

音声のスペクトルのローカルピークによる認識

東北工業大学
環境情報工学科
古賀秀昭

I. 単母音のスペクトルのローカルピークと傾斜を用いた話者認識

1.1 はじめに

ンるして化犯てフ術コ使実
コあとしれルタシル技とにを
、で法とさバーにアの間利会
し声用ド表一ユかや識人便社
出の利一代ロピい字認、り報
抽誰のワにグンを数者はよ情
をし識ストのコ報の話とをな
徴析認パツそて情来、コンか
特解者のネはし人従てるヨ豊
のを話人一でと個。べれシで
のを話人一でと個。べれシで
のの特る。個夕会償。るらら一全うで
手の特る。個夕会償。るらら一全うで
しそでのに報さ大題号でユてるの己研とカ知力の本るは一
話てとる特情利が問暗面ミしき識自本斜一が入誰があでデ
はっこれ。の便性なる場コにで認や、傾口とははら声が験の
とよるらる代、能要よなののが者ムが3、こにか音合実音
識にすえあ現い可重にうタもと話ラたたはい識中力照た母
認タ別考で、なくはトよ一いこのトれめれ強認の入者っ単
者一判に用にも招かッのユする来スら求こに者声、話行の
ピュを主応うとをるるべこピヤす従プいら。音話るとる回分
ピかてのよに罪守アがンい現ケ用かる騒あ別す今人る。

1.2 実験方法

1-2-1. 実験の流れ

一た調べのつと
ピしを述べよこ
ルと数で一にる
カ眼波後ピ識べ
一主周はル認調
口に割てかるか
と特分い一よの
斜でのつ口にる
傾斜にと用れ
3 実験細識併さ
のの3 詳認の善
ル回なのの方改
ト今適数み両に
ク、最波の、う
ペがに周斜しよ
スる識割傾較の
のはい認分3比ど
でて①(②)をが
1-2-1. 実験の流れ
本をのること認識る。
クもべる)の認識る。

1-2-2. 音声資料

実験に使用した音声データについて述べる。本実験では単母音を対象としているので、日本語の5母音 /a/, /o/, /u/, /i/, /e/ を実験資料とする。被験者（日本人男女）に防音室で、各母音をランダムに10回ずつ発声してもらい録音する。録音にはDATを用いた。録音した音声を4.5kHzのローパスフィルタに通した後、サンプリング周波数12kHz、分解能12bitでAD変換した。これを視察で1536点切り出した。さらに512点を1フレームとして3フレームに重複なしで分割した。従って、被験者一人あたりのデータの数は5（5母音）×10（10回発声）×3（フレーム）=150個である。（表1-1）

1-2-3. 特徴抽出

前述した音声資料を1フレーム（512点）毎に以下の処理を行う。（図1-1）

1. 時間窓をかける。（ハミング窓）
2. 14次線形予測係数（LPC係数）を求める。
3. LPC係数を高速フーリエ変換（FFT）し、LPCスペクトルを求める。
4. スペクトルに最小二乗近似直線を引いてスペクトルの3傾斜を求める。
5. LPCスペクトルからローカルピークを抽出し1/6 Oct.帯域幅のバンドパスフィルタに割り当て、25chのローカルピーク系列に圧縮する。

この処理を傾斜の分割周波数を500Hzから1/3 Oct.毎に変化させ、4000Hzまで繰り返した。分割周波数を固定させずに変化させたのは、傾斜の引き方によって認識にどのような影響があるのか、どこに最適な分割周波数が存在するのかどうかを調べるためである。

1-2-4. 傾斜の引き方

LPCスペクトルの全帯域（235Hzから4.5kHzまで）の最小二乗近似直線の傾斜（AT:Total 図1-2）、235Hzから分割周波数（500Hzから4.0kHzまで変化させる）までの傾斜（AL:Low 図1-3）、分割周波数から4.5kHzまでの傾斜（AH:High 図1-3）とした。また、分割周波数は500Hzから1/3 Oct.毎に変化させたので

$$f_n = 500 \cdot 2^{\frac{(n-1)}{3}}$$

より求まる周波数 f_n をその分割周波数とした。
ただし、 $(n=1, 2, \dots, 10)$ である。

1-2-5. 認識方法

LPC スペクトルの 3 傾斜による BAYES 判定を行い、その対数尤度 (SL) を求める。SL は次のように求める。入力データを X とする。話者をカテゴリとする k を用いて

$$SL_k = -\frac{1}{2} \left\{ n \cdot \log(2\pi) + \log \left| \sum_k \right| + (X - \mu k)' \sum_k^{-1} (X - \mu k) \right\}$$

として、この SL_k が最大であるカテゴリ k を BAYES 判定の結果とする。ただし、 μ は平均値、 \sum は共分散行列、 n は特徴の次元である。この場合は 3 傾斜なので $n=3$ である。

続いてローカルピーク系列 $\{0, 1\}$ で、2ch 同時確率の対数尤度 (LP) を求める。LP は次のように求める。認識用の入力音声データ系列を X_m とする。

$$X_m = (x_1, x_2, \dots, x_i, \dots, x_j, \dots, x_n)$$

x_i, x_j は 1 または 0 であり、それぞれピークの有無を示す。話者をカテゴリとする k を用い、 $x_i = x_j = 1$ となるときに確率を P_{ijk} とする。もし、 $x_i = x_j = 1$ となった場合に標準パターン P_{ijk} として、そうでない場合は $(1 - P_{ijk})$ より次式を計算し対数尤度 LP_k を求める。

$$\begin{aligned} LP_k &= \log \left(\prod_{j=1}^n \prod_{i=1}^n P_{ijk}^{x_i x_j} \cdot (1 - P_{ijk})^{1 - x_i x_j} \right) \\ &= \sum_{j=1}^n \sum_{i=1}^n \log(1 - P_{ijk}) + \sum_{j=1}^n \sum_{i=1}^n x_i x_j \{ \log P_{ijk} - \log(1 - P_{ijk}) \} \end{aligned}$$

この対数尤度 LP_k が最大であるカテゴリ k が 2ch 同時確率の認識結果となる。
 以上のように求められた SL と LP の線形結合 NL を

$$NL = \alpha \cdot SL + (1.0 - \alpha) \cdot LP$$

で計算する。このときの NL が最大となるものを認識結果とする。ここで α は線形結合係数で $\alpha = 0.0$ を 0.1 刻みで 1.0 まで変化させて認識率の変化を見る。
 従って $\alpha = 0.0$ ではローカルピークのみ、 $\alpha = 1.0$ では 3 傾斜のみの認識となる。

1-3. 認識結果

/a/ の場合では線形結合によって認識率が最大となったのは分割周波数 4kHz、 $\alpha = 0.6$ のときで 97.00% を得た。/o/ の場合では線形結合によって認識率が最大となったのは分割周波数 2519.84Hz、 $\alpha = 0.7$ のときで 94.33% を得た。/u/ の場合では線形結合によって認識率が最大となったのは分割周波数 2kHz、 $\alpha = 0.5$ のときで 98.67% を得た。/i/ の場合では線形結合によって認識率が最大となったのは分割周波数 2kHz、 $\alpha = 0.6$ のときで 95.33% を得た。/e/ の場合で線形結合によって認識率が最大となったのは分割周波数 2kHz、 $\alpha = 0.6$ のときで 97.67% を得た。最後に全単母音の認識では線形結合によって認識率が最大となったのは分割周波数 4kHz、 $\alpha = 0.5$ のときで 80.20% を得た。以上の結果を(表 1-2)に簡単にまとめた。これは各単母音の認識率が最大をとる分割周波数と線形結合係数 α と、その認識率である。結果として /u/ の認識率 98.67% が最も高い認識率であった。

1-4. あとがき

今回の実験では、傾斜とローカルピークの線形結合をすることにより認識率が相対的に改善されたことがわかった。
 また本研究に引き続き、音声に雑音を加えた場合のロバストネスの検討、および伝達特性の変化へのロバストネスについても検討してみる必要がある。

言語障害者自身の単母音を標準パターンとした言語障害者自身の単母音の音声認識を行い、その認識率を比較した。さらに、その結果も合わせて報告する。

2-2. 音声資料

実験で使用する単母音の音声資料は、表2-1に示すように、言語障害者は、男4名、女2名の合計6名に、また、健常者は、男60名、女60名の合計120名に、*/a/*, */o/*, */u/*, */i/*, */e/*の5母音をランダムな順序で、各10回ずつ発声して、その音声データをDATに録音した。表2-2に、発声者の言語障害者の原因と障害の種類を示す。言語障害の種類は、マヒ性構音障害である。構音障害の程度は、多少聞き取りやすい人から聞き取りが困難な人まで、様々である。録音した音声資料から認識用の音声データへの加工は、サンプリング定理を満足させるように、遮断周波数4.5kHz、遮断特性-80db/octのローパスフィルタを通して、サンプリング周波数12kHz、分解能12bitでAD変換した。AD変換した音声波形から母音中心付近の1536点を視察で切り出した。分析の際には、このデータを重複なしで1フレームを512点として3フレーム用いた。したがって、母音数は、各人150個となる。

2-3. 認識実験

音声認識は、まず初めに、男女別に60人の健常者の単母音を標準パターンとして、言語障害者6名の単母音を一人ずつ認識した。したがって、こ者名のはオプン認識実験となる。次に、言語障害者各人の単母音を標準パターンとして、言語障害者6名の単母音を一人ずつ認識する。これは、クロス認識実験となる。

2-3-1. 健常者の単母音を標準パターンとした認識結果

図2-1は、男女別に60人の健常者の単母音を標準パターンとして、言語障害者6名の単母音を一人ずつ認識し、各母音についての母音と認識したかを示すグラフである。健常者の単母音

を、認識した場合、我々の実験では、各母音について、約90%前後の高い認識率を得ているが、今回の言語障害者の実験では、16%から43%と全体的に低い認識率となり、各母音にわたって分散して認識されていることがわかった。このことは、言語障害者の発声の仕方が、健常者と大きく違うためであると考えられる。

2-3-2. 言語障害者各人の単母音を標準パターンとした認識結果

図2-2は、言語障害者各人の単母音を標準パターンとして、言語障害者6名の単母音を一人ずつ認識し、各母音についてどの母音として認識したかを示すグラフである。認識率は各母音で、84%から94%台と高く、安定して発声することが難しい/u/と/i/についても、84.4%と85.0%の高い認識率となった。この結果から、我々には聴き取りにくい音声でも、発声者自身は、個人の中で実際に発声の区別をしていることがわかった。

2-3-3. 各発声被験者の認識結果

表2-3は、発声被験者ごとに、オープンとクローズの2つの認識実験で得られた認識率を比較したものである。この結果から、どの被験者もオープン認識率が18%から43%台と低いのに対して、クローズの認識率が75.0%から97.3%と高いことがわかる。

特に、被験者4と被験者6のオープンの認識率が18%台と低く、健常者とは、発声法が大きく違うと考えられるにもかかわらず、クローズでは、79.3%と90.0%の高い認識率を示している。このことから、各被験者は個人の中では安定に発声の区別をしていることがわかる。

2-4. 聴取実験

聴取実験では、認識実験で得られた認識率と聞き取りによる単母音の認識率を比較した。各人の聴取実験用の音声は、認識実験に用いた各人の単母音音声データを、1母音につき10個、合計50個を持続時間128msとしてDA変換し、DATに録音した。この際、音声信号の立ち上がりや下りのクリック音などを除去するために、データの両側に10msのテーパーを付けた。

にくいが言語障害者自身は安定な発声をしており、個人の中で発声の区別をしていくことがわかった。このことから、言語障害者個人の特性に合わせて音声認識装置の開発に向け、さらに実験を進める予定である。

謝 辞

本研究を進めるにあたり、音声収集にご協力頂いた宮城県拓桃医療療育センターの関係者の皆様に深く感謝し、心からお礼を申し上げます。また、音声データの収集とデータ処理に協力して頂いた本学卒業生（平成12年度）の高田真一君と若生悟君に心から感謝いたします。

全体の謝辞

これらの研究を進めるにあたり、貴財団からの助成により多くの成果を挙げる事ができました。ここに記して厚く御礼を申し上げます。ありがとうございました。

単母音の数	5 個
遮断周波数	4.5kHz
サンプリング周波数	12kHz
分解能	12bit
切り出した点の数	512 点×3=1536 点
一人当たりのデータの総数	5×10×3=150 個

表 1 - 1 話者認識用音声資料

単母音	分割周波数 Hz	線形結合係数	認識率
/a/	4000.00	0.6	97.00
/o/	2519.84	0.7	94.33
/u/	2000.00	0.5	98.67
/i/	2000.00	0.6	95.33
/e/	2000.00	0.6	97.67
全単母音	4000.00	0.5	80.20

表 1 - 2 10 人の話者認識率

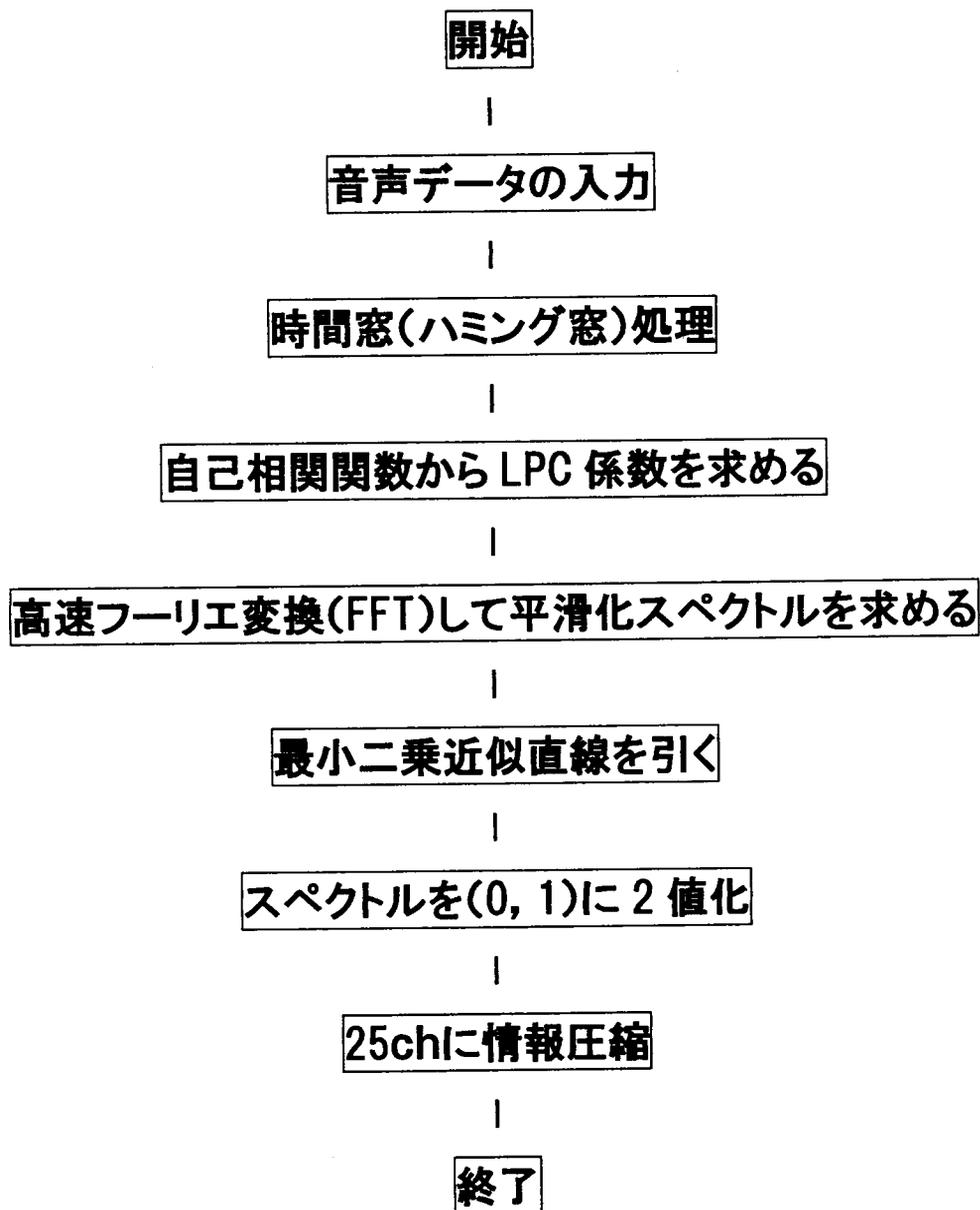


図 1-1 ローカルピークの抽出法

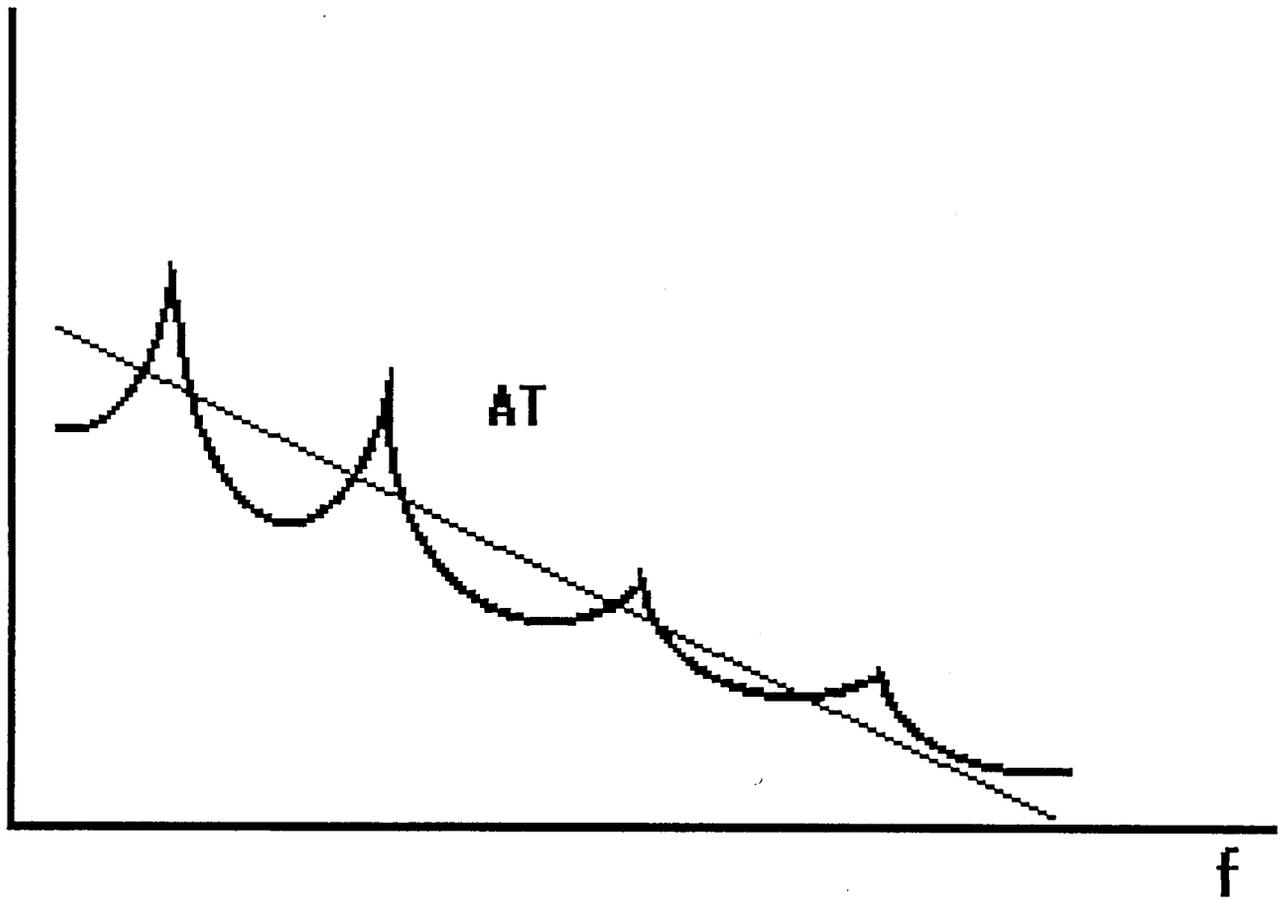


図 1 - 2 AT の抽出法

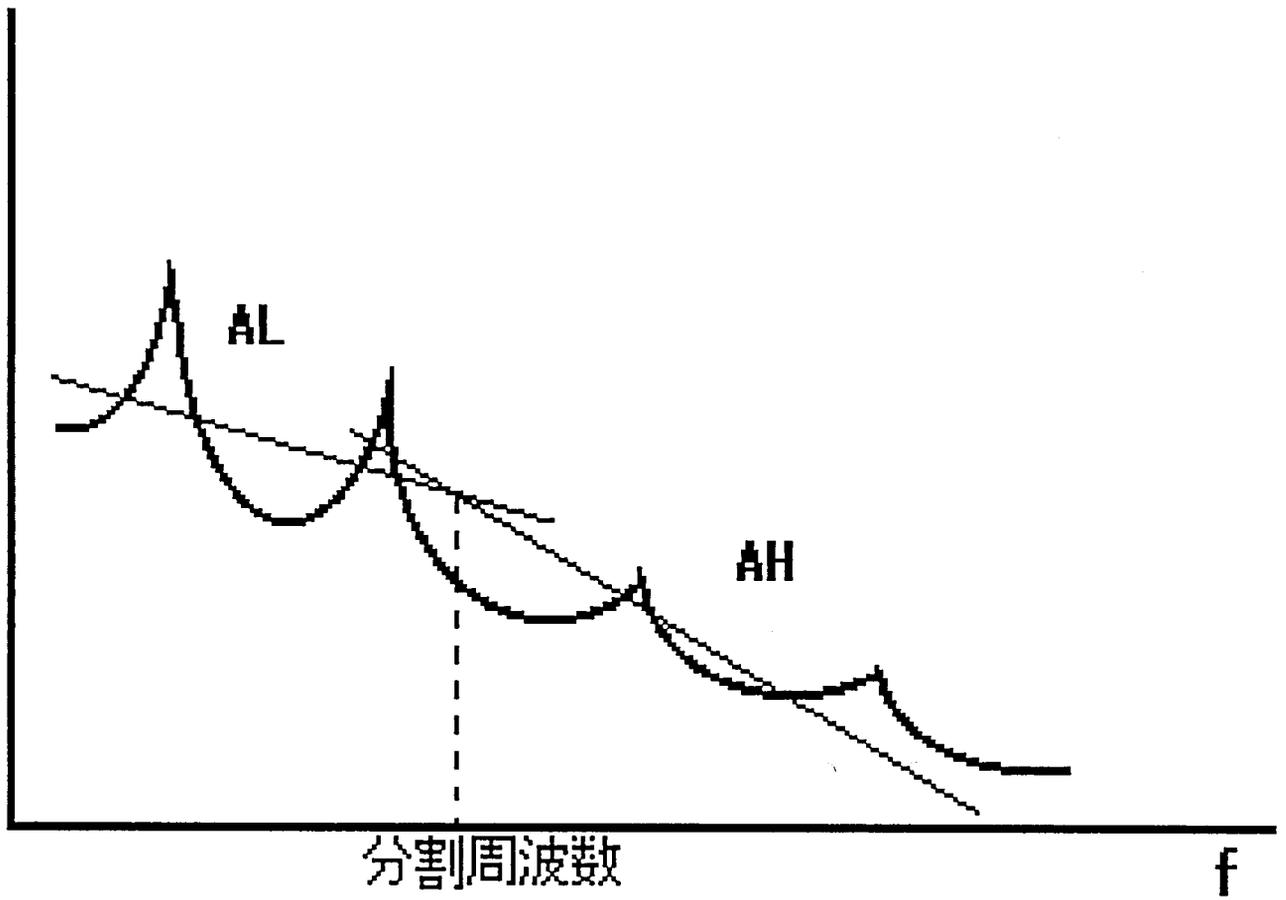


図1-3 AL、AHの抽出法

	言語障害者	健常者
話者数	男 4名 女 2名 計 6名	男 60名 女 60名 計 120名
障害の種類	マヒ性構音障害	—
LPF	遮断周波数 4.5kHz, 遮断特性-80dB/oct	
サンプリング周波数	12 kHz	
分解能	12 bit	
音声波形の切り出し	1,536 点	
1フレーム	512 点 (データ重複なし)	
母音数	各人 150 個	

表 2-1 単母音音声資料

発声被験者	診断名	言語障害の種類
被験者 1(女)	多発奇形	口蓋裂構音障害
被験者 2(女)	脳性マヒ	マヒ性構音障害
被験者 3(男)	頸直型両手マヒ	マヒ性構音障害
被験者 4(男)	頸性四肢マヒ	マヒ性構音障害
被験者 5(男)	脳性マヒ	マヒ性構音障害
被験者 6(男)	頸性四肢マヒ	マヒ性構音障害

表 2 - 2 言語障害の原因と障害の種類

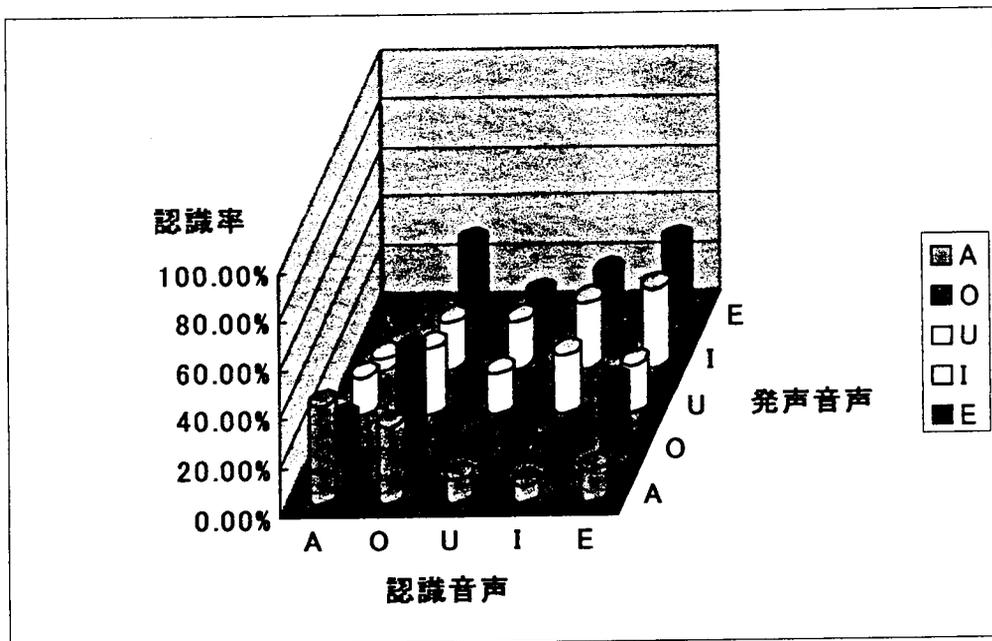


図2-1 健常者の単母音を標準パターンとした
言語障害者の認識結果

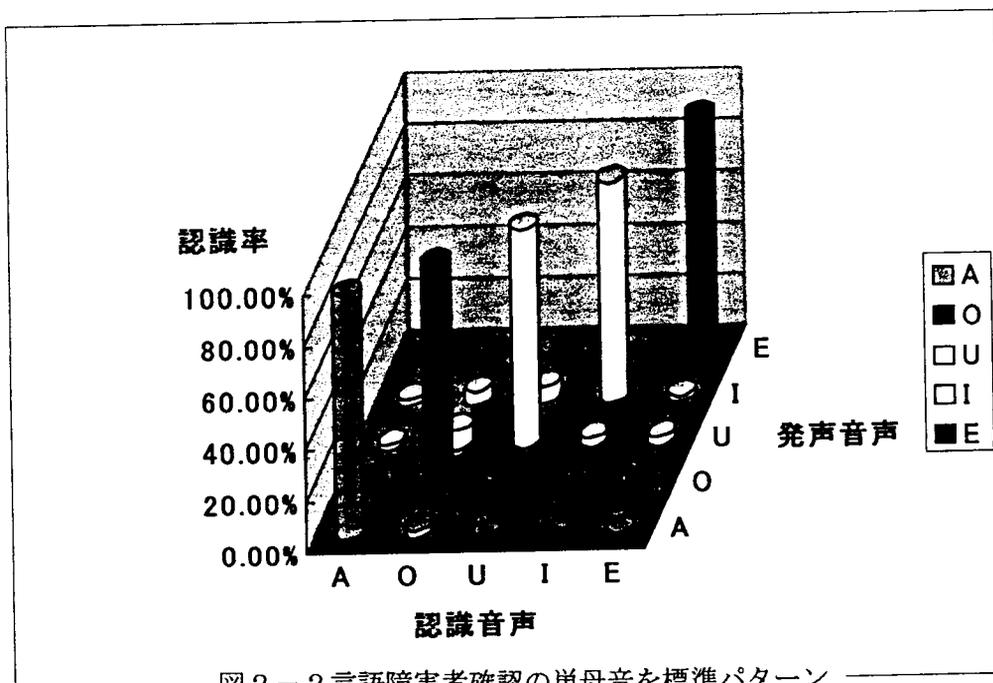


図2-2 言語障害者確認の単母音を標準パターンとした認識結果

発声被験者	(1) オープン	(2) クローズ
被験者 1(女)	37.3%	97.3%
被験者 2(女)	32.7%	75.3%
被験者 3(男)	39.3%	95.3%
被験者 4(男)	18.0%	79.3%
被験者 5(男)	43.3%	96.0%
被験者 6(男)	18.7%	90.0%
平均値	31.6%	88.9%

表 2 - 3 認識実験結果

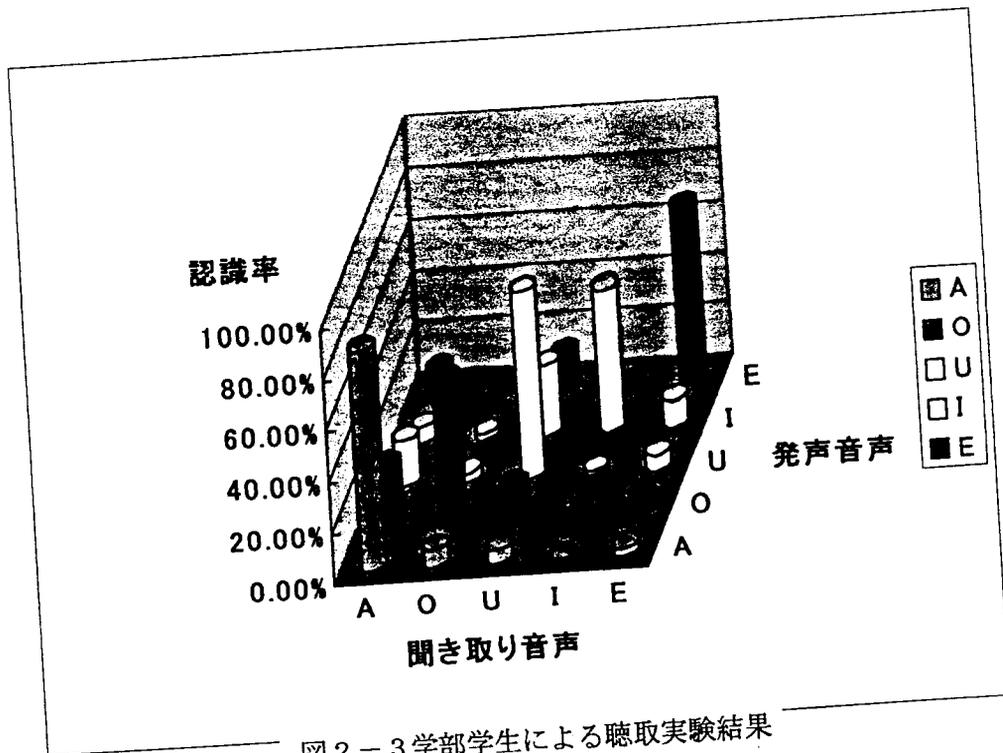


図2-3 学部学生による聴取実験結果

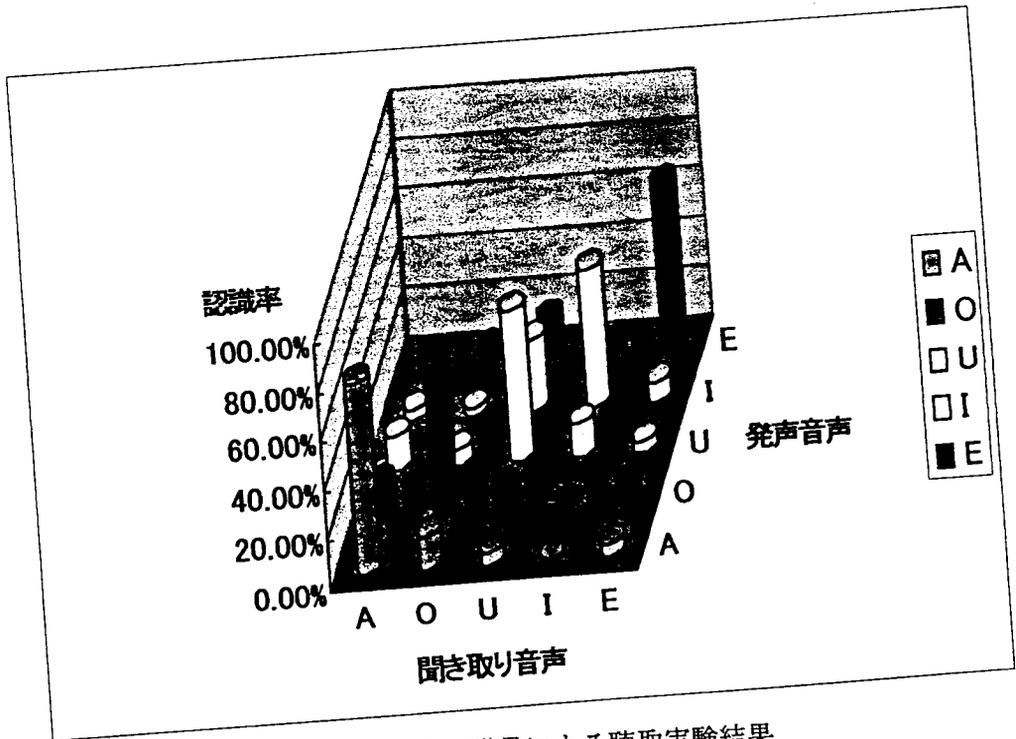


図 2 - 4 介護職員による聴取実験結果

発声被験者	学部学生	介護職員
被験者 1(女)	76.0%	63.6%
被験者 2(女)	41.6%	38.8%
被験者 3(男)	79.2%	76.0%
被験者 4(男)	63.2%	67.2%
被験者 5(男)	96.4%	90.8%
被験者 6(男)	48.0%	56.4%
平均値	67.4%	65.5%

表 2 - 4 聴取実験結果