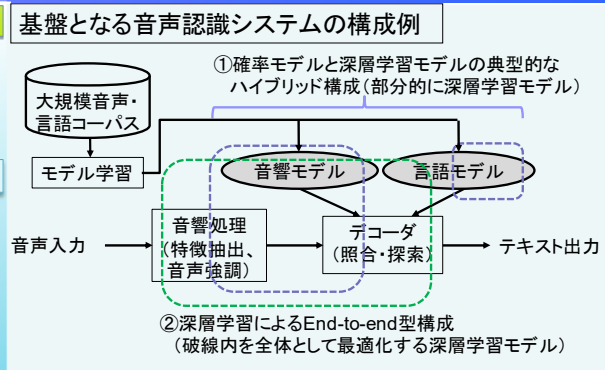
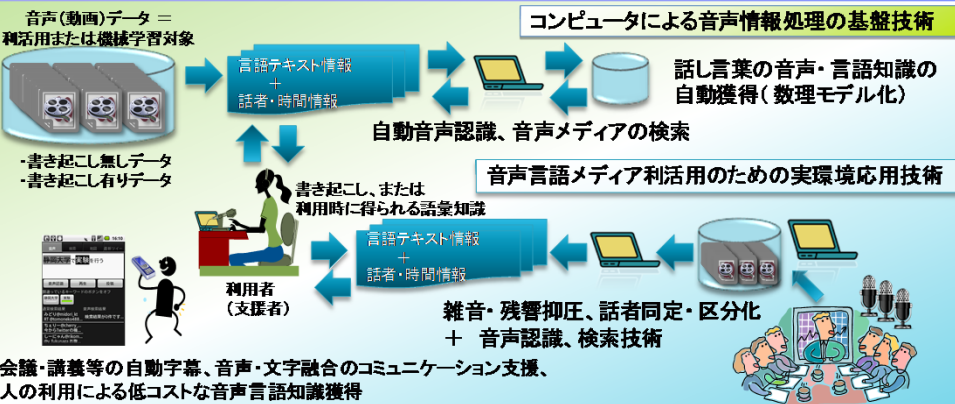
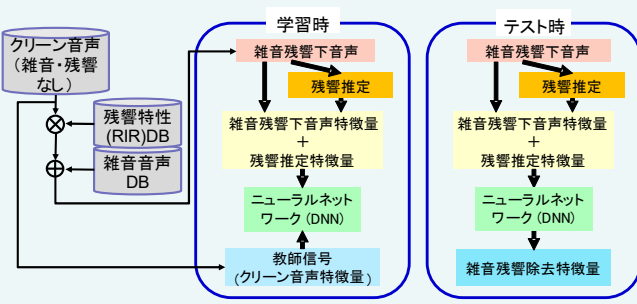


# 実環境向けの音声認識とその応用

工学領域 数理システム工学系列 准教授 甲斐充彦



## 深層学習を用いた単一マイク遠隔音声認識技術



要素技術  
 ・音声・言語処理のモデル(確率モデル、深層学習モデル)  
 ・機械学習  
 ・パターン認識  
 など

## 「話し言葉コーパス(CSJ)」によるモデル構築例

CSJ評価用データセット(10講演)での認識精度

音響モデル	単語誤り率(文字誤り率)
従来の確率モデル(GMM-HMM) (話者適応学習、話者適応の特徴変換)	15.1%
深層学習ハイブリッドモデル(DNN-HMM) (話者適応の特徴変換なし)	12.4%
深層学習ハイブリッドモデル(TDNN-HMM) (話者特徴の補助入力あり)	10.9% (9.0%)
End-to-endモデル(CTC/Attention-hybrid)	(8.2%)

✓ 多種の人工データ生成による雑音残響抑圧モデルの学習 (音声認識システム向けの汎用的モデルとして最適化)  
 ✓ 未知の残響環境に適応する仕組みを学習

・Y. Ueda, L. Wang, A. Kai, EURASIP Journal on Advances in Signal Processing (2015)  
 ・上田, 王, 甲斐, 電子情報通信学会技術報告(SP) (2015.12)

■学習データ  
 □クリーン音声: 英語読み上げ音声(WSJCAM0) 約18時間、話者92名  
 □人工雑音・残響用: 室内インパルス応答 ( $T_{60} = 0.1s \sim 0.8s$ ), 雑音(6種類)

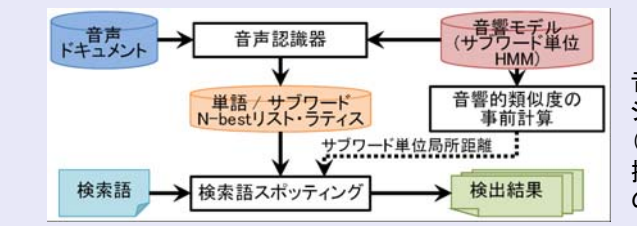
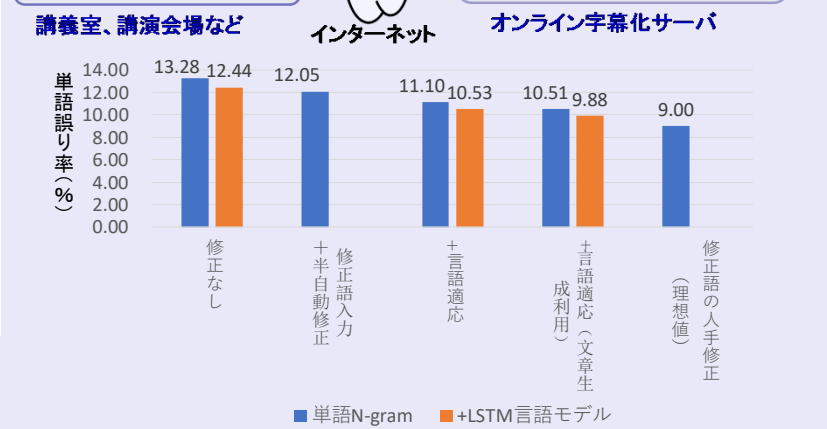
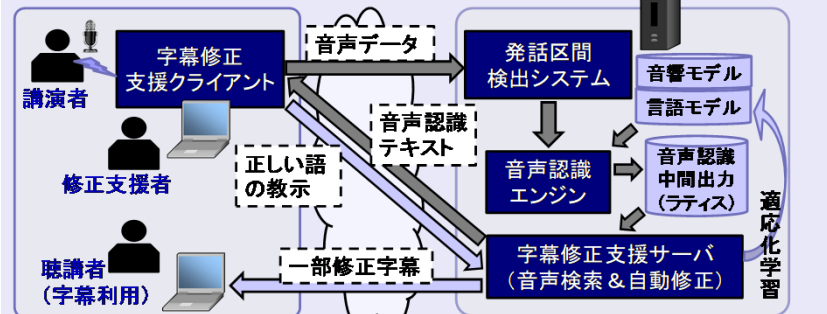
■テストデータ

	残響時間			雑音 SN比	話者-マイク間距離	
	room1	room2	room3		near	far
人工環境	0.25s	0.5s	0.7s	20dB	50cm	200cm
実環境	0.7s	-	-	-	100cm	250cm

■実験結果  
 ■音響モデル: DNN-HMM & SGMM, 評価指標: 単語誤り率(WER) [%]

フロントエンド処理	人工環境						実環境			
	room1		room2		room3		平均	room1		平均
	near	far	near	far	near	far		near	far	
CMVN	4.90	5.39	6.33	11.77	7.68	15.57	8.61	43.40	42.98	43.19
MSLP	4.71	5.18	5.95	9.95	7.32	14.45	7.93	35.52	36.09	35.81
従来DAE	4.79	5.40	5.64	9.00	7.06	10.85	7.12	30.02	31.09	30.56
Reverberation-aware DAE	<b>4.54</b>	<b>5.05</b>	<b>5.37</b>	<b>7.62</b>	<b>6.50</b>	<b>9.40</b>	<b>6.41</b>	<b>26.38</b>	<b>27.28</b>	<b>26.83</b>

## オンライン字幕修正支援システムと要素技術



音声検索システム (字幕修正支援システムの要素技術)  
 ・Y. Terada, K. Tamiya, and A. Kai, "Investigation of Efficient Semi-automatic Correction Method Using STD for Automatic Captioning," Proc. IEEE GCCE 2017, Oct 2017.  
 ・寺田, 塚本, 甲斐, "講演音声認識の修正語のオンライン教示による半自動的な修正手法と語彙適応の併用の効果", 日本音響学会講演論文集, 2019.9

keyword: 話し言葉音声認識、雑音・残響下音声認識、音声検索、深層学習、適応学習