

論文 / 著書情報
Article / Book Information

論題(和文)	対談音声を対象とした音声認識の検討
Title(English)	Investigation into Recognizing Interview Speech
著者(和文)	斎藤洋平, 古井 貞熙
Authors(English)	Yihei Saito, SADAOKI FURUI
出典(和文)	日本音響学会2001年春季講演論文集, Vol. , No. 1-3-17, pp. 37-38
Citation(English)	, Vol. , No. 1-3-17, pp. 37-38
発行日 / Pub. date	2001, 3

対談音声を対象とした音声認識の検討*

齋藤 洋平 古井 貞熙 (東工大)

1. はじめに

NHKのTV番組「クローズアップ現代」を用いて、対談音声の認識を試みた。97年7月29日および30日に放送された音声をテストセットとして用いた[1,2]。内容はそれぞれ「新人官僚の研修」と「ヨットレース事故」で、いずれも、男女の対談部分および女性キャスターの単独発話部分を含んでいる。認識タスクとしては、対談部分の女性発話を用いた。また、比較対照用として、女性キャスターの単独部分の認識も行った。テストセットの大きさを表1に示す。

表1 テストセットのサイズ

Test Set	文数	形態素数
対談	19	494
キャスター	22	450

2. 言語モデル

以下の言語モデルを作成した。各学習データの形態素解析にはJTAGを用いた。

LM1: 放送ニュース原稿5年分(92年7月から96年5月まで)から作成。総学習文数0.4M。語彙サイズ20K。

LM2: 「クローズアップ現代」書き起こしテキスト(97年6月2日から99年6月30日放送分まで合計159日分(テストセットの2日間を除く))。総形態素数0.8M。語彙サイズ20K。

LM3: LM2に品詞N-gramを用いて未知語の追加を行った[3]。追加する未知語の選択は、ニュース原稿の書き起こしから、テストセットの内容に近いものを自動的に選び出し、その書き起こし中に含まれる単語の中で、頻度の上位M単語を選択した。その単語の品詞(主品詞および細品詞)情報に基づくクラスN-gramを用いて、未知語と既知語間の単語N-gramを計算した。なお、本実験ではM=50とした。

LM4: テストセットと同日に放送された番組中のVTRの書き起こし(約6.5K形態素)を、学習テキストに重み付けして足し合わせたモデル。

3. 音響モデル

音響モデルは、以下の3種類のモデルを用いた。

AM1: IPAの男女別不特定話者モデル。状態数

2000、混合数16。

AM2: 話者に適応化したモデル。対談の女声部分は、女性キャスターと同一人物であるため、女性キャスターの単独発話部分を用いて、AM1をもとにML推定を行った。

また、対談音声認識の重要な問題点であるクロストーク部分に関しては、音響的なback-off(以下、AB)を行った[4]。すなわち、複数話者の発声が同時に行われている場合、その部分は音響スコアの計算を行わず、その話者に対するフレームごとの平均尤度を割り当てるという手法を用いた(クロストーク区間は、何らかの方法で検出できていると仮定した)。

4. 実験結果

表2と3に、各言語モデルのパープレキシティを示す。ニュース原稿を用いて学習したLM1に比べ、テストセットと同じ番組の書き起こしを用いることにより、パープレキシティを大幅に削減することができる。LM4に関しては、AM1,2それぞれにおいて、認識率が最も高くなる重みでVTR書き起こしを加えたときの平均値を示す。

表2 Test Set Perplexity (Bigram)

Test Set	LM1	LM2	LM3	LM4
対談	539.99	69.16	82.24	94.61
キャスター	322.48	113.72	128.17	162.9

表3 Test Set Perplexity (Trigram)

Test Set	LM1	LM2	LM3	LM4
対談	510.79	56.56	66.42	85.53
キャスター	306.70	101.49	115.24	198.7

図1に、各言語モデルにおける未知語率を示す。LM1とLM2を比較すると、キャスター部分ではそれほどの違いは見られないが、対談部分に関しては、LM2で未知語率が大幅に削減されている。これは、話し言葉独特の言い回しがLM2で学習されることによる。また、関連ニュース原稿を用いて未知語を追加することにより、未知語率を大幅に削減できることが分かる(LM3)。また、放送されたVTRを用いることにより、LM3よりも容易に未知語率を削減できることが分かった(LM4)。

*Investigation into recognizing interview speech. By Yohei Saito and Sadaoki Furui(Tokyo Institute of Technology)

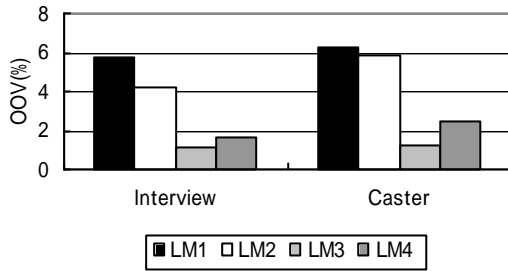


図1 未知語率

各言語モデル別の認識率を図2に示す。認識実験には Julius2.1 を用いた。音響モデルは AM1 を用いている。また、各音響モデルと音響 back-off を使用した場合の認識率を図3に示す。言語モデルは LM2 を用いている。この図中の「Interview」では、クロストーク区間だけでなく、対談の相手話者の(重なっていない)相槌などに対しても適用している。そこで、クロストーク部分のみに AB を適用した結果を「Crosstalk Only」に示す。クロストークを含むテストセットは対談19文中9文、クロストーク区間長は全体の約5%であった。これらの結果から、ここで検討した種々の方法の有効性が分かった。

また、LM4においてVTR書き起こしに対する重みを変えたときの認識率の変化を図4に示す。言語モデル重みや挿入ペナルティは一定値を用いたが、音響モデルの違いによって認識率の飽和点が大きく異なることが分かった。

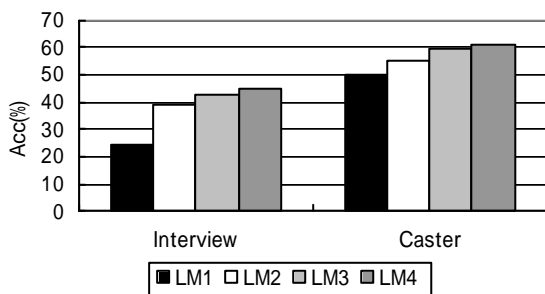


図2 言語モデルの性能評価

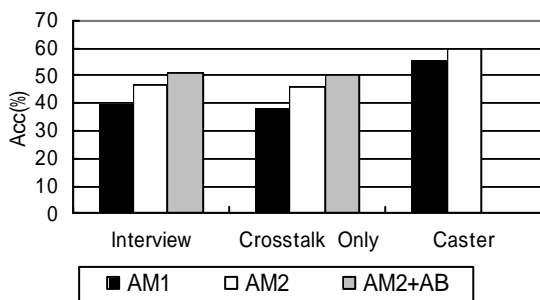


図3 音響モデルの性能評価

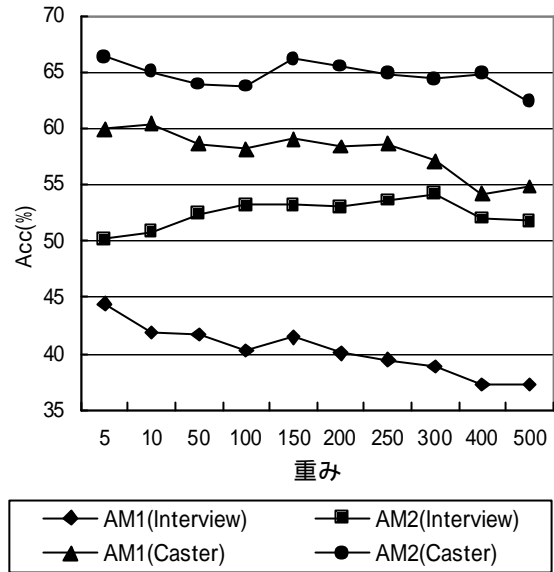


図4 重みによる認識率の変化

5. 考察

話し言葉を忠実に書き起こしたデータから学習した言語モデルの有効性が確認できた。また、対談の話題に沿った単語を言語モデルに追加することにより、未知語率を削減することができた。また、少量のテキストでも、対談の話題に近いものを学習コーパスに足し合わせるにより、認識率を向上させることが分かった。

音響的には、音響モデルの話者適応や再学習により認識率を向上させることができた。また、対談の大きな問題点であるクロストークに関しても、音響的な back-off を適用することが有効であることがわかった。

謝辞

放送音声や書き起こしコーパスを提供いただいた NHK 放送技術研究所、また、古井研究室の皆様の協力に感謝する。

参考文献

- [1] 本間、今井、安藤：“対談番組を対象にした音声認識の検討”、春季音学講論、3-Q-24 (1999)
- [2] 篠崎、斎藤、堀、古井：“話し言葉音声の認識を目指して”、信学技報、SP2000-96 (2000)
- [3] 松井、加藤、小林、今井、田中、安藤：“ニュース音声認識における直前原稿を利用した認識性能の改善”、信学技報、SP99-128 (1999)
- [4] J.de Veth, B.Cranen & L.Boves, “Acoustic backing-off in the local distance computation for robust automatic speech recognition”, Proc.ICSLP-98, Sydney, pp. 1427-1430 (1998)