

Main contribution

- ▶ We propose a novel speech enhancement algorithm for removing reverberation and noise from recorded speech data.
- ▶ Compared to conventional methods, our approach results in:
 - ▷ Substantial improvement in PESQ and other objective metrics.
 - ▷ Fewer artifacts in informal listening.
- ▶ Our method effectively increases the analysis window duration that can be used for voiced speech.
 - ▷ We extend the *coherence time*, which is the duration over which an analysis method is coherent with the signal.
 - ▷ Conventional methods assume a speech coherence time of 10-30 ms; we extend this time to 128 ms.

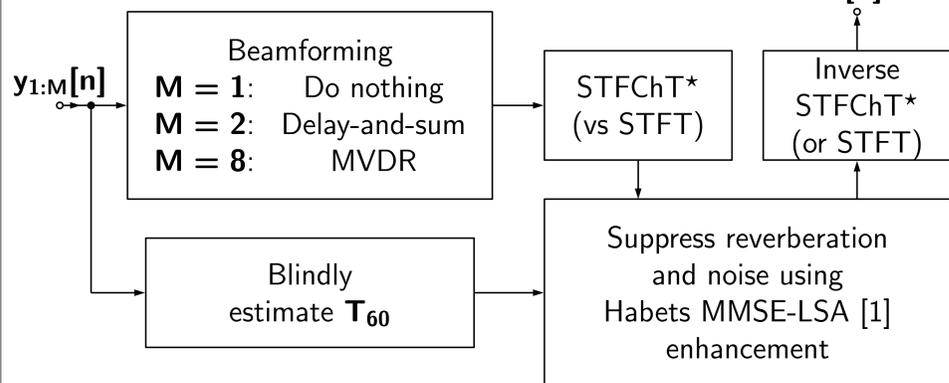
System block diagram

Input is M channels of reverberant, noisy speech:

$$\mathbf{y}_{1:M}[n] = \mathbf{h}_m[n] * \mathbf{s}[n] + \mathbf{v}[n].$$

Output is estimate of clean speech:

$$\hat{\mathbf{s}}[n]$$



* "Short-time fan-chirp transform", see panels B and C.

A. Blindly estimate \mathbf{T}_{60}

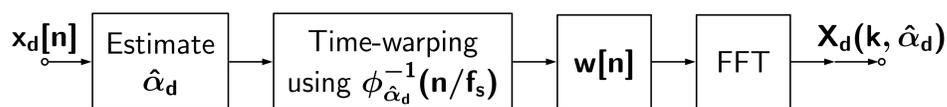
- ▶ Suppress additive noise in each channel using Ephraim and Malah MMSE-LSA [2] and concatenate enhanced channels.
- ▶ Use maximum-likelihood blind \mathbf{T}_{60} estimator by Löllmann et al. [3].
- ▶ \mathbf{T}_{60} estimation improves with more data (i.e., more channels).

B. Short-time fan-chirp transform (STFChT) [4]

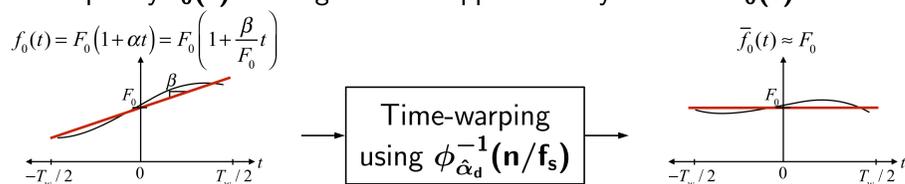
The STFChT is used to implement our method of extending the coherence time of analysis. Definition of STFChT:

$$\mathbf{X}_d(\mathbf{f}, \hat{\alpha}_d) = \int_{-T_w/2}^{T_w/2} \mathbf{w}(\tau) \mathbf{x}_d(\phi_{\hat{\alpha}_d}^{-1}(\tau)) e^{-j2\pi\mathbf{f}\tau} d\tau \quad (1)$$

- ▶ $\mathbf{w}(\mathbf{t})$ is an analysis window of duration \mathbf{T}_w .
- ▶ $\mathbf{x}_d(\mathbf{t}) = \mathbf{x}(\mathbf{t} - d\mathbf{T}_{hop})$, $\mathbf{t} \in [0, \mathbf{T}_w]$, is a short frame of a time-domain signal.
- ▶ $\phi_{\hat{\alpha}_d}(\mathbf{t})$ is a linear phase trajectory
- ▶ $\phi_{\hat{\alpha}_d}^{-1}(\mathbf{t})$ is a time-warping function.
- ▶ $\hat{\alpha}_d$ is an estimated chirp rate α for the d th frame.

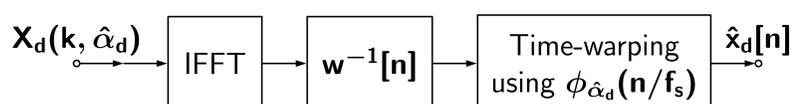


- ▶ Time-warping converts signals with linearly time-varying fundamental frequency $\mathbf{f}_0(\mathbf{t})$ into signals with approximately constant $\mathbf{f}_0(\mathbf{t})$.

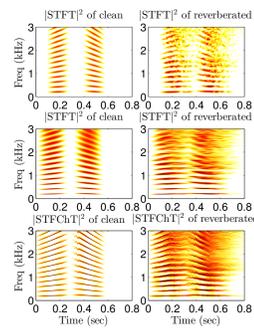


C. Inverse STFChT

- ▶ Time-warping implemented as combination of oversampling and interpolation, which achieves almost perfect STFChT reconstruction.

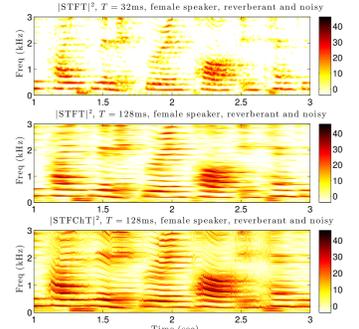


D. STFT versus STFChT (32 ms and 128 ms window durations)



Synthetic harmonics.

A direct demonstration of the advantage of extending coherence: examples of STFT with 32 ms and 128 ms windows versus STFChT with 128 ms window.

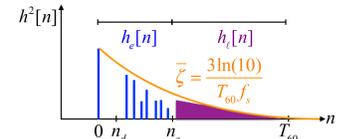


Reverberated speech utterance.

E. Suppress reverberation and noise using Habet's MMSE-LSA [1]

- ▶ Employs a statistical model of reverberation:

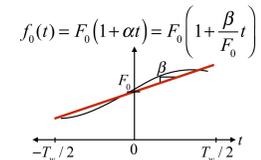
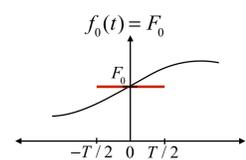
$$\mathbf{E} [h^2[n]] = \begin{cases} \sigma_d^2 e^{-2\zeta n}, & 0 \leq n < n_d \\ \sigma_r^2 e^{-2\zeta n}, & n \geq n_d \\ 0 & \text{otherwise.} \end{cases}$$



- ▶ Estimate complex time-frequency coefficients $\hat{\mathbf{X}}_e(\mathbf{d}, \mathbf{k})$ of early reverberant component using Habet's MMSE-LSA gains [1]:

$$\hat{\mathbf{X}}_e(\mathbf{d}, \mathbf{k}) = \mathbf{G}_{\text{MMSE-LSA}}(\mathbf{d}, \mathbf{k}) \cdot \mathbf{Y}(\mathbf{d}, \mathbf{k}) \quad (2)$$

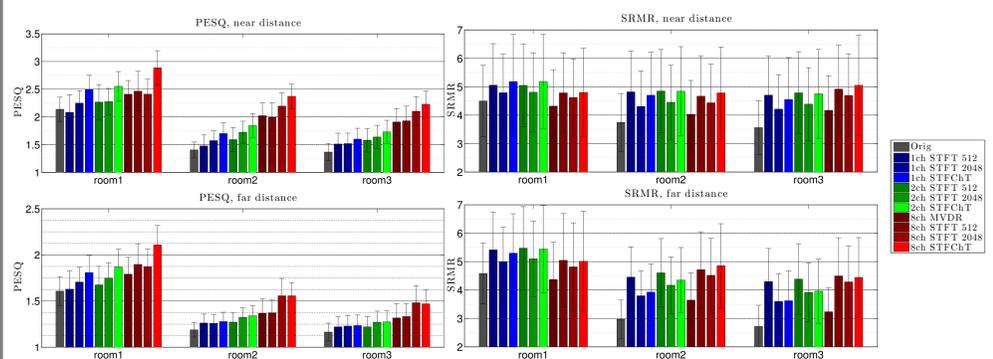
- ▶ Assumes stationary signal with constant $\mathbf{f}_0(\mathbf{t})$ over analysis duration.



- ▶ This assumption limits STFT window length, which limits data record length for statistical estimators.

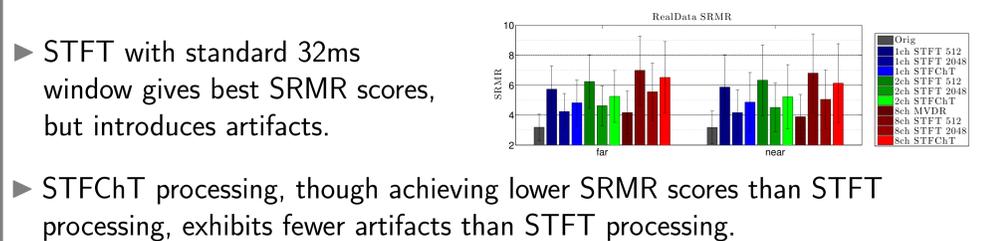
- ▶ The STFChT is coherent with speech over longer durations, which allows longer data records and thus provides higher SNR.

Results for SimData test set



STFChT processing achieves substantially better PESQ scores while maintaining roughly equivalent SRMR scores. STFT 512 uses a 32 ms window; STFT 2048 and STFChT use 128 ms windows.

Results for RealData test set



- ▶ STFT with standard 32ms window gives best SRMR scores, but introduces artifacts.

- ▶ STFChT processing, though achieving lower SRMR scores than STFT processing, exhibits fewer artifacts than STFT processing.

References

- [1] E. A. P. Habet, "Speech dereverberation using statistical reverberation models," in *Speech Dereverberation*, Patrick A. Naylor and Nikolay D. Gaubitch, Eds. Springer, July 2010.
- [2] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error log-spectral amplitude estimator," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 33, no. 2, pp. 443–445, 1985.
- [3] H. Löllmann, E. Yilmaz, M. Jeub, and P. Vary, "An improved algorithm for blind reverberation time estimation," in *Proc. IWAENC*, Tel Aviv, Israel, 2010, p. 1–4.
- [4] M. Képesi and L. Weruaga, "Adaptive chirp-based time-frequency analysis of speech signals," *Speech Communication*, vol. 48, no. 5, pp. 474–492, May 2006.