



Towards a speaker diarization system

for the CHiME 2020 dinner party transcription

<u>Christoph Boeddeker</u>¹, Tobias Cord-Landwehr¹, Jens Heitkaemper¹, Cătălin Zorilă², Daichi Hayakawa³, Mohan Li², Min Liu⁴, Rama Doddipatla², Reinhold Haeb-Umbach¹

¹Paderborn University, Department of Communications Engineering, Paderborn, Germany
²Toshiba Cambridge Research Laboratory, Cambridge, United Kingdom
³ Toshiba Corporation Corporate R&D Center, Kawasaki, Japan
⁴ Toshiba China R&D Center, Beijing, China









- Guided Source Separation (GSS)
 - Designed for CHiME-5
 - Baseline in CHiME-6 Track 1
 - ▶ Here: From Track 1 to Track 2 (i.e. drop human annotations)
- Acoustic models
- Results on Track 2
- Ongoing work
 - PIT Neural Speaker Diarization
 - Spatial features from Spatial Mixture Model (SMM)



CHiME-6 Dataset



For a complete description of the dataset see presentation:

S. Watanabe, M. Mandel, J. Barker, E. Vincent

"Overview of the 6th CHiME Challenge"

or publication:

S. Watanabe, M. Mandel, J. Barker, E. Vincent "CHiME-6 Challenge:Tackling Multispeaker Speech Recognition for Unsegmented Recordings" arXiv preprint arXiv:2004.09249 (2020)













CHiMe WORKSHOP Spatial Mixture Model





CHiMe WORKSHOP Spatial Mixture Model





CHiMe WORKSHOP Spatial Mixture Model











Guided Source Separation (GSS)









CHiMe WORKSHOP Implementation



- Guide replacement: Already implemented for refinement of human annotations
- Utterance boundary replacement: Missing
 - ► When forgotten, may increase or decrease WER
- After the Workshop: Publish code to support RTTM files for both





Acoustic model and results in Track 2 with baseline diarization system

	Acoustic model			Enh. in test	WER in $\%$			
		Training Data	NN layers			•	DEV	EVAL
			CNN	TDNNF	•			
le	→ AM ₀	unproc. & reverb.		15		WPE+BFIt	81.92	76.37
	AM ₀	unproc. & reverb.		15		GSS	78.12	73.06
ell	AM_1	worn & WPE	10	9		GSS	76.44	72.04
Dag	AM_2	worn & GSS	10	9		GSS	74.74	71.27
V	AM_3	worn & GSS	10	9	DT	GSS	74.67	70.55
ack	AM_4	worn & GSS	40 RESNET	9		GSS	74.05	70.47
	AM_5	worn & GSS	40 RESNET	9	DT	GSS	74.73	70.14
	AM_{0-5}					GSS	73.50	68.96

Simular to C. Zorila, M. Li, D. Hayakawa, M. Liu, N. Ding and R. Doddipatla "Toshiba's speech recognition system for the CHiME 2020 Challenge" but used on Track 2 here. DT: Discriminative Training, TDNNF: Factorized Time Delay Neural Network





Acoustic model and results in Track 2 with baseline diarization system

Acoustic model						Enh. in test	WER in $\%$	
		Training Data	NN layers				DEV	EVAL
			CNN	TDNNF				
ack z daseline	→ AM ₀	unproc. & reverb.		15		$WPE{+}BFIt$	81.92	76.37
	AM ₀	unproc. & reverb.		15		GSS	78.12	73.06
	AM_1	worn & WPE	10	9		GSS	76.44	72.04
	AM_2	worn & GSS	10	9		GSS	74.74	71.27
	AM_3	worn & GSS	10	9 DT		GSS	74.67	70.55
	AM_4	worn & GSS	40 RESNET	9		GSS	74.05	70.47
	AM_5	worn & GSS	40 RESNET	9	DT	GSS	74.73	70.14
	AM_{0-5}					GSS	73.50	68.96

Simular to C. Zorila, M. Li, D. Hayakawa, M. Liu, N. Ding and R. Doddipatla "Toshiba's speech recognition system for the CHiME 2020 Challenge" but used on Track 2 here. DT: Discriminative Training, TDNNF: Factorized Time Delay Neural Network



CHiMe WORKSHOP Ongoing work



Track 2: Ongoing work

• Not conforming to challenge regulations



Proposed in: Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, S. Watanabe "End-to-end neural speaker diarization with permutation-free objectives.", Interspeech (2019)



Onset and offset from simple Viterbi decoding on an HMM as used in the SAD baseline system



Christoph Boeddeker

Towards a speaker diarization system for the CHiME 2020 dinner party transcription



Towards a speaker diarization system for the CHiME 2020 dinner party transcription



Towards a speaker diarization system for the CHiME 2020 dinner party transcription



PIT Neural Speaker Diarization (NSD) results

CHiM

	Mel bins	SMM spat. features	out	D	ER	W		
Diarization			Dropo	TRAIN	DEV+EVAL	DEV	EVAL	
Baseline ¹	-	-	-	-	60.87	77.49	71.92	
2 BLSTM	80	No	No	34.28	60.09	72.68	71.25	
2 BLSTM	80	Yes	No	2.52	53.75	70.05	69.50	
1 BLSTM	24	Yes	0.25	7.03	48.89	63.82	63.47	5
Oracle	-	-	-	-	0	47.67	49.54 ←	—

Enhancement in test: GSS

Track 1 baseline







However: Not conforming to challenge regulations

- Applied to 40 s segments
- Intra segment speaker permutation solved by oracle

Why didn't we address this?

- Focused on reduce overfitting first
- Overfitting source:
 - Training dataset has only 32 speakers
 - ► X-Vector baseline system: more than 7000 speakers from VoxCeleb







- Guided Source Separation: Replaced human annotations with diarization output
- Submitted to Track 2: Baseline Diarization System + GSS + 6 AMs
 - ► Category A: DEV: 73.50 %, EVAL 68.96 %
 - Category B: DEV: 73.05 %, EVAL: 68.45 %
- Neural Speaker Diarization (NSD) trained on CHiME-6 data
 - Proposed spatial feature from spatial mixture model
 - Overfitting problem
 - ToDo: Segment permutation problem solver

Thank you for your attention! Questions?

Contact me via email: boeddeker@nt.upb.de







- Guided Source Separation: Replaced human annotations with diarization output
- Submitted to Track 2: Baseline Diarization System + GSS + 6 AMs
 - ► Category A: DEV: 73.50 %, EVAL 68.96 %
 - Category B: DEV: 73.05 %, EVAL: 68.45 %
- Neural Speaker Diarization (NSD) trained on CHiME-6 data
 - Proposed spatial feature from spatial mixture model
 - Overfitting problem
 - ToDo: Segment permutation problem solver

Thank you for your attention! Questions?

Contact me via email: boeddeker@nt.upb.de