */$$

$$G(k, l) = \max \left\{ \frac{\hat{\xi}(k, l)}{1 + \hat{\xi}(k, l)}, G_{\min} \right\}; \quad /* \text{ spectral flooring } */$$

$$\hat{X}(k, l) = G(k, l) Y(k, l);$$

end

end

Tikhonov Regularization

- useful for dealing with multicollinearity in linear regression problems
- obtain Tikhonov-regularized matrix $\tilde{\mathbf{A}}$ from $\mathbf{A} \in \mathbb{C}^{N \times N}$ as

$$\tilde{\mathbf{A}} = \mathbf{A} + \frac{\delta}{N} \text{trace}\{|\mathbf{A}|\} \mathbf{I}_N \quad (7)$$