

DA-IICT/IITV SYSTEM FOR THE 5th CHiME 2018 CHALLENGE

Ankur T. Patil¹, Maddala V. Siva Krishna², Mehak Piplani², Pulikonda Aditya Sai², Hardik B. Sailor¹, Hemant A. Patil¹

¹Speech Research Lab, Dhirubhai Ambani Institute of Information and Communication Technology (DA-IICT), Gandhinagar, Gujarat, India.

²Indian Institute of Information Technology (IIIT), Vadodra, Gujarat, India.

ankur_patil@daiict.ac.in, 201551045@iiitvadodara.ac.in, 201551072@iiitvadodara.ac.in, 201551013@iiitvadodara.ac.in, sailor_hardik@daiict.ac.in,

hemant_patil@daiict.ac.in

The 5th International Workshop on Speech Processing in Everyday Environments (CHiME 2018).



EVOLUTION OF CHiME CHALLENGES

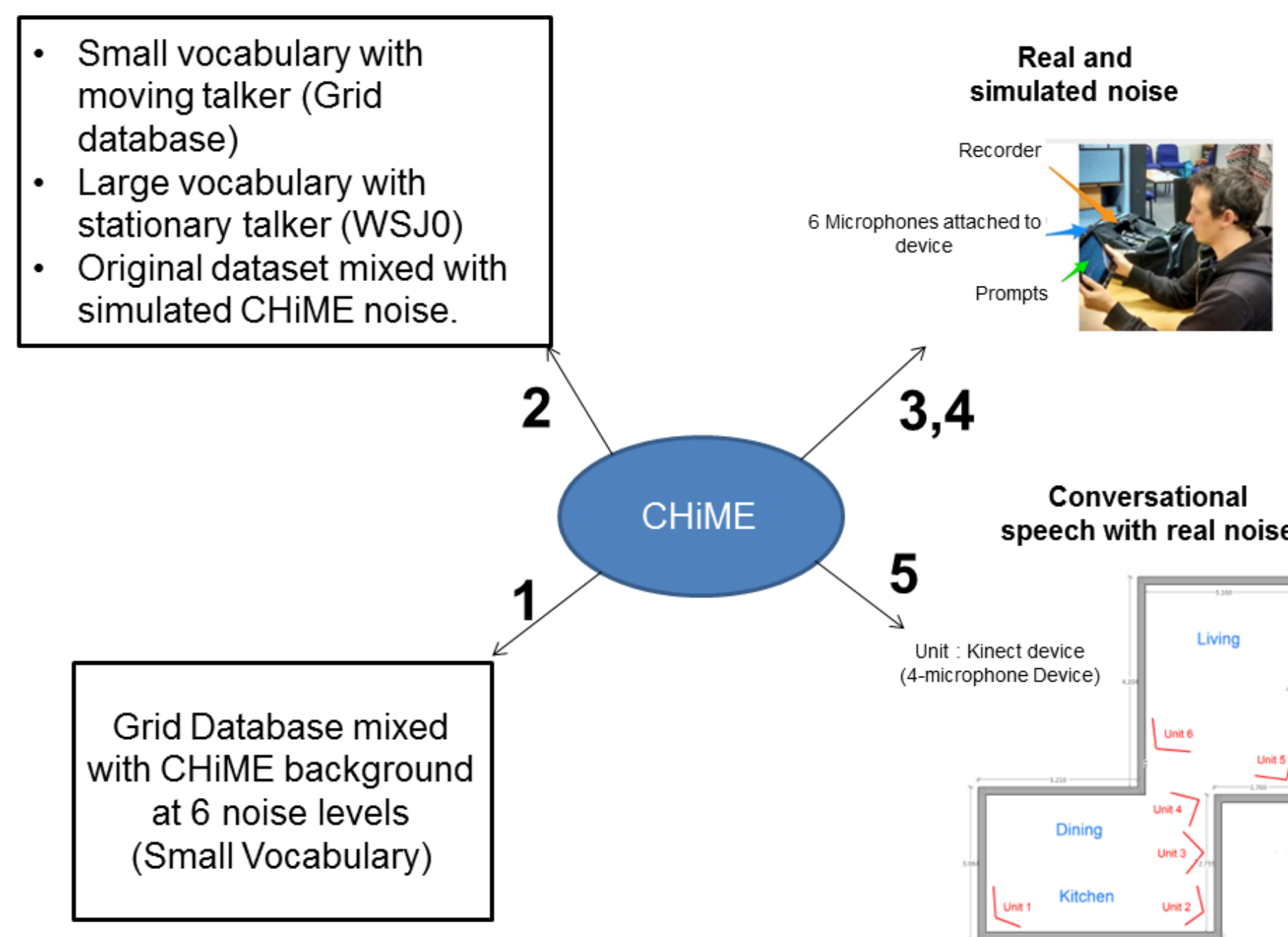


Figure 1: Evolution of database from CHiME-1 to CHiME-5.

- Top CHiME-1 and CHiME-2 systems came close to human performance.
- CHiME-3 best system : 5.8%WER
- CHiME-4 best system :2.2%WER
- CHiME-5 : The task of distant multi-microphone conversational speech recognition.
- Database of CHiME-5 Challenge [1]:
 - The scenario
 - * Recording of twenty separate dinner parties taking place in real homes. (natural conversational speech).
 - * Each party is of minimum 2 hours composed of 3 phases, corresponding to different locations namely, kitchen, dining, and living.
 - * Natural movement from one location to other location is allowed.
 - Audio
 - * 6 Kinect devices are strategically placed for each party.
 - * Each Kinect device has linear array of 4 microphone.
 - * Binaural microphone for each participant.

Table 1: CHiME-5 dataset

Dataset	Parties	Speakers	Hours	Utterances
Train	16	32	40:33	79,980
Dev	2	8	4:27	7,440
Eval	2	8	5:12	11,028

- Observations about database
 - Multiple speakers speaking at the same time.
 - Multisource noise is present.
 - Training text is not in context as it is conversational speech.

BLOCK DIAGRAM OF PROPOSED ASR SYSTEM

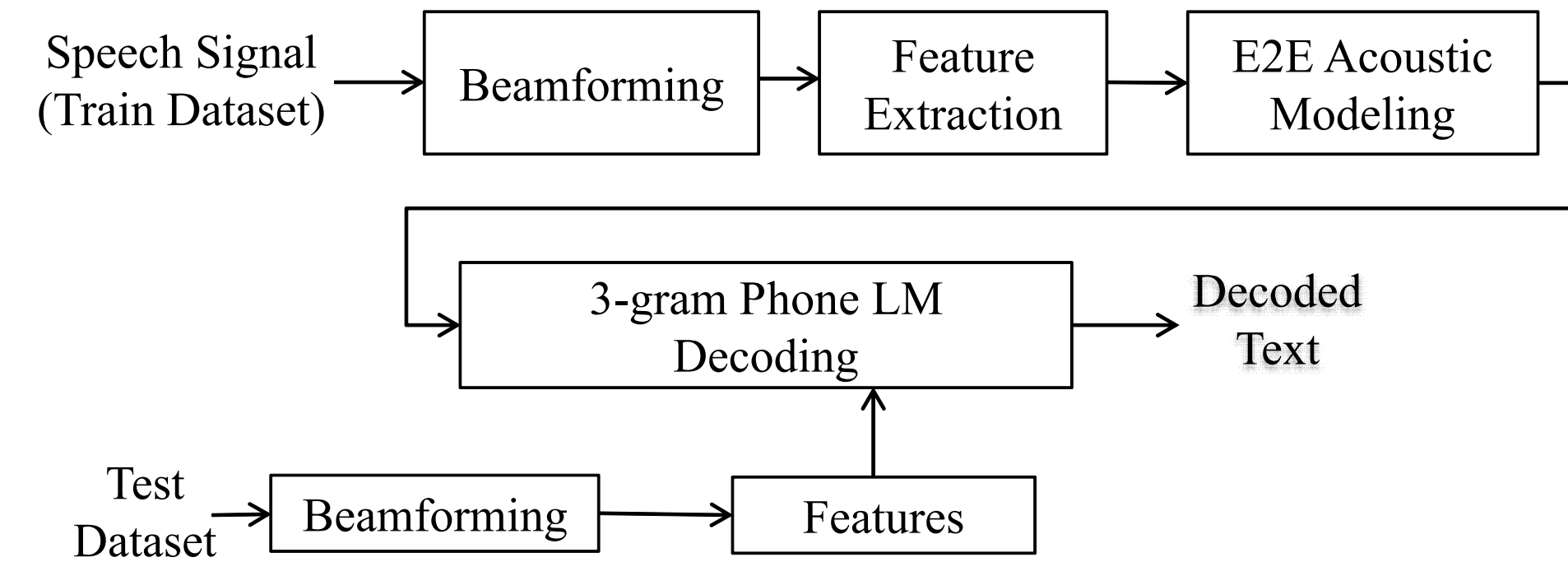


Figure 2: The comparison of ConvRBM filterbank scales with standard auditory frequency scales.

PROPOSED SYSTEM ARCHITECTURE DETAILS

- E2E using LF-MMI [2] :
 - E2E : No prior alignments required from initial models, hence called flatstart.
 - Context dependency trees are not required.
 - Monophones or full left biphones are used for language modeling.
 - Entire training process in one stage.
- Speech enhancement : delay-and-sum beamformer. Applied on training and testing data [3].
- Feature extraction : Mel-frequency spectral coefficient (MFSC) and Power normalised spectral coefficient (PNSC) [4].
- Decoding is performed using 3-gram LM followed by RNNLM.
- Acoustic modeling using DNN:
 - Time Delayed Neural Network (TDNN): 8 hidden layers (2048 neurons per layer) [5].
 - Long-Short Term Memory along with TDNN (TDNN-LSTM) : 7-TDNN layers and 3-LSTM layers (1024 neurons per layer) [6].

SYSTEM SPECIFICATIONS

Table 2: E2E LF-MMI ASR System Specification

System	DNN Model	Features	Training Data
S1	TDNN	MFSC	Full Data *
S2	TDNN	MFSC	Enhanced Speech **
S3	TDNN-LSTM	MFSC	Enhanced Speech
S4	TDNN	PNSC	Enhanced Speech
S5	TDNN-LSTM	PNSC	Enhanced Speech

* Full data : Audio data from all microphone channels and from all devices are used for training.

** Enhanced Speech : Audio data from each device is applied to beamformer. This enhanced speech is used for training.

- S1 : MFSC is 40-D (filterbank coefficients)
- S2-S5 : MFSC is 120-D (filterbank coefficients + Δ + $\Delta\Delta$)

RESULTS-1

Table 3: Results of various E2E system and their combinations [7] using 3-gram LM per session and location together with the overall % WER on development set.

System	Session	Kitchen	Dining	Living	Overall
S1	S02	88.30	83.22	80.79	83.85
	S09	84.92	85.15	81.10	
S2	S02	88.62	83.01	80.54	83.75
	S09	84.42	85.51	80.88	
S3	S02	90.23	84.94	82.49	84.79
	S09	84.50	83.69	81.46	
S4	S02	89.13	85.65	83.93	85.17
	S09	84.82	83.95	81.97	
S5	S02	93.99	91.44	87.31	89.30
	S09	87.35	88.55	86.14	
SC-1	S02	85.89	79.88	77.49	80.14
	S09	79.79	79.44	77.06	
SC-2	S02	84.35	77.60	75.29	78.69
	S09	79.24	78.87	76.49	
SC-3	S02	84.21	78.46	75.64	78.63
	S09	78.23	78.15	76.27	

- System combinations:

- SC-1 : S2 \oplus S3
- SC-2 : S1 \oplus S2 \oplus S3
- SC-3 : S1 \oplus S2 \oplus S3 \oplus S4

Table 4: Results of various E2E system and their combinations using RNNLM per session and location together with the overall % WER on development set.

System	Session	Kitchen	Dining	Living	Overall
S1	S02	88.07	83.02	80.50	83.61
	S09	84.53	84.61	81.32	
S2	S02	88.56	83.09	80.10	83.40
	S09	83.38	84.68	80.94	
S3	S02	89.97	85.19	82.52	84.65
	S09	83.75	83.71	81.34	
S4	S02	89.50	86.65	84.21	85.64
	S09	85.04	84.75	82.47	
S5	S02	94.51	92.12	88.47	90.11
	S09	87.99	89.75	86.77	
SC-1	S02	86.47	80.43	77.97	80.64
	S09	79.69	80.41	77.76	
SC-2	S02	84.91	78.33	76.19	79.04
	S09	78.93	78.35	76.37	
SC-3	S02	84.79	79.27	76.54	79.36
	S09	79.12	78.68	76.83	

Table 5: Results of our best system (SC-3) per session and location together with the overall %WER on evaluation set.

System	Session	Kitchen	Dining	Living	Overall
SC-3	S01	82.65	73.38	84.68	76.42
	S21	79.49	72.55	69.82	

RESULT-2

Table 6: Comparison of proposed system combination with baseline systems on development set.

System	Session	Kitchen	Dining	Living	Overall
LF-MMI TDNN*	S02	87.3	79.5	79	81.3
	S09	81.6	80.6	77.6	
ESPnet E2E **	S02	-	-	-	94.7
	S09	-	-	-	
S2	S02	88.62	83.01	80.54	83.75
	S09	84.42	85.51	80.88	
SC-3	S02	84.21	78.46	75.64	78.63
	S09	78.23	78.15	76.27	

SUMMARY AND CONCLUSIONS

- Developed E2E system using LF-MMI as objective function.
- Performance in Kitchen is very poor due to presence of more multi-source noise.
- RNNLM rescoring do not show any improvement because of conversational speech.

ACKNOWLEDGMENTS

- CHiME-5 challenge organizers for providing database.
- NVIDIA for the hardware grant of TITAN X GPU.

REFERENCES

1. J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth CHiME Speech Separation and Recognition Challenge: Dataset, task and baselines," in INTERSPEECH 2018, Hyderabad, India, Sep. 2018.
2. H. Hadian, H. Sameti, D. Povey, and S. Khudanpur, "End-to-end speech recognition using lattice-free MMI," in INTERSPEECH 2018, Hyderabad, India, Sep. 2018.
3. M. Wolfel and J. McDonough, Distant Speech Recognition. John Wiley & Sons, 2009.
4. C. Kim and R. M. Stern, "Power-normalized cepstral coefficients (PNCC) for robust speech recognition," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 24, no. 7, pp. 1315-1329, July 2016.
5. V. Peddinti, D. Povey, and S. Khudanpur, "A time delay neural network architecture for efficient modeling of long temporal contexts," in INTERSPEECH, Dresden, Germany, 2015, pp. 2440-2444.
6. V. Peddinti, Y. Wang, D. Povey, and S. Khudanpur, "Low latency acoustic modeling using temporal convolution and LSTMs," IEEE Signal Processing Letters, vol. 25, no. 3, pp. 373-377, March 2018.
7. H. Xu, D. Povey, L. Mangu, and J. Zhu, "Minimum Bayes risk decoding and system combination based on a recursion for edit distance," Computer Speech & Language, vol. 25, no. 4, pp. 802-828, 2011.