

雑音下音声認識評価ワーキンググループ活動報告： 認識に影響する要因の個別評価環境

北岡 教英¹ 山田 武志² 滝口 哲也³ 柘植 覚⁴ 山本 一公⁵
宮島千代美¹ 西浦 敬信⁶ 中山 雅人⁶ 傳田 遊亀⁶ 藤本 雅清⁷
田村 哲嗣⁸ 黒岩 眞吾⁴ 武田 一哉¹ 中村 哲⁹

¹ 名古屋大学 ² 筑波大学 ³ 神戸大学 ⁴ 徳島大学 ⁵ 信州大学

⁶ 立命館大学 ⁷ NTT CS 基礎研 ⁸ 岐阜大学 ⁹ NiCT/ATR

E-mail: kitaoka@slp.ics.tut.ac.jp

あらまし 雑音下の音声認識の性能向上は音声認識実用化のために急務である。これまでに数多くの研究が行われてきており、これらの手法を客観的に比較評価できる標準評価基盤の構築を目的として、2001年10月、情報処理学会音声言語情報処理研究会の下に雑音下音声認識評価ワーキンググループを組織した。本稿ではこれまでの標準評価基盤 CENSREC シリーズを振り返り、今年度新たに配付した CENSREC-1-C の概要と位置づけを述べる。さらに、今後どのような方針で新たな評価基盤を設計・構築・配付するののかについての考えを述べる。

キーワード 音声認識, 雑音, 評価基盤, CENSREC

Progress Report of SLP Noisy Speech Recognition Evaluation WG: Individual evaluation framework for each factor affecting recognition performance

Norihide KITAOKA¹, Takeshi YAMADA², Tetsuya TAKIGUCHI³, Satoru TSUGE⁴,
Kazumasa YAMAMOTO⁵, Chiyomi MIYAJIMA¹, Takanobu NISHIURA⁶,
Masato NAKAYAMA⁶, Yuki DENDA⁶, Masakiyo FUJIMOTO⁷, Satoshi TAMURA⁸,
Shingo KUROIWA⁴, Kazuya TAKEDA¹, and Satoshi NAKAMURA⁹

¹ Nagoya Univ. ² Univ. of Tsukuba ³ Kobe Univ. ⁴ Univ. of Tokushima ⁵ Shinshu Univ.

⁶ Ritsumeikan Univ. ⁷ NTT CS research Lab. ⁸ Gifu Univ. ⁹ NiCT/ATR

E-mail: kitaoka@slp.ics.tut.ac.jp

Abstract Performance improvement of noisy speech recognition is urgent for practical use of speech recognition. Many researchers have been trying to overcome this problem. We organized a working group under Special Interest Group of Spoken Language Processing in Information Processing Society of Japan, to develop evaluation frameworks of noisy speech recognition to compare many methods for processing of noisy speech. In this paper, we first review the series of CENSREC (Corpus and Environment of Noisy Speech REcognition) and then introduce the CENSREC-1-C, the newest CENSREC. Finally we describe the road-maps of future CENSRECs.

Key words Speech recognition, noise, evaluation frameworks, CENSREC.

1. はじめに

音声認識の性能が飛躍的に向上し実用化が進められている中、

実環境のさまざまな雑音下で頑健に動作するシステムはまだまだ実現されていない。雑音下の音声認識の研究は多くなされてきており、個々の実験条件の下では十分な性能を得ることができ

ているものも多い。そこで、これらの多くの手法を客観的に比較評価できれば、例えば開発者が意中の目的にふさわしい手法を見付けることが可能となると考えられる。

これまでに、こうした問題に対して欧州の AURORA 研究プロジェクトでは、雑音抑圧手法の標準評価コーパスと HTK を利用した標準認識スクリプト、および標準評価用 Excel シートをセットとした評価基盤を配布した。これらは、TI digit に雑音を付与した連続数字音声認識タスクの AURORA-2、自動車内で実際に発声した連続数字認識タスクの AURORA-3、Wall Street Journal タスクをベースとした雑音下大語彙連続音声認識タスクの AURORA-4 からなる。[1], [2]。この中でも AURORA-2 はコンパクトで認識評価が容易なことから国際会議論文でもよく利用されている。

我々も上記の活動に追従するように 2001 年 10 月に雑音下音声認識評価ワーキンググループを組織し、日本語における同様の評価基盤の構築を開始した。さらにそれを拡充し、雑音下音声認識手法をさまざまな側面から評価できる評価基盤群を構築すべく活動してきた [3], [4]。その一環として、2006 年 9 月に雑音下音声区間検出評価基盤である CENSREC-1-C の配布を開始した。

本稿では、これまでのワーキンググループの活動を振り返るとともに、基盤群構築のための考え方を述べ、今後の展開について述べる。

2. 情報処理学会 音声言語情報処理研究会 雑音下音声認識評価ワーキンググループ

本ワーキンググループ (WG) では、雑音下音声認識の標準評価基盤として、どのようなものを提供すればよいかを課題として議論を重ねてきた。そして、AURORA プロジェクトの配布形態にならない、共通学習データ/テストデータ、標準学習/認識スクリプトおよび評価用 Excel シートを含むいくつかの評価基盤を配布している。これらは CENSREC (Corpus and Environment for Noisy Speech REcognition) と呼び、CENSREC-1.2.3 の作成・配布を行ってきた。

CENSREC-1 は AURORA-2 の日本語版であり、連続数字発声に種々の雑音を付加したデータによる実験環境である。しかし一般に、雑音付加によるシミュレーション環境は実際の雑音環境をある程度模擬しているものの、実環境下ではそれ以外の要因も作用していると考えられている。CENSREC-2 では、この違いを検討できるように、自動車内の実環境下で同様の内容を実際に発声して収録したデータを用いている。さらに、タスクによる影響も考えられるため、その検討のために、CENSREC-2 を孤立単語認識タスクに変更したものが CENSREC-3 である。すなわち、雑音下音声認識の困難さを「雑音の違い」「シミュレーション/実環境の違い」および「タスクの違い」の軸上で比較できる環境群となっている。以下に、それぞれの説明を行う。

2.1 CENSREC-1/AURORA-2J

欧州 AURORA プロジェクトの作成した雑音環境下連続英語数字音声認識タスクの共通評価フレームワークである AURORA-2 の日本語版である。2003 年 7 月の配布開始から、100 部を越

表 1 CENSREC-1 の雑音環境

	加法性雑音	フィルタ特性
Set A	Subway, Babble, Car, Exhibition	G.712
Set B	Restaurant, Street, Airport, Station	G.712
Set C	Subway, Street	MIRS

える部数を配布し、各種研究発表で利用されている。

学習/テストデータとともに学習認識スクリプトとベースライン性能を提供し手軽に実験を開始できる。タスクが連続数字認識とシンプルで、加算性雑音の影響とその対策に絞っており、逆にこれからの発展がさまざまな考えられるという意味で、本 WG の評価基盤の原点であるといえる。

2.1.1 データ構成

1 から 7 桁の連続数字であり、話者数、男女構成などは AURORA-2 と同一で、発声リストは同じものを日本語に翻訳して使用した。10 種類の数字のうち、英語の“ゼロ”と“オー”に対応して 0 のみが“Z(ぜろ)”, “O(まる)”の 2 種類がある。

学習データは、clean training (クリーン音声によるモデル学習)、multicondition training (雑音重畳音声による学習) 用にそれぞれ 110 名、8,140 発話 (男女 55 名、4,220 発話ずつ) 用意されている。Multicondition training の場合のみ 4 種類の雑音 (Subway, Babble, Car, Exhibition) を 5 種類の SNR レベル (clean, 20dB, 15dB, 10dB, 5dB) で重畳した音声 (各雑音・SNR で 422 発話ずつの学習データとなる) を用いて学習を行なう。チャンネルフィルタとして G.712 の中で規定されているフィルタを用いている。

テストデータは男女 104 名による 4,004 発話を 4 セットに等分した 1001 発話に表 1 のいずれかの雑音を 7 種類の SNR レベル (clean, 20dB, 15dB, 10dB, 5dB, 0dB, -5dB) で重畳したものである。

2.1.2 評価方法

HTK を用いた HMM の学習および認識スクリプトが添付されている。HMM は 10 数字 (11 種類) の単語モデルに無音 (sil, sp) を加えた 13 モデルで、数字 HMM は 18 状態 (2 状態はダミー)、出力確率は対角正規分布の 20 混合 (無音モデルは 36 混合) である。特徴パラメータは、MFCC (12 次元 + Δ + $\Delta\Delta$ + 対数パワー (1 次元) + Δ パワー + $\Delta\Delta$ パワー) の計 39 次元である。

この認識スクリプトを動作させることによりベースライン性能が得られる。これは既に入力された Excel シートが用意されており、独自の手法で雑音対策した結果を入力することによりベースライン性能からの相対的な性能改善が計算される。

詳細は文献 [5] を参照されたい。

2.2 CENSREC-2

CENSREC-1 と同様の連続数字データを自動車内で発話したデータに基づいており、CENSREC-1 の雑音加算シミュレーションに対し実収録データへの効果を確認できる。本評価環境は、2005 年 12 月より配布を開始している。

2.2.1 収録データ

名古屋大学の実験車両を用いて収録を行った。収録時には 5

本(接話1本, ダッシュボード上2本, 天井3本)のマイクロフォン(SONY ECM77B)を用いており, CENSREC-2では, 接話マイクロホンと天井の1本(遠隔マイクロホンと呼ぶ)で収録された音声を用いた[10]. 収録条件は, 3種類の走行速度(アイドリング, 低速(市街地)走行, 高速走行)と4種類の車内環境, エアコン On, オーディオ On, 窓あけ)を組み合わせた11種類の環境である. 総発話者数は104名(学習データ: 73名, 評価データ: 31名)で, 収録音声の総数は17,651発話である.

2.2.2 評価環境の設計

前節のような種々の条件で収録された音声データを組み合わせ, そのときの学習データと評価データの収録条件(マイク種別, 収録環境)の一致度合により, 以下の4種類の音声認識評価条件(Condition 1~4)が設定されている.

Condition 1: マイク種別, 収録環境がともに一致

Condition 2: マイク種別が一致, 収録環境が相違

Condition 3: マイク種別が相違, 収録環境が一致

Condition 4: マイク種別, 収録環境がともに相違

CENSREC-2に関する詳細は, 文献[6]を参照されたい.

2.3 CENSREC-3

CENSREC-3は, CENSREC-1に続いて, 実音響環境での発話を対象として設計されたものであり, 実走行車内での孤立単語音声認識の評価環境を提供する. このデータおよび評価環境は2005年2月より配布を行っている

2.3.1 収録データ

収録環境はCENSREC-2と同じで, 接話マイクロホンと遠隔マイクロホンの2種類を用いて3種類の走行速度(アイドリング, 低速(市街地)走行, 高速走行)と, 6種類の車内環境(通常走行, ハザード On, エアコン(Low), エアコン(High), オーディオ On, 窓開)を組み合わせた16種類の環境[11]である. 内容は, カーナビゲーション等で使用されることを想定した50個のコマンドワードである. これらの18名(男性8名, 女性10名)による14,216発話が収録されている

2.3.2 評価環境の設計

学習データと評価データの環境の組み合わせにより, 欧州AURORAプロジェクトのAURORA-3のWell-matched condition(WM), Moderate-mismatched condition(MM), High-mismatched condition(HM)に準じた6種類の評価環境(Condition 1~6)を設定しており, 以下のような対応となっている.

Condition 1, 2, 3 マイク種別, 収録環境がともに一致(WM)

Condition 4 マイク種別が一致, 収録環境が相違(MM)

Condition 5, 6 マイク種別, 収録環境(の一部)が相違(HM)

これらの学習環境条件の元で収録された音素バランス文を用いて学習したWord Internal Triphone HMMで認識を行うスクリプトと, ベースラインからの性能向上が評価できるExcelシートを用意している. また評価時のバックエンド部分の変更レベルにより他手法との比較が可能となるように評価カテゴリを設定している.

詳細は, 文献[7]を参照されたい.

3. 雑音下音声区間検出評価基盤 CENSREC-1-C [8]

これまでは認識対象データに含まれる種々の要因の音声認識性能に及ぼす影響を比較検討できるという観点で基盤群を構築してきた. しかし実環境での利用の状況をより広範囲に鑑みて, それ以外の評価すべき要因が存在することがわかってきた.

そのひとつに音声区間検出手法(Voice Activity Detection; VAD)がある. 近年, 音声認識の性能向上を目的として雑音下のVADの研究が盛んである. 正確なVADは音声区間のみを音声認識することを可能にし, 非音声区間からの湧き出し誤りや音声区間からの脱落誤りを減少させて結果的に音声認識性能を向上させることができる. また, VADは音声強調や音声符号化などの音声処理においてもその精度向上に大きな役割を果たす. ここでは雑音下のVAD評価環境として構築したCENSREC-1-C(CENSREC-1-Concatenated)について述べる. 雑音下連続数字発声データと評価ツールからなり, VADを行った結果を2種類の評価基準に基づいて評価できるものである.

3.1 データ構成

連続数字を間隔をあけて発声したもので, 個々の発声はCENSREC-1[5]に準じている. 雑音加算によるシミュレーションデータと実際の雑音環境下で発声された実環境データでの評価の両方が行えるよう, 以下のデータ群を, 目視により作成された音声区間の始端・終端データとともに用意した.

3.1.1 雑音加算によるシミュレーションデータ

CENSREC-1[5]の音声データを複数接続することにより, 連続的な発話のデータを作成している. CENSREC-1の1つの雑音環境で同一話者により発話されている9ないし10の発話を1秒の無音を挟んで接続した. CENSREC-1の1つの雑音環境での話者数は104名(男女各52名)であり, 話者毎に接続データを作成するため, 1環境におけるデータ数は104となる. 本評価環境における雑音環境はCENSREC-1のSet A, Bと同様である(表1参照). ただし, Set Cは提供しない.

3.2 実環境データ

実環境データの収録は, 2つの雑音環境(学生食堂, 高速道路付近)および2つのSNR環境(低SNR, 高SNR)にて行った. マイクロホンは近接位置と遠隔位置を同期収録した(近接マイクのデータは評価対象ではない). 被験者は, 男女5名, 合わせて10名(被験者の年齢内訳は, 20歳前後男女各3名, 30歳前後男女各1名, 40歳以上男女各1名)とした. 1名につき各雑音環境および各SNR環境につき1~12桁の連続数字を8~10回, 約2秒間の間隔で発声した音声を1つのファイルとして, 計4ファイル(総発話数: 38~39発話)である(ただし1名分は意図通りに発話できていないデータとして評価から除外).

シミュレーションデータおよび実環境データの内容を表2にまとめる.

表 2 CENSREC-1-C 収録データの全容

	シミュレーションデータ	実環境データ
各ファイル中の発声	CENSREC-1 の 9~10 発話を無音を挟んで接続	連続数字を約 2 秒間隔で 8~10 発話
雑音	CENSREC-1 と同じ雑音を -5~20dB(5dB 刻み) で加算, またはクリーン	学生食堂 (低 SNR 69.7dB, 高 SNR 53.4dB) および高速道路付近 (低 SNR 69.2dB, 高 SNR 58.4dB)
データ量	各条件 104 ファイル, 全 5824 ファイル	各条件 4 ファイル, 全 160 ファイル (評価対象は 144 ファイル)

3.3 評価尺度

3.3.1 フレーム単位の評価尺度

音声区間検出の性能を測るためのフレーム単位の評価尺度として, FRR (False Rejection Rate) と FAR (False Acceptance Rate) を用いる。

$$FRR = \frac{N_{FR}}{N_s} \times 100 [\%], \quad FAR = \frac{N_{FA}}{N_{ns}} \times 100 [\%]$$

ここで, N_s は音声フレーム数, N_{FR} は音声を非音声と検出したフレーム数, N_{ns} は非音声フレーム数, N_{FA} は非音声を音声と検出したフレーム数である。複数のデータを対象とする場合には, データ毎に FRR と FAR を求め, その平均値を用いて評価する。手法によってフレーム長は異なるので, サンプル単位の正解データをフレーム長に換算して評価する。

一般に FRR と FAR はトレードオフの関係にあり, どちらの性能を重視するかは対象とするアプリケーションによって決まる。そこで, 閾値によって FRR と FAR を調整することで, ROC 曲線 (x 軸: FAR, y 軸: $100 - FRR$) を示すことを推奨している。以下の 2 つの ROC 曲線を描くための Excel シートを用意している。

- (1) SNR 別 ROC 曲線: 雑音レベルによる性能変化の評価
- (2) 雑音別 ROC 曲線: 雑音種類による性能変化の評価

3.3.2 発話単位の評価尺度

音声認識のための音声区間検出は, 一般に発話単位 (単語や文など) で行う。発話単位の音声区間検出性能を評価する尺度には, 発話区間検出正解率 Corr と発話区間検出正解精度 Acc を用いる。

$$Corr = \frac{N_c}{N} \times 100 [\%], \quad Acc = \frac{N_c - N_f}{N} \times 100 [\%]$$

ここで, N は総発話区間数, N_c は正解発話区間検出数, N_f は誤発話区間検出数である。Corr は, 発話区間をどれだけ多く検出できるかを評価する尺度であるのに対し, Acc は, 発話区間をどれだけ過不足なく検出できるかを評価する尺度である。複数のデータを対象とする場合には, データ毎に Corr と Acc を求め, その平均値を用いて評価する。

音声認識では, 発話区間を短かめに検出すると認識誤りを起こしやすいが, 前後にやや長めに検出しても認識結果への影響は比較的少ない。そこで, 検出区間の中に, ある 1 つの発話区間全体が含まれ, その前後の発話区間と重なりがなければ, 正解検出とし, その他の検出区間は全て誤検出とする。

3.4 ベースライン評価

音声パワーに基づくベースライン VAD を設定し, その結果を公開している。雑音区間と音声区間のフレームごとに音声/非音声の 2 クラスに分類する分類問題と考えると, クラス分類に

表 3 ベースラインの実環境データに対する発話単位の評価結果 (上: 正解率, 下: 精度)

Real Data	Correct Rate [%]		
	Remote Microphone	Restaurant	Street
Correct Rate [%]	High SNR: 74.20	39.42	56.91
	Low SNR: 56.52	41.45	48.99
Average	65.36	40.44	52.90

Real Data	Accuracy [%]		
	Remote Microphone	Restaurant	Street
Accuracy [%]	High SNR: 21.45	-15.65	2.90
	Low SNR: -43.48	-33.91	18.70
Average	-11.02	-24.78	-17.90

最適な閾値を判別基準

$$\eta(s) = \frac{\sigma_B^2(s)}{\sigma_T^2}, \quad \sigma_T^2 = \sigma_W^2(s) + \sigma_B^2(s)$$

を最大とするしきい値 (s) を求めることにより決定する方法である [9]。ここで, $\sigma_W^2(s)$ は平均クラス内分散, $\sigma_B^2(s)$ は平均クラス間分散を示す。この手法で求めた最適なしきい値を初期しきい値 (THR_{int}) とし, 以下の式で $k = -40, -39, \dots, 39, 40$ の 81 通りに変動させて評価した。

$$THR = THR_{int} + k \cdot \alpha, \quad \alpha = \frac{(POW_h - POW_l)}{K}$$

ここで, POW_l は初期しきい値未満のフレームの平均パワー, POW_h は初期しきい値以上のフレームの平均パワーを示す。

フレームのパワーがしきい値以上となった場合, 当該フレームを音声開始フレーム候補とし, 次の音声終了フレーム検出を行う。音声開始フレームの候補が決定された後, フレームのパワーがしきい値未満となりかつそのフレーム以降のある一定区間 (ベースライン実験では 500ms) 以上のフレームのパワーがしきい値未満の場合, 音声終了フレームとする。ただし, 一定区間以上連続してフレームのパワーがしきい値未満でない場合は棄却し, 再度音声終了フレームを検出する。さらに, 音声終了フレームが確定しても, 音声区間が一定区間未満の場合にはその音声区間は棄却する。

フレームベース評価における SNR 別, 雑音別の ROC 曲線を図 1 に示す。また, 発話単位の性能評価結果を表 3 に示す。これは, 最も発話区間検出正解率 (Corr) が高かったしきい値の結果である。ここでの評価では, 上記手法で得られた発話区間を前後に 300ms 延長している。

シミュレーションデータに対しても同様の評価が行える。利川者は自身の手法で定められたフォーマットで VAD 結果を出力すれば容易に評価結果が得られ, ベースラインと比較できる。詳細は文献 [8] を参照されたい。

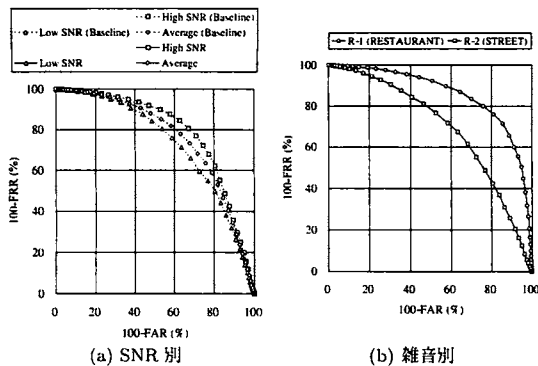


図1 ベースラインの実環境データに対するVADのROC曲線

4. 今後の展開

さきにも述べた通り、これまで認識対象データに含まれる種々の要因の音声認識性能に及ぼす影響を比較検討できる基盤群を構築してきた。ここでは、それを含めより広い観点から評価すべき要因を挙げ、今後どのような評価基盤を構築するかを考える。

4.1 環境

まず、認識対象データのバリエーションについて考える。

4.1.1 加法的雑音

これまで、実在する雑音の種類を集めて重畳したり、自動車の走行条件を変えるなどしてきた。これは実在の雑音のカバレッジを向上させるという目的がある。しかし、手当たり次第に集めることで十分なカバレッジを得ることは容易ではないし、得たとしても膨大な種類の雑音での評価が必要となってしまう。

そこで、音声認識の性能への影響という視点から、雑音の種類を客観的に評価・分類し、雑音の種類を効率的にサンプリングする方法を探求していきたいと考えている。例えば、WGメンバーにより音声の歪尺度 PESQ [12] と擬似音声 [13] を用いた簡便な認識率推定手法を提案されている [14]。また、雑音はよく SNR や非正常性で分類されるが、同じ SNR でもスペクトルが異なれば音声への影響は異なり、非正常性にもさまざまなある。平均スペクトルから非正常性、突発性などを複合的に用いて雑音の姿を記述できる量に関しても議論していく。そして、COSMOS マップ [15] のように俯瞰的・可視的にみること雑音収集の方向が見えてくるものと考えている。

4.1.2 残響

これまで実環境の大きな要素として注目され研究もされながら、音声認識の評価基盤として確立されていないのが残響である。残響もインパルス応答が測定できれば畳み込みにより残響環境下音声シミュレーション可能である。しかし、加算性の場合と同様にシミュレーションと実環境下の差異が問題となる。そこで、現在 CENSREC-1 と同様な畳み込みによるシミュレーション残響音声と、実際に残響の存在する場所において収録した音声とを含んだ残響音声認識評価基盤を設計している。

しかし残響も、加算性雑音と同様にそのカバレッジをいかに確保するかの問題が生じる。これに対して、メンバーにより、

初期残響と音声認識率の関係を明らかにするために第1次反射音の遅れ時間と音声認識率の関係についての調査 [16] が行われており、残響が一概に音声認識率の低下を引き起こすばかりではないことが確認されている。このように残響を記述する量を探求する研究結果を基盤構築に生かしていく。

4.2 雑音下音声認識のためのアプローチ

これまで雑音下の音声認識へのアプローチの一般的な分類として、(1) 雑音に強い特徴量, (2) 雑音の除去・抑圧, (3) 音響モデル・適応があった。これらはいずれも与えられた音声があり、それを HMM とマッチングする際のミスマッチを軽減するための方法であり、これまでに提供してきた CENSREC シリーズで評価が可能であった。

しかし最近、より異なるアプローチも多数提案されてきている。そのひとつに、音声区間検出 (VAD) がある。音声と HMM とのミスマッチ軽減という直接的な方法ではなく、不要なマッチングを除外し、必要十分な部分のみを認識することで、大幅な湧き出し誤りと脱落誤りの削減効果がある。この方法は認識そのものによる対処法と独立に開発し、組み合わせることができる。このアプローチの独立した評価基盤として CENSREC-1-C を配布している。

また、フロントエンドの処理としてマルチマイクロフォンを用いる手法が数多く提案され、遅延和アレイによる SN の向上や適応的に特定の到来方向の雑音を抑圧する方法、さらに空間的に音声を除く音成分を取り出して減算することによる空間的スペクトルサブトラクション [17] などが効果的であると言われている。このように、(比較的少数の) 複数マイクロフォンを用いた手法はいくつも提案されてきており、各手法の能力を客観的に比較できる基盤は必要であろう。ただし、マイクロフォンの配置にはそれぞれの手法や研究者の意図によってさまざまな工夫が施されており、共通に利用可能な収録条件が設定できるかどうか難しい面がある。今後、慎重に議論していきたい。

さらに、認識段階における新たなアプローチも見られる。音響的な雑音に強い映像情報を音声とハイブリッドに用いて頑健に音声認識を実現するパイモダル手法も数多く研究されている。そして、大語彙連続音声認識アルゴリズムにおける耐雑音性の向上も考えられる。まだ研究例は少ないが、雑音レベルに応じて音響尤度と言語尤度の重みを動的に変化させて、高 SNR 時には音響モデルに、低 SNR 時には言語モデルに頼るなどデコーディングの工夫による研究例もある (例えば [18])。英語の AURORA-4 は大語彙連続音声認識タスクにおける評価基盤であり、大規模さのために動作に時間がかかるなど問題があり現在のところあまり利用されているとは言えない。しかしこのようなアプローチにも十分な効果があることが確認されれば必要とされよう。研究動向を注視しながら議論を重ねていく。(注1)

このように雑音に対するアプローチも多様化している。そ

(注1): さらに上位のアプローチとしてアプリケーションレベルもありえよう。すなわち対話的に誤認識を解消したり他の入力デバイスとの併用により認識結果の複数候補から正解を選択するなどがある。しかしこれらは本 WG の対象外としたい。

れに見合った評価基盤の提供が必要である。次節では、その一例として顔映像と音声のバイモーダル情報による認識のための評価基盤として配付準備を進めている CENSREC-1-AV, CENSREC-2-AV を紹介する。

4.3 バイモーダル音声認識評価用基盤 CENSREC-1-AV, CENSREC-2-AV

4.3.1 概要

CENSREC-1-AV は, CENSREC-1 と同様に, 室内で収録したデータに事後的に雑音を重畳したデータである。音声には, CENSREC-1 と同じ雑音を, 映像には車内でダッシュボードを撮影した映像から得た γ 値の変化を加えている。発話内容も CENSREC-1 と同じ連続数字単語であり, 約 90 名のデータの公開を予定している。

一方, CENSREC-2-AV は, CENSREC-2 と同様に, 実車内でドライバの音声と映像を収録したデータである。名古屋大学のデータ収録車を用いて, アイドリング・走行中に収録した約 50 名のドライバの音声と顔映像を収録している。近年, 安全運転支援のための居眠りや脇見の検知, 防犯用の顔撮影のためのカメラが実用化され始めており, 今後車内カメラの普及が予想され, バイモーダル音声認識にも利用できるかと期待できる。

音声は 2 つのマイク, 顔映像はカラーカメラと可視光カットフィルタ装着した近赤外カメラで収録しており, メディアコンバータでステレオ音声とそれぞれの映像を統合している。[19]

4.3.2 配布方法

バイモーダル音声データベースの映像に関しては, それぞれ 2 つの形式で公開する予定である。一つは, 顔全体を含む収録した映像のデータそのものであり, もう一つは, 顔映像から唇付近のみを抽出した口唇映像のデータである。口唇映像は, 唇位置検出アルゴリズムや追跡アルゴリズムを用いて唇部分を抽出し, 各フレームに分割した画像を, 音声・評価用スクリプト等と併せて DVD に収めて配布する。顔映像のデータは容量が大きいため, HDD での配布を予定している。音声と顔映像は, DV フォーマットの AVI ファイルとして保存しているが, AVI ファイルから映像をステレオ音声と映像に分割するためのソフトウェア等も併せて配布する予定である。

音声は MFCC, 映像は主成分スコアを特徴量とした認識結果を, 比較用の認識結果として掲載する。音声と映像を統合するためのストリーム重みは, 学習時は等重みで固定し, 認識時にもみ重みの割合を変化させて, 各 SNR において認識精度が最も高くなる重みを選択できるようなスクリプトを配布する。

評価用スクリプトや認識結果を集計するための Excel シートの形式については, CENSREC-1, CENSREC-2 に準拠する。これらの他に, 映像付加による認識性能の改善率を計算するための Excel シートも追加する。

5. まとめ

本稿では, 雑音下音声認識評価ワーキンググループのこれまでの標準評価基盤 CENSREC シリーズを振り返り, 今年度新たに構築した CENSREC-1-C の概要と, これまでおよびこれからの CENSREC シリーズ内での位置づけを述べた。今後は,

CENSREC-1,2,3 と同様に認識対象データの雑音環境の幅を広げていくとともに, CENSREC-1-C のように雑音下音声認識に対する音声認識そのものだけではなく, 違ったアプローチに対しても個別評価・比較検討できる環境を提案していく。

文 献

- [1] D. Pearce, "Developing the ETSI AURORA advanced distributed speech recognition front-end & What next." Proc. Eurospeech 2001, 2001.
- [2] Aurora document no. AU/345/01, "Large vocabulary evaluation of front-ends- baseline recognition system description." Mississippi State University, Aug. 2001.
- [3] 中村 哲, 武田一哉, 黒岩真吾, 山田武志, 北岡教英, 山本一公, 西浦敬信, 藤本雅清, 水町光徳, "SLP 雑音下音声認識評価ワーキンググループ活動報告." 情報処理学会研究報告, 2002-SLP-42-11, pp.65-70, July 2002.
- [4] 中村 哲, 武田一哉, 黒岩真吾, 山田武志, 北岡教英, 山本一公, 西浦敬信, 藤本雅清, 水町光徳, "SLP 雑音下音声認識評価のための WG: 評価データ収集について." 情報処理学会研究報告, 2002-SLP-45-9, pp.51-56, Feb. 2003.
- [5] S. Nakamura et al., "AURORA-2J: An Evaluation Framework for Japanese Noisy Speech Recognition." *IEICE Transactions on Information and Systems*, Vol. E88-D, No. 3, pp. 535-544, 2005.
- [6] 藤本 雅清, 武田一哉, 中村 哲, "CENSREC-2: 実走行車内における連続数字音声データベースと評価環境の構築." 情報処理学会研究報告, SLP-60-3, pp. 13- 18, 2006.
- [7] M. Fujimoto, K. Takeda, and S. Nakamura, "CENSREC-3: An Evaluation Framework for Japanese Speech Recognition in Real Driving-Car Environments." *IEICE Transactions on Information and Systems*, (accepted).
- [8] 北岡教英ら, "CENSREC-1-C: 雑音下音声区間検出手法評価基盤の構築." 情報処理学会研究報告, 2006-SLP-63-1, pp. 1-6, 2006.
- [9] 大津展之, "判別および最小 2 乗基準に基づく自動しきい値選定法." *信学論*, Vol. J63-D, No. 4, pp. 349-356, 1980.
- [10] K. Takeda et al., "Construction and Evaluation a Large In-Car Speech Corpus," *IEICE Transactions on Information and Systems*, Vol. E88-D, No.3, Mar. 2005.
- [11] 武田一哉他 "走行状況別車内音声データベースとその予備評価," 音講論集, 3-P-10, pp. 185-186, Mar. 2002.
- [12] ITU-T Recommendation P.862, "Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs." Feb. 2001.
- [13] ITU-T Recommendation P.50, "Artificial voices," Sept. 1999.
- [14] T. Yamada, M. Kumakura, N. Kitawaki, "Performance estimation of speech recognition system under noise conditions using objective quality measures and artificial voice," *IEEE Transactions on Audio, Speech and Language Processing*, Vol. 14, No. 6, pp. 2006-2013, 2006.
- [15] 庄境 誠 "複数音声コーパスの俯瞰的分析." 音声言語シンポジウム, 情報処理学会研究報告, 2005-SLP-59-23, 2005.
- [16] 西浦敬信, 傳田遊也 "音声認識における初期反射音の影響についての検討." 日本音響学会 2006 年度春期全国大会講演論文集, 3-1-17, pp. 141-142, Mar. 2006.
- [17] Yu Takahashi, Tomoya Takatani, Hiroshi Saruwatari, Kiyohiro Shikano, "Noise reduction using spatial subtraction array based on independent component analysis". 電子情報通信学会技術研究報告, EA2006-22, vol.106, no.125, pp.13-18, June 2006.
- [18] N. B. Yoma, F. R. McInnes, M. A. Jack, S. D. Stump, and L. L. Ling, "On including temporal constraints in viterbi alignment for speech recognition in noise." *IEEE Transactions on Speech and Audio Processing*, Vol. 9, pp. 179-182, 2001.
- [19] 根本 大輔, 前野 俊希, 北坂 孝幸, 森 健策, 末永 康仁, 宮崎 千代美, 伊藤 克己, 武田一哉, 板倉 文忠, 佐野 昌己, 二宮 芳樹, "バイモーダル車内音声認識評価用データベースの構築", 情報処理学会研究報告, 2005-SLP-55-7, Feb. 2005.