

Introduction

- Filters for single-microphone speech enhancement can be applied to suppress background noise while preventing speech distortion [1, 2, 3, 4]
- Single- & multi-frame trade-off filters (also known as speech-distortion weighted inter-frame Wiener filters) can increase **noise reduction** while limiting **speech distortion**
- Single-frame filters **disregard correlation** across time frames and/or frequency bins
- Multi-frame filters can **exploit temporal correlation** across present and past frames

IN THIS POSTER

Real- and complex-valued **trade-off** filters are analysed using **low-delay** STFT filterbanks for **speech enhancement** with a single-microphone.

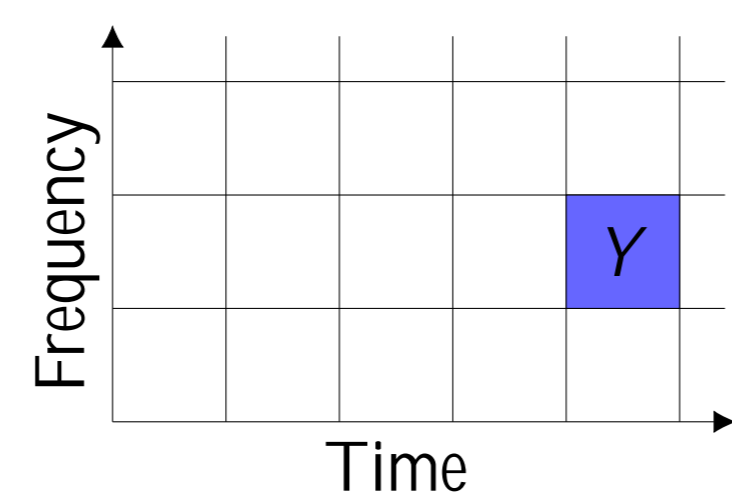
Signal Models

Single-Frame Signal Model

- Assumption: no correlation over time frames and/or frequency bins
- Signal model:

$$Y = S + N$$

noisy speech noise



- Estimate speech S by applying gain G to each noisy speech STFT coefficient Y independently:

$$\hat{S} = GY$$

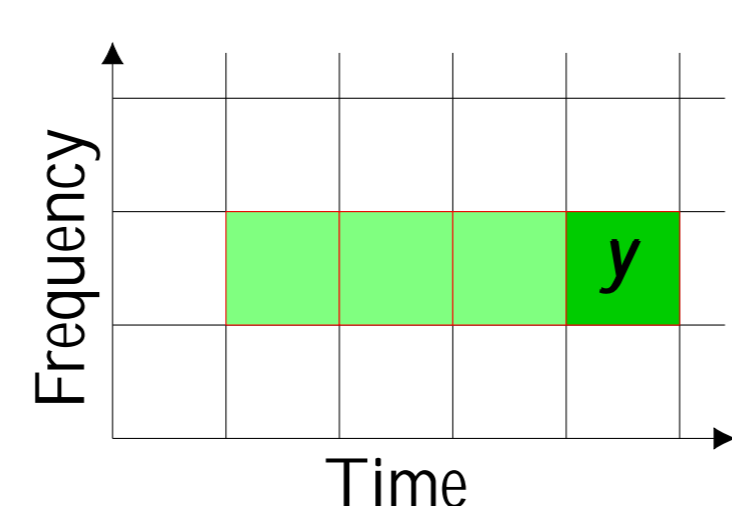
Multi-Frame Signal Model

- Assumption: speech correlated over consecutive time frames
- Speech correlation vector:

$$s = \frac{E[sS]}{E[S^2]} = \frac{R_s e}{s}$$

- Signal model of noisy speech y in terms of speech s and noise n or correlated speech sS and undesired u :

$$y = s + n = sS + u$$



- Estimate speech S by applying filter h to noisy speech vector:

$$\hat{S} = h^H y$$

Trade-off Filters

Less speech distortion \leftarrow Trade-off μ \rightarrow More noise reduction

Single-Frame Trade-off Filter:

$$G = \underset{G}{\operatorname{argmin}} \frac{E[|GS - S|^2]}{\text{Filtered speech}} + \mu \frac{E[|GN|^2]}{\text{Filtered noise}}$$

Real-Valued Trade-off Filter (R-TF)

- Implemented in STFT filterbank with **high frequency resolution** & **high time resolution**
- Parameters are estimated using [5], [6]

$$G = \frac{s}{\mu N + s}$$

- For $\mu = 0$ applies **no noise reduction**
- For $\mu = 1$ is a **Wiener gain**

Multi-Frame Trade-off Filter:

$$h^{\text{MFTF}} = \underset{h}{\operatorname{argmin}} \frac{E[|h^H sS - S|^2]}{\text{Filtered correlated speech}} + \mu \frac{E[|h^H u|^2]}{\text{Filtered undesired}}$$

Complex-Valued Multi-Frame Trade-off Filter (C-MFTF)

- Implemented in STFT filterbank with **low frequency resolution** & **high time resolution**
- Parameters are estimated as in [1]

$$h^{\text{C-MFTF}} = \frac{R_y^{-1} s}{s^H R_y^{-1} s} \frac{s}{\mu \frac{\text{out}}{s} + s}$$

- For $\mu = 0$ is a **C-MFMPDR filter**
- For $\mu = 1$ is a **C-MFWF**

Real-Valued Multi-Frame Trade-off Filter (R-MFTF)

- Implemented in STFT filterbank with **high frequency resolution** & **high time resolution**
- Parameters are estimated as in [2]
- Solution is obtained by applying DFT to cost function: $g = Dh$ (D : DFT matrix)

$$g^{\text{R-MFTF}} = \frac{y^{-1} s \mathbf{1}}{\mu + (1 - \mu) \mathbf{1}^T s y^{-1} s \mathbf{1} A}, \quad A = \frac{1}{\mathbf{1}^T s \mathbf{1}}$$

- Real-valued filter vector $g^{\text{R-MFTF}}$ can be overlapped into high-resolution scalar gain: $G^{\text{R-MFTF}}$
- For $\mu = 0$ is a **R-MFMPDR filter**
- For $\mu = 1$ is a **R-MFWF**

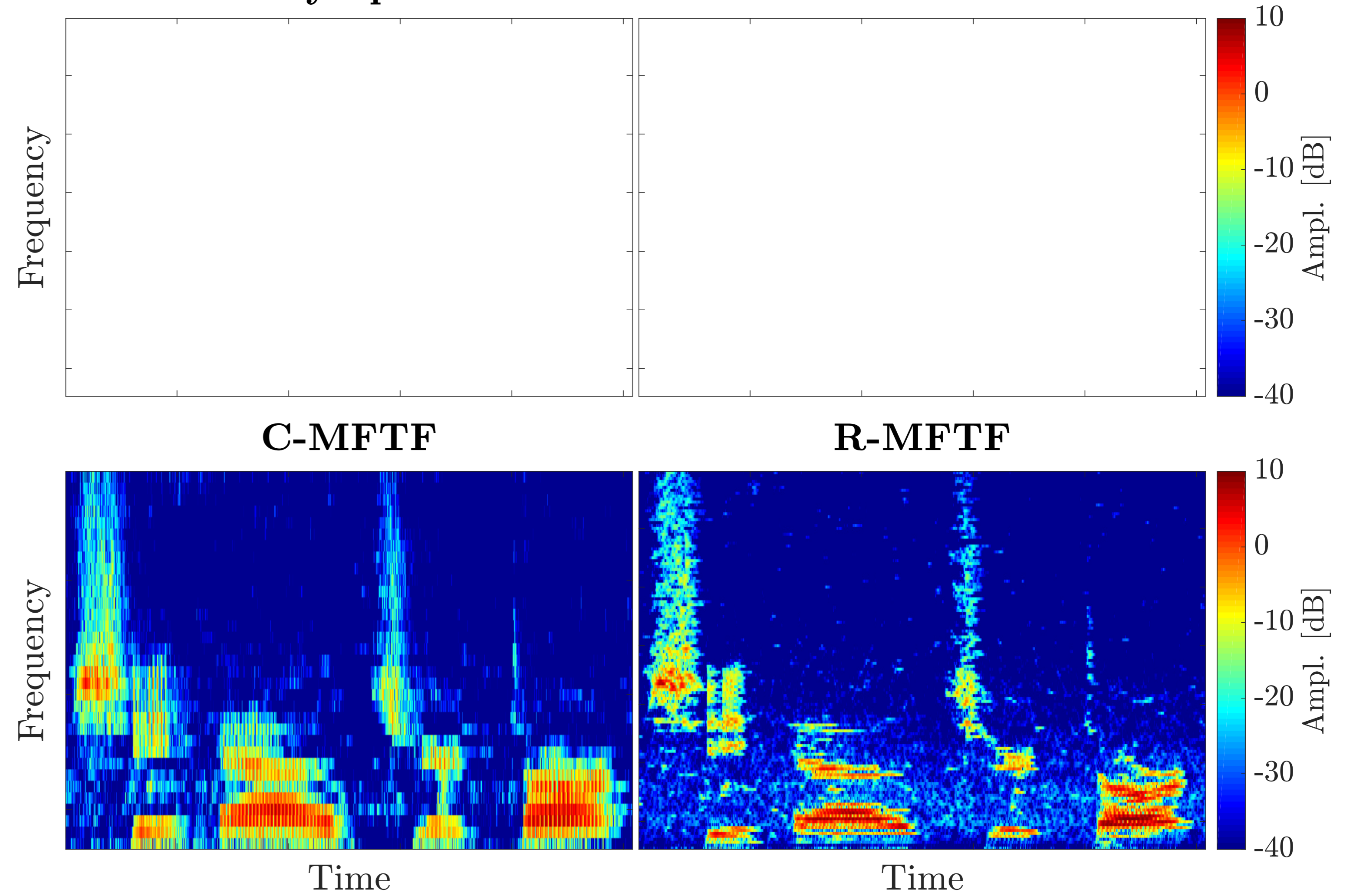
Framework

Filterbank	Analysis window length	Synthesis window length
Low resolution STFT	64 (4ms)	64 (4ms)
High resolution STFT	256 (16ms)	64 (4ms)

- Sampling frequency: 16 kHz
- Both filterbanks use a frame-shift of 16 (1 ms) and have a synthesis delay of 3 ms

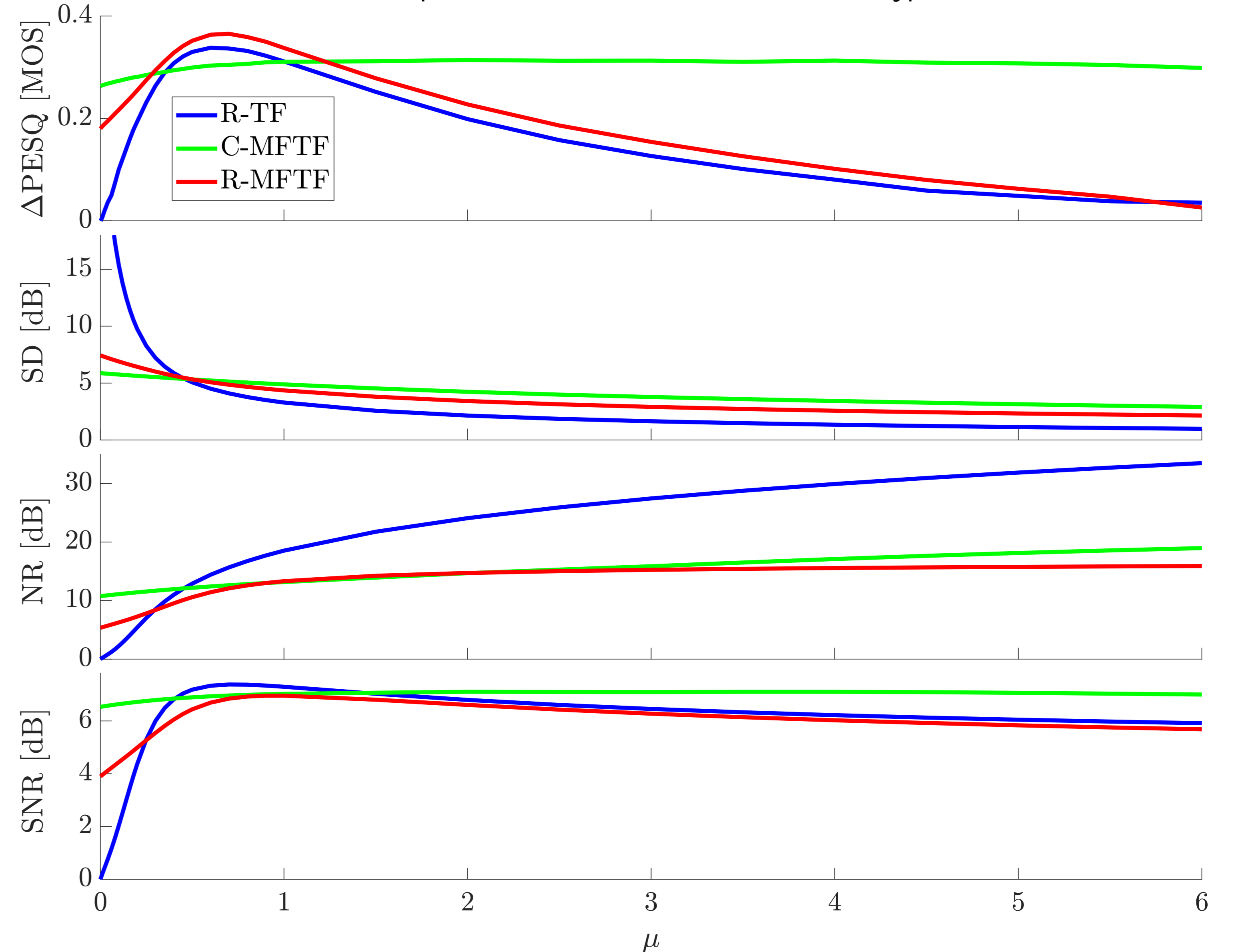
Evaluation

Speech with crossroad noise at SNR: 0 dB. Trade-off parameter $\mu = 0.8$



Results

Objective measures: Perceptual speech quality (PESQ), segmental speech distortion (SD), segmental noise reduction (NR), and segmental signal to noise ratio (SNR). Evaluated over 105 s speech material for 5 different noise types at SNR: 0 dB



- At $\mu = 0$, the C-MFTF (i.e. C-MFMPDR filter) is the most effective filter
- R-MFTF applies best overall PESQ improvement at $\mu = 0.7$
- R-TF can apply highest noise reduction but with high speech distortion

Conclusions

- Complex- and real-valued trade-off filters are effective for speech enhancement
- Increasing trade-off parameter μ increases both noise reduction and speech distortion
- C-MFMPDR filter performs better than R-MFMPDR filter (MFTFs for $\mu = 0$)

References

- D. Fischer, K. Brümnn, and S. Doclo, "Comparison of parameter estimation methods for Single-Microphone Multi-Frame Wiener filtering," in 27th European Signal Processing Conference (EUSIPCO), (Submitted).
- K. T. Andersen and M. Moonen, "Robust speech-distortion weighted interframe Wiener filters for single-channel noise reduction," IEEE Trans. Audio, Speech, Language Process., vol. 26, no. 1, pp. 97–107, Jan. 2018.
- Y. Huang and J. Benesty, "A multi-frame approach to the frequency-domain single-channel noise reduction problem," IEEE Trans. Audio, Speech, Language Process., vol. 20, no. 4, pp. 1256–1269, May 2012.
- A. Schasse and R. Martin, "Estimation of subband speech correlations for noise reduction via MVDR processing," IEEE Trans. Audio, Speech, Language Process., vol. 22, no. 9, pp. 1355–1365, Sep. 2014.
- Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error log-spectral amplitude estimator," IEEE Trans. Acoust., Speech, Signal Process., vol. 33, no. 2, pp. 443–445, Apr. 1985.
- T. Gerkmann and R. C. Hendriks, "Unbiased MMSE-based noise power estimation with low complexity and low tracking delay," IEEE Trans. Audio, Speech, Language Process., vol. 20, no. 4, pp. 1383–1393, May 2012.