

Speech Enhancement using End-to-End Speech Recognition Objectives

Aswin Shanmugam Subramanian¹, Xiaofei Wang¹, Murali Karthick Baskar^{1,2}, Shinji Watanabe¹, Toru Taniguchi³, Dung Tran³, Yuya Fujita³



¹Center for Language and Speech Processing, Johns Hopkins University, Baltimore, MD, USA

²Brno University of Technology, Brno, Czech Republic

³Yahoo Japan Corporation, Tokyo, Japan



Abstract

- Denoising and Dereverberation systems usually optimized with **signal reconstruction objective**.
 - Issue 1 - Can be trained only on simulated data
 - Issue 2 - Not application oriented
- Alternative - Optimize with automatic **speech recognition (ASR) objective**.
- Contributions of the paper:
 - Check how joint optimization of far-field denoising and dereverberation with ASR objective as a single network performs in terms of enhancement objectives
 - See which enhancement metric correlates well with ASR metric
 - Learn to predict important hyper-parameters using the data

Dereverberation Subnetwork

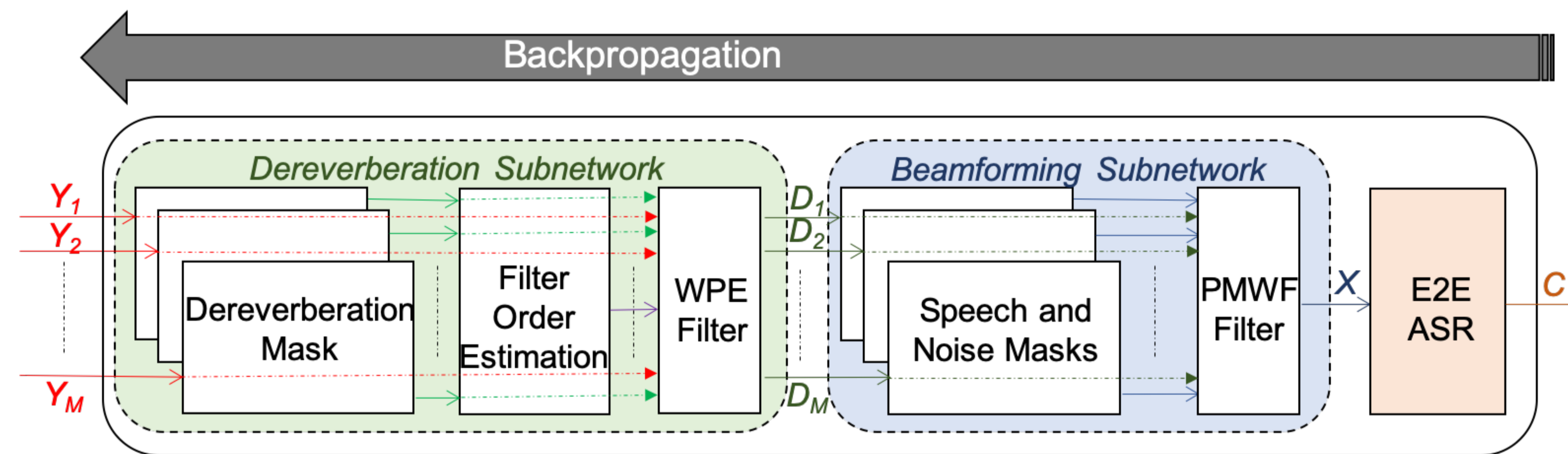
- Based on weighted prediction error (WPE) method.
- $\mathcal{Y} = (\{\mathbf{y}(t, b)\}_{b=1}^B \in \mathbb{C}^M | t = 1, \dots, T)$: sequence of T -length M -channel noisy speech spectrum
- dereverberated signal:

$$\mathbf{d}(t, b) = \mathbf{y}(t, b) - (\mathbf{R}(b)^{-1} \mathbf{P}(b))^H \tilde{\mathbf{y}}(t - \Delta, b), \Delta$$
: prediction delay, $\tilde{\mathbf{y}}(t - \Delta, b) \in \mathbb{C}^{ML}$: delayed and stacked input with filter order L .
- $\mathbf{R}(b) = \sum_t \frac{\tilde{\mathbf{y}}(t-\Delta, b) \tilde{\mathbf{y}}^H(t-\Delta, b)}{\sum_m |\tilde{d}(t, b, m; \theta_{dry})|^2 / M} \in \mathbb{C}^{ML \times ML}$
- $\mathbf{P}(b) = \sum_t \frac{\tilde{\mathbf{y}}(t-\Delta, b) \mathbf{y}^H(t, b)}{\sum_m |\tilde{d}(t, b, m; \theta_{dry})|^2 / M} \in \mathbb{C}^{ML \times M}$
- **Neural network** with learnable parameter θ_{dry} used to **predict** $\tilde{d}(t, b, m; \theta_{dry})$
- $\mathcal{D} = \text{Dry}(\mathcal{Y}; \theta_{dry})$

Beamforming Subnetwork

- Beamformed signal: $x(t, b) = \mathbf{f}^H(b) \mathbf{d}(t, b)$
- **Parametrized multi-channel Wiener filter**:

$$\mathbf{f}(b) = \frac{\Phi_N(b)^{-1} \Phi_S(b)}{\beta(b) + \text{Trace}(\Phi_N(b)^{-1} \Phi_S(b))} \mathbf{u} \in \mathbb{C}^M$$



Beamforming Subnetwork Cnt'd

- $\Phi_v(b) = \sum_{t=1}^T w_v(t, b; \theta_{fcs}) \mathbf{d}(t, b) \mathbf{d}^H(t, b)$ where $v \in \{S, N\}$
- $w_S(t, b; \theta_{fcs}) \in [0, 1]$ and $w_N(t, b; \theta_{fcs}) \in [0, 1]$ obtained from a neural network with learnable parameter θ_{fcs}
- $\beta(b) \in \mathbb{R}_{\geq 0}$ is the trade-off factor between speech distortion and noise reduction: **predicted from the network** as a function of $\Phi_v(b)$
- $\mathbf{X} = \text{Fcs}(\mathcal{D}; \theta_{fcs})$

Filter Order Estimation

- Attention based ASR: $p(C|\mathcal{Y}) = \text{Trn}(X; \theta_{trn})$, C is the character sequence.
- $p(L|\mathcal{Y})$ is obtained as a **softmax function of a network output** with learnable parameter θ_{flt}
- **Reinforcement learning** to update θ_{flt} :

$$\nabla_{\theta_{flt}} \mathcal{L}_{flt} = \sum_L \text{CE}(C_{ref}, p_L(C|\mathcal{Y})) \nabla_{\theta_{flt}} \log p(L|\mathcal{Y})$$

Experimental Setup

- Training Data: 2ch simulation data from REVERB
- Evaluation Data: REVERB 8ch & DIRHA Living Room Array - 6ch real data
- E2E ASR system - ESPnet REVERB recipe

Results

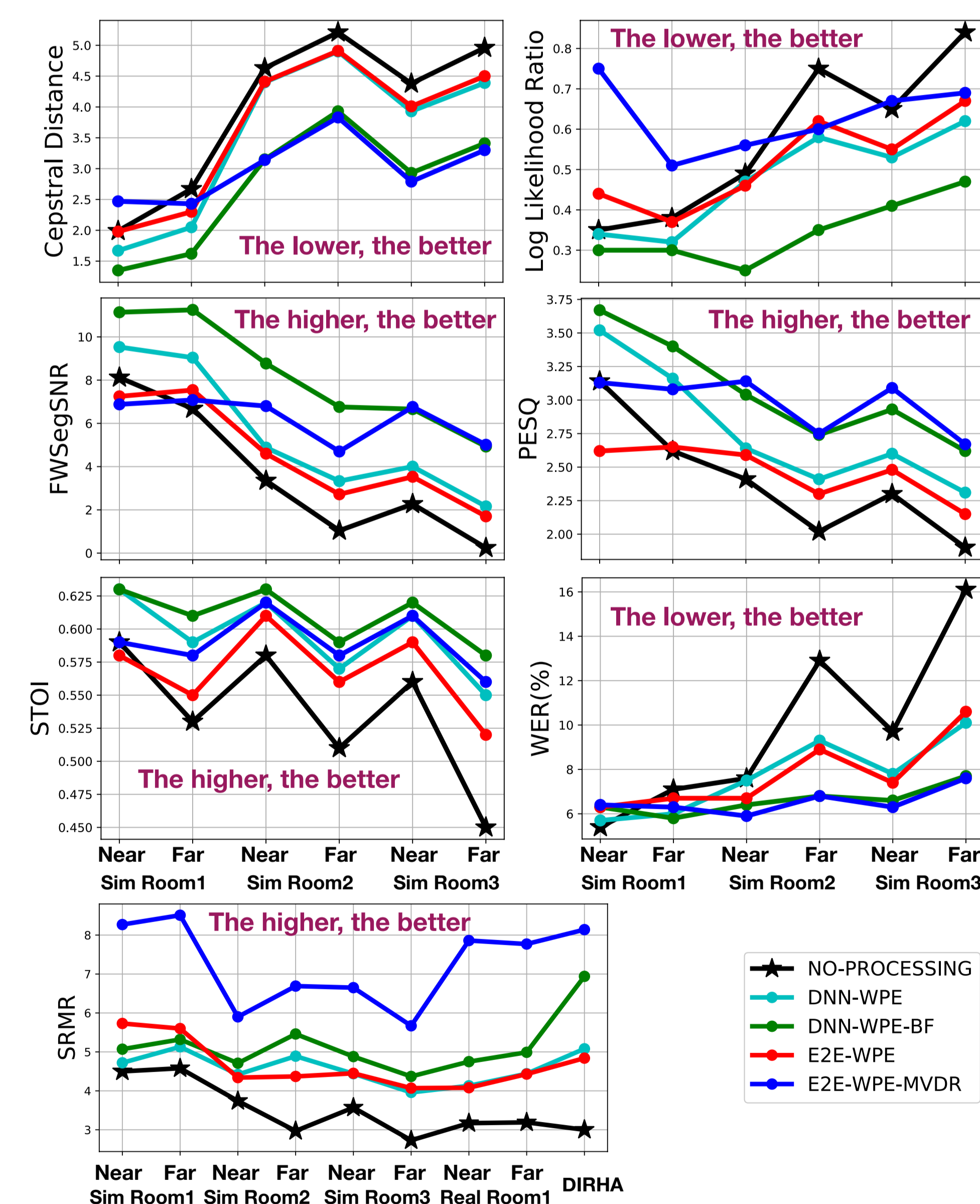


Figure: Objective measures of various methods. (1) cepstral distance (CD), (2) log-likelihood ratio (LLR), (3) frequency-weighted segmental SNR (FWSegSNR), (4) PESQ, (5) STOI, (6) WER and (7) SRMR.

Correlation coefficients	SRMR	CD	LLR	FWSegSNR	PESQ	STOI
WAR (= 100 - WER)	0.48	-0.57	-0.57	0.71	0.78	0.77
MUSHRA: PAR	0.59	-0.76	-0.42	0.74	0.84	-
MUSHRA: OQ	0.06	-0.38	-0.39	0.49	0.67	-

ASR Performance

Frontend Type	Dereverberation		Beamformer Method	REVERB Real Room 1		DIRHA LA Array
	Filter Order Estimation	Method		Near	Far	
-	-	-	-	23.9	26.8	55.3
Pipeline	N	DNN-WPE	-	16.4	18.5	41.3
	N	DNN-WPE	BeamformIt	11.0	10.8	31.3
E2E	N	WPE	-	18.0	19.8	42.3
	Y	WPE	-	15.1	16.9	36.9
	N	WPE	MVDR	8.7	12.4	29.1
	N	WPE	PMWF	9.7	11.8	27.9

Discussion

- **ASR objective** method considerably **degrades LLR** but significantly **improves SRMR**.
- **PESQ and STOI** - well **correlated with WER**
- E2E approach - significantly improves ASR performance on challenging DIRHA data.
- Filter order prediction seems to be based on the room size and not the microphone position
- Higher filter order chosen for real data

Mode	REVERB Simulated			REVERB Real Room 1		DIRHA LA Array
	Near	Far	Room 3	Near	Far	
Order L	9	9	4	4	4	9
Percentage	87.1	82.6	44.4	50.7	93.1	92.5
				71.0	70.4	70.4

Conclusion

- Speech enhancement using **ASR objective** - **effective on most enhancement metrics**.
- Proposed distortion weight estimation performs well on the DIRHA ASR task.
- Future work: (1) apply on more realistic and challenging environments like CHiME5, (2) Subjective evaluation.

Acknowledgement

Our dereverberation subnetwork implementation is based on DNN WPE module from NTT-CS Labs https://github.com/nttcs-lab-sp/dnn_wpe