

# The Hitachi/JHU CHiME-5 system: Advances in speech recognition for everyday home environments using multiple microphone arrays

*Naoyuki Kanda<sup>1</sup>, Rintaro Ikeshita<sup>1</sup>, Shota Horiguchi<sup>1</sup>, Yusuke Fujita<sup>1</sup>, Kenji Nagamatsu<sup>1</sup>,  
Xiaofei Wang<sup>2</sup>, Vimal Manohar<sup>2</sup>, Nelson Enrique Yalta Soplín<sup>2</sup>, Matthew Maciejewski<sup>2</sup>,  
Szu-Jui Chen<sup>2</sup>, Aswin Shanmugam Subramanian<sup>2</sup>, Ruizhi Li<sup>2</sup>, Zhiqi Wang<sup>2</sup>, Jason Naradowsky<sup>2</sup>,  
L. Paola Garcia-Perera<sup>2</sup>, Gregory Sell<sup>2</sup>*

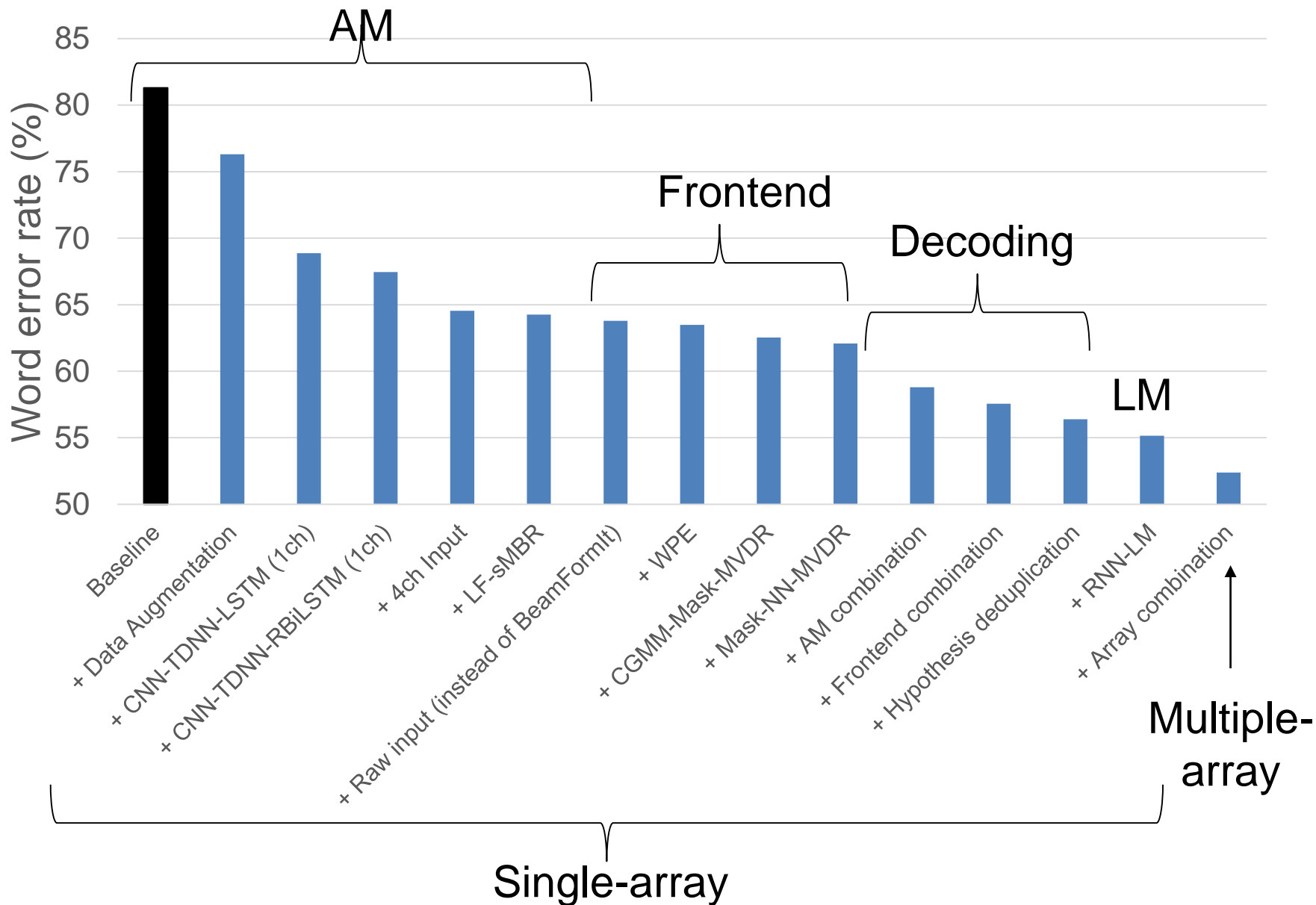
1

**HITACHI**  
Inspire the Next

2

 **JOHNS HOPKINS**  
WHITING SCHOOL  
of ENGINEERING

# Step-by-Step Improvements for Dev



# Acoustic Model Training Pipeline

**Step 1.**  
GMM-AM

1ch worn L  
1ch worn R  
1ch worn L+R



**Step 2.**  
Alignment

1ch worn L+R



“Cleaned” 1ch worn L+R  
& phone-state alignment

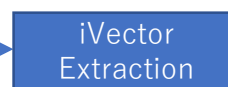
Full set



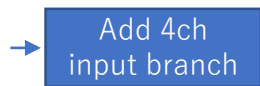
“Cleaned” full set  
& phone-state alignment

**Step 3.**  
1ch AM

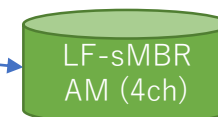
“Cleaned” full set  
& phone-state alignment



**Step 4.**  
4ch AM



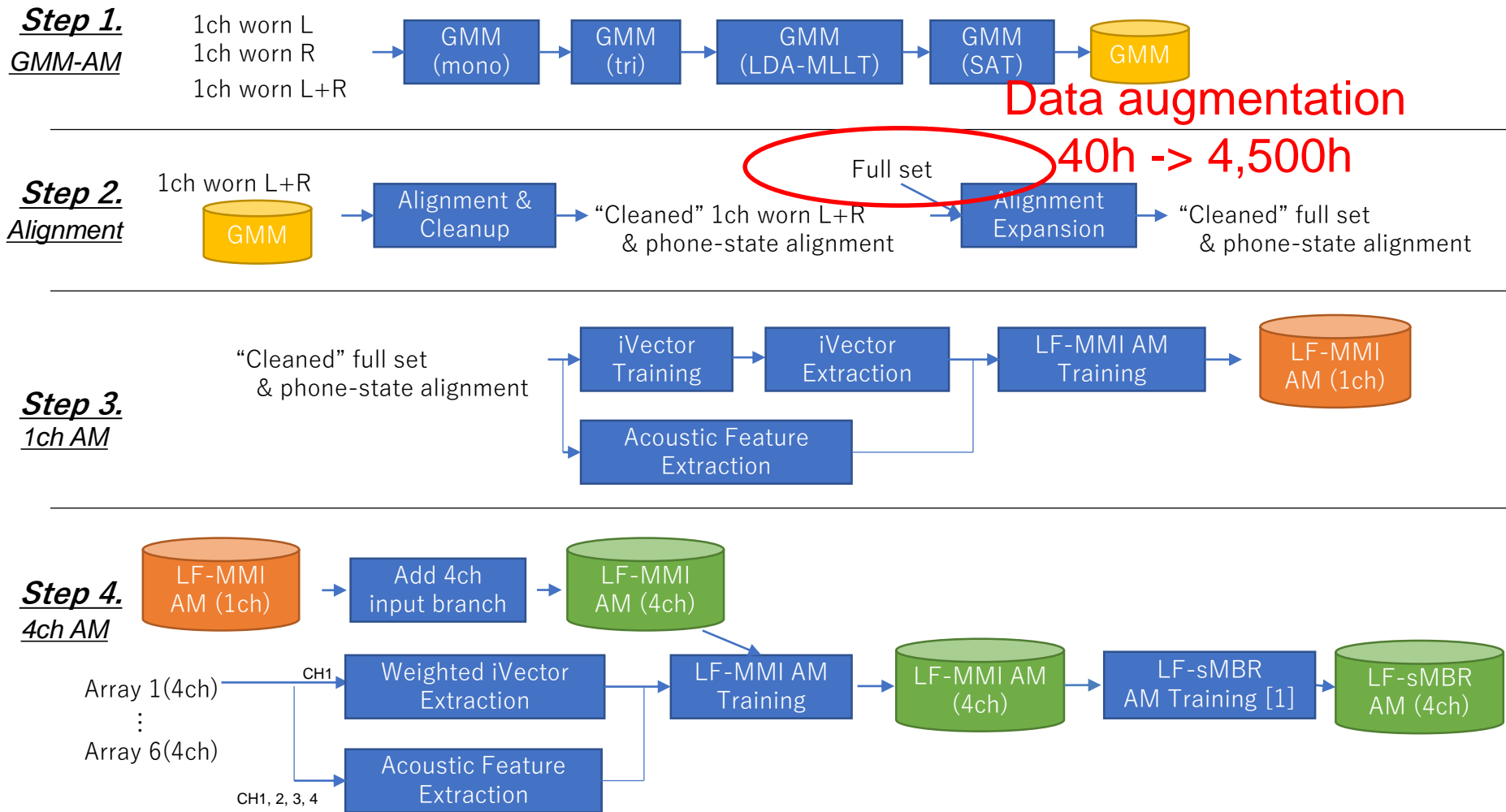
Array 1(4ch)  
⋮  
Array 6(4ch)



CH1, 2, 3, 4



# Acoustic Model Training Pipeline



[1] Naoyuki Kanda, Yusuke Fujita, Kenji Nagamatsu, Lattice-free state-level minimum Bayes risk training of acoustic models, Interspeech 2018.

## Effect of data augmentation with baseline AM

Data			Data Augmentation			Training Epoch	Worn-Dev	Ref-Array-Dev
Worn	Array (Raw, CH1)	Array (BeamFormIt)	Speed & Volume	Reverb. & Noise(*)	Bandpass			
L, R, L+R	1		✓			4	44.05	79.65
L, R, L+R	1	1	✓			4	44.49	78.72
L, R, L+R	1 ... 6	1 ... 6	✓			4	48.92	78.51
L, R, L+R	1 .... 6	1 ... 6	✓	✓		2	45.82	77.26
L, R, L+R	1 ... 6	1 ... 6	✓	✓	✓	1	45.37	76.31

(\*) Reverb. & noise perturbation was applied only for worn microphone data.

Speed: 0.9, 1.0, 1.1

Volume: 0.125 – 2.0

Reverberation: Generate impulse responses of simulated rooms by image method.

Follow the settings of {small, medium}-size rooms in [1].

Noise: Add non-speech region of array data with SNR of {20,15,10,5, 0}

Bandpass: Randomly-selected frequency band was cut off.

(leave at least 1,000 Hz band within the range of less than 4,000 Hz)

[1] T. Ko, et al.: A study on data augmentation of reverberant speech for robust speech recognition, Proc. ICASSP, pp. 5220—5224, 2017.

# Acoustic Model Training Pipeline

**Step 1.**  
GMM-AM

1ch worn L  
1ch worn R  
1ch worn L+R



**Step 2.**  
Alignment

1ch worn L+R



“Cleaned” 1ch worn L+R  
& phone-state alignment

Full set



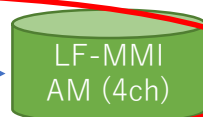
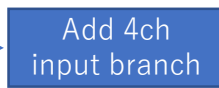
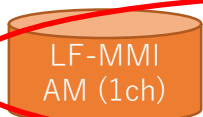
“Cleaned” full set  
& phone-state alignment

**Step 3.**  
1ch AM

“Cleaned” full set  
& phone-state alignment

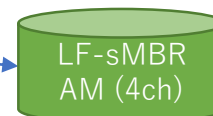


**Step 4.**  
4ch AM



4ch AM

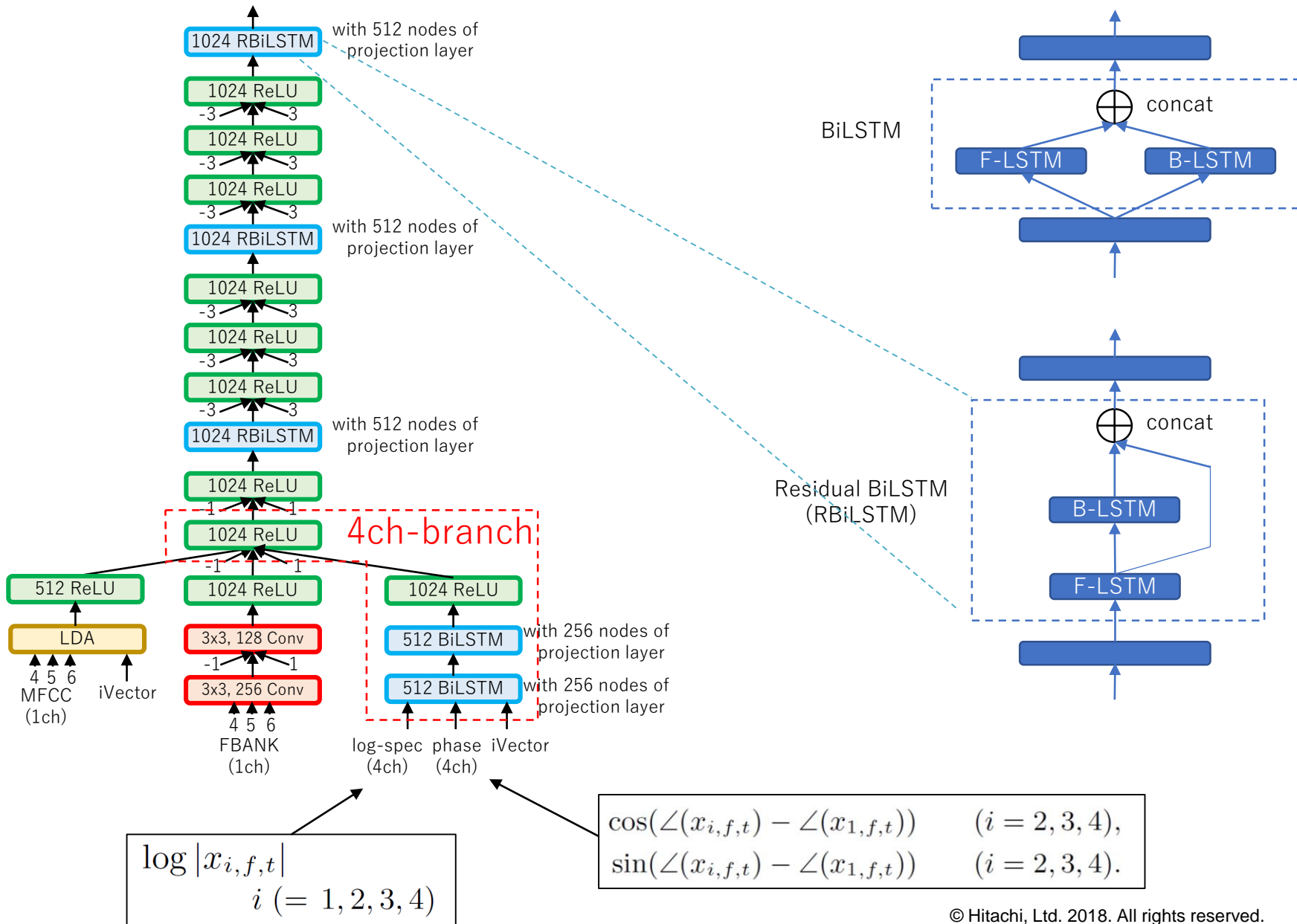
Array 1(4ch)  
⋮  
Array 6(4ch)



CH1, 2, 3, 4

[1] Naoyuki Kanda, Yusuke Fujita, Kenji Nagamatsu, Lattice-free state-level minimum Bayes risk training of acoustic models, Interspeech 2018.

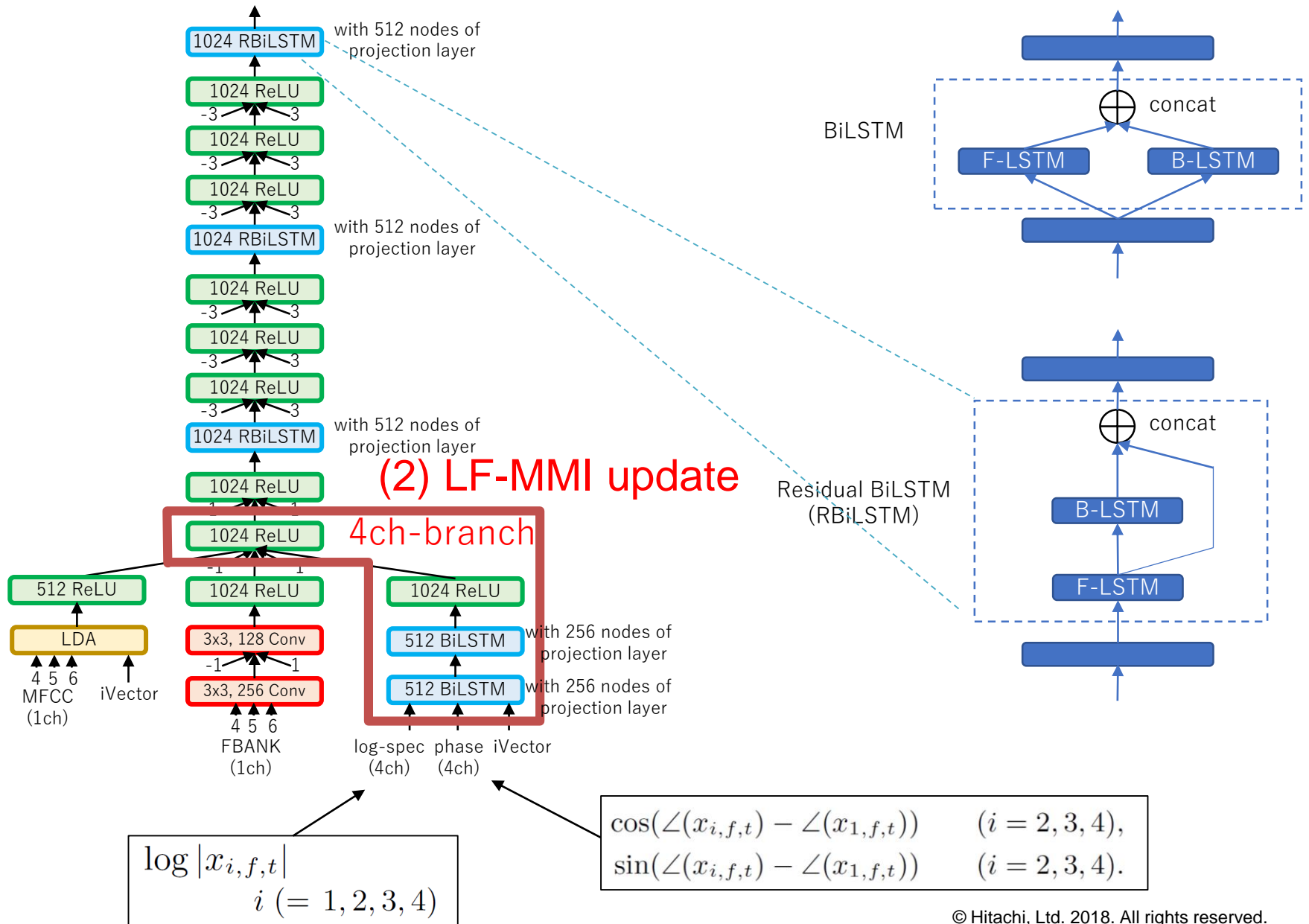
# 4ch CNN-TDNN-RBiLSTM



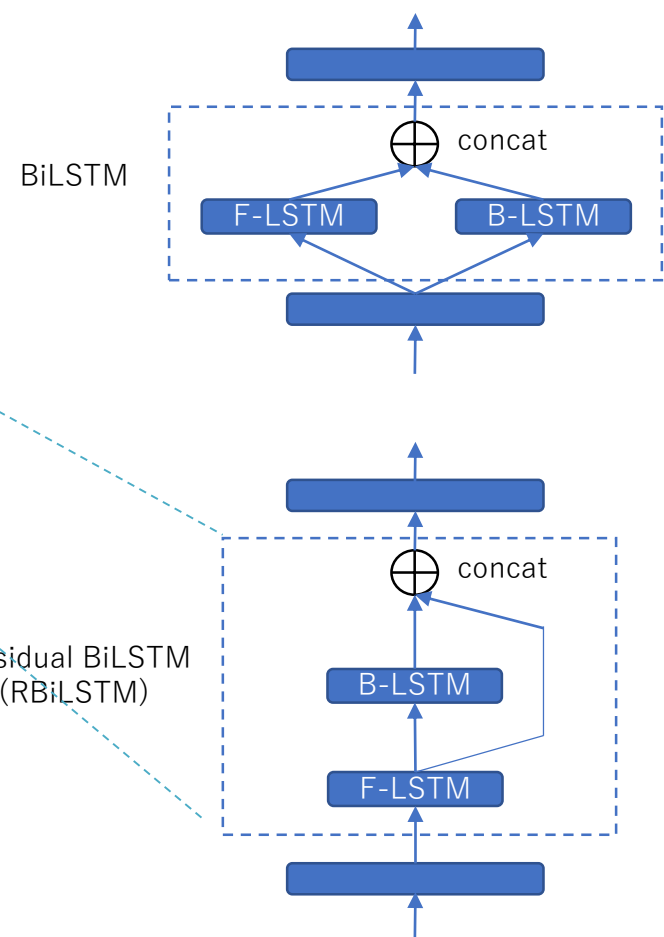
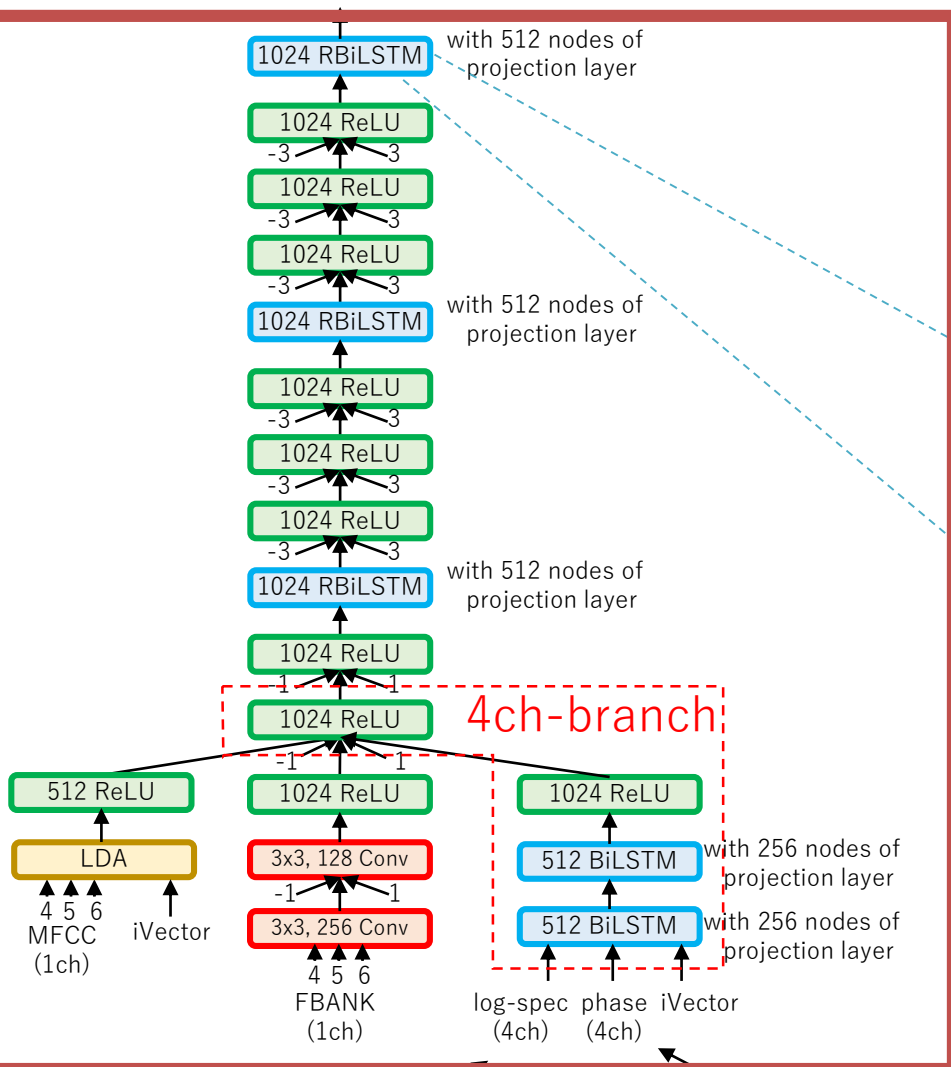




# 4ch CNN-TDNN-RBiLSTM



# 4ch CNN-TDNN-RBiLSTM



(3) LF-sMBR update

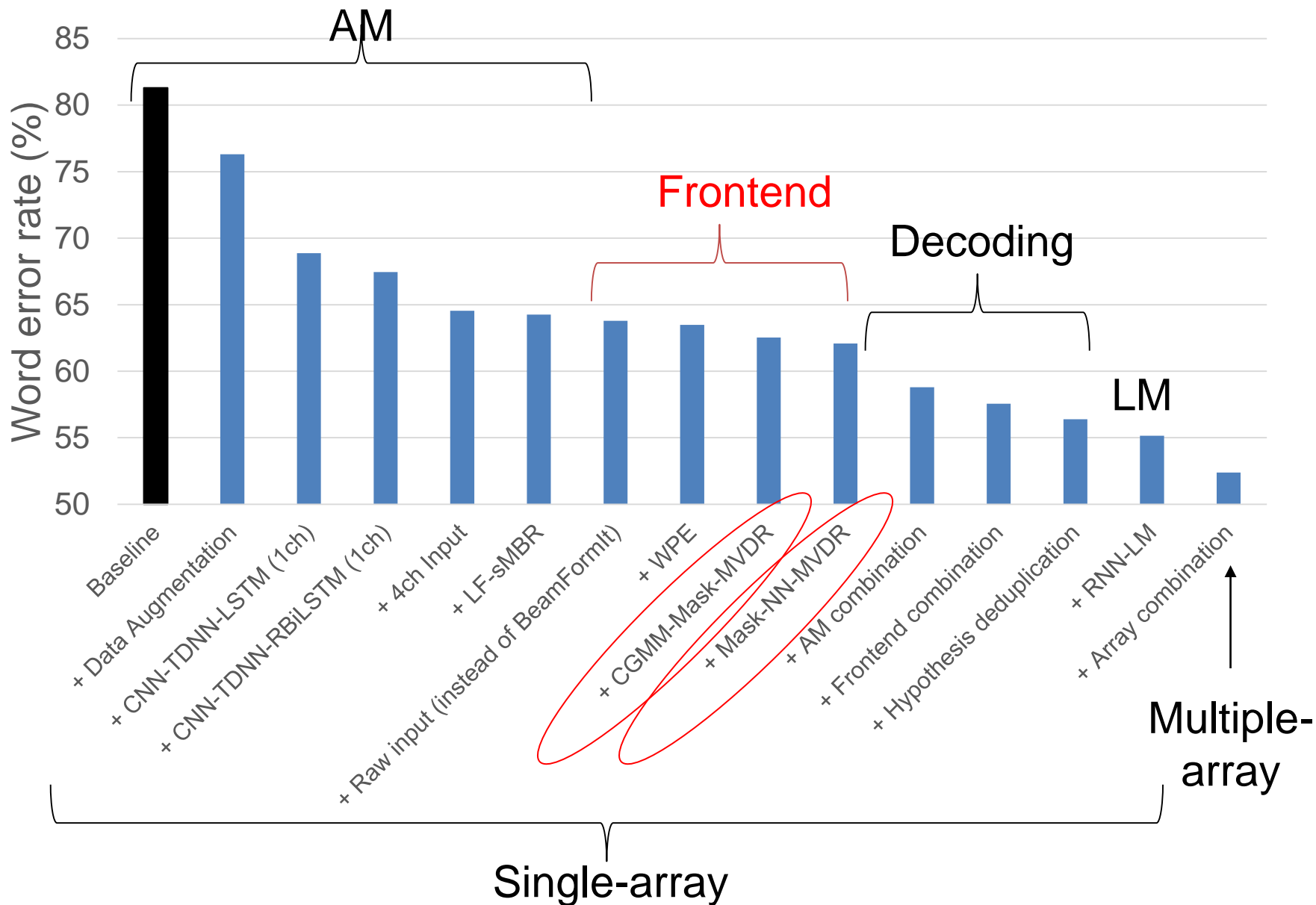
$$\log |x_{i,f,t}|$$

$$i (= 1, 2, 3, 4)$$

$$\cos(\angle(x_{i,f,t}) - \angle(x_{1,f,t})) \quad (i = 2, 3, 4),$$

$$\sin(\angle(x_{i,f,t}) - \angle(x_{1,f,t})) \quad (i = 2, 3, 4).$$

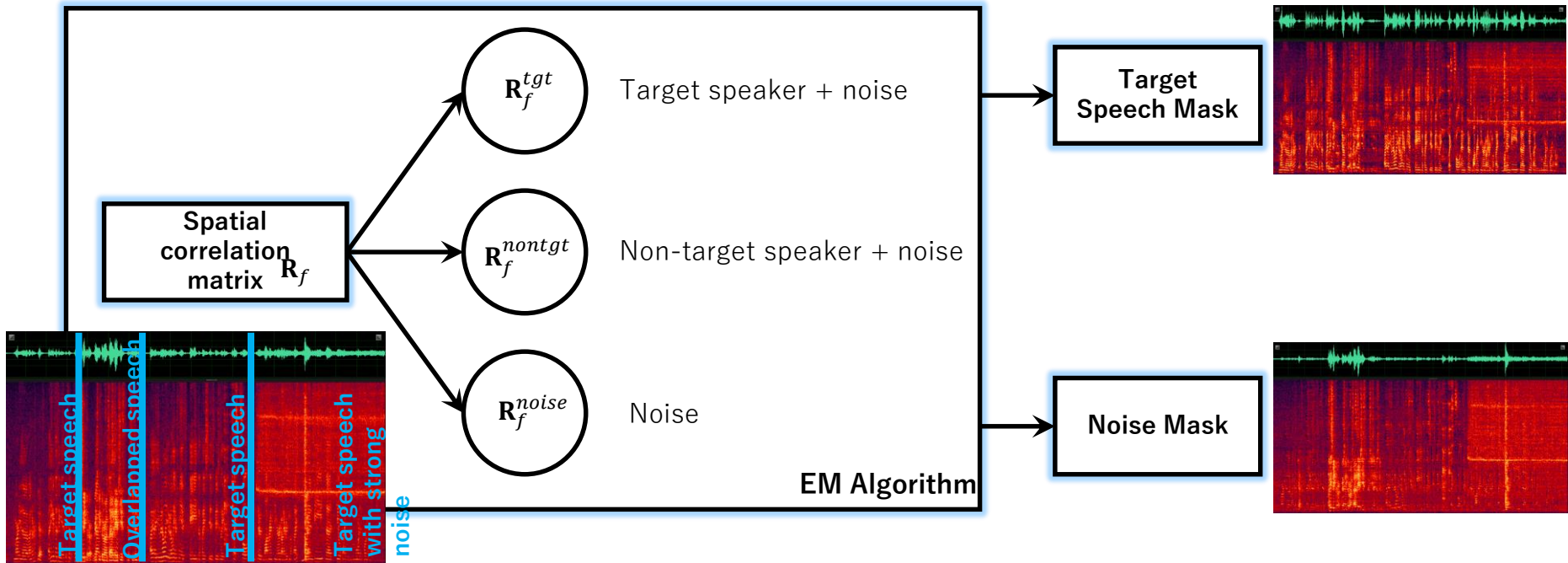
# Step-by-Step Improvements for Dev



- 3-class mixture: target, non-target, and noise

$$y_f(t) \sim \alpha_{tgt} \mathcal{N}_{\mathbb{C}}(0, v_f^{tgt}(t) R_f^{tgt}) + \alpha_{nontgt} \mathcal{N}_{\mathbb{C}}(0, v_f^{nontgt}(t) R_f^{nontgt}) + \alpha_{noise} \mathcal{N}_{\mathbb{C}}(0, v_f^{noise}(t) R_f^{noise})$$

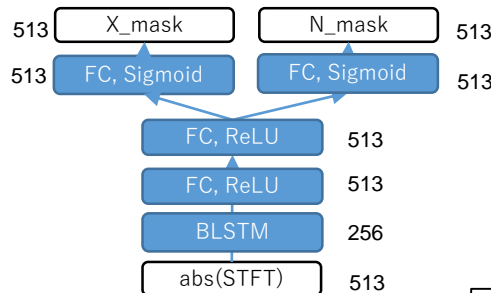
- Mask estimation using EM Algorithm  $\rightarrow$  MVDR-based Beamformer



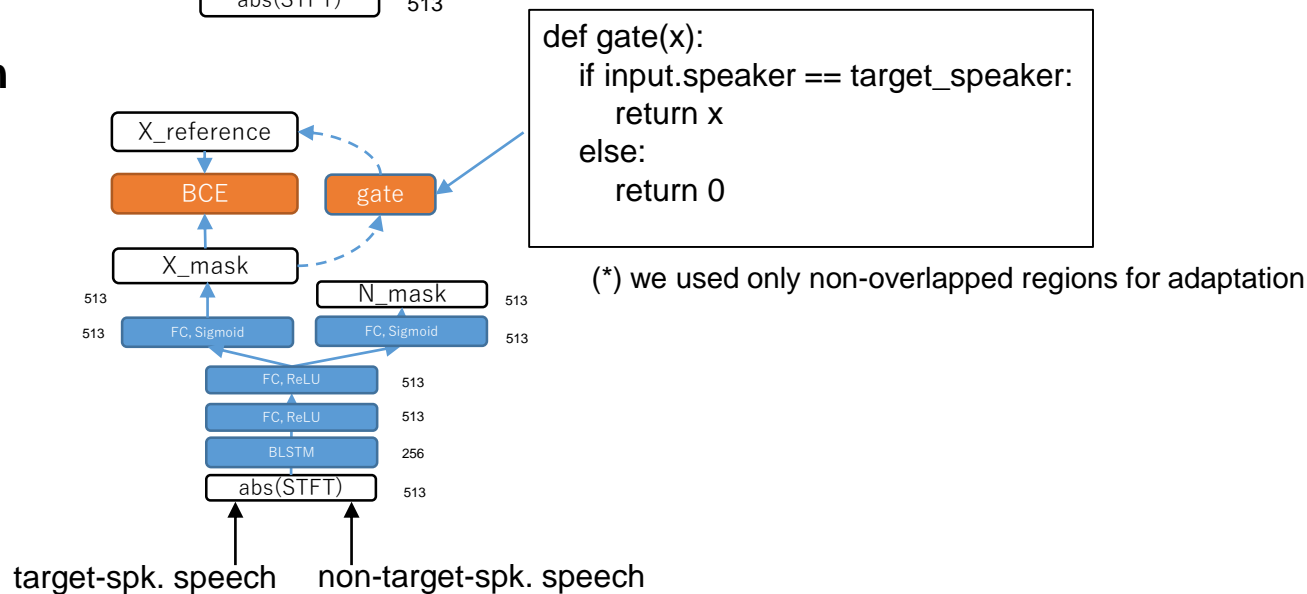
Higuchi, Takuya, et al. "Online MVDR beamformer based on complex Gaussian mixture model with spatial prior for noise robust ASR." *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 25.4 (2017): 780-793.

## 1. Train mask estimation (ME) network [1][2]

by using mixture of speech (worn non-speaker-overlapped region) and noise (array non-speech region) in the training set



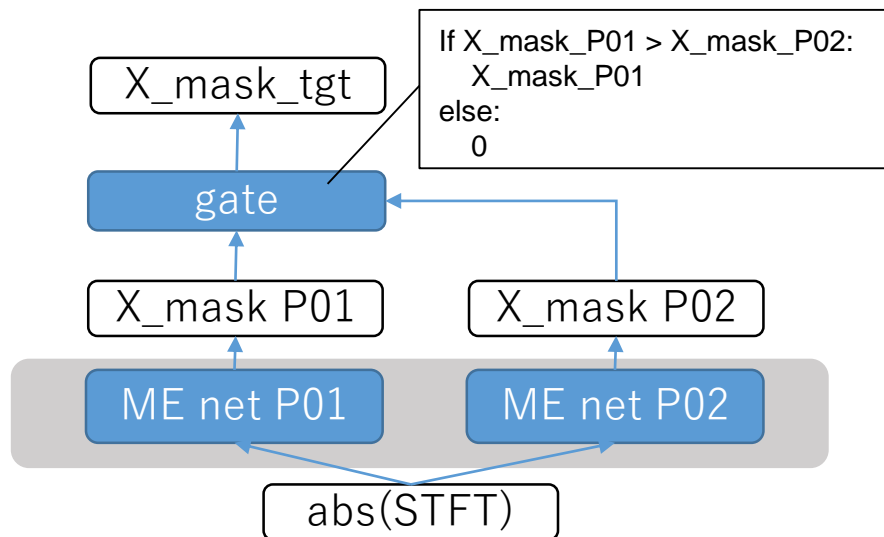
## 2. Speaker adaptation



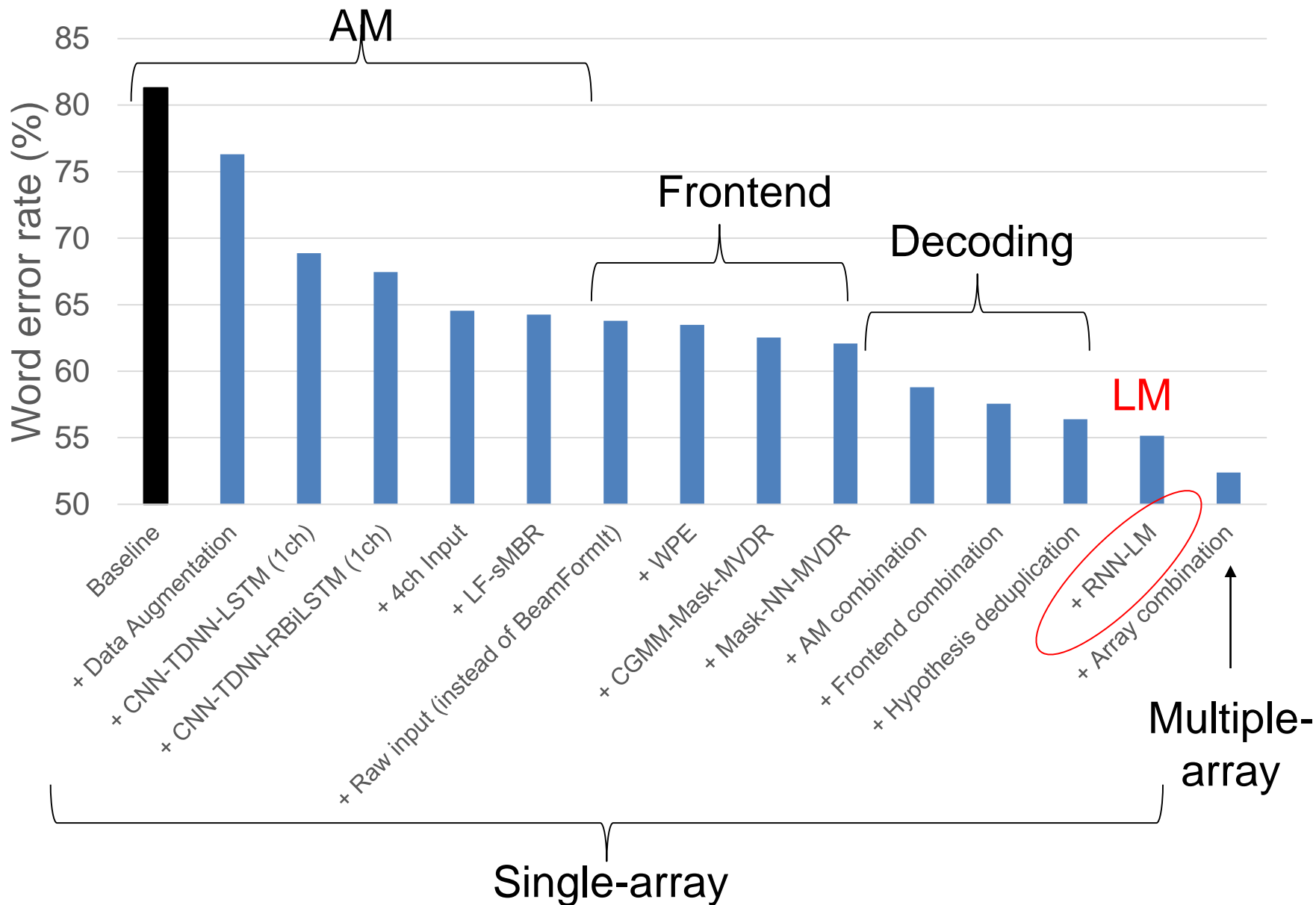
## 3. Mask inference

Target speaker's mask is selected only if target speaker's output value is higher than all other non-targets values.

*Example: P01(target) and P02(non-target)*



# Step-by-Step Improvements for Dev



- Recurrent neural network based word-LM
  - 2 layer LSTM with 512 nodes, 50% dropout
  - 512 dim embeddings
  - PyTorch implementation
- Official-LM: forward-RNN-LM: backward-RNN-LM  
= 0.5 : 0.25 : 0.25

WER (%) for Dev set

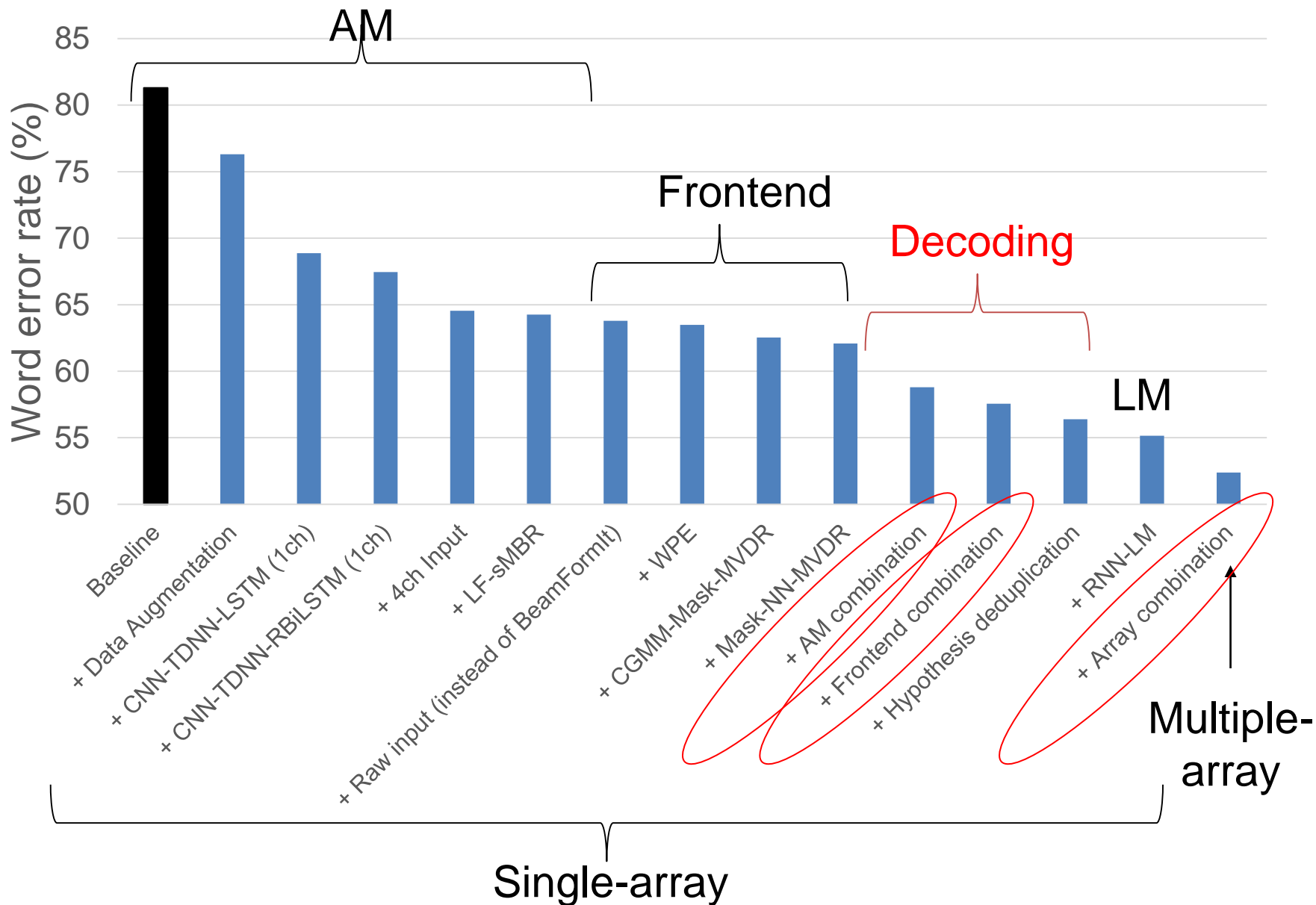
	without RNN-LM	with RNN-LM
Single-array	56.40	55.15
Multiple-array	54.00	52.38

Annotations: Red boxes around the improvement percentages (1.3% and 1.6%) and red arrows pointing from the 'without RNN-LM' column to the 'with RNN-LM' column.

(\*) Results with model combination and hypothesis deduplication



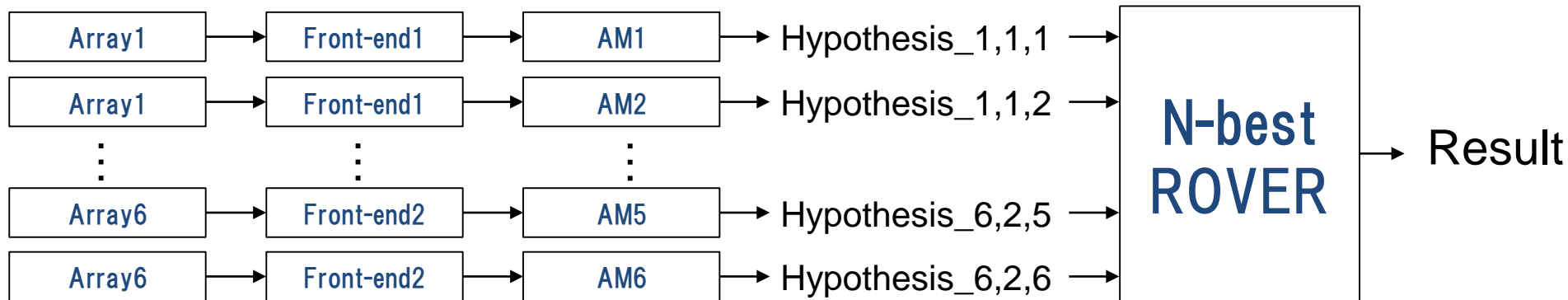
# Step-by-Step Improvements for Dev



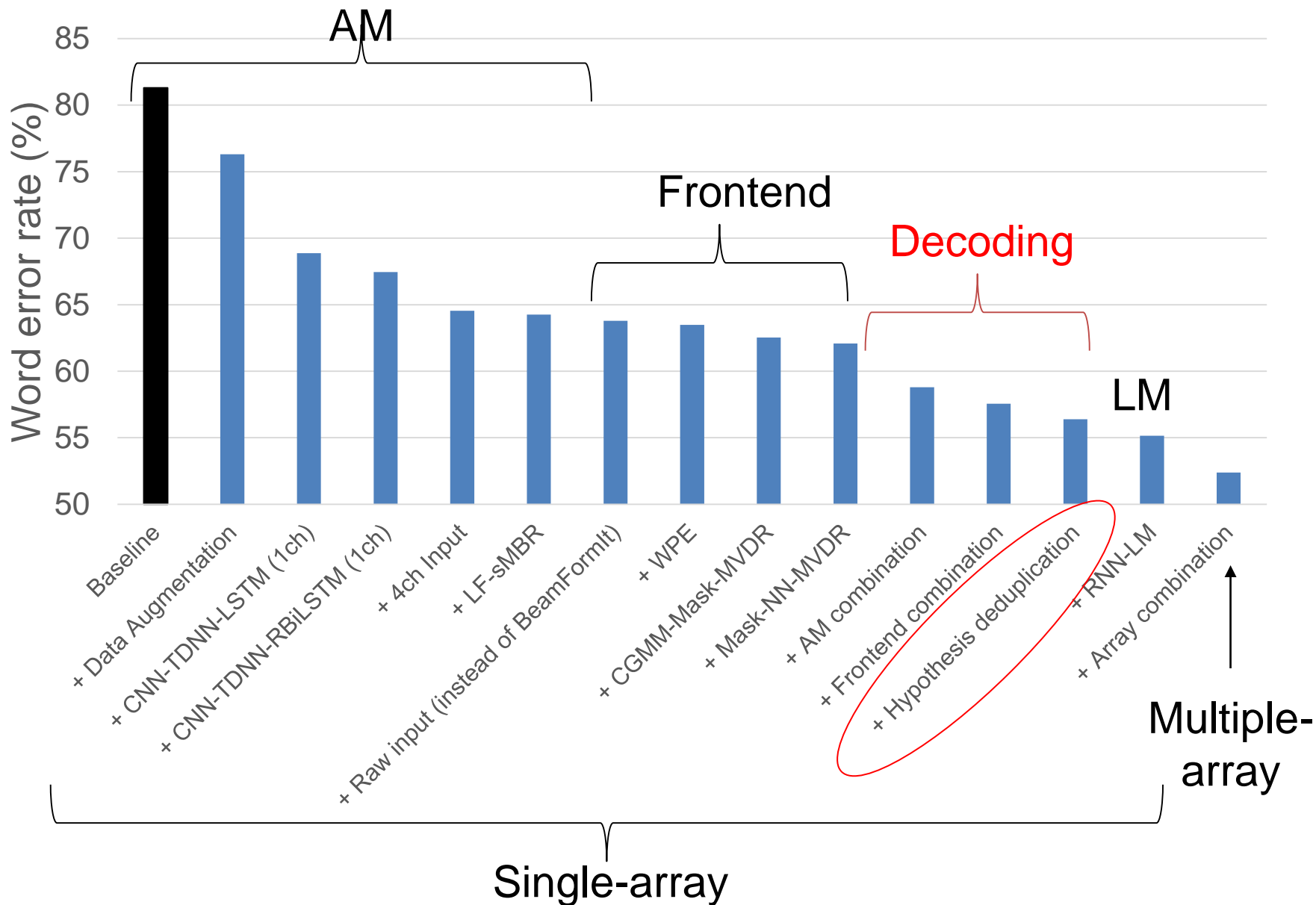
## ■ Hypotheses combination by N-best ROVER

- 6 AMs := CNN-TDNN- $\{LSTM, BiLSTM, RBiLSTM\}$   $\times$   $\{3500, 7000\}$  senones
- 2 Front-ends := Mask Network, CGMM
- 6 Arrays

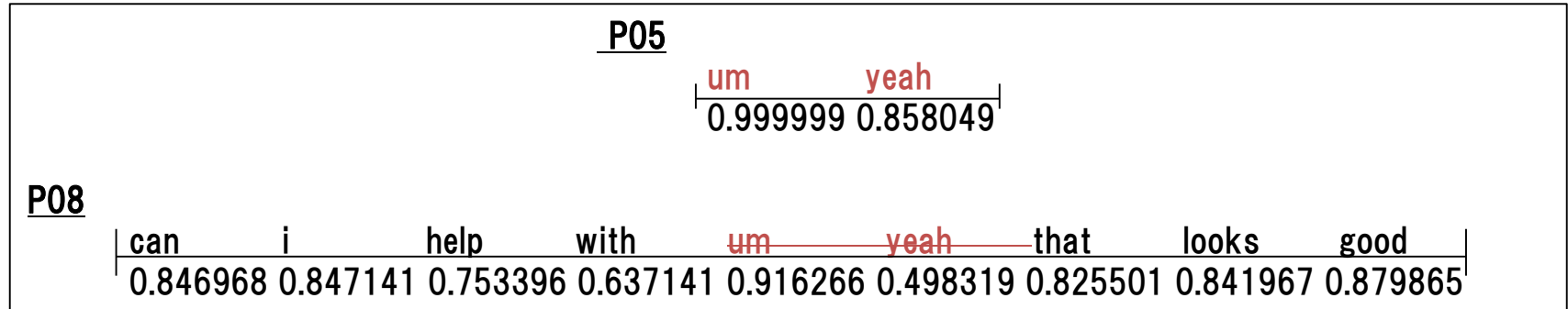
We didn't select array, instead combined hypotheses from each array.



# Step-by-Step Improvements for Dev



- Same words were sometimes recognized for overlapped utterances



- Duplicated words with lower confidence were excluded from the hypothesis.
  - HD was applied after ROVER, so precise time boundary could not be used. Minimum edit distance-based word alignment was used to detect word duplication.

WER (%) for Dev set

	without HD	with HD
Single-array	56.44	55.15
Multiple-array	53.69	52.38

(\* ) Results with RNN-LM

---

## Final results & Conclusion

WER (%)  
without RNN-LM / with RNN-LM

Our best eval result

Track	Session		Kitchen	Dining	Living	Overall
Single-array	Dev	S02	66.37 / 65.13	56.79 / 55.42	50.89 / 49.54	<b>56.40 / 55.15</b>
		S09	55.89 / 55.24	55.94 / 54.37	51.57 / 50.15	
	Eval	S01	59.42 / 57.62	44.18 / 41.81	63.85 / 62.33	<b>50.36 / 48.20</b>
		S21	52.11 / 49.68	42.14 / 39.78	46.71 / 44.59	
Multiple-array	Dev	S02	61.05 / 59.31	54.56 / 52.96	50.47 / 48.95	<b>54.00 / 52.38</b>
		S09	51.87 / 50.64	52.46 / 50.69	52.48 / 50.46	
	Eval	S01	59.82 / 57.01	43.59 / 41.22	62.28 / 60.67	<b>50.59 / 48.24</b>
		S21	54.70 / 51.59	44.12 / 42.17	45.95 / 43.82	

## WER (%) without RNN-LM / with RNN-LM

Track	Session		Kitchen	Dining	Living	Overall
Single-array	Dev	S02	66.37 / 65.13	56.79 / 55.42	50.89 / 49.54	<b>56.40 / 55.15</b>
		S09	55.89 / 55.24	55.94 / 54.37	51.57 / 50.15	
	Eval	S01	59.42 / 57.62	44.18 / 41.81	63.85 / 62.33	<b>50.36 / 48.20</b>
		S21	52.11 / 49.68	42.14 / 39.78	46.71 / 44.59	
Multiple-array	Dev	S02	61.05 / 59.31	54.56 / 52.96	50.47 / 48.95	<b>54.00 / 52.38</b>
		S09	51.87 / 50.64	52.46 / 50.69	52.48 / 50.46	
	Eval	S01	59.82 / 57.01	43.59 / 41.22	62.28 / 60.67	<b>50.59 / 48.24</b>
		S21	54.70 / 51.59	44.12 / 42.17	45.95 / 43.82	

- Array combination by ROVER worked well for dev, but not effective for eval set.
  - Why? Different types of rooms? Speaker-array distance?
- Anyway, better array combination methods should be pursued.

## ■ Our contributions

- Multiple data augmentation
- 4-ch AM with Residual BiLSTM
- Speaker adaptive mask estimation network / CGMM-based beamformer
- Hypothesis Deduplication
- Array combination by ROVER (found not effective for evaluation set)

## ■ Our results

- 48.2% WER for evaluation set
- 2<sup>nd</sup> ranked, with only 2.1 point difference to the best result

Thank you for your attention!



---

# Appendix

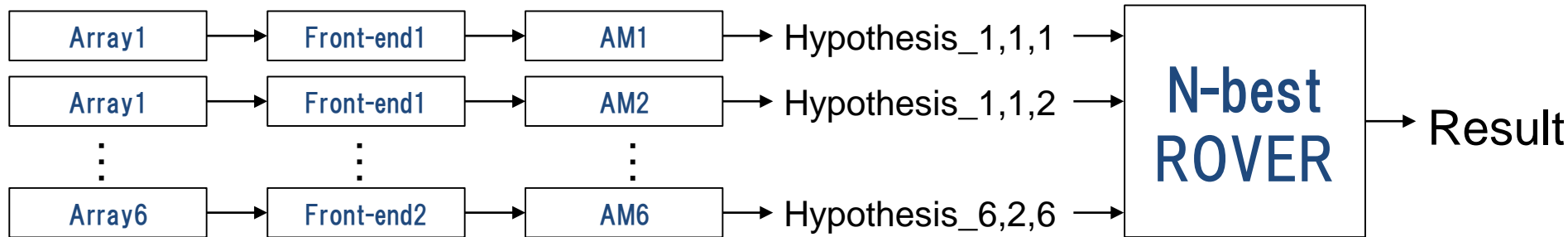
Model	Input	Training	Worn-Dev	Ref-Array-Dev
Baseline	1ch	LF-MMI	45.37	76.31
CNN-TDNN-LSTM	1ch	LF-MMI	39.22	68.87
CNN-TDNN-BiLSTM	1ch	LF-MMI	40.04	68.42
CNN-TDNN-RBiLSTM	1ch	LF-MMI	39.21	67.46
CNN-TDNN-RBiLSTM	4ch	LF-MMI	n/a	64.54
CNN-TDNN-RBiLSTM	4ch	LF-sMBR [1]	n/a	64.25

[1] Naoyuki Kanda, Yusuke Fujita, Kenji Nagamatsu, Lattice-free state-level minimum Bayes risk training of acoustic models, Interspeech 2018.

Front-end for 1ch input	Front-end for 4ch input	Dev
BeamFormIt (= Baseline)	Raw	64.28
Raw	Raw	63.79
WPE	WPE	63.49
CGMM-MVDR	WPE	62.53
Speaker adaptive mask NN-MVDR	WPE	62.09

## ■ Hypotheses combination by N-best ROVER

- 6 AMs := CNN-TDNN- $\{LSTM, BiLSTM, RBiLSTM\}$  x  $\{3500, 7000\}$  senones
- 2 Front-ends := Mask Network, CGMM
- 6 Arrays We didn't select array. Instead we combined hypotheses from each array.



WER (%) for Dev set

	AM	Array	Frontend	Dev
Single-array	1	1	MaskNet	62.09
	6	1	MaskNet	58.79
	6	1	MaskNet, CGMM	57.55
Multiple-array	6	6	MaskNet, CGMM	55.08

(\*) Results w/o RNN-LM

---

**Thank you**