

平成 13 年度
卒業研究論文

ロボットの聴覚機能実現の研究

98 k c 090 中島 一晃

98 k c 149 吉田 友輔

東京電機大学情報通信工学科

音響信号処理研究室

指導教授：金田 豊

目次

1 .	はじめに	．．． 2
2 .	研究目的	．．． 3
3 .	音源方向の検出	．．． 4
3 - 1 .	音源方向の検出原理	．．． 4
3 - 2 .	音源方向検出のプログラム	．．． 7
3 - 3 .	音源方向検出の例	．．． 10
3 - 4 .	サンプルと検出角度の関係	．．． 11
4 .	実音場における音源方向の検出実験	．．． 13
4 - 1 .	近距離における方向検出	．．． 13
4 - 2 .	周期音(母音)による距離をつけた実験	．．． 16
4 - 3 .	不規則な波形の音による距離をつけた実験	．．． 22
4 - 4 .	高域を強調する	．．． 27
4 - 5 .	立体的に考える	．．． 32
5 .	音声認識判断	．．． 34
5 - 1 .	音声認識判断の研究方法	．．． 35
5 - 2 .	音声特徴量の分析	．．． 36
5 - 3 .	分析結果	．．． 38
6 .	ロボットの動作制御の方法	．．． 44
7 .	まとめ	．．． 48
	謝辞	．．． 49
	参考文献	．．． 49

1 . はじめに

われわれの周辺には音が充満し、音の存在は生活の背景として至極当然のこととなっているのである。そして音の中で生き、耳を持つことの利便は、全くはかり知れないものがある。特に、言葉という音は人と人とを結ぶ架け橋であり、日常生活の中で、音は人とその環境との間の最も重要な橋渡し役を担っているのである。情報伝達に、危険防止に、方向探知に、リクリエーションなどに欠くことのできない存在となっている。

人間はもちろん、あらゆる動物はまずその生命を維持するために環境に適応した行動をとらねばならない。生物の行動は変化に富んでおり、一見したところでは、いったいどのくらい行動の種類があるのか見当もつかない。しかし十分長い期間にわたって観察すると生物は周期的な生活をしており、ある行動のパターンはたびたび繰り返されている事が分かる。生物は幾種類かの型にはまった行動 威嚇、攻撃、逃避、食餌、求愛等などを“持ち駒”として常時準備しており、環境からくる刺激の中にある特定の合図の現われを認めると、持ち駒の中の最も適したものを選んだものを選びその場面に適応しようとする。すなわちその場合の刺激（音声など）はある行動を起こす信号の役目となるのである。

「2000年からの10年間は間違いなくロボットの時代」と言われるように、近い将来に人間の生活空間で活動する人間のパートナー型ロボットが多く誕生するであろう。

そこでロボットにも聴覚機能を搭載し、刺激である音声を発すると、音源方向や音声をロボット自身で認識し、その音声に適応した行動をとれることができるようにすることを目指す。

2 . 研究目的

ロボットの聴覚機能の実現において音声を取り込み、ロボットの動作制御までの一連の流れの中には多くの課題がある。まず色々な妨害音の中からの欲しい音声だけを取り込むこと。またその音声がどの方向から聞こえているかを判断する音源方向の検出。何を言われたかを判断する音声認識判断。その次にやっとその音声による命令に従うロボットの動作制御となる。今回のロボットの聴覚機能実現の研究においては方向検出と音声認識判断に重点をおいて研究を進めた。

2つのマイクロホンにより音を取り込みそれぞれに入ってくる音の時間の差をMATLABにおいて作成したプログラムによって計算し音源の方向を角度によって割り出す。

最終的にそのプログラムを市販のロボットに移行し、そのロボットに何らかのアクションを起こさせる。またそのプログラムにおいて音声認識判断を行い、それらの違いでロボットのアクションも異なるようにする。

3 . 音源方向の検出

3 - 1 . 音源方向の検出原理

距離が d だけ離れた 2 つのマイクロホン M 1、M 2 に到達する音波の距離差、即ち時間差 s によって角度 θ を割り出す。音波の録音状況を図 1 に示す。

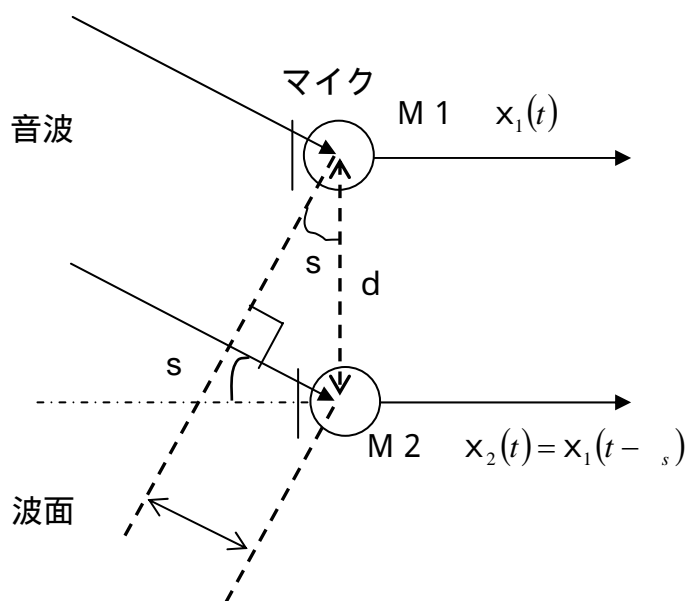


図 1 : 音波とマイクロホンの位置関係

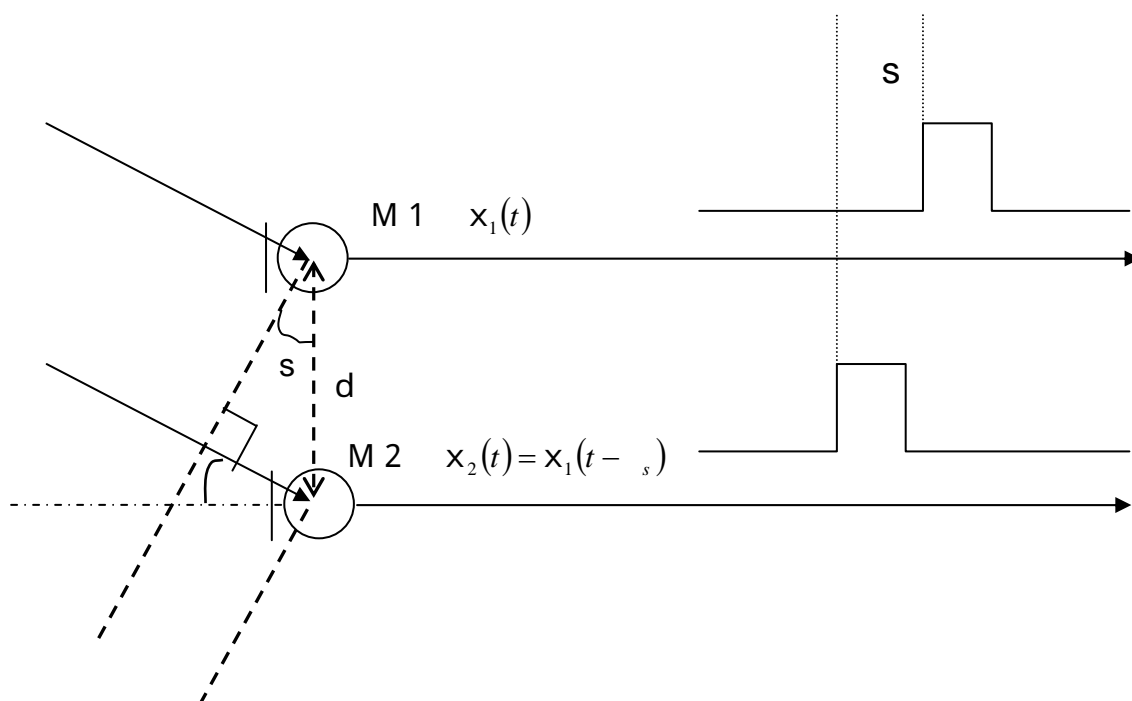


図 2 : 受信信号と遅延時間の関係

まず s から到来した音波はマイクロホンM 1 において録音される。次に距離を進んだぶん時間 s 遅れてマイクロホンM 2 に到達し録音される。この距離差は

$$= d \sin \theta_s \quad \dots$$

と表すことができる。図2からも分かるようにM 2 での受信信号 $x_2(t)$ はM 1 での受信信号 $x_1(t)$ と比べて音波が距離 $d \sin \theta_s$ だけ進むのにかかった時間 s だけ遅れた信号になっている。この時間差 s は音速を c とすると

$$s = (d \sin \theta_s) / c \quad \dots$$

となる。式と式より

$$\theta_s = \sin^{-1}(c \cdot s / d) \quad \dots$$

の関係式が成り立つ。したがって録音される時間差 s が分かれば音の到来方向 s つまり音源方向を割り出すことができる。

時間差 s は受信信号 $x_1(t)$ 、 $x_2(t)$ との相互相関関数 $r_{12}(\tau)$ から求めることができる。

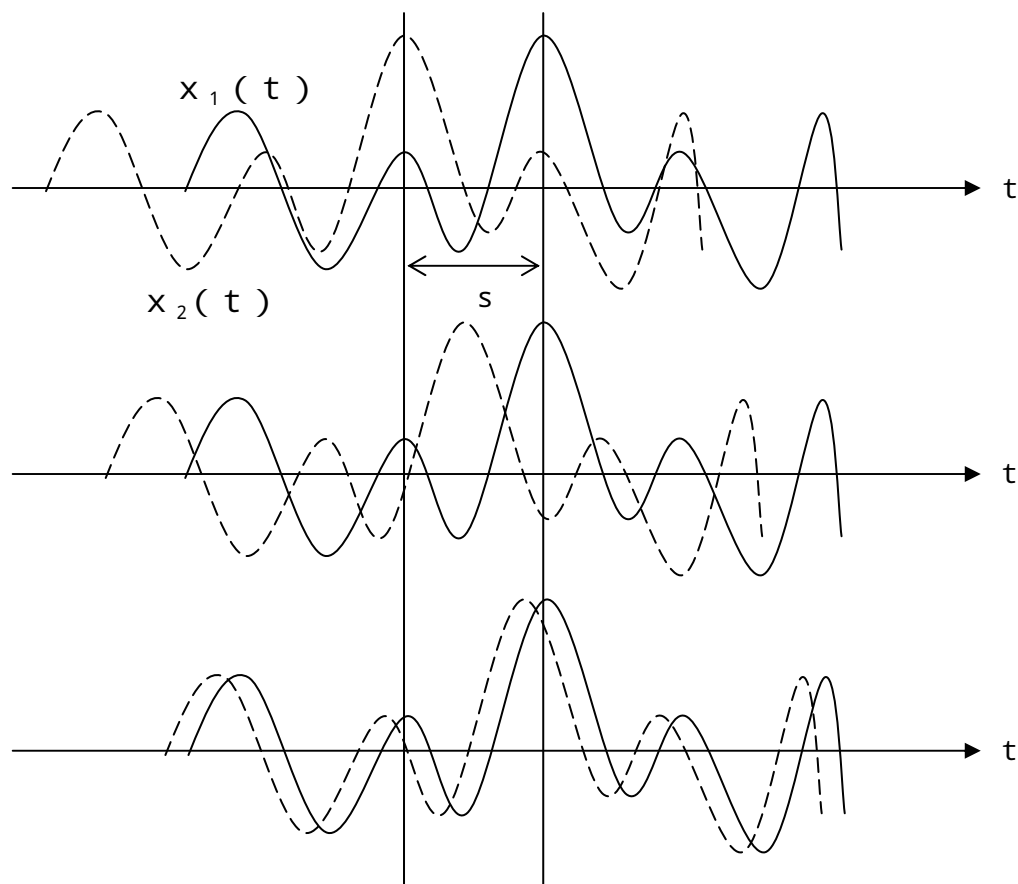


図3：相互相関関数の考えた

$\phi_{12}(\tau)$ の定義式

$$\begin{aligned}\phi_{12}(\tau) &= (1/N) \sum \{x_1(t) \cdot x_2(t + \tau)\} \\ &= (1/N) \sum \{x_1(t) \cdot x_1(t + \tau - \tau_s)\} \quad \dots\end{aligned}$$

において $\tau = \tau_s$ となる時に2つの受信信号 $x_1(t)$ 、 $x_2(t)$ は同一波形となり、最大となるので $\phi_{12}(\tau)$ の最大値を与える τ を求めれば時間差 τ_s が得られる。

3 - 2 . 音源方向検出のプログラム

実験条件としてマイクロホンM1、M2間の距離を17[cm]とし、角度はマイクロホンより人に向かい右側が-、左側を+としマイクロホンM1、M2を回転させることによって $\pm 90[^\circ]$ の範囲で変化させ測定を行う。

音源方向検出のプログラムを以下に示す。

```
%パラメータの設定
fs=48000          %サンプリング周波数
c=340            %音速
d=0.17          %マイクロホン M1,M2 間の距離

[y,fs]=wavread('d:\検出\record1\r45-2.wav');
          %音声を読み込んでくる

Ctsum=zeros(100,1); %100×1の行列のCtsumを作る

ysou=y(:,1);
xt=y(:,2);

taumax=30          % taumax を 30 とする(1)

%相互相関関数の計算(2)

for tau= - taumax : taumax

    Ctsum(tau+31)=xt(taumax+1+tau : taumax+1+10000+tau)
    *ysou(taumax+1 : taumax+1+10000)/10000 ;

end

figure(1)
plot(- taumax : taumax, Ctsum(1 : taumax*2+1));
[yjiku xjiku]=max(Ctsum) %表示されたグラフの最大値を示す
          x 軸の値を xjiku,
          y 軸の値を yjiku とする

xjiku=xjiku - (taumax+1) % x 軸の値 0 が tmax+1 番目であるため
```



```
xlabel('時間 [サンプル]') % x 軸のラベル表示
ylabel('相互相関関数') % y 軸のラベル表示
```

```
whos
```

```
hoko=asin(c*xjiku/Fs/d); %時間差から角度[ラジアン]を求める
```

```
'kaku=' ,hoko/pi*180 %単位を[ラジアン]から[°]に変換
```

(1) t_{\max} を 30 とした理由

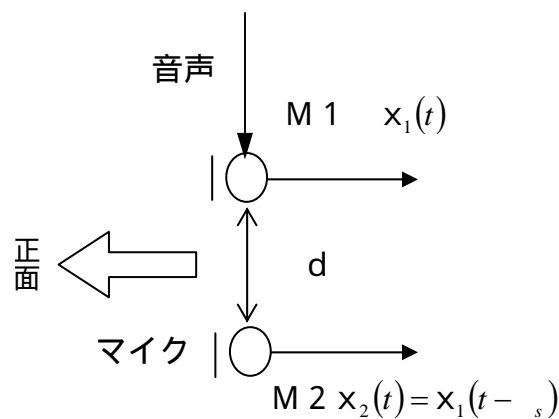


図 4 : が最大になる時の音源とマイクロホンの位置関係

マイクロホン M1 での受信信号 $x_1(t)$ とマイクロホン M2 での受信信号 $x_2(t)$ の時間差 s が最大値 (s_{\max}) をとるときは、図 4 にあるように音波が正面方向に対して $s = 90[^\circ]$ つまり直角に到来してきた場合である。

$$s_{\max} = d/c$$

$$= d/c \quad [\text{秒}] \quad \dots$$

上式の単位は[秒]となるので、[サンプル]にするにはサンプリング周波数 f_s をかければよいので

$$s_{\max} = s_{\max} \times f_s$$

$$= d/c \times fs \quad [\text{サンプル}] \quad \dots$$

実験条件と式より

$$\max = (0.17/340) \times 48000 = 24 \quad [\text{サンプル}]$$

となる。式の相互相関関数 $r_{12}(\tau)$ は、 τ で最大値をとることが分かっているため、 $r_{12}(\tau)$ が最大になる値は、 \max より大きくなることはない。したがってプログラム中の変数 τ_{\max} は 24 を少し上回る値である 30 に設定した。

(2) 相互相関関数とは

相互相関関数とは異なる 2 つの信号波形に対して相関処理を行い、2 つの信号波形の類似度を調べることである。N 個のサンプリング値からなる信号波形 $x_1(t)$ と遅れ時間 τ だけずらした信号波形 $x_1(t+\tau)$ においてそれぞれ対応する部分をかけ合わせて累積し、平均化する処理を行う。

3 - 3 . 音源方向検出の例

先に示したプログラムを実行すると次に示すグラフが得られる。

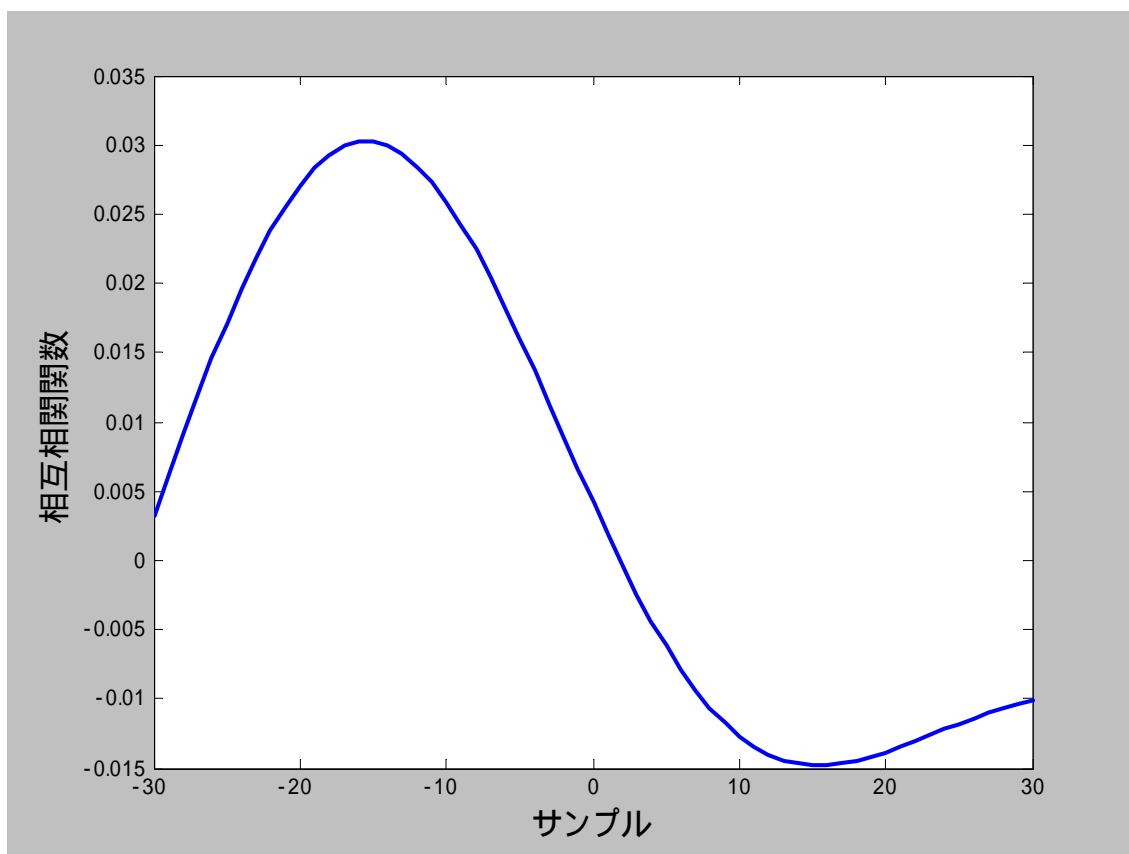


図 5 : 相互相関関数のグラフ

相互相関関数が最大となる時のx軸の値が - 16 となっているので音源方向はマイクロホンより人に向かい正面を $0 [^\circ]$ として、左側に $41.81 [^\circ]$ となった。

この実験において使用した音声は左側 $45 [^\circ]$ より録音したものを使用した。表 1 より分かるように $45 [^\circ]$ はサンプルが整数で表すことができない。そこで $45 [^\circ]$ に一番近い値は $17 [サンプル]$ となるが次に近い値が $16 [サンプル]$ の $41.81 [^\circ]$ となるので、評価としてはまずまずの結果と言えるだろう。

3 - 4 . サンプルと検出角度の関係

サンプルと検出角度の関係を計算式により求める。

$$\theta_s = \sin^{-1}(c \cdot \tau_s / d)$$

τ_s はアナログ表示[秒]なので、デジタル表示である(taus)[サンプル]にするにはサンプリング周波数 fs をかければよいので

$$\begin{aligned} \text{taus} &= \tau_s \cdot fs \\ \tau_s &= \text{taus} / fs \end{aligned}$$

したがって

$$\theta_s = \sin^{-1}(c \cdot \text{taus} / fs / d) \quad \dots$$

ここで θ_s はラジアン計算している。よって求めたい θ は

$$\theta = \theta_s / \pi 180$$

となるので

$$\theta = \sin^{-1}(c \cdot \text{taus} / fs / d) / \pi 180 \quad \dots$$

で求めることができる。

実験条件であるパラメータを

fs=48000	(サンプリング周波数)
c=340	(音速)
d=0.17	(マイクロホン M1,M2 間の距離)

とし、求めた値をグラフにし次ページに示す。

表 1 : サンプルと検出角度の関係

サンプル	角度[°]
0.00	0.00
1.00	2.39
2.00	4.78
3.00	7.018
4.00	9.59
5.00	12.02
6.00	14.48
7.00	16.96
8.00	19.47
9.00	22.02
10.00	24.62
11.00	27.28
12.00	30.00
13.00	32.80
14.00	35.69
15.00	38.68
16.00	41.81
16.96	45.00
17.00	45.10
18.00	48.60
19.00	52.34
20.00	56.44
20.78	60.00
21.00	61.05
22.00	66.44
23.00	73.40
24.00	90.00

サンプルと角度の関係は表 1 に示したようになったが、ここでは 90[°]を 24 等分(180[°]を 47 等分)にしたため、細かい角度まで求めることができなかつた。また、角度が ±90[°]に近くなると 1[サンプル]あがるごとに角度の開きが大きくなるので、人間の感覚と同じように正面付近はほぼ正確に方向検出をできだが ±90[°]付近は誤差が大きく出ることがおおかつた。

4 . 実音場における音源方向検出実験

4 - 1 . 近距離における方向検出

録音した音声の方向を検索する上でプログラムの精度を検証する必要がある。そこでマイクロホンから距離をおかずに近距離で録音し部屋の反射波やノイズ等の影響があまり出ないような状況で方向検出をした。

実験条件：マイクロホンと話者の位置関係を図6に示す。2つのマイクロホンM1、M2の距離dはそれぞれ17[cm]とし、角度 θ はマイクロホンより人に向かい右側が-、左を+としマイクロホンM1、M2を回転させることで θ を $\pm 90[^\circ]$ の範囲で変化させて測定した。録音したデータは“あ”という母音のみとし、データの長さは1[秒]、1回目と2回目、3回目の録音は同一人物が行う事とする。

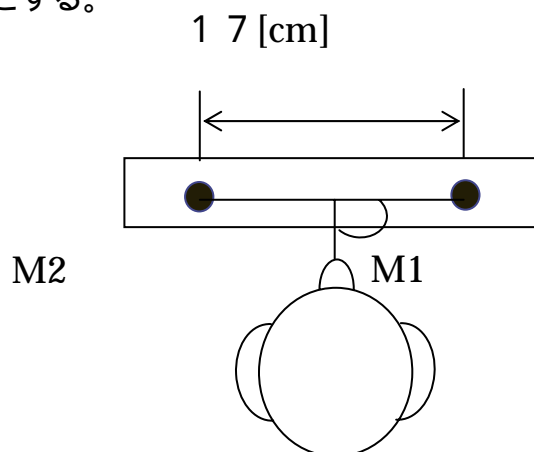


図6：音源とマイクロホンの関係

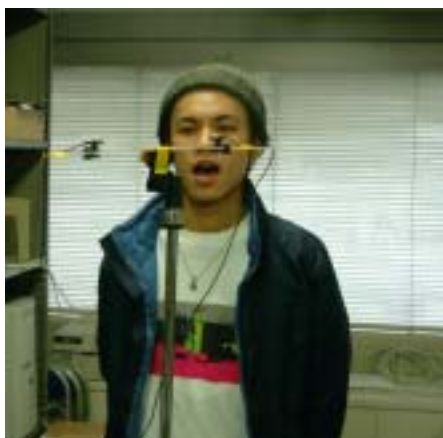


図7：録音風景

表 2 : 音源方向と検出方向の関係

音源方向 [°]	検出方向 [°]
0 1回目	4.7802
右 30 1回目	-24.6243
右 30 2回目	-27.2796
左 30 1回目	30
左 30 2回目	32.7972
右 45 1回目	-35.6853
右 45 2回目	-45.0995
左 45 1回目	45.0995
左 45 2回目	41.8103
右 60 1回目	-52.3415
右 60 2回目	-90
左 60 1回目	56.4427
左 60 2回目	56.4427
右 90 1回目	-90
右 90 2回目	-73.4022
右 90 3回目	-56.4427
左 90 1回目	90
左 90 2回目	90
左 90 3回目	90+39.7144i

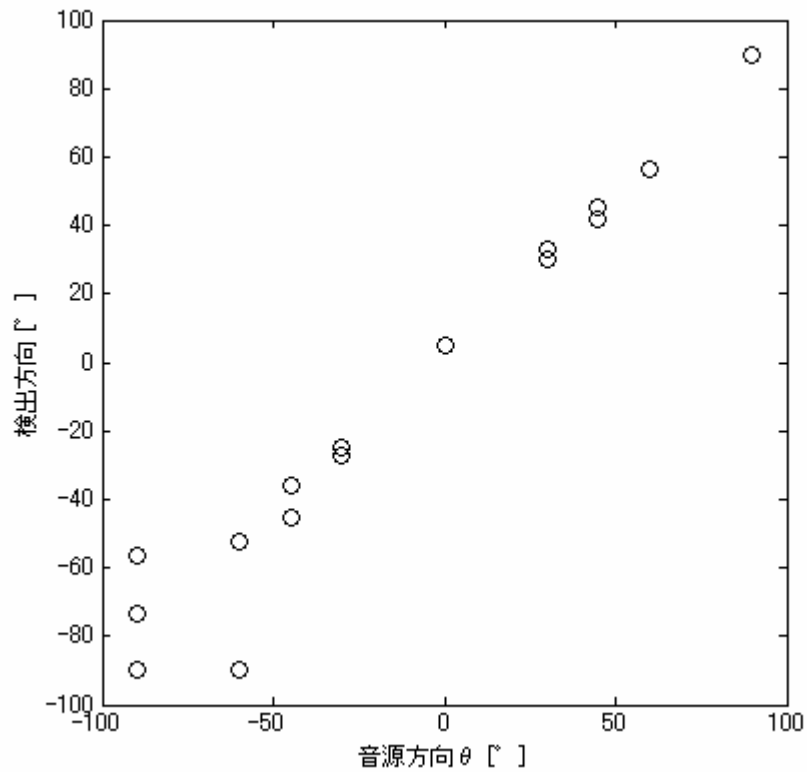


図 8 : 音源方向と検出方向の対照図

図 6 に示すように近距離において録音した時の結果を表 2 と図 8 に示す。評価としては表 2 からも見取れるようにさすがに近距離からの録音では、音源方向とプログラムで検出された方向が近い値をとっている。図 8 においてもプロットされた場所がほぼ直線に近い形になった。

4 - 2 . 周期音 (母音) による距離をつけた実験

マイクロホンに対していくつかの音源方向 をつけそれに距離をつけて録音する。

実験条件 : マイクロホンと話者の位置関係を図 9 に示す。2 つのマイクロホン M 1、M 2 の距離 d は 17[cm] とし、角度 θ はマイクロホンより人に向かい右側が -、左を + としマイクロホン M 1、M 2 を回転させることで θ を $\pm 90[^\circ]$ の範囲で変化させて測定した。マイクロホンから人までの距離 r をそれぞれ 1 [m] 2 [m] 3 [m] として録音を行った。録音したデータは“あ”という母音のみとし、データの長さは 1 [秒]、1 回目と 2 回目の録音は同一人物が行う事とする。

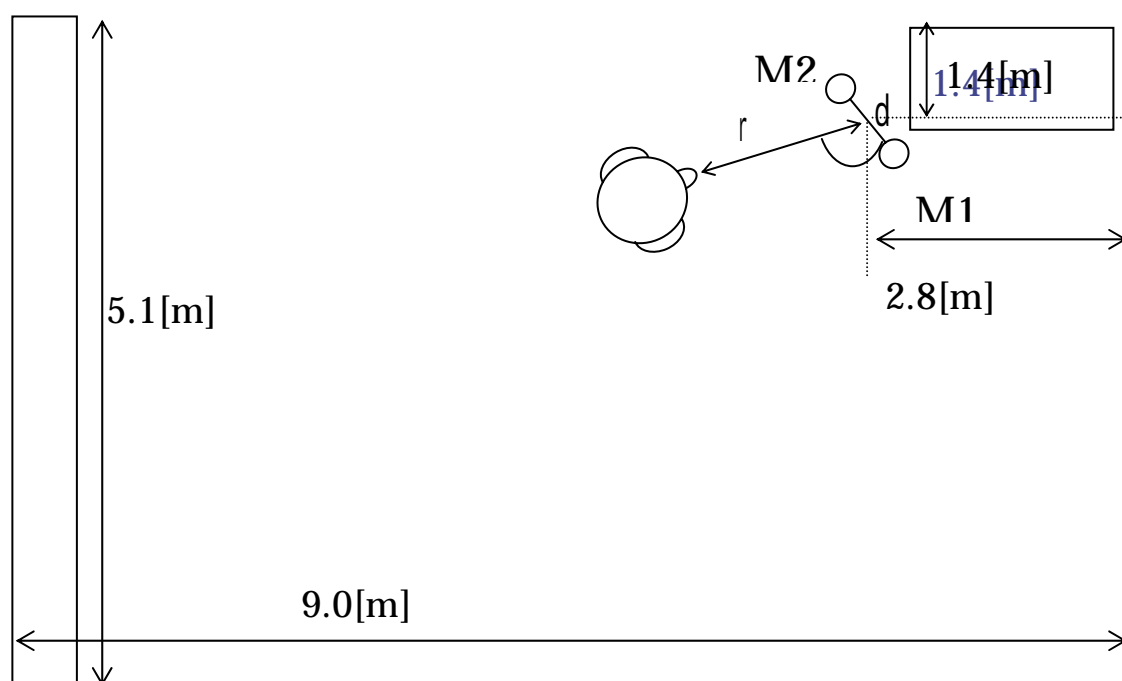


図 9 : 部屋の端における録音状況

表 3 : 音源方向 と検出方向の関係 (部屋の端)

音源方向 [°]	マイクロホンからの距離		
	1[m]	2[m]	3[m]
	検出角度 [°]		
0 1回目	2.388	4.7802	14.4775
0 2回目	4.7802	4.7802	-12.0247
右 30 1回目	-4.7802	-12.0247	-4.7802
右 30 2回目	-22.0243	-16.9578	-7.1808
左 30 1回目	16.9578	22.0243	22.0243
左 30 2回目	19.4712	14.4775	16.9578
右 45 1回目	-48.5904	-45.0995	-14.4775
右 45 2回目	-45.0995	-56.4427	-56.4427
左 45 1回目	24.6243	30	90
左 45 2回目	45.0995	41.8103	90+16.48i
右 60 1回目	-56.4427	-56.4427	-12.0247
右 60 2回目	-52.3415	-45.0995	-66.4435
左 60 1回目	35.6853	48.5904	30
左 60 2回目	35.6853	45.0995	48.5904
右 90 1回目	-56.4427	-48.5904	-32.7972
右 90 2回目	-66.4435	-66.4435	-9.5941
左 90 1回目	73.4022	56.4427	24.6243
左 90 2回目	90+28.35i	90+23.23i	14.4775

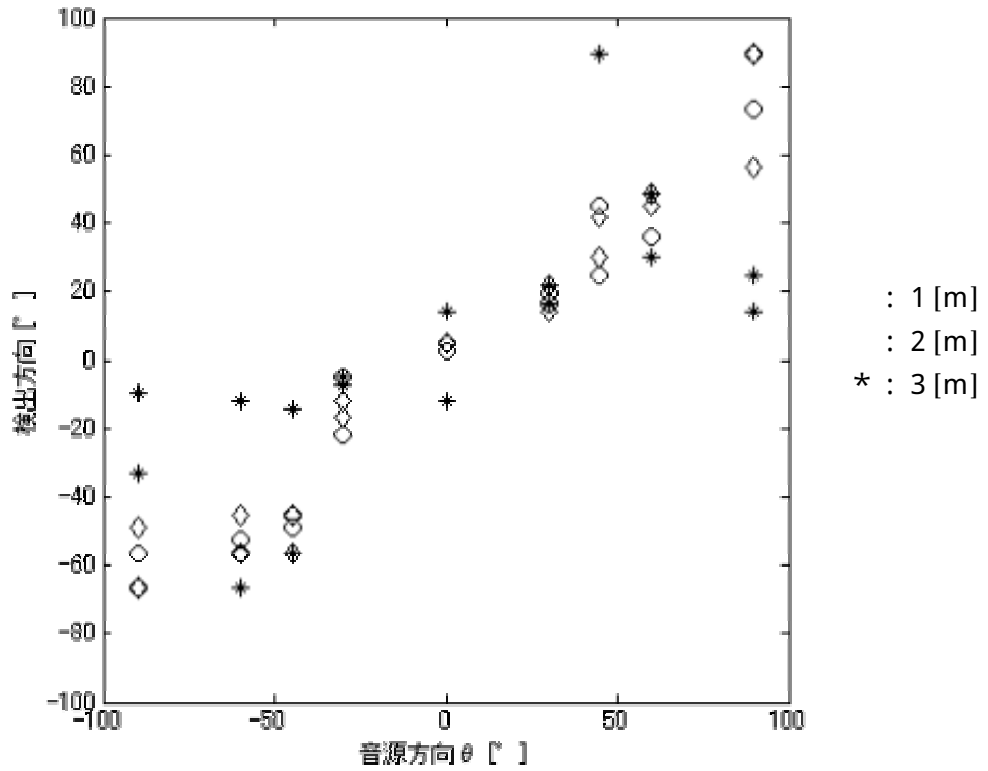


図 10：音源方向 と検出方向の対照図（部屋の端）

図 9 に示すようにマイクロホンを部屋の端に置いて録音した時の結果を表 3 と図 10 に示す。図 10 から分かるように距離 r が離れるにしたがって検出方向がずれてくる。ここで分かるのは 3 [m] の場合が一番分かりやすいが、 $\pm 90[^\circ]$ に近づくと検出方向が $0[^\circ]$ に近づいている。

マイクロホンの場所を変えて同様に録音する。

実験条件：マイクロホンと話者の位置関係は図 11 に示す通りである。このときは1回目と2回目の録音は別の人物が録音した。その他の条件は実験と同様。

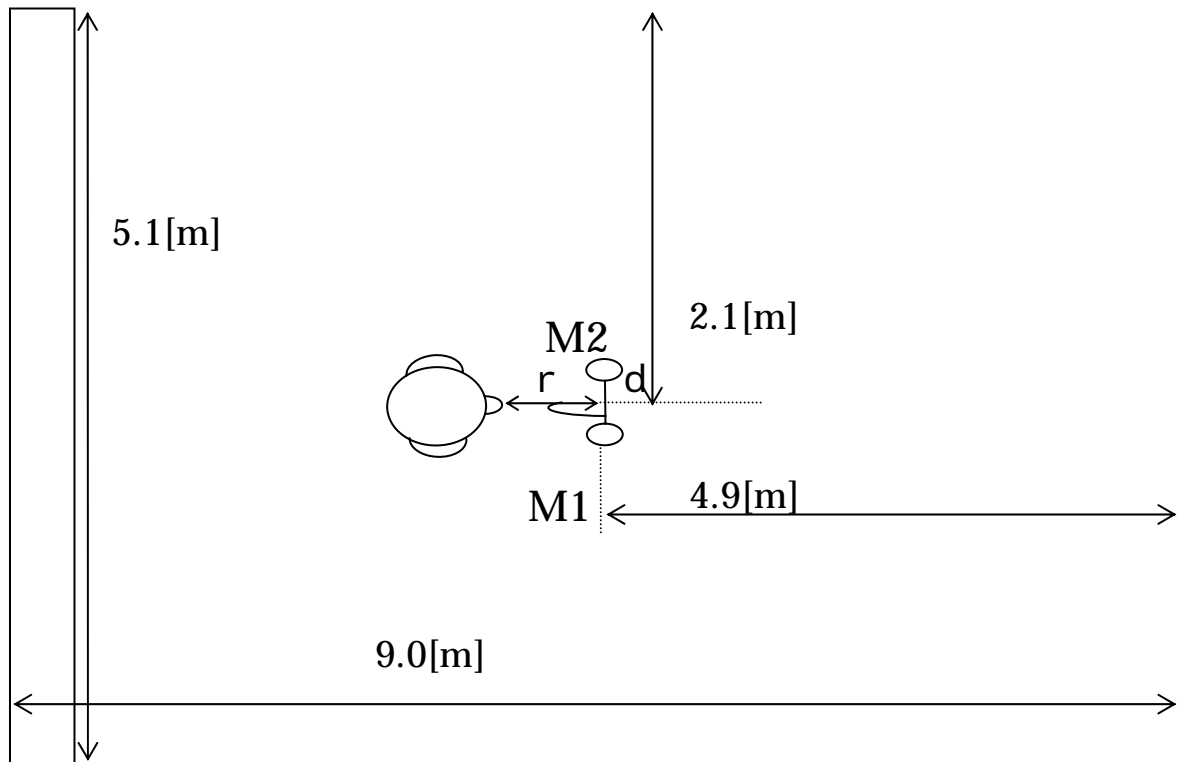


図 11：部屋の中心における録音状況

表 4 : 音源方向 と検出方向の関係 (部屋の中心)

音源方向 [°]	マイクロホンからの距離		
	1[m]	2[m]	3[m]
	検出角度 [°]		
0 1回目	4.7802	2.388	14.4775
0 2回目	2.388	2.388	0
左 30 1回目	22.0243	19.4712	90+42.76i
左 30 2回目	14.4775	16.9578	12.0247
右 30 1回目	-19.4712	-14.4775	-7.1803
右 30 2回目	-4.7802	-2.388	90+39.71i
左 45 1回目	27.2796	27.2796	73.4022
左 45 2回目	32.7972	24.6243	14.4775
右 45 1回目	-19.4712	-30	-30
右 45 2回目	-14.4775	-9.5941	-30
左 60 1回目	38.6822	32.7972	14.4775
左 60 2回目	24.6243	35.6853	4.7802
右 60 1回目	-35.6653	-27.2796	-19.4712
右 60 2回目	-38.6822	-24.6243	-12.0247
左 90 1回目	48.5904	16.9578	27.2796
左 90 2回目	41.8103	14.4775	7.1808
右 90 1回目	-22.0243	-24.6243	-9.5941
右 90 2回目	-48.5904	-45.0995	90+42.76i

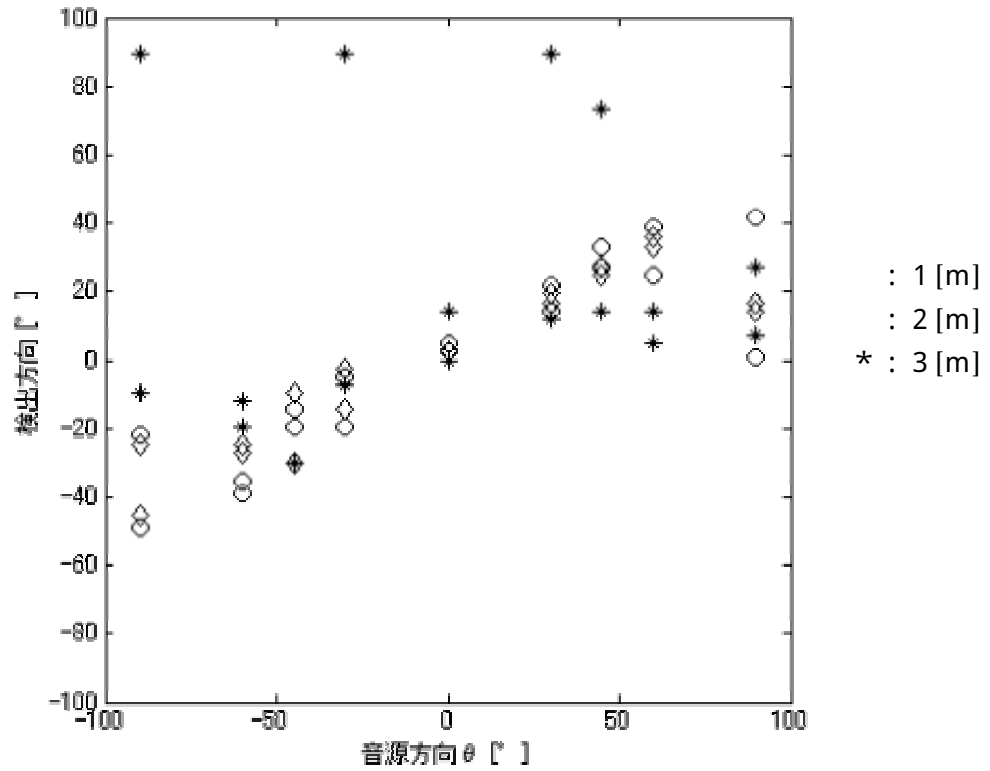


図 12：音源方向 と検出方向の対照図（部屋の中心）

図 11 に示すように部屋の中心にマイクロホンを置いた時の結果を表 4 と図 12 に示す。部屋の端に置いた時（図 10）と比べると全体的に $0 [^\circ]$ に近い検出結果となっている。

4 - 3 . 不規則な波形の音を使用し距離をつけた実験

今までの録音では発声する音を「あ」等の母音（周期音）にしていた。しかし周期音で録音すると1周期ずれていてもそのずれを認識することができないという曖昧さがある。そこで濁音等の波形が不規則な音を使用し周期音の時と同様に部屋の端と中心の2通りの実験結果を図13、14に示す。

実験条件：マイクロホンと話者の位置関係は図9，12に示す通りである。
その他の条件は実験と同様。

表5：濁音時の音源方向 と検出方向の関係（部屋の端）

音源方向 [°]	マイクロホンからの距離		
	1[m]	2[m]	3[m]
	検出角度 [°]		
0 1回目	-4.7802	-4.7802	4.7802
0 2回目	-7.1808	-7.1808	-4.7802
右 30 1回目	-45.0995	-19.4712	24.6243
右 30 2回目	-41.8103	-32.7972	14.4775
左 30 1回目	61.045	41.8103	45.0995
左 30 2回目	56.4427	32.7972	-35.6853
右 45 1回目	-14.4775	-27.2796	-30
右 45 2回目	-24.6243	-27.2796	-35.6853
左 45 1回目	27.2796	24.6243	24.6243
左 45 2回目	32.7972	22.0243	35.6853
右 60 1回目	-22.0243	-16.9578	-30
右 60 2回目	-19.4712	-16.9578	-27.2796
左 60 1回目	16.9578	22.0243	16.9578
左 60 2回目	19.4712	22.0243	16.9578
右 90 1回目	-52.3415	-16.9578	-48.5904
右 90 2回目	-52.3415	-27.2796	-16.9578
左 90 1回目	56.4427	56.4427	19.4712
左 90 2回目	32.7972	41.8103	32.7972

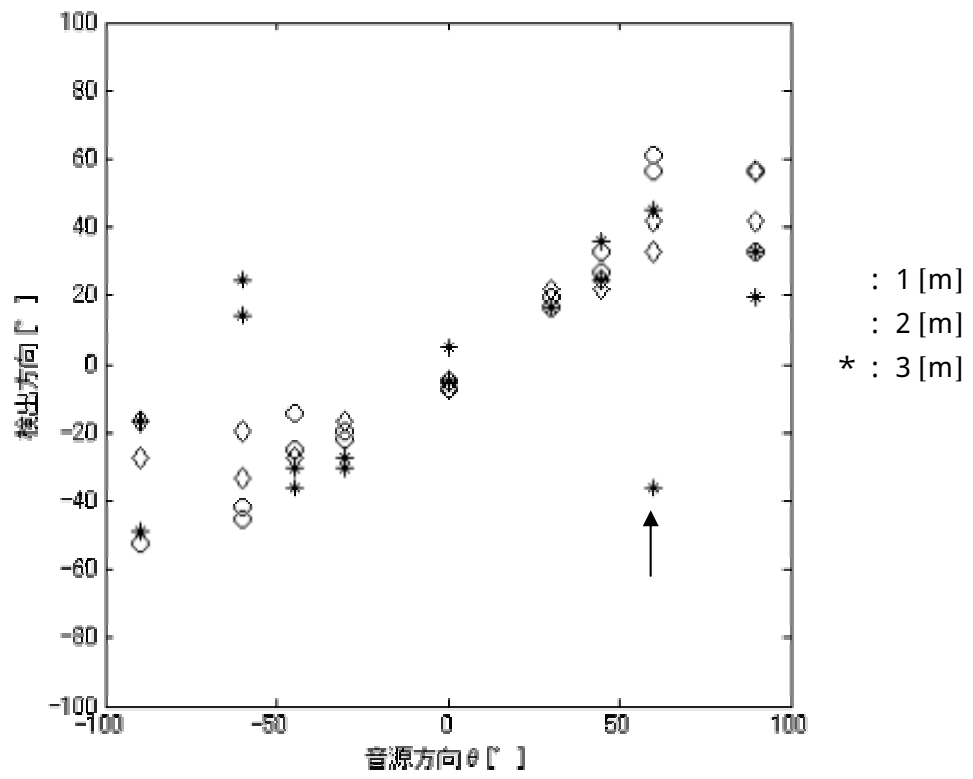


図 13 : 濁音時の音源方向 と検出方向の関係 (部屋の端)

表 6 : 濁音時の音源方向 と検出方向の関係 (部屋の中心)

音源方向 [°]	マイクロホンからの距離		
	1[m]	2[m]	3[m]
	検出角度 [°]		
0 1回目	4.7802	4.7802	4.7802
0 2回目	4.7802	2.388	7.1808
右 30 1回目	-16.9578	-12.0247	-14.4775
右 30 2回目	-14.4775	-9.5941	-9.5941
左 30 1回目	35.6853	32.7972	27.2796
左 30 2回目	32.7972	19.4712	24.6243
右 45 1回目	-24.6243	-27.2796	-19.4712
右 45 2回目	-19.4712	-14.4775	-14.4775
左 45 1回目	41.8103	38.6822	32.7972
左 45 2回目	48.5904	35.6853	30
右 60 1回目	-38.6822	-35.6853	-30
右 60 2回目	-38.6822	-30	-24.6243
左 60 1回目	56.4427	45.0995	32.7972
左 60 2回目	45.0995	48.5904	45.0995
右 90 1回目	-52.3415	-38.6822	-27.2796
右 90 2回目	-61.045	-38.6822	-32.7972
左 90 1回目	61.045	52.3415	35.6853
左 90 2回目	66.4435	48.5904	45.0995

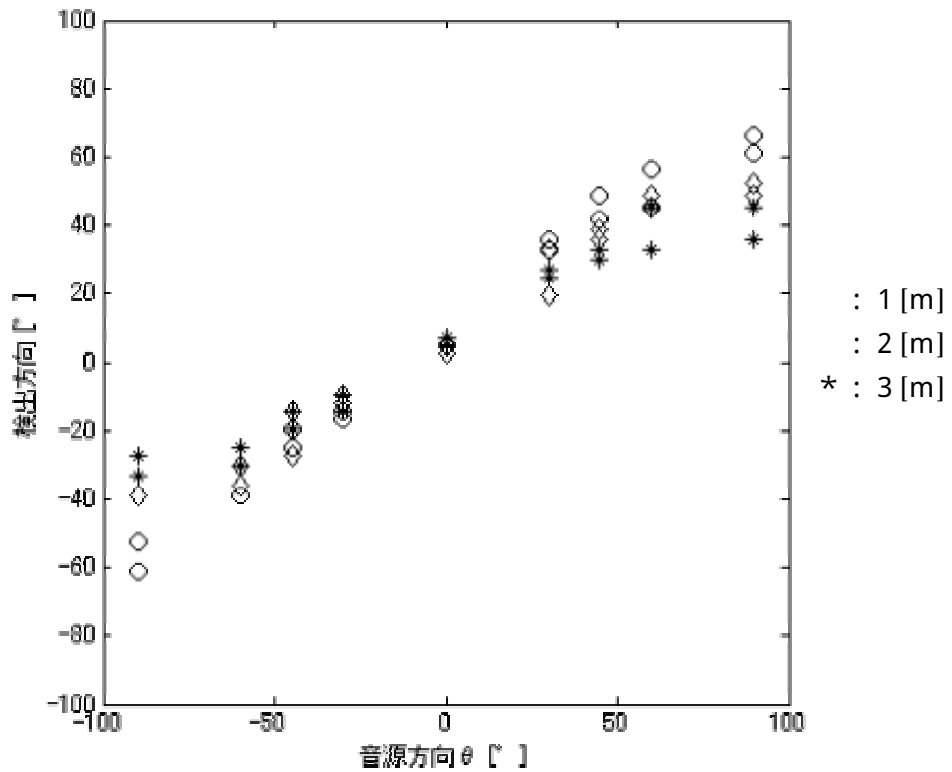


図 14：濁音時の音源方向 と検出方向の関係（部屋の中心）

図 10 と図 13、図 12 と図 14 を比較すると「あ」等の周期音と濁音等の不特定の波形ではやはり周期音の方が検出方向に散らばりがみられる。やはりこれは相互相関関数において計算される際、周期音のような信号波形より不特定の信号波形のほうが波形のずれを正確に捉える事ができるからではないかと考えられる。さらに図 13 と 14 を比較すると部屋の反射の影響が少ない部屋の中心で録音をおこなった図 14 の方がよりきれいなグラフとなった。

ここで図 13 において期待する値と違う点が数点ある。一例として図においての矢印で示した点の相関関数とサンプル数の関係を表した図が図 15 である。図 15 をみると本来右の山が高くなるはずだが左の山のほうに反射波が加わりエネルギー的に大きくなってしまいそのために正確な方向検出ができなかったと考えられる。

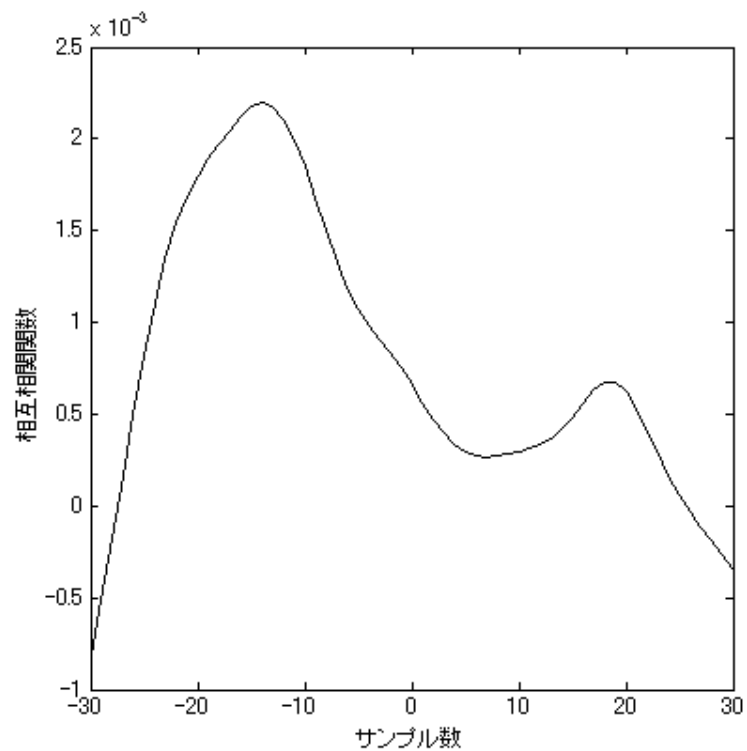


図 15 : 相関関数とサンプル数の関係

4 - 4 . 高域を強調する

今まで検出してきた周期音または濁音等の不規則な音波では、図 16 のように低域であるエネルギーの強い部分だけで検出してきた可能性がある。そこで図 17 のように高域のエネルギーの弱い部分を強調し、録音した音波のすべての部分において方向検出をするようにする白色化を試みた。

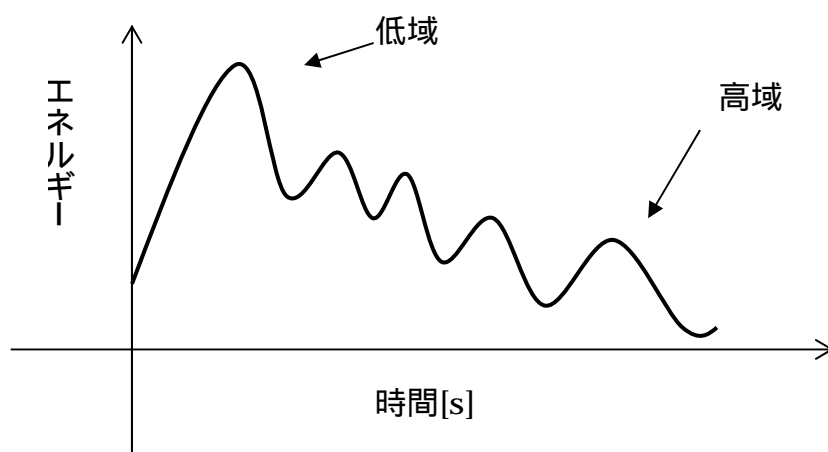


図 16 : 一般的な波形

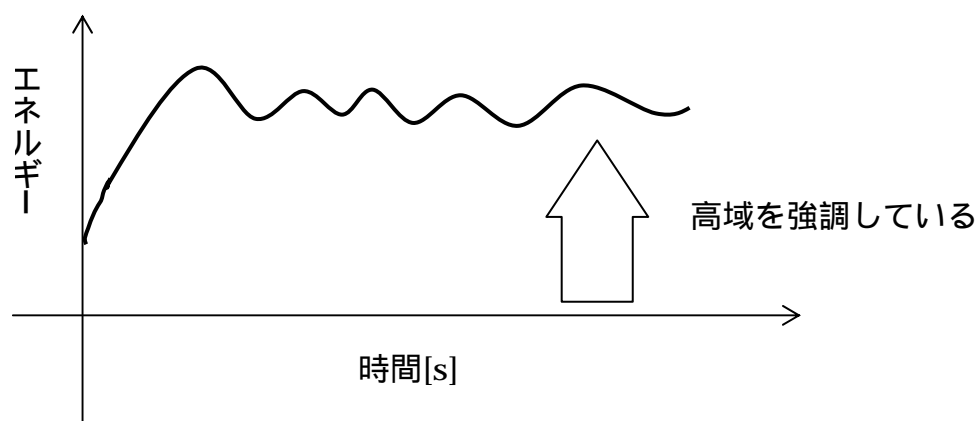


図 17 : 高域を強調した時の波形

以上のことをプログラムで表すと以下のようなになる。

%パラメータの設定

```
fs=48000          %サンプリング周波数
c=340             %音速
d=0.17           %マイクロホン M1,M2 間の距離
```

```
[ y , fs]=wavread('d:¥検出¥record2¥032.wav') ;
                %音声を読み込んでくる
```

%白色化

```
a=0.7
NN=length(y)
y0=y(2:NN,:) - a*y(1:NN-1,:);
```

```
Ctsum=zeros(100,1);    %100×1の行列のCtsumを作る
```

```
ysou=y0(:,1);          %ysouにy0の1列目を代入
xt=y0(:,2);            %xtにy0の2列目を代入
```

```
taumax=30              % taumax を 30 とする
```

%相互相関関数の計算

```
for tau= - taumax : taumax
```

```
    Ctsum(tau+31)=xt(taumax+1+tau : taumax+1+10000+tau)
    *ysou(taumax+1 : taumax+1+10000)/10000 ;
```

```
end
```

```
figure(1)
```

```
plot( - taumax : taumax, Ctsum(1 : taumax*2+1));
```

```
[yjiku xjiku]=max(Ctsum) %表示されたグラフの最大値を示す
                        x軸の値を xjiku,
                        y軸の値を yjiku とする
```

```
xjiku=xjiku - (taumax+1) % x 軸の値 0 が tmax+1 番目であるため
```

```
xlabel('時間 [サンプル]') % x 軸のラベル表示
```

```
ylabel('相互相関関数') % y 軸のラベル表示
```

```
whos
```

```
hoko=asin(c*xjiku/Fs/d); %時間差から角度[ラジアン]を求める
```

```
'kaku=',hoko/pi*180 %単位を[ラジアン]から[°]に変換
```

結果としては検出方向には直接変化は現れなかった。ただ相互相関関数とサンプル数のグラフにおいて図 18 と図 19 を比べた時図に示した矢印の部分に微妙な差が出た。(以下の図は表 3 における $0 [^\circ]$ 、 $3[m]$ 、2 回目のものである。)

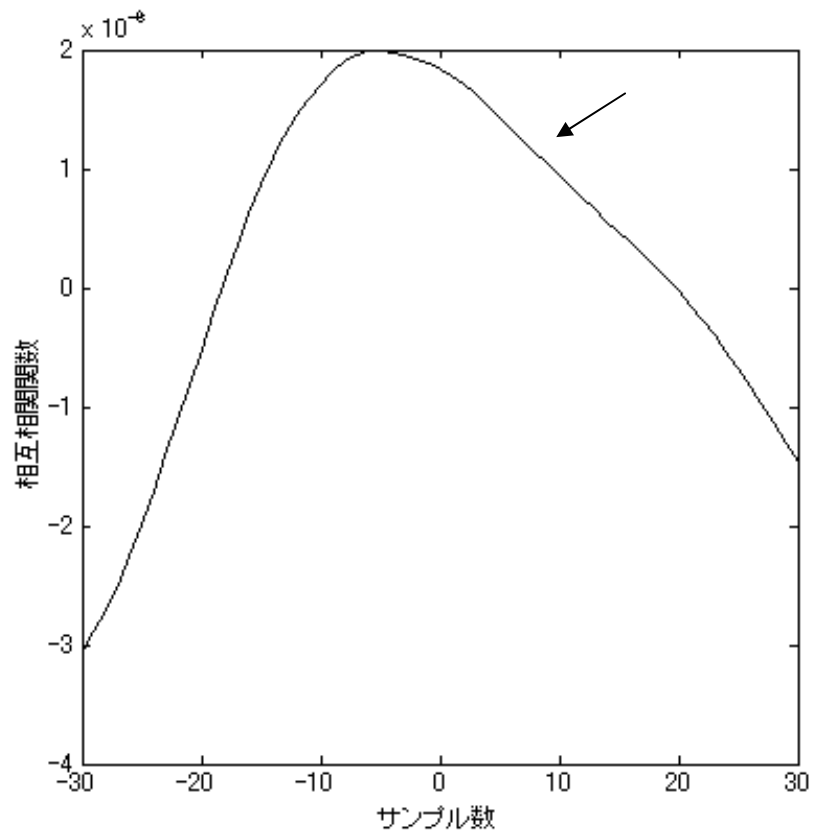


図 18：白色化していない相関関数とサンプル数の関係

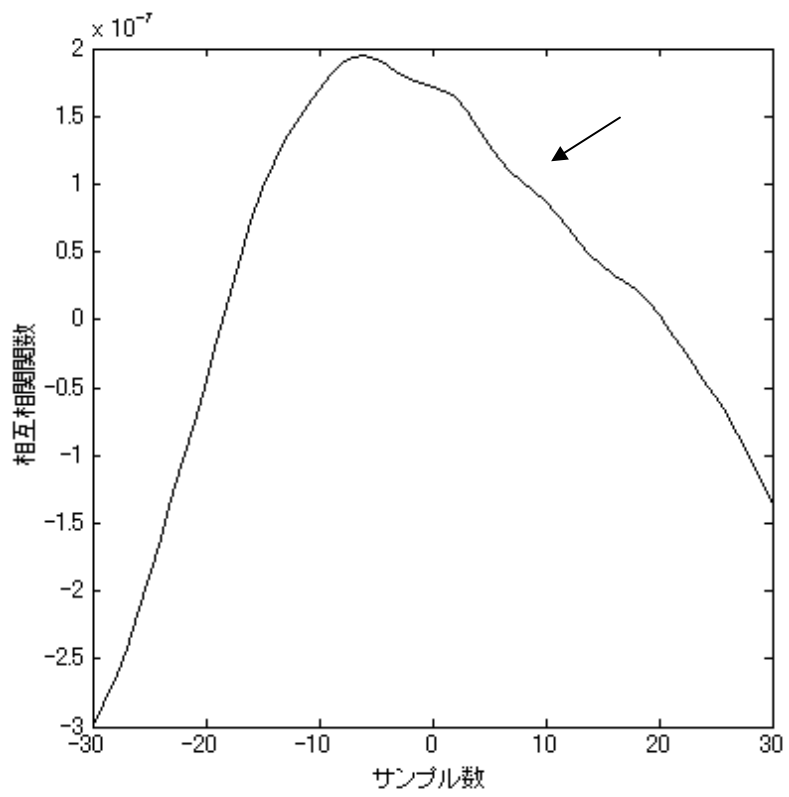


図 19：白色化した相関関数とサンプル数の関係

4 - 5 . 立体的に考える

今までの研究では、録音を平面的に行っていたが、 $90[^\circ]$ 付近では平面的に数度ずれていても検出方向に大きく影響を与えていたが、それは立体的(3次元空間的)に考えても同様のことがいえる。

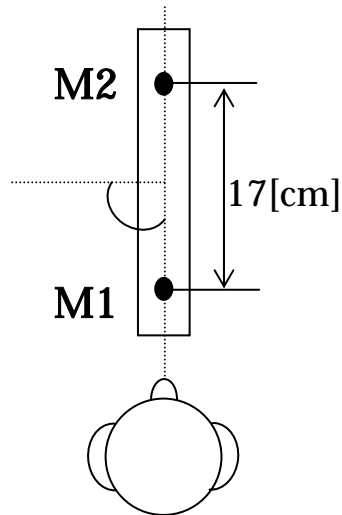


図 20 : 平面的に考えた場合

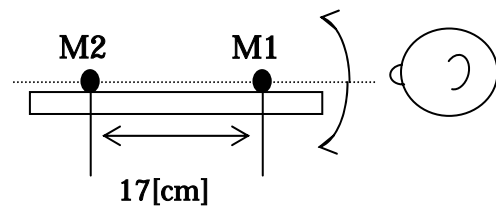


図 21 : 立体的に考えた場合

図 20,21 から分かるように平面的に見た $\theta = 90[^\circ]$ において(図 20) 音源位置が立体的に(上下に)数度ずれた場合(図 21) 今回の判定方法では、平面方向のずれと判断してしまう。

そこで立体的に見た音源方向と検出できる信号の時間差について考察してみた。具体的には図 22 のように平面上での音源から 2 つのマイクロホンへの距離の差を立体的に考察してみた。

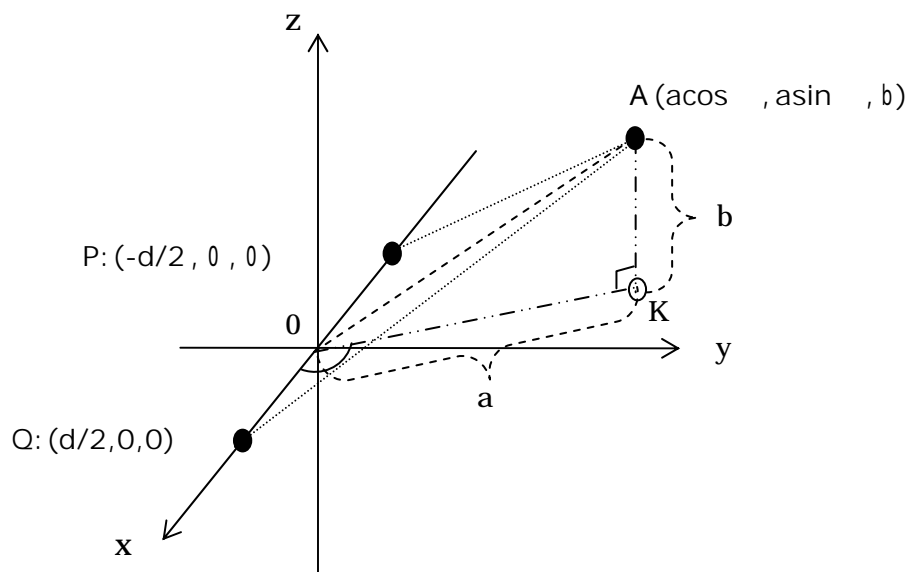


図 22 : 音源とマイクロホンの位置関係

A : 音源
P : マイクロホン M1
Q : マイクロホン M2
 : x-y 平面においての原点から音源への角度
d : 2つのマイクロホン間の距離
a : A - K間の距離
b : Aの高さ

音源 A から 2つのマイクロホン M1、M2 の距離の差を α とすると求める距離差 は以下の式で表せる。

$$\begin{aligned}\alpha &= |A-P| - |A-Q| \\ &= \sqrt{(a \cos \phi - d/2)^2 + (a \sin \phi)^2 + b^2} - \sqrt{(a \cos \phi + d/2)^2 + (a \sin \phi)^2 + b^2}\end{aligned}$$

この式をプログラムの中に使用するのだが、本研究では実際にはそこまで到達しなかった。

5 . 音声認識判断

「音声に対してロボットを反応させる時にただ音に反応するだけではなく、音声をきちんと認識することができ、それに対応させて音声別に行動をとれるようなプログラムを組んでいきたい」という目標で研究を進めた。

本研究では“あ・い・う・え・お”といった母音のみの音声認識について重点をおいて検討をした。

5 - 1 . 音声認識判断の研究方法

まずはじめに“あ・い・う・え・お”といった母音のみをそれぞれ数回録音し、MATLABに取り込み周波数スペクトルで表す。次にその周波数スペクトルを逆FFT変換してケプストラムを求める。“あ”の音声において求めたケプストラムを図に示す。

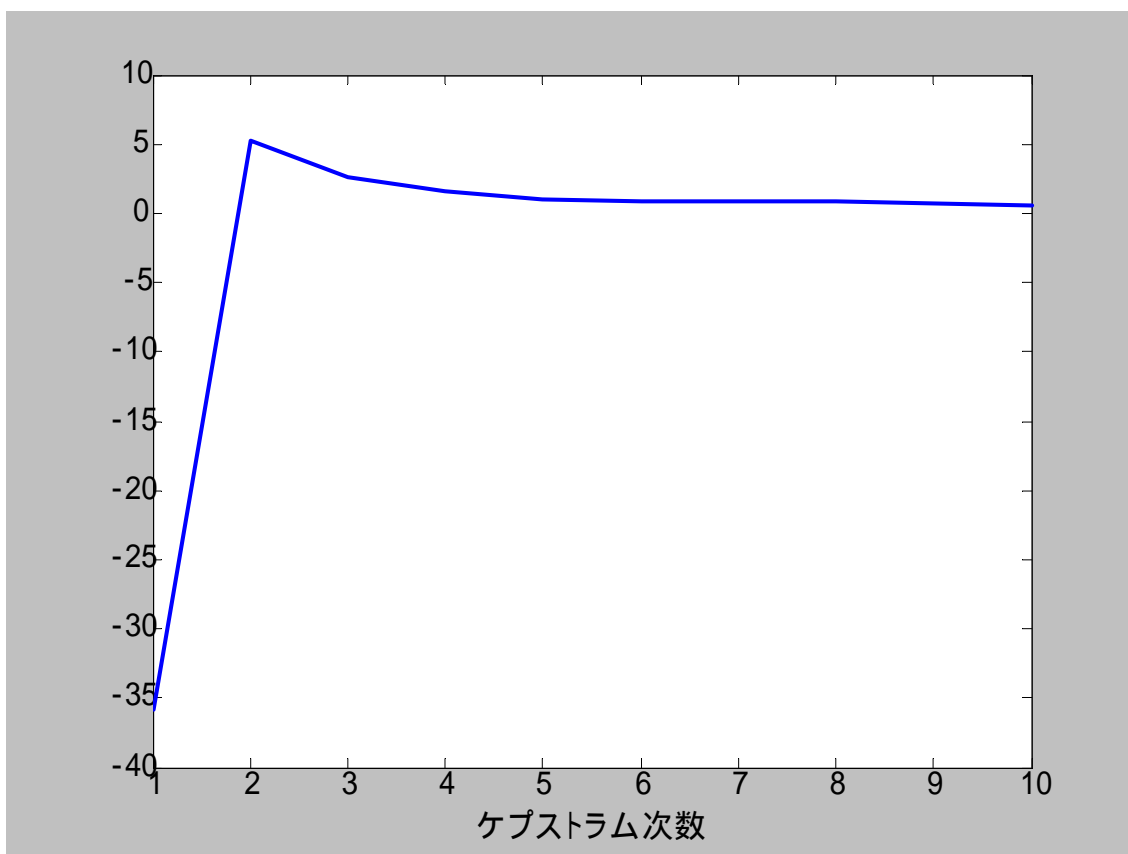


図 23 : “あ” の音声ケプストラム

ここで2次ケプストラムは傾斜なので無視をして、3次ケプストラムをx軸、4次ケプストラムをy軸としてグラフを作成し、それぞれの特徴を捉え比較する。

5 - 2 . 音声特徴量の分析

音声特徴量の分析のプログラミングを以下に示す。

```
%パラメータの設定
fs=48000;           %サンプリング周波数

[y, fs]=wavread('d:¥検出¥record6¥a5.wav');
                    %音声を読み込んでくる

Ctsum=zeros(100,1); %100×1の行列のCtsumを作る

ysou=y(:,1);       %ysouにyの1列目を代入

%FFT変換する
ip=1;
nfft0=2048;

for iii=1:20

    yfsum=zeros(nfft0,1);
    Nrep=1

    for i=1:Nrep
        y1=y(ip:ip+nfft0-1);
        yf=abs(fft(y1)).*abs(fft(y1));
        yfsum=yfsum+yfsum;
        ip=ip+nfft0;
    end

    yfsum=yfsum/Nrep;

%逆FFT変換する
NFFT=length(yfsum);
x=1:nfft0;
fx=(x-1)*(fs/2)/(NFFT/2);
```

```
ylog=10*log10(yf);  
cep1=ifft(ylog);
```

```
figure(1);  
plot(cep1(3),cep1(4),'o'); %3 次ケプストラムの値を x 軸,  
4 次ケプストラムの値を y 軸としてプロット
```

```
hold on  
xlabel('3 次ケプストラム') % x 軸のラベル表示  
ylabel('4 次ケプストラム') % y 軸のラベル表示  
axis([0 7 0 3]); % グラフの範囲指定
```

```
end
```

5 - 3 . 分析結果

図 24、25、26、27、28 は“あ・い・う・え・お”といった母音のみ 1 秒間録音したものをプログラムを実行させることにより音声を 20 個に区切りそれぞれの音声特徴を分析し、3 次ケプストラムの値と 4 次ケプストラムの値をプロットしたものである。

図 29 と図 30 は“あ・い・う・え・お”それぞれの母音を繰り返し重ねてプロットしたものである。見て分かるように 2 次ケプストラムと 3 次ケプストラムを使用して判別するより、3 次ケプストラムと 4 次ケプストラムを使用した方が分解度はあがった。

図 31 と図 32 はそれぞれ“あ・お”の音声特徴と“あ・い”の音声特徴を表したものである。2 つの図より、比較してみると“あ・い”の音声特徴の差ははっきりと見て取れるが、“あ・お”の音声特徴においてはほとんど似た音声特徴を持っており判別することは難しかった。

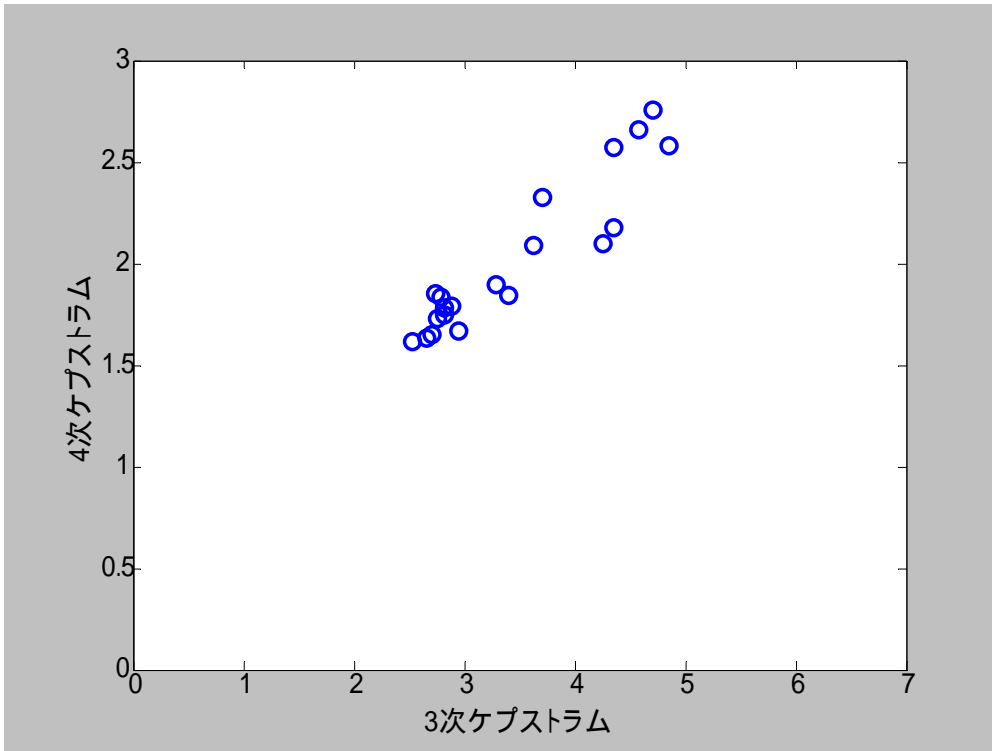


図 24 : “ あ ” の音声特徴

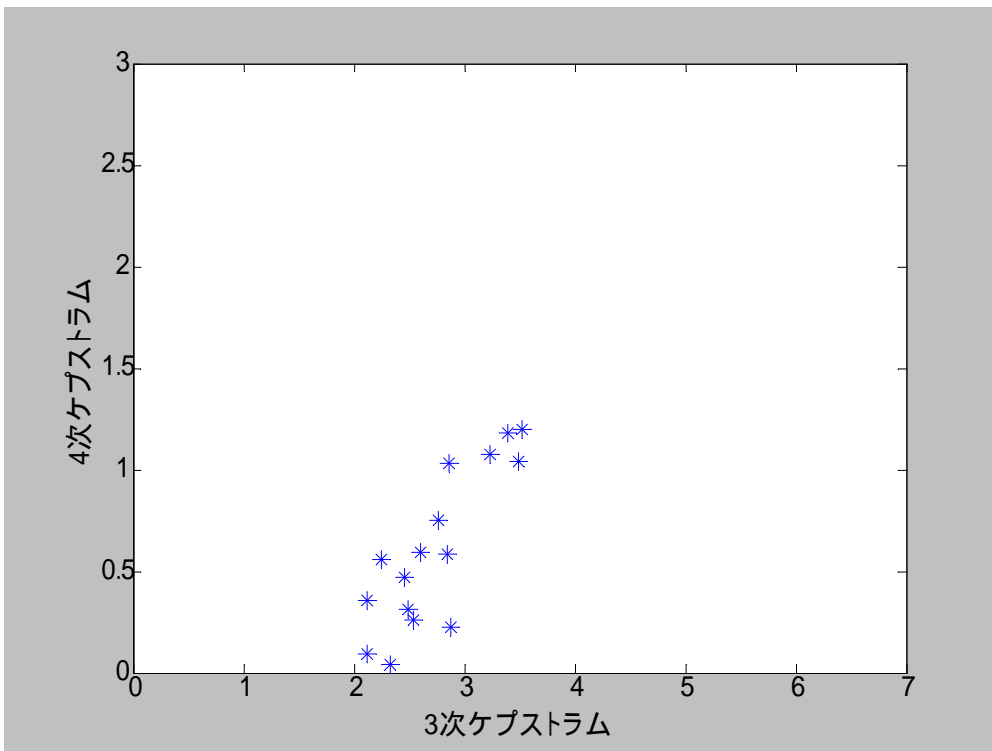


図 25 : “ い ” の音声特徴

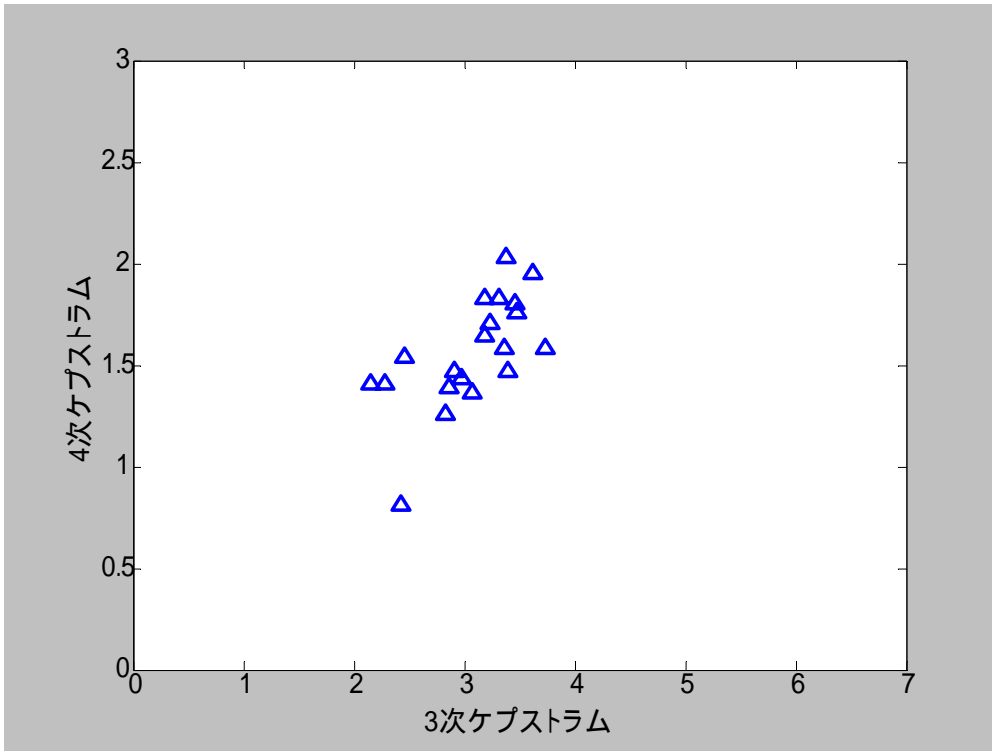


図 26：“う”の音声特徴

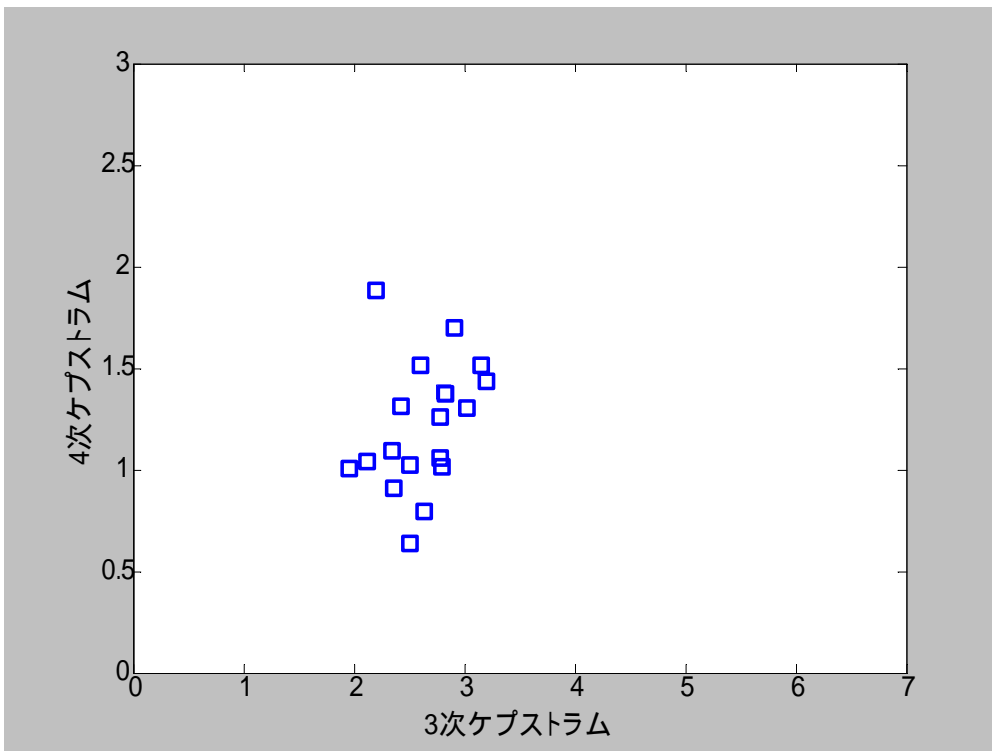


図 27：“え”の音声特徴

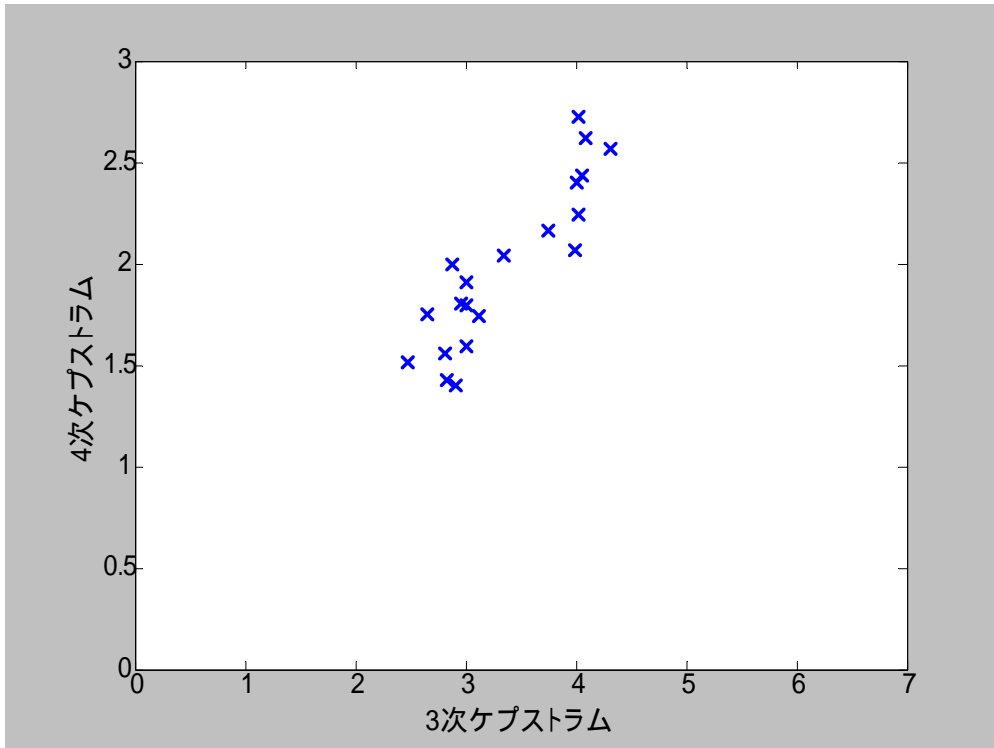


図 28：“お”の音声特徴

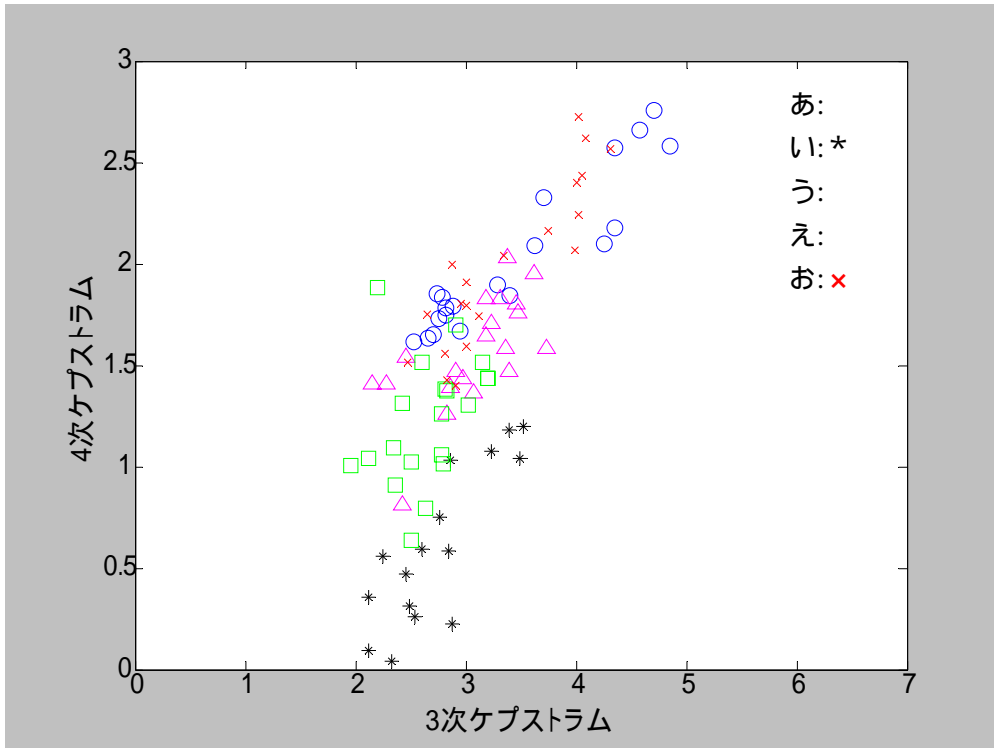


図 29 : それぞれの音声特徴(3 次,4 次ケプストラムを使用)

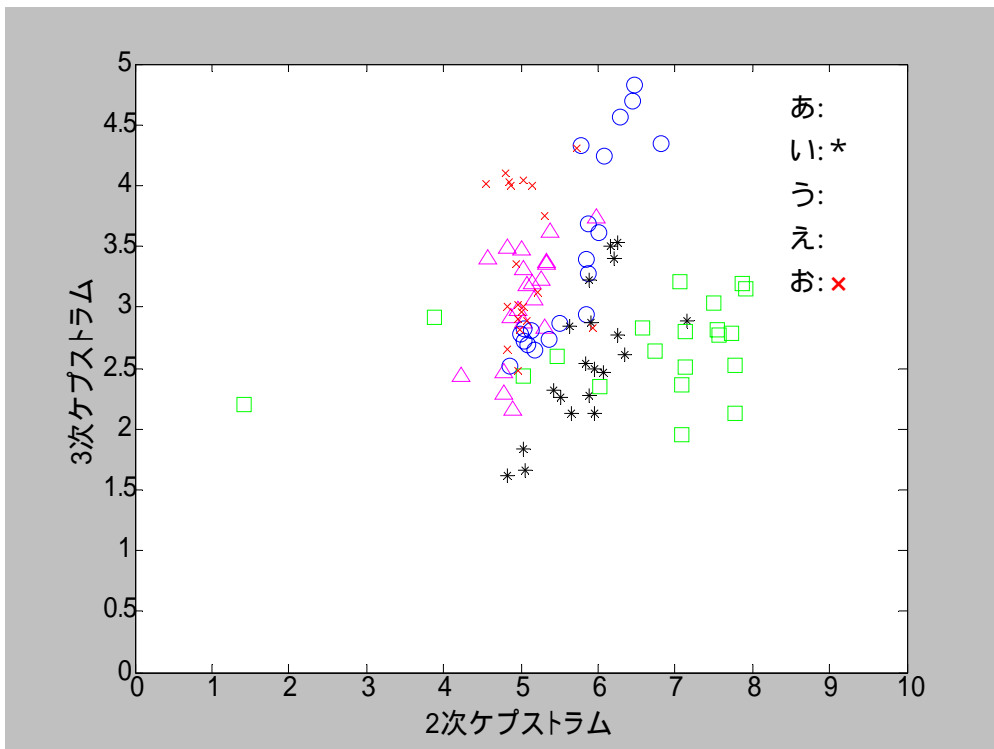


図 30 : それぞれの音声特徴(2 次,3 次ケプストラムを使用)

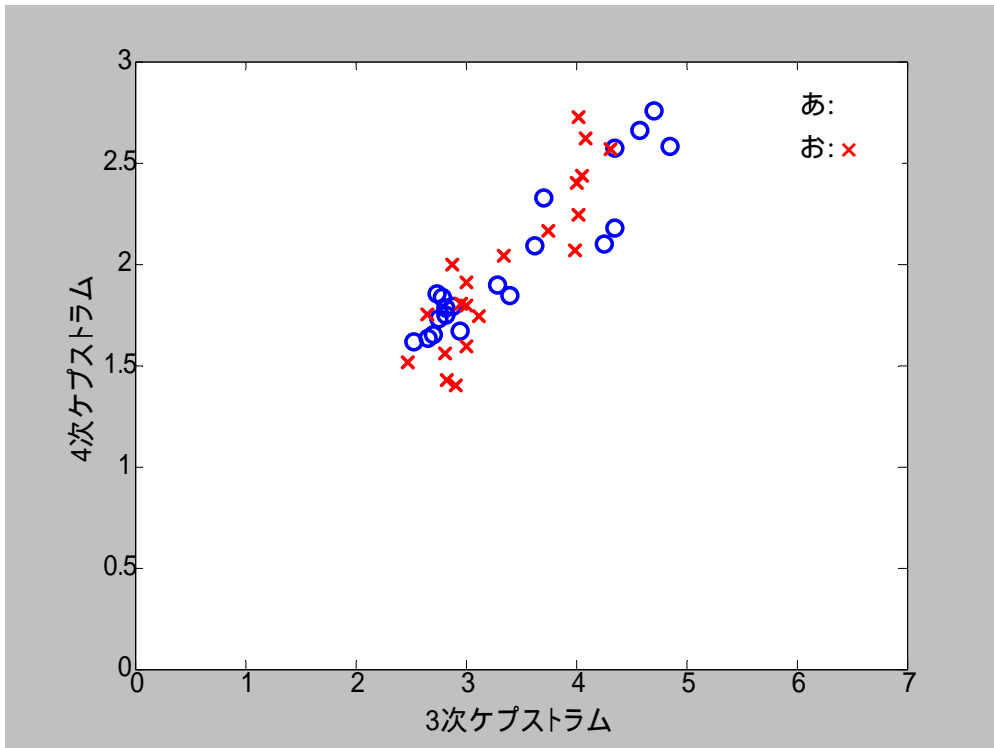


図 31 : “ あ・お ” の音声特徴

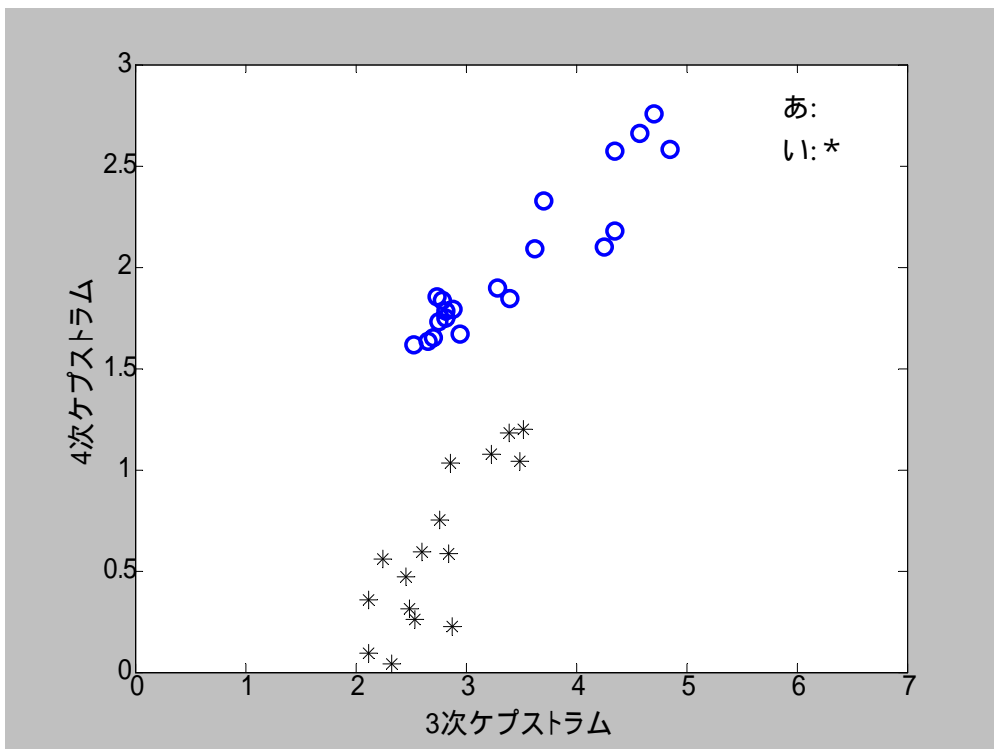


図 32 : “ あ・い ” の音声特徴

6 . ロボットの動作制御の方法*1

音声認識判断を行い命令に従うロボットはパソコンで制御可能なLEGOロボットを使用した。使用したLEGOロボットの標準ファームウェア⁽³⁾を図に示す。



図 33 : LEGOロボットのプログラミング画面

このプログラミング画面は左上のビッグブロックでロボットの動作を指定することができる。その下にあるスモールブロックではモーターのパワー調節、光やタッチセンサー機能、Loopやif文といったプログラムの流れを指定することができる。しかし肝心の音声センサーがなかった。この標準ファームウェアでは、32個の変数と10個のタスク、そして8個のサブルーチンを利用できるようになっている。また、RCX同士で1[Byte]のメッセージでやり取りすることも可能。

そこで標準ファームウェアを使用すると音声センサーがないので、一度音声をPCに取り込んでMATLABにおいて音源方向の角度を計算する。次に音声認識判断を行いそれからその角度や音声に対応した動きを考え、標準のファームウェアに手動で打ち込みロボットを動作させるという手間のかかる事になってしまった。

そこでロボットを動かすファームウェアも標準のものではなく、Visual Basic⁽⁴⁾やLEGOロボットのためのプログラミングソフトのNQC/RcxCc⁽⁵⁾などを使用して、自分たちでファームウェアを作成していきたいと考えたがVisual BasicではLEGOがバージョンアップされていたので動作させることができなかった。NQC/RcxCcにおいてもSpirit.ocx⁽⁶⁾というRCXのコントロールやプログラムのダウンロードをするソフトウェアがなく動作させることができなかった。動作制御においては課題が残ってしまった。

(3)ファームウェア

LEGO MindStorms「Robotics Invention System(以下RIS)」の標準的なプログラミング環境であるRCXコードではパソコンで作ったプログラムをRCX(ロボットの頭脳)にダウンロードした後、私たちが[Run]ボタンを押すと実行が開始される。このようなダウンロードされたプログラムを[Prgn]ボタンが押されたら切り替えたり、[Run]ボタンを押されたら実行したりというような処理を行うソフトウェアをRISでは「ファームウェア」と呼ぶ。

ファームウェアとは、パソコンのOSに似た基本的な機能を提供する為のソフトウェアのこと。ファームウェアは一般的にハードウェアとソフトウェアの中間的な位置に存在するソフトウェアで、マイコンを搭載している家電などで利用されている。例えば、エアコンでは温度が20[]になるように冷気を調節して部屋の中を送るといったような処理を行っている。

ただ家電などではほとんどの場合ファームウェアはROMに焼かれていて、私たちが簡単に交換できるようにはなっていない。それにたいしてRCXのファームウェアは私たちがダウンロードできるようになっているという大きな特徴がある。また、私たちが作成したRCXコードを実行することができるという特徴もある。すなわち、RCXのファームウェアは、RCXコードを解釈しながら実行するインプリタ(命令文を1つずつ解釈しながらプログラムを実行する方法)を搭載している。

(4) Visual Basicでのプログラミング

RISの標準のプログラミング環境(RCXコード)は、Spirit.ocxというActiveXコントロールをしている。ActiveXコントロールとはいくつかの機能をまとめたプログラミング部品で、Spirit.ocxにはRCXをコントロールする機能がまとめられている。そして、Visual BasicなどのActiveXコントロールに対応したプログラミング言語からこれを利用すれば、ハードウェアに関する知識がほとんどなくてもRCXをコントロールプログラムが作れるようになっている。つまりSpirit.ocxはRCXを操作するためのもっとも重要な部分だといえる。

Visual Basicを使ってRCXをプログラムすることのメリットは、標準ファームウェアのもっている能力を使いきれんというメリットもあるが、パソコン周辺機器のようにRCXをコントロールできる点に1番の魅力を感じる。例えば、パソコンの画面からマウスを使ってモーターをon/offにしたり、RCXの接続されたセンサーの値を読んでパソコンでいろいろな処理をしながらRCXに命令を送るような場合、Visual Basicを使ったプログラミングは非常に有効だといえるだろう。

(5) NQC/RcxCC

NQCとは「Not Quite C : 完全なCではない」ということで、プログラミング言語Cに類似した構文をもつ簡単な言語である。NQCはLEGO社とはまったく無関係に開発されていて、どんな形でもLEGO社には属さない。伝統的なプログラミング方法、つまりテキストファイルとして記述し、コンパイルしRCXにダウンロードして使う。グラフィックツールのように簡単には始められないが、プログラミングに親しんでいる人にとっては手ごろであり文字を使った言語はとても効率のよい開発ツールと言える。変数は32個、サブルーチンは10個、タスクは8個定義することができ標準ファームウェアの性能をさいだいげんに引き出すことができる。NQCの利点の1つは速度である。いくつかある他のツールは最速のコンピュータハードウェアでも遅くて不便である。NQCは旧式のハードウェアでもコンパイルできるうえ、高速にプログラムのダウンロードができる。Windows、MacOS、Linuxの各環境用にプログラムが用意されている。短所としては学習が少し難しいところだろう。

RcxCCとは「Rcx Command Center」の略。これはNQC向けの総合開発環境で、リアルタイムにRCXをコントロールすることもできる。こちらはWindows用のみとなる。

(6) Spirit.ocx

Spirit.ocxは開発ツールではない。VisualBasicなどの他のプログラミング言語で使用するActiveXコントロールで、RCXのコントロールやプログラムのダウンロードをするものである。Spirit.ocxはRISのソフトウェアでWindowsのみで動作する。Spirit.ocxのドキュメントはLEGO社からMindStorms用のソフトウェア開発キット(SDK)として入手ができる。VisualBasicでRCXのプログラミングができるのは素晴らしいが、Spirit.ocxは完全ではない。Spirit.ocxはVisualBasicのプログラミングが働いているホストコンピュータからRCXへの問い合わせと制御を行う。また、VisualBasicプログラムでRCXのプログラムをダウンロードすることができる。しかし、作成されたRCXプログラムSpirit.ocxの提供する分だけしか使えないという制約がある。VisualBasicで使い慣れた関数や機能は実際にRCXプログラムには使用できない。

まとめ

本研究ではロボットの聴覚機能実現を目的とした研究の第1歩として音源方法の検出、音声認識判断、ロボットの動作制御の検討をおこなった。

音源の方向検出においては人間の感覚と同じように正面付近はほぼ正確に方向を検出する事ができたが、角度が開いていくにつれ検出の結果に誤差が大きくなるが多かった。つまり、音源が真横に近づくにつれて方向検出が不正確になった。これは人間にも同様のことが言える。また、発声する音の高低やその声の質、そしてマイクロホンを置く場所によってその部屋のノイズや反射波等による影響のために検出する音源方向に大きな変化があった。

音声認識判断においては、はっきりと判別可能な音声もあったが“あ”や“お”といった発声するときの口の開き具合や、口の中の形状が似ている音声では多くのプロットが重なり、音声特徴が似たものを持っていたため判別する事が難しかった。今回の研究では単音(母音)のみの音声認識判断の検討だけだったが、単語からちょっとした文章を認識し、判断していくことができるようになればロボットの聴覚機能実現に向けて大きく前進することだろう。

ロボットの動作制御においてはパソコンで制御可能な市販のLEGO Mind Stormsを購入し作動させることを目的とした。LEGO Mind Stormsには音声センサーが搭載されていなかったので標準ファームウェアを使用せずにVisual BasicやNQC/RcxCcなどを使用して、自分たちでファームウェアを作成していきたいと考えた。しかし、Visual BasicではLEGOがバージョンアップされていたので動作させることができずNQC/RcxCcにおいてもSpirit.ocxというRCXのコントロールやプログラムのダウンロードをするソフトウェアがなく動作させることができなかった。残念ながら動作制御においては多くの課題が残ってしまった。

本研究においては結果を出せなかった研究もいくつかあったが現実的にロボットの聴覚機能実現のためにはこのほかにまだ数多くの研究が必要になる。それらに対応できるようにしていくことがロボットの聴覚機能実現につながっていくことだろう。

謝辞

本研究を進めるにあたりご指導下さいました金田豊教授ならびに、同研究室の研究生の方々に深く感謝申し上げます。

参考文献

* 1

- (1) Jin Sato・白川裕記・牧瀬哲郎・倉林大輔・衛藤仁郎，
“ LEGO MINDSTORMS ハ°-フヱクトガ°イト° ”
pp.79～91，翔泳社，1999．
<http://www.legomindstorms.com/sdk/>

- (2) Davit Baum
“ LEGO MINDSTORMS ラ°-ニング° &° フ°ロク°ラミング°ガ°イト° ”
pp.38～，シュ°リンガ°ー°フヱアラ°ク，2001．