



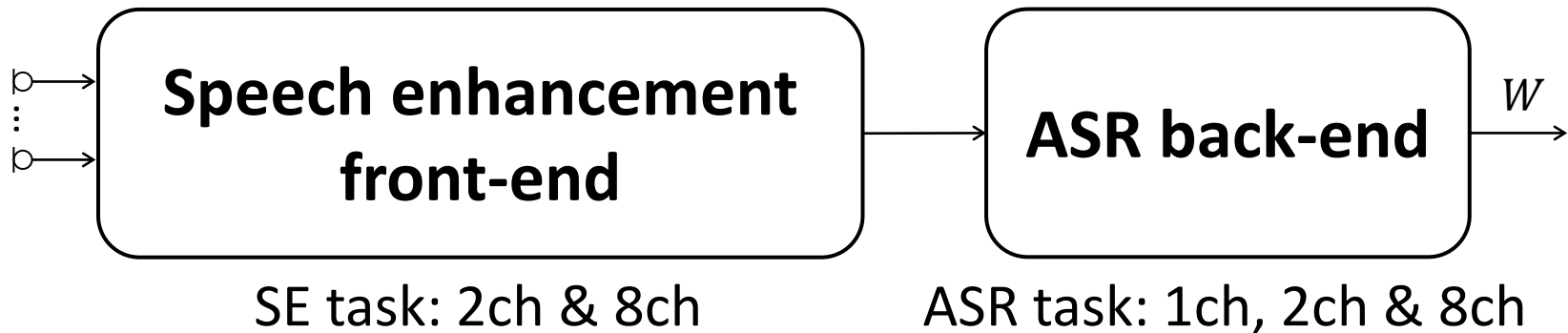
Innovative R&D by NTT

# Linear prediction-based dereverberation with advanced speech enhancement and recognition technologies for the REVERB challenge

M. Delcroix, T. Yoshioka, A. Ogawa, Y. Kubo, M. Fujimoto, N. Ito, K. Kinoshita, M. Espi, T. Hori, T. Nakatani, A. Nakamura

NTT Communication Science Laboratories

- Combine state-of-the-art SE and ASR techniques



- Significant performance improvement:

SE 8ch (FWSegSNR) : 3.62 dB → 10.31 dB

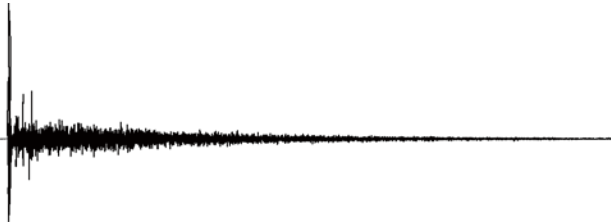
ASR 8ch (WER) : 4.2 % (SimData)

**9.0 %** (RealData)

# The challenges of the challenge...

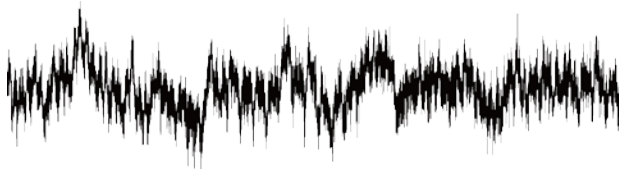


## 1. Severe reverberation



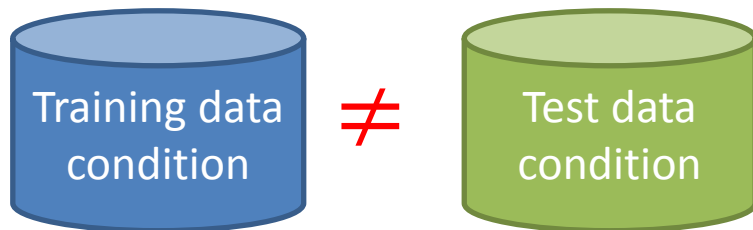
→ Powerful dereverberation

## 2. Significant amount of noise



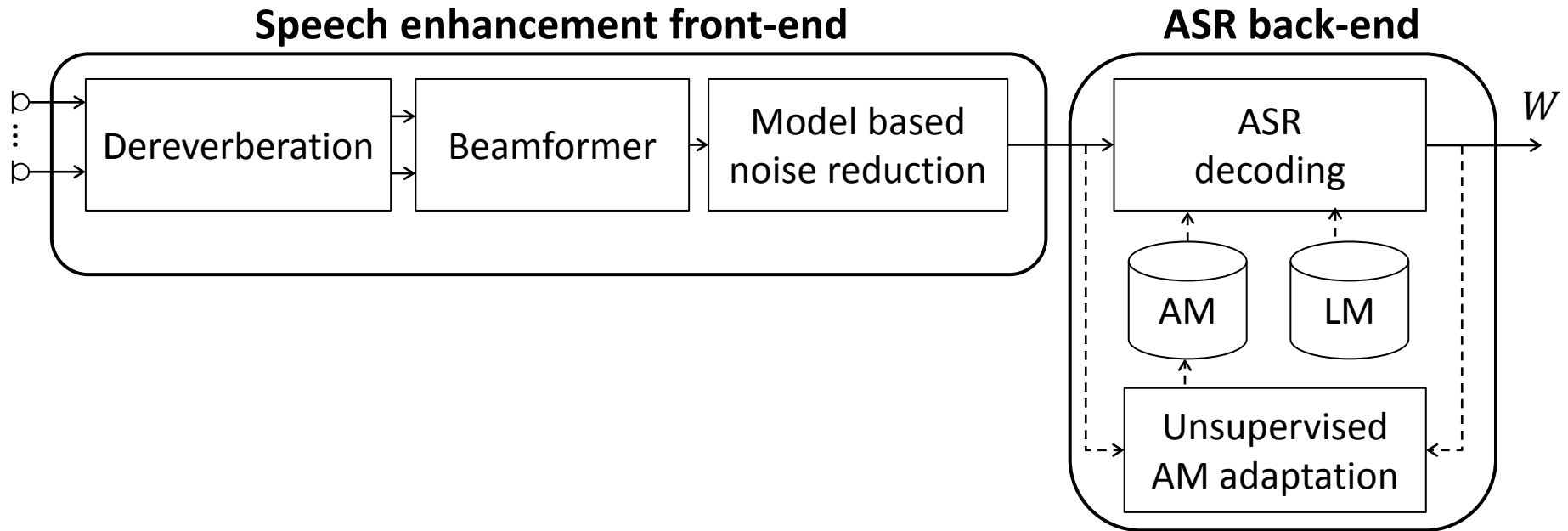
→ Noise robust dereverberation  
+ additional noise reduction

## 3. Mismatch between training and testing



→ High performance w/o overfitting  
*Extended training data*  
*+ unsupervised environmental adaptation*

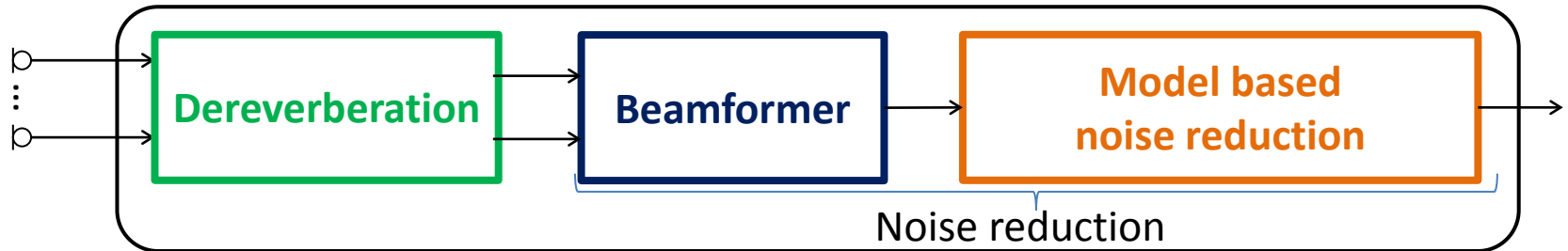
# Proposed system



# Characteristics of the SE front-end

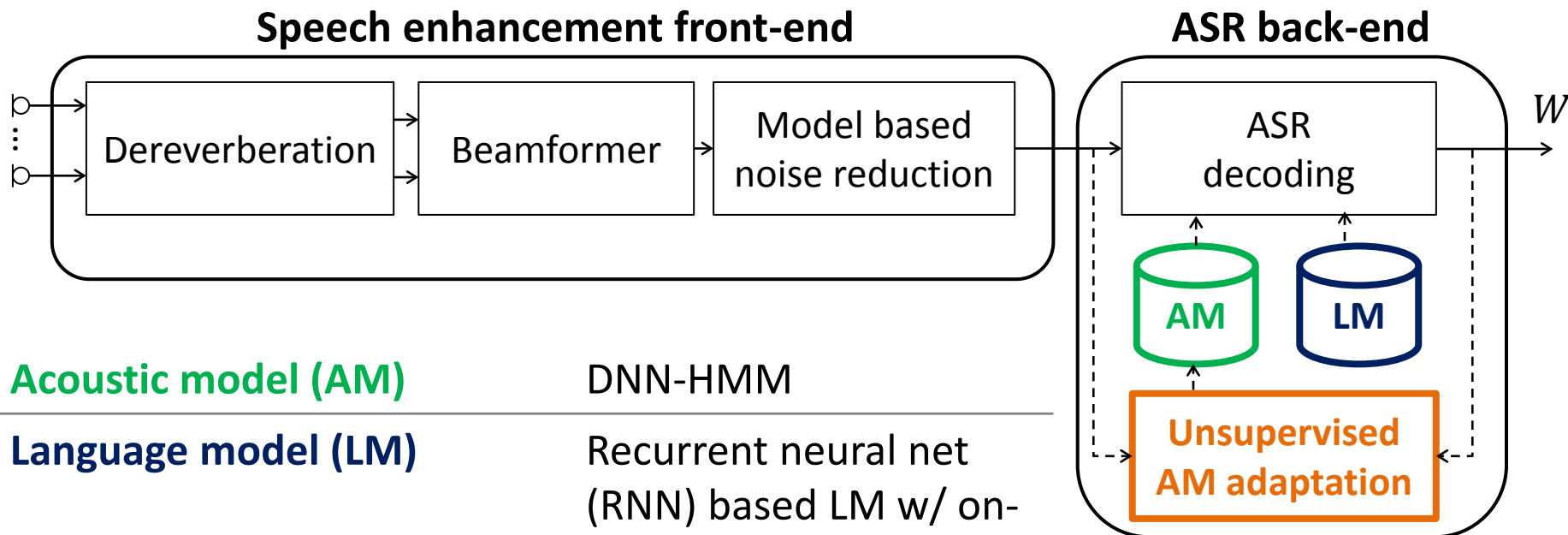


## Speech enhancement front-end



<i>Approach</i>	<b>Linear-prediction based derev. [Yoshioka, 2012]</b>	<b>MVDR beamformer [Souden, 2010]</b>	<b>Modified VTS [Fujimoto, 2012]</b>	<b>DOLPHIN [Nakatani, 2013]</b>
<i>Processing</i>	<b>Linear filtering</b>	<b>Linear filtering</b>	Model-based SE	Model-based SE
<i>Scenario</i>	1ch/2ch/8ch SE/ASR	2ch/8ch SE/ASR	<b>2ch/8ch</b> <b>SE</b>	<b>2ch/8ch</b> <b>ASR</b>
<i>Output</i>	<b>1ch/2ch/8ch</b>	1ch	1ch	1ch
<i>Proc. mode</i>	Utterance batch	Utterance batch	Utterance batch	Full batch
<i>RTF</i>	0.2-2.8	0.03	0.5	6-10

# Characteristics of the ASR back-end



**Acoustic model (AM)**

DNN-HMM

**Language model (LM)**

Recurrent neural net (RNN) based LM w/ on-the-fly rescoring [Hori, 2014]

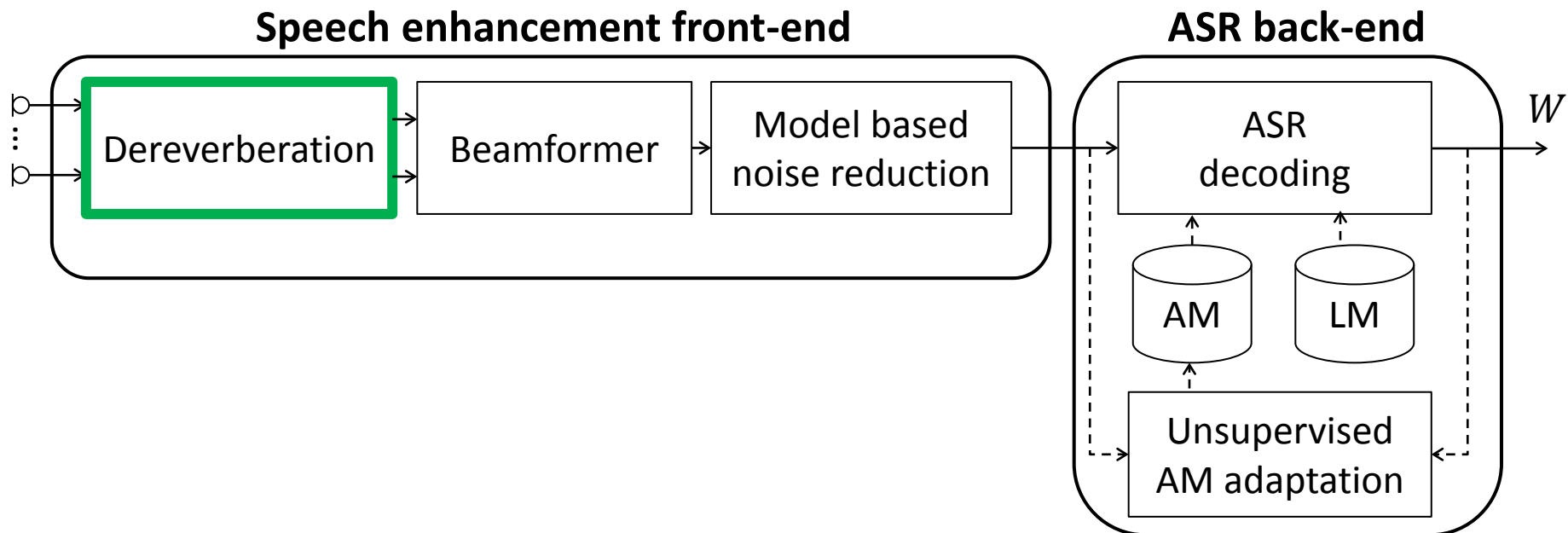
**Unsupervised environmental adaptation**

Retraining of the input layer of DNN-HMM

# Dereverberation



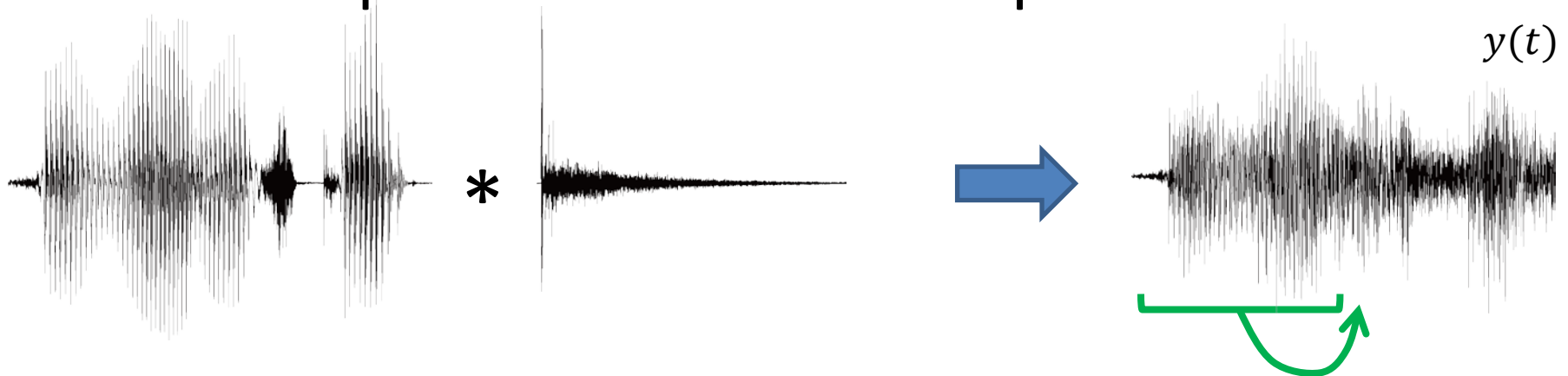
Innovative R&D by NTT



# Dereverberation based on linear prediction (1/2)



- Linear prediction for blind equalization



- Prediction residual:

$$e(t) = y(t) - \sum_{n=1} g_n y(t - n)$$

$g_n$ : prediction filter

Predict reverberation in current observation from past obs.

- Assuming  $e(t)$  is Gaussian,  $g_n$  can be obtained by minimizing:

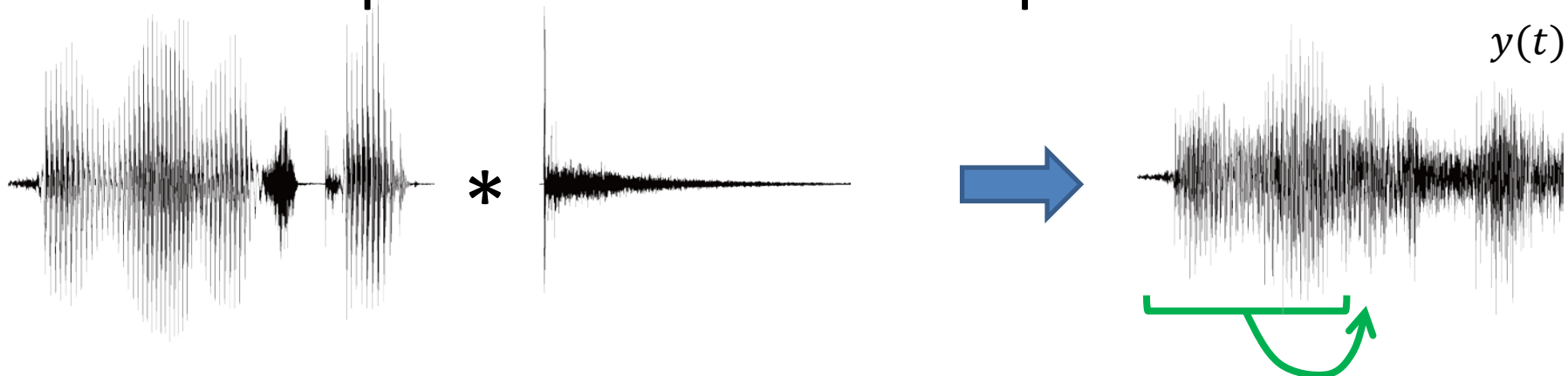
$$F = \sum_t \left| y(t) - \sum_{n=1}^N g_n y(t - n) \right|^2$$



# Dereverberation based on linear prediction (1/2)



- Linear prediction for blind equalization



- Prediction residual:

$$e(t) = y(t) - \sum_{n=1} g_n y(t - n)$$

$g_n$ : prediction filter

- Assuming  $e(t)$  is Gaussian,  $g_n$  can be obtained by

⊖ Not well suited for speech dereverberation

$$F = \sum_t \left| y(t) - \sum_{n=1}^N g_n y(t - n) \right|^2$$

Predict reverberation in current observation from past obs.

# Dereverberation based on linear prediction (2/2)



- For speech dereverberation [Yoshioka, 2012]

- **Introduce time delay  $\tau$**

*Do not equalize speech generative process*

→ Focus on late reverberation

- **Better modeling of speech**

*Short time Gaussianity of speech w/ time varying variance  $\sigma_t^2$*

→ Weighted prediction error (WPE)

- $g_n, \sigma_t^2$  can be obtained by minimizing:

$$F = \sum_t \frac{|y(t) - \sum_{n=\tau}^N g_n y(t-n)|^2}{\sigma_t^2} + \sum_t \log(\sigma_t^2)$$

# Dereverberation based on linear prediction (2/2)



- For speech dereverberation [Yoshioka, 2012]

- **Introduce time delay  $\tau$**

*Do not equalize speech generative process*

→ Focus on late reverberation

- **Better modeling of speech**

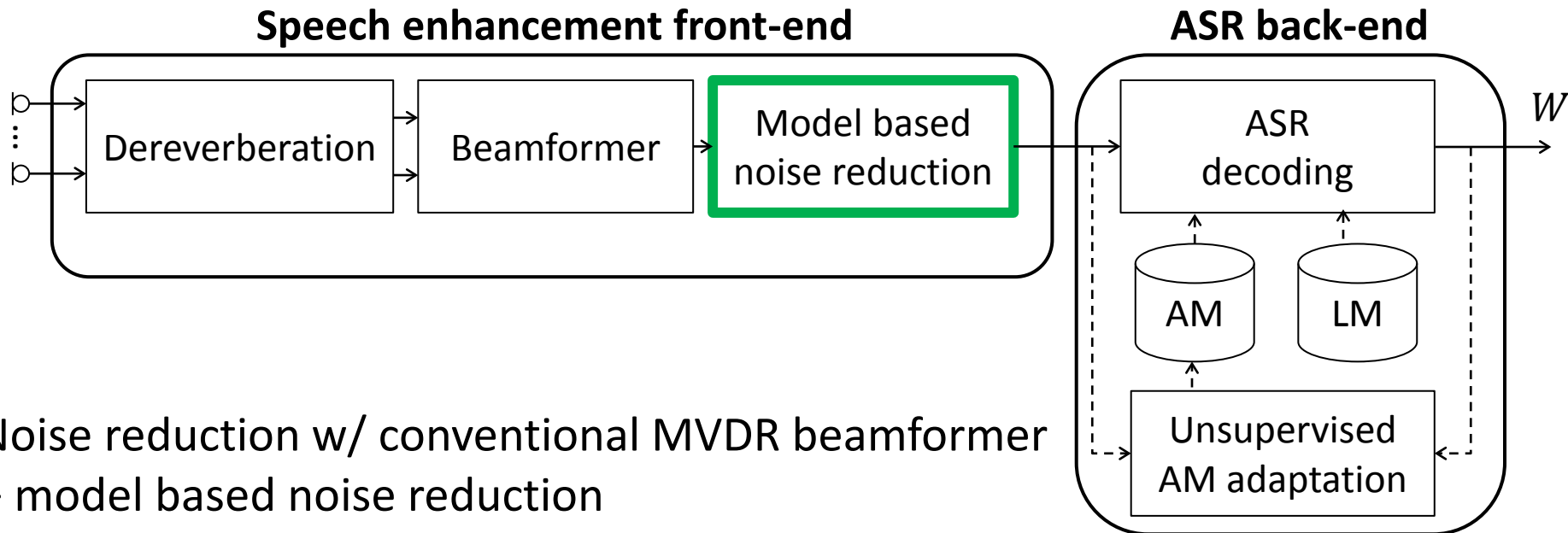
*Short time Gaussianity of speech w/ time varying variance  $\sigma_t^2$*

→ Weighted

- **Precise speech dereverberation**  
**Relatively robust to noise**  
**Implemented in STFT domain → efficient**

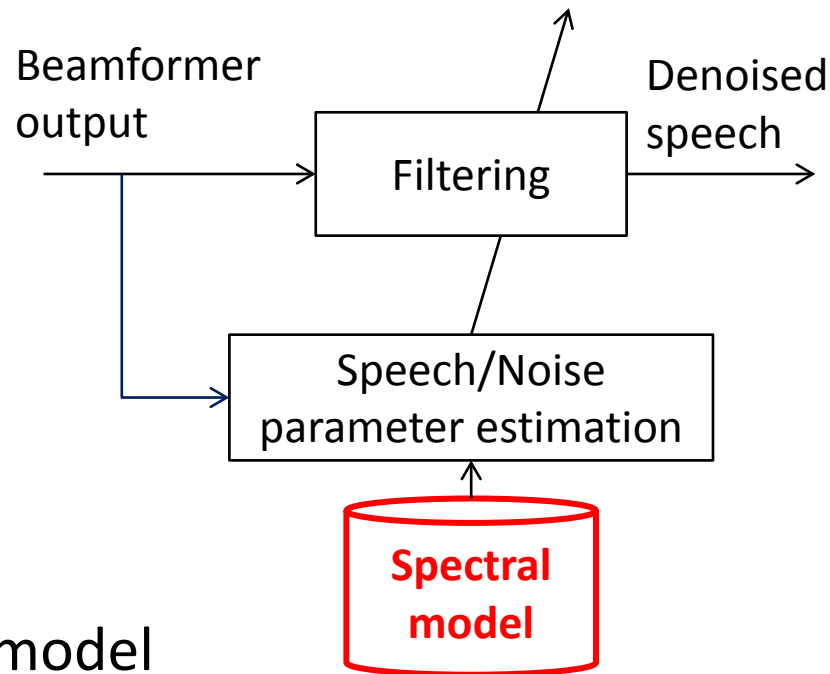
$$F = \sum_t \frac{|y(t) - \sum_{n=1}^N \hat{y}_n(t-n)|^2}{\sigma_t^2} + \sum_t \log(\sigma_t^2)$$

# Model based Noise reduction



Noise reduction w/ conventional MVDR beamformer  
+ model based noise reduction

## VTS-like model-based noise reduction [Fujimoto, 2012]

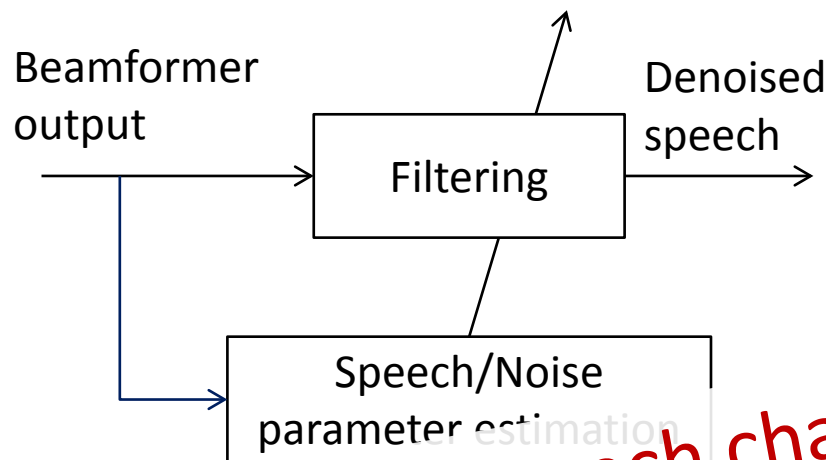


- Spectral model
  - Speech model: pre-trained using clean training data + adaptation (per utterance)
  - Noise model: estimated per utterance

# Modified VTS (M-VTS)



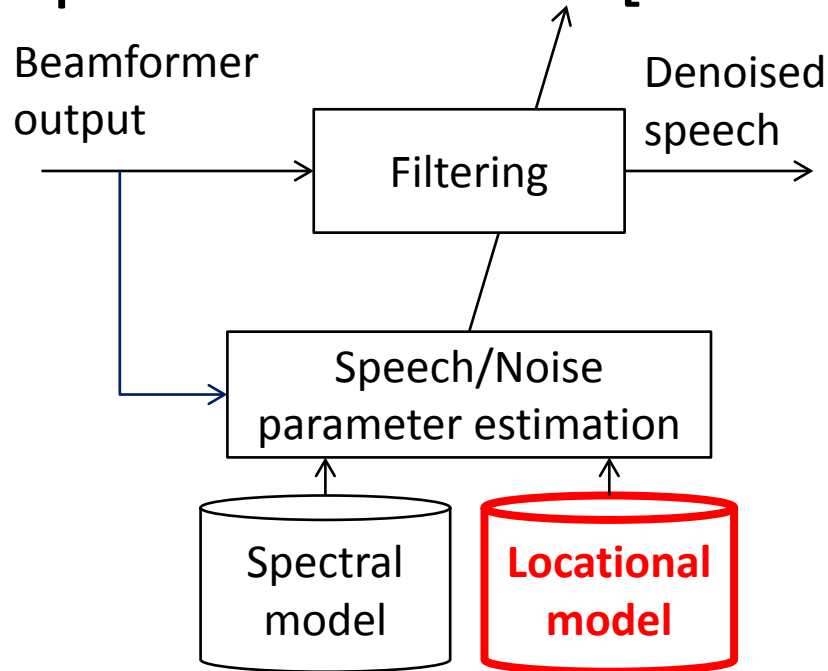
## VTS-like model-based noise reduction [Fujimoto, 2012]



😊 Denoised speech characteristics close to that of clean speech  
Used for 2ch/8ch SE tasks  
Utterance batch

- Spectral model-based noise reduction
  - Speech model: pre-trained using clean training data + adaptation (per utterance)
  - Noise model: estimated per utterance

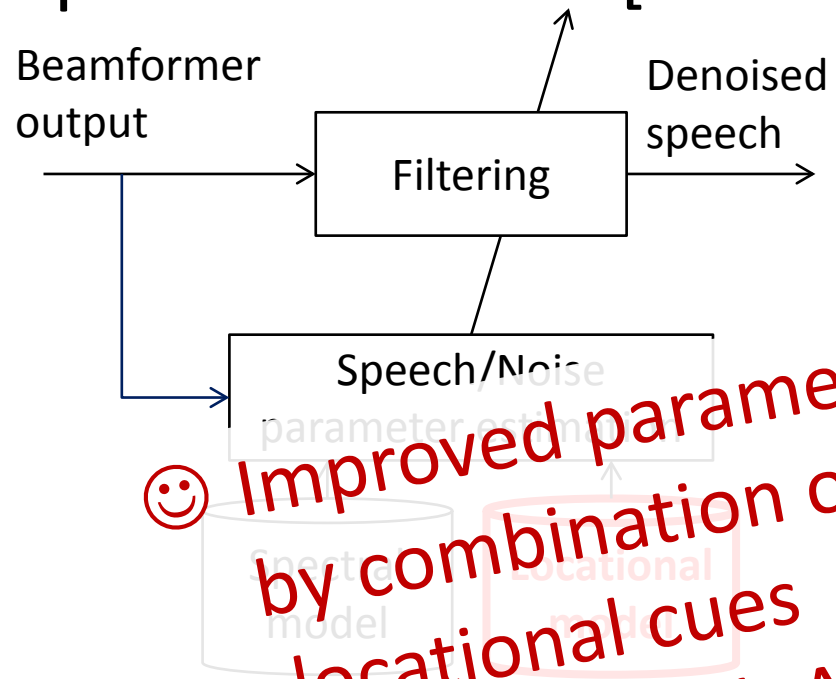
## Combines spectral and locational (DOA feature) models of speech and noise [Nakatani, 2013]



- Locational model estimated from multi-channel dereverberation output



Combines spectral and locational (DOA feature) models of speech and noise [Nakatani, 2013]



😊 Improved parameter estimation by combination of spectral and locational cues  
Used for 2ch/8ch ASR tasks  
Full batch

- Locational model dereverberation



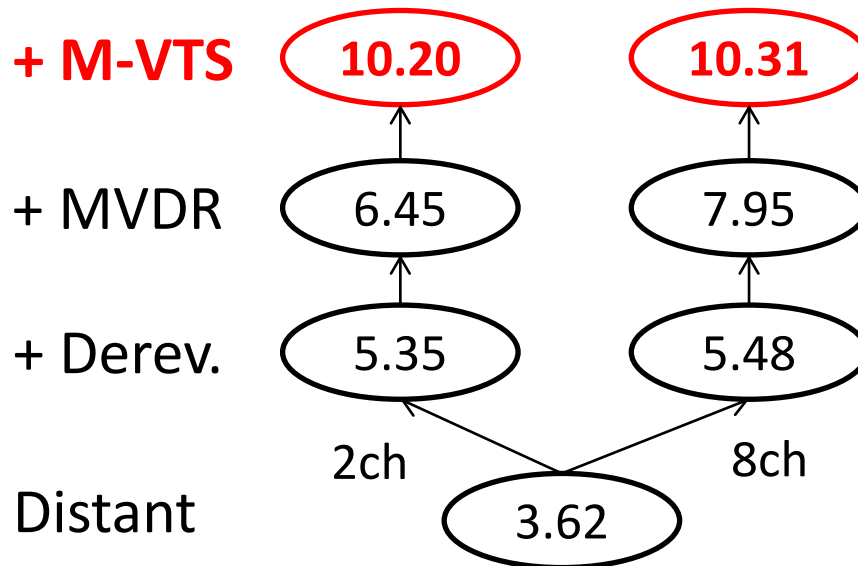
# SE results: Objective evaluation (Eval set)



	CD	SRMR	LLR	FWSegSNR	PESQ
w/o front-end	3.97	3.68	0.58	3.62	1.48
Proposed (2ch)	2.34	5.06	0.41	10.20	2.43
<b>Proposed (8ch)</b>	<b>2.25</b>	<b>5.39</b>	0.43	<b>10.31</b>	<b>2.82</b>

**Improvement for all measures**

Avg. FWSegSNR (dB) for *SimData*



**Consistent improvement**

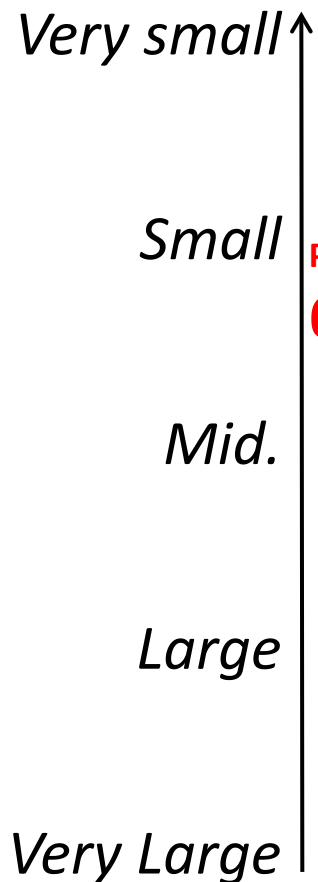
# SE results: Subjective evaluation (Eval set)



For RealData (far) 8ch

Perceived reverb.

Overall Quality



**Proposed**  
**62.4 (20.5)**

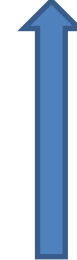


Distant  
**16.6 (12.0)**

**Significant improvement**  
*Same for all conditions*



**Proposed**  
**70.5 (10.6)**

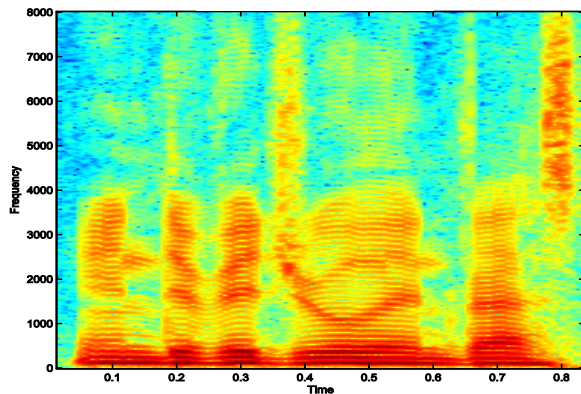


Distant  
**25.3 (7.5)**

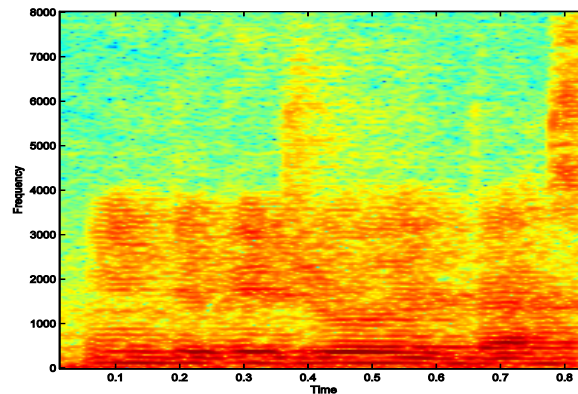
# SE results: Spectrograms RealData (Eval)



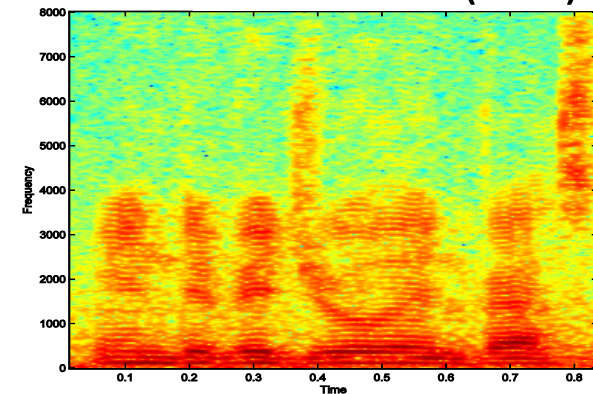
Headset



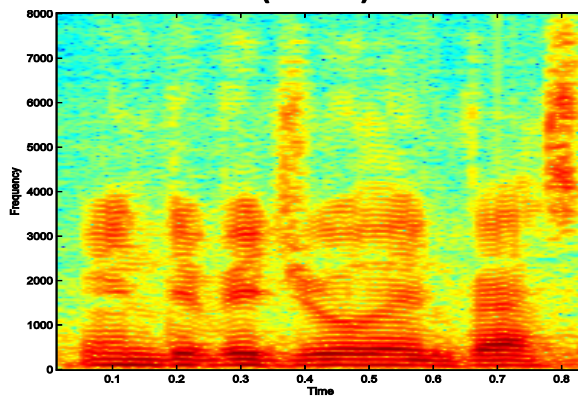
Distant



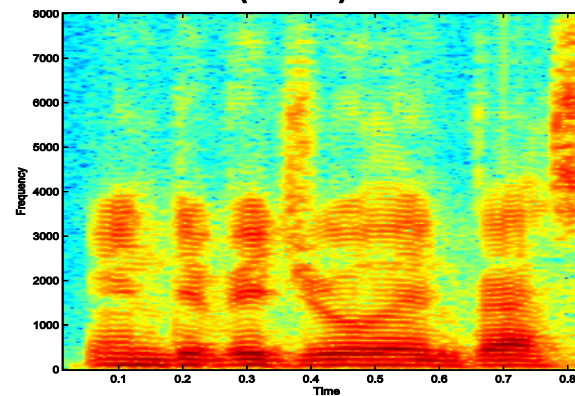
Dereverb. (8ch)



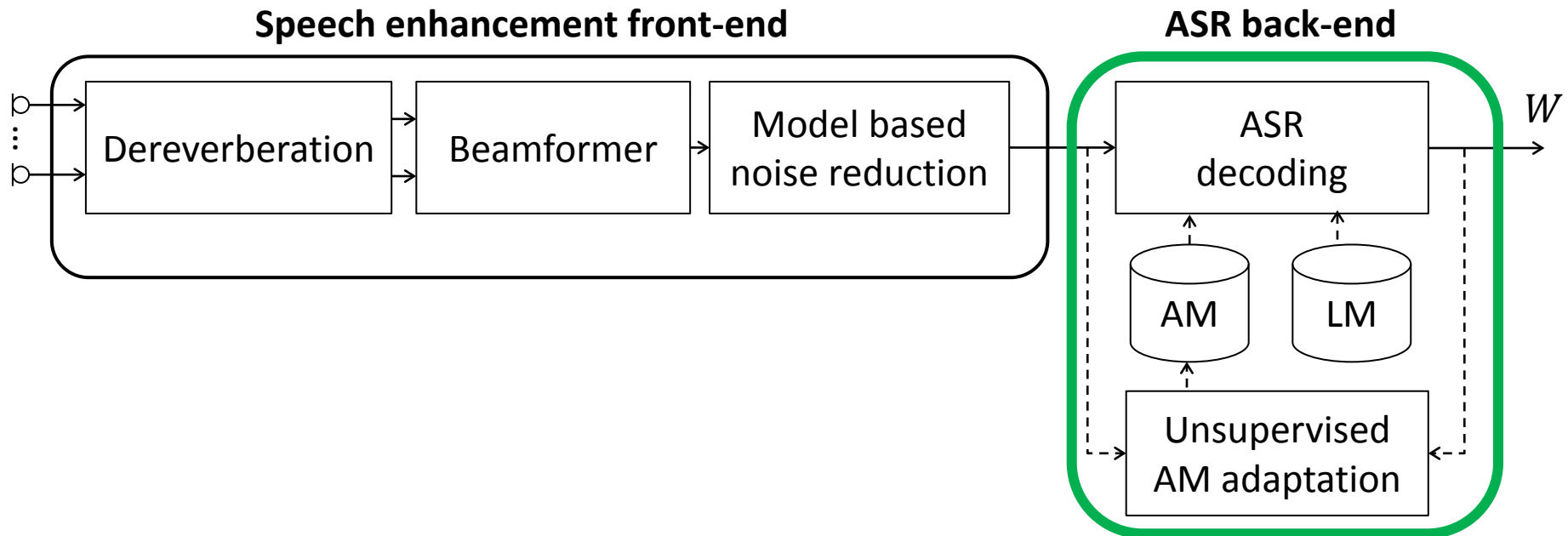
Dereverb.(8ch) + MVDR



Dereverb.(8ch) + MVDR + M-VTS



# System overview



- DNN-HMM based acoustic model
  - *layer-wise RBM pre-training + fine tuning w/ stochastic gradient descent*
  - *3129 HMM states/ State alignment from clean training data*
  - *7 hidden layers (2048 units)*
- Features
  - *40 Log mel filter-bank coefficients +  $\Delta$  +  $\Delta\Delta$  (120)*
  - *5 left+5 right context (11 frames)*
- Language model LM
  - *Trigram*
  - *RNN with fast decoding using on-the-fly rescoring [Hori, 2014]*
- Unsupervised environmental adaptation
  - *Retrain 1<sup>st</sup> layer of DNN-HMM w/ small learning rate using labels obtained from a 1<sup>st</sup> recognition pass*

# Extended training data



- Combat mismatch between training/SimData and RealData

→ Add acoustic variety during DNN training

Challenge  
multi-cond.  
Data  
**@ 20 dB**

+

multi-cond.  
Data  
**@ 10 dB**

+

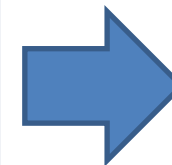
multi-cond.  
Data  
**@ 15 dB**

+

WSJCAMO  
Clean training

+

WSJCAMO  
**Desktop mic**  
training data



Extended  
training Data

**Only used data available  
to all participants**

# Summary of ASR results (Eval)

	SimData						RealData			
	Room 1		Room 2		Room 3		Ave.	Room1		Ave.
	Near	Far	Near	Far	Near	Far		Near	Far	
w/o front-end	3.8	4.4	5.3	8.5	5.8	9.5	<b>6.2</b>	21.1	23.3	<b>22.2</b>
<b>Best system (8ch)</b>	<b>3.7</b>	<b>4.0</b>	<b>4.0</b>	<b>4.5</b>	<b>4.4</b>	<b>4.8</b>	<b>4.2</b>	<b>8.8</b>	<b>9.3</b>	<b>9.0</b>

+ DOLPHIN

12.7

**9.0**

+ MVDR

13.9

10.0

+ Derev.

17.4

14.9

13.8

+ Adap

1ch

**22.2**

2ch

8ch

+ RNN LM

26.3

+ Extended data

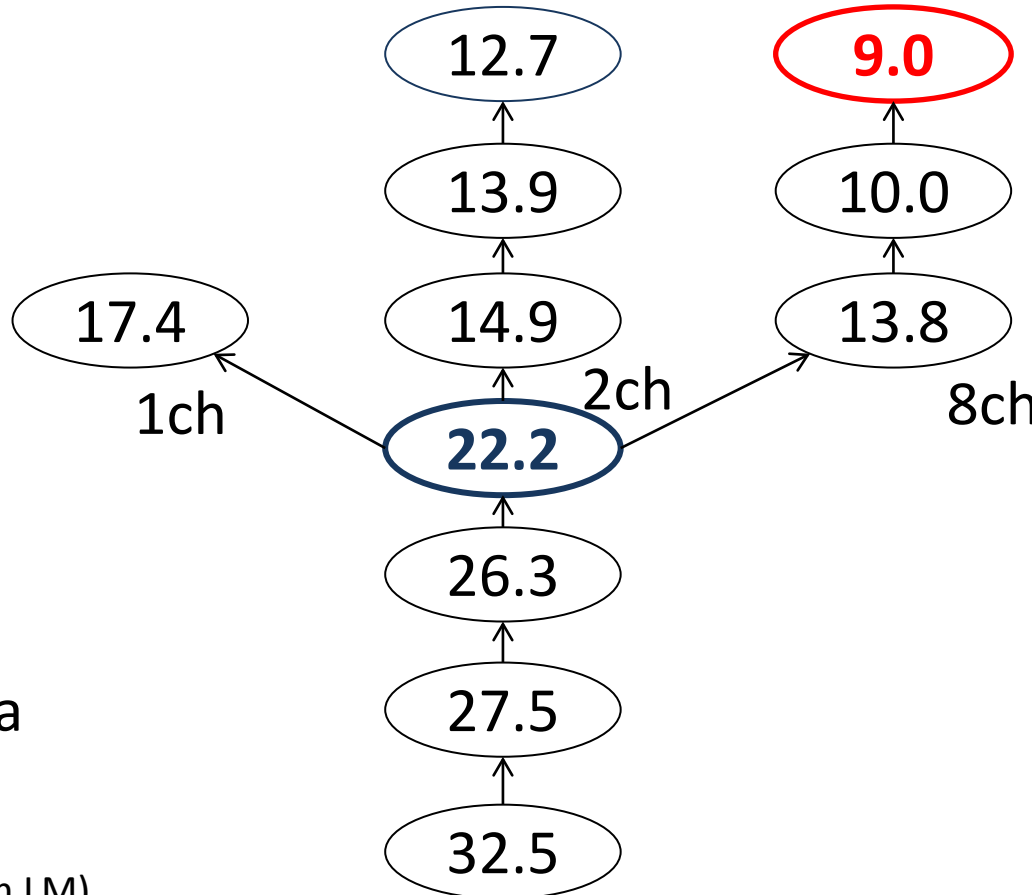
27.5

Distant

32.5

(w/ DNN AM & Trigram LM)

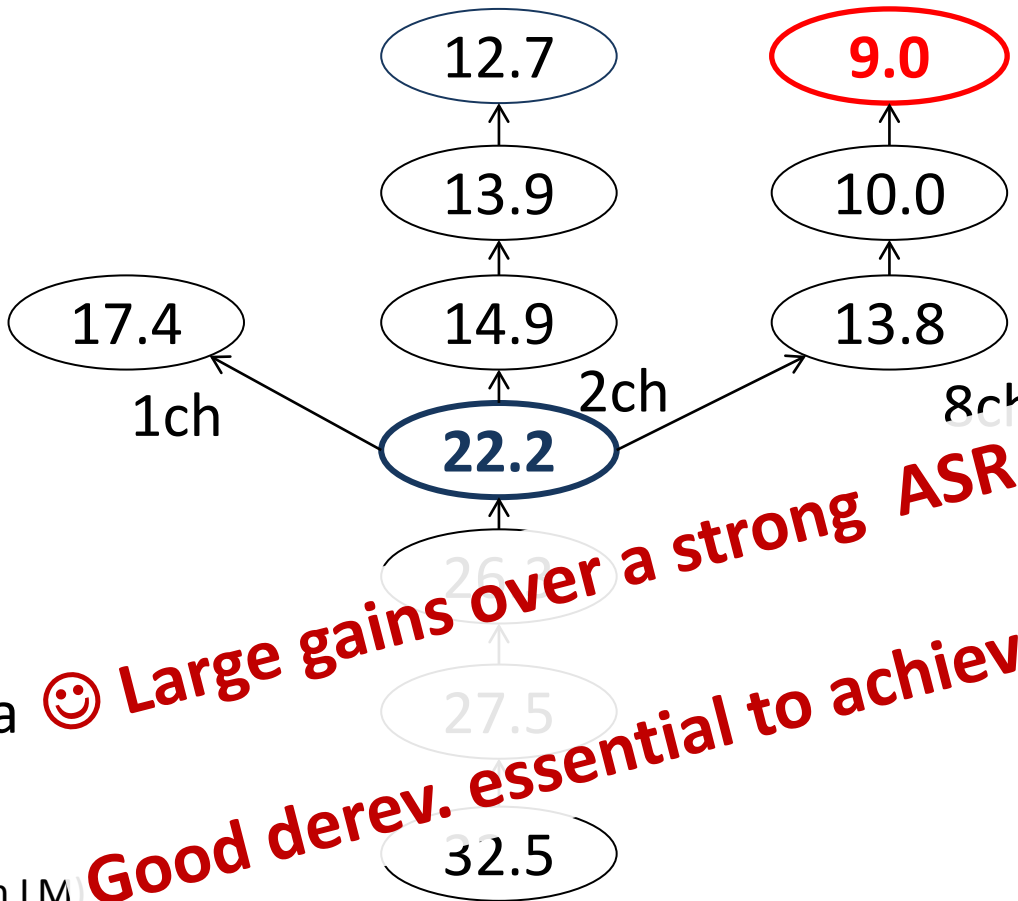
*For RealData*



# Summary of ASR results (Eval)

	SimData						RealData			
	Room 1		Room 2		Room 3		Ave.	Room1		Ave.
	Near	Far	Near	Far	Near	Far		Near	Far	
w/o front-end	3.8	4.4	5.3	8.5	5.8	9.5	6.2	21.1	23.3	22.2
<b>Best system (8ch)</b>	3.7	4.0	4.0	4.5	4.4	4.8	4.2	8.8	9.3	9.0

- + DOLPHIN
- + MVDR
- + Derev.
- + Adap
- + RNN LM
- + Extended data
- Distant  
(w/ DNN AM & Trigram LM)



**Large gains over a strong ASR backend**  
**Good derev. essential to achieve large gains**  
*For RealData*



# Conclusion



- Combined advanced SE and ASR techniques for reverberant speech
- Large performance improvement for both SE and ASR tasks

## Where are we?

<b>Proposed (1ch)</b>	<b>Proposed (2ch)</b>	<b>Proposed (8ch)</b>	Lapel Headset (MC-WSJ)	Clean (WSJCAM0)	
<b>17.4 %</b>	<b>12.7 %</b>	<b>9.0 %</b>	8.3 %	5.9 %	3.6 %

← Room for improvement for robust ASR → Non-environmental mismatch between SimData and RealData →



Innovative R&D by NTT

**Thank you!**