

## ヘッドマウンテドカメラを用いた ポインティングジェスチャと音声による移動ロボットの誘導

貴島 茂雄 西川 敦 宮崎 文夫

大阪大学大学院 基礎工学研究科 システム人間系専攻 機械科学分野

本研究では、人間の頭部に装着されたカメラを用い、人間の人差指によるポインティングジェスチャと音声によってロボットを指示者が意図した地点まで誘導するシステムを提案し、誘導実験によりその有効性を検証する。提案するシステムでは、ヘッドマウンテドカメラと人間のポインティングに関する考察を巧く組み合わせる事によって、1台のカメラによるポインティングシステムを実現する。また、カメラ座標とロボット座標の座標変換に平面の幾何学的不変量を用いることにより、キャリブレーションの負荷を軽減するとともに、計算コストをより低くすることが可能となる。さらに、人間にとって自然なインタフェースである音声を積極的に用いることにより、より直感的で扱いやすいシステムをめざす。

## Navigation of a Mobile Robot by Pointing Gesture and Voice Based on the Use of Head Mounted Camera

Shigeo Kijima, Atsushi Nishikawa and Fumio Miyazaki

Division of Mechanical Science, Department of Systems and Human Science  
Graduate School of Engineering Science, Osaka University

In this paper, we propose a mobile robot navigation system by pointing gestures and voice based on the use of a head mounted camera. The proposed system is based on a hypothesis on the pointing gestures human naturally makes. A simple result from projective geometry (projective invariants for five coplanar points) is applied to estimation of the operator's instruction points, which contributes to the reduction of the calculation cost and in the load of the camera calibration because all calculation takes place in the two-dimensional image and ground planes. Furthermore, we also utilize voice inputs towards more simple and intuitive system for an operator to use.

### 1 はじめに

近年、パーソナルコンピュータの一般社会への普及等により、人間が日常的にコンピュータやこれを搭載した機器に触れる機会が増加して来ている。この様な状況において、それらは、一部の操作に熟練した者だけではなく、子どもからお年寄りまで誰にでも容易に扱えることが要求される。

このような背景から、近年マウスやキーボードに代わる、より直感的で人間にやさしいマンマシンインタフェースに関する研究が

行われてきている [1]。

人間にとって自然なインタフェースの1つとして、音声挙げられる。我々が普段社会のなかで他人とのコミュニケーションをとる時、その情報の大部分は音声によるものである。この事を考慮すると、音声とは人間にとって最も身近で直感的な情報の伝達手段であると言える。

しかし、音声は人間にとって最も直感的な情報の伝達手段である一方で、音声で伝達することが困難な情報も存在する。その一例と

して、意図した地点や物体など、対象の位置に関するものが挙げられる。これらを相手に的確に伝達する時、我々は、しばしば、「ここ」「これ」などの指示語とともに、対象を指さす「ポインティングジェスチャ」を用いる。

人間の意図した対象を検出するポインティングシステムとして、Cipollaら [2] は、キャリブレーションの負荷を軽減したステレオカメラを用いて、ポインティングジェスチャによりオペレータ(人間)が意図した平面上の点を特定するシステムを提案している。このシステムは実際使用する時にはカメラを台座に固定しているため、ポインティングできる方向が限定され、作業平面の領域内しか扱う事ができない。

これに対し、渡辺ら [3] は、マルチカメラを用いて、室内でどの方向をポインティングしてもポインティングジェスチャを検出できるシステムを提案している。

さて、これらのポインティングシステムでは、いずれもカメラを固定して使用することを前提としているため、システムを適用できる範囲が限定される。また、複数個のカメラを使用しているため、大がかりなシステムになってしまう。さらに、ポインティング動作をシステムに伝えるタイミングをどうするかなどの、実際にポインティングシステムをインタフェースとして使用する場合に起こる問題をあまり考慮していない。

本研究では、ポインティングシステムのこれらの問題点を解決するために以下に示す事項を採用したポインティングシステムを開発し、実際に移動ロボットを誘導する。

- カメラをユーザの頭部に装着する(カメラは人間と共に自由に移動できるため、画像を処理するPCを小型化できれば、開かれた世界での使用が可能)。
- 頭部に装着した1台のCCDカメラのみを用いてポインティングジェスチャを認識する。(システムの簡素化)
- ポインティングジェスチャを相手に知らせる手段(トリガ)として、またシステムへのコマンド入力として、人間にとって

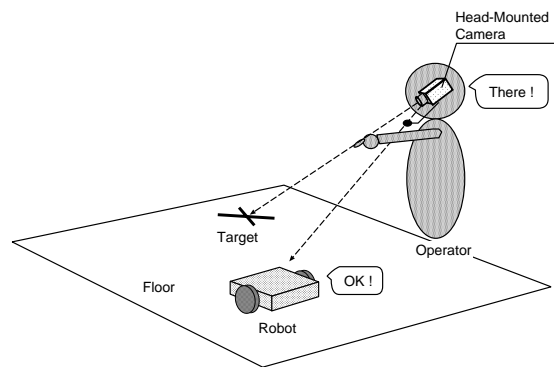


Fig.1 The Pointing System

最も直感的なインタフェースである音声を利用する。

## 2 システム概要

本システムは、人間の頭部に装着された単眼のカメラを用いて、ポインティングジェスチャと、音声により、フィールド上を移動するロボットを、指示者が意図した地点まで誘導することを目的としたシステムである。本システムの概念図をFig.1に示す。

指示者はロボットの方向を向き、人差指で自分の意図した目標地点をポインティングジェスチャにより指し示し、「ここ」等の音声を入力することで、ロボットに目標地点を示してやり、「走れ」と入力することによりロボットを移動させる。「走れ」と入力する前に、ロボットに数点の目標地点を示してやることで、ロボットの軌道を与えることも可能である。

画像処理とロボットへのコマンドの送信は、Linux2.2をOSとし、CPUにPentiumIII(700MHz)、そしてビデオキャプチャデバイス(GV-VCP2/PCI)を搭載した、標準的なIBM-PC互換アーキテクチャのコンピュータで行う。カメラからの画像の取得はLinux上で利用可能なビデオキャプチャ用のAPIであるVideo for Linux [4]を利用した。一方、音声処理は、WindowsMeをOSとし、CPUにCeleron(500MHz)を搭載したIBM社製のノートPC、ThinkPadで行う。音声認識にはIBM社製のViaVoicePro[5]を使用した。

画像の入力装置であるヘッドマウンテドカ



Fig.2 Head Mounted Camera

カメラは、Fig.2に示すものであり、ミノトロン社製 CCD カメラ MTV-7366 を側頭部左側に位置するように、ヘルメットに固定したものである。また、カメラのレンズは焦点距離 4mm、水平画角 60 度の CS マウントレンズである。一方、音声の入力装置には、IBM 社製 ViaVoice 付属の標準的なヘッドセットマイクを使用した。

操作対象であるロボットは、200mm×200mm の四角形のボディに、NEC 製の CPU、V25 を使用した日本システムデザイン社製の CPU ボードと制御用 OS、JSDOS を搭載している [6]。ロボットは、独立したモータによって駆動される 2 つの車輪によって移動する。これらにはロータリエンコーダが取り付けられており、これを利用してデッドレコニングが可能である。また、外界との通信が可能なシリアルポートを持っており、外部の PC 等からプログラムのダウンロードや、実行中のプログラムの通信の用途として利用可能である。ロボットシステム全体は 7.2V の Ni-Cd バッテリーで機能する。ロボット本体の外観を Fig.3(a) に示す。実際の実験では、同図 (b) に示すように、ロボット上に複数の色マーカを設置する。ロボットの周囲に配置されている 5 つの円形マーカは、ポインティングジェスチャによって指し示された点をロボット座標に変換するのに利用する (3 章参照)。また、ロボット中央に設置されているマーカは実験データ (移動軌跡) の記録用である (6 章参照)。

これらの接続形態を Fig.4 に示す。

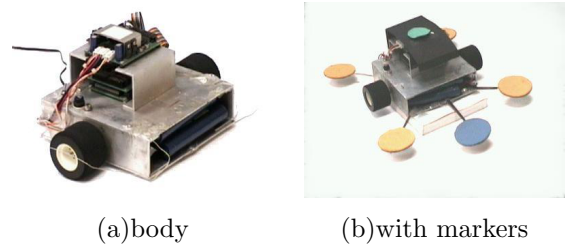


Fig.3 Moving Robot Overview

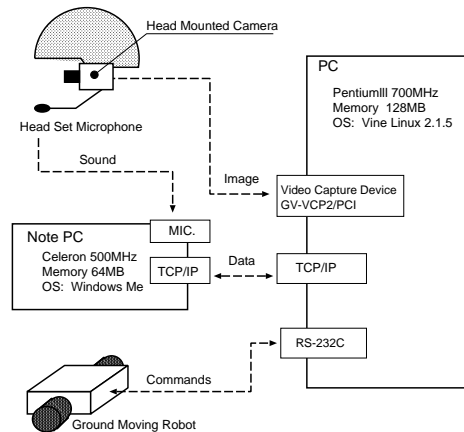


Fig.4 Hardware Components

### 3 平面 5 点の幾何学的不変量を利用したカメラ座標からロボット座標への座標変換

本システムでは、頭部に装着した一台の CCD カメラのみから画像情報を入力する。カメラから得られた画像平面上での点を、ロボットが走行する床平面上での点に対応づける処理は本システムを考えるうえで重要である。このような座標変換の手法として、1 フレーム毎に、ロボット座標からのカメラ座標の相対的な位置姿勢を常に計算する方法 [7] が考えられる。しかし、この手法では、カメラ座標、ロボット座標共に 3 次元の座標情報を扱うので、繁雑であり計算コストが高い。

これに対して、本システムではロボットは常にある平面を移動することを想定していることに着目し、比較的計算コストの低い「平面の幾何学的不変量」を用いて、カメラ座標、ロボット座標共に 2 次元の座標情報のみを扱い、カメラ座標からロボット座標への

座標変換を行う。

本手法はカメラパラメータ(焦点距離や画像中心など)を陽に扱わないので,カメラキャリブレーションの負担が大幅に軽減されるといふ利点もある<sup>1</sup>。

### 3.1 平面5点の幾何学的不変量

ある3次元直交座標空間と,その内部において,視点から透視投影して得られた画像平面を考える。幾何学的不変量とは,視点が3次元空間内を動くとき,画像面上で物体を特徴づける,点の座標や,線分あるいは曲線などの方程式の係数を変数として,視点の動きに依存しない値をとる関数である。同一平面上に存在する5点(ただしどの3点も同一直線上にない)に対して,次の2つの不変量 $I_1, I_2$ が存在する[9]。すなわち,平面上の5点を $x_i(i=1,2,3,4,5)$ とすると

$$I_1 = \frac{\det P_{431} \cdot \det P_{521}}{\det P_{421} \cdot \det P_{531}} \quad (1)$$

$$I_2 = \frac{\det P_{421} \cdot \det P_{532}}{\det P_{432} \cdot \det P_{521}} \quad (2)$$

ただし,平面上の点は $x_i = [X_i \ Y_i \ 1]^T$ といったような同次座標で表し, $P_{ijk} = [x_i \ x_j \ x_k]$ なる $3 \times 3$ の行列である。

### 3.2 任意の点のカメラ座標からロボット座標への変換

ロボットには,円形のマーカが五角形の頂点をなすように配置してある(Fig.3(b)参照)。ここで,ロボット座標系でのマーカの重心座標を ${}^rM_i(i=1,2,3,4,5)[\text{mm}]$ ,カメラ画像座標系でのマーカの重心座標を ${}^cM_i(i=1,2,3,4,5)[\text{pixel}]$ で表す。ただしロボット座標系とは,ロボットの中心を原点とし,そこからロボットの進行方向前向きにx座標をとり,それに垂直方向左向きにy軸をとったx-y直交座標系の事であり,カメラ画像座標系とはキャブチャした画

<sup>1</sup>本研究で使用しているカメラレンズは焦点距離4mmであり,画角が広い。それゆえ,座標変換の計算の際に仮定している透視変換による画像平面に対し,カメラから得られる画像には光学的な歪みが生じている。したがって,現状ではこの歪みを補正するためのカメラキャリブレーションプロセスは必要である。この処理は, Camera Calibration Toolbox for Matlab[8]を利用して行っている。

像の左上を原点とし,右向きにx軸,下向きにy軸をとったx-y直交座標系のことである。

今,人間が指示した床上の点をカメラ画像座標系で表したものを ${}^cp[\text{pixel}]$ とし,既知(5章で述べる画像処理によって推定可能)とする。 ${}^cM_i$ も画像処理によって計測可能である[7]。また, ${}^rM_i$ は先見的知識としてロボットに与えられているものとする。このとき ${}^rM_i, {}^cM_i, (i=1,2,3,4,5)$ と ${}^cp$ を用いて,指示点のロボット座標系上での点 ${}^rp$ を導出する方法について述べる。

まず始めに ${}^cM_i(i=1,2,3,4,5)$ の5点から以下の様な適当な4点の組合せを選び,番号の若い順に ${}^cx_j(j=1,2,3,4)$ とする。適当な4点とは,その4点と ${}^cp$ の5点のうち,どの3点も同一直線上に無い様な4点の組合せである。

次にその5点から,以下に示す式で幾何学的不変量 $I_1, I_2$ を計算する。

$$I_1 = \frac{\det[{}^cx_4 \ {}^cx_3 \ {}^cx_1] \cdot \det[{}^cp \ {}^cx_2 \ {}^cx_1]}{\det[{}^cx_4 \ {}^cx_2 \ {}^cx_1] \cdot \det[{}^cp \ {}^cx_3 \ {}^cx_1]} \quad (3)$$

$$I_2 = \frac{\det[{}^cx_4 \ {}^cx_2 \ {}^cx_1] \cdot \det[{}^cp \ {}^cx_3 \ {}^cx_2]}{\det[{}^cx_4 \ {}^cx_3 \ {}^cx_2] \cdot \det[{}^cp \ {}^cx_2 \ {}^cx_1]} \quad (4)$$

ここで ${}^cx_j(j=1,2,3,4)$ に対応するロボット座標系上での点の座標を ${}^rx_j(j=1,2,3,4)$ とすると, $I_1, I_2$ は記述する座標系によらず不変であるから以下の式が成り立つ。

$$\frac{\det[{}^rx_4 \ {}^rx_3 \ {}^rx_1] \cdot \det[{}^rp \ {}^rx_2 \ {}^rx_1]}{\det[{}^rx_4 \ {}^rx_2 \ {}^rx_1] \cdot \det[{}^rp \ {}^rx_3 \ {}^rx_1]} = I_1 \quad (5)$$

$$\frac{\det[{}^rx_4 \ {}^rx_2 \ {}^rx_1] \cdot \det[{}^rp \ {}^rx_3 \ {}^rx_2]}{\det[{}^rx_4 \ {}^rx_3 \ {}^rx_2] \cdot \det[{}^rp \ {}^rx_2 \ {}^rx_1]} = I_2 \quad (6)$$

従って式(5),(6)より ${}^rp[\text{mm}]$ が求まる。

## 4 ポインティングジェスチャについて

ここで,人間がポインティングジェスチャを行うときの,目線,指先,対象の位置関係を考えてみる。それは,対象までの距離,又は伝えようと思う対象の精度によって変わってくる(例えば,対象が或る一点ならば高い精度が要求されるが,対象が比較的大きな物体などであればそれ程高い精度は要求されない)。Fig.5に示すように,対象が遠くにある場合または高い精度が要求される場合,人差し指の指先,指のつけね,そして対象を結ん

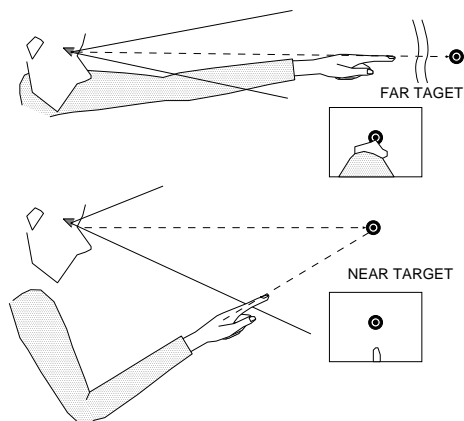


Fig.5 Difference between Pointing Action at Far Target and at Near Target

だ直線は、人間の視線と一致する。これに対し、対象が比較的近くにある場合、またはそれほど高い精度が要求されない場合には、人間は視界の中に対象と人差し指が入るように調節し、さらに指のつけねから指先へ向かう直線上に対象がくる位置で指をとめる。

このポインティングジェスチャの違いの原因としては、そもそもポインティングジェスチャ自体が、相手に対象を伝える為に用いることなどが挙げられる。すなわち、対象が近くにあるか、又はポインティングジェスチャに高い精度が要求されない場合には、相手は指示者の指の方向だけをたどって対象を認識することができる。しかし対象が遠くにあったり、ポインティングジェスチャに高い精度が要求される場合は、相手は単に人差し指の方向だけではなく、指示者の視線と人差し指の関係まで、注意して見る必要が出てくる。

本システムでは左手でポインティングジェスチャを行う。また、本システムにおいて、カメラは両眼の真横、左肩の真上に位置するように頭部に装着されている。

従って指示者が対象への自分の視線に人差し指を重ねるようにポインティングジェスチャを行うならば、ヘッドマウントドカメラから観測される手と目標の映像は、人間の目から観測される映像とほぼ同じものとなる。実際にヘッドマウントドカメラからとらえた画像を Fig.6 に示す。

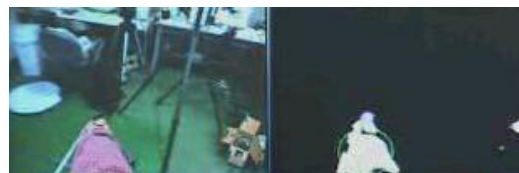


Fig.6 Pointing Gesture taken by Head Mounted Camera

さて、本システムは、ロボットを床平面上のある目標地点へ、ポインティングジェスチャによって誘導することを目的とする。目標座標を指示するので、指示者が想定しているポインティングジェスチャの精度は高いものである。従って、ポインティングジェスチャとしては前述した、対象への自分の視線に人差し指を重ねる方法が指示者にとってより自然であろう。

さらに、カメラが頭部に水平に取り付けられているため、ポインティングジェスチャを自然に行ったならば、それを捕らえたカメラ画像のなかで、人差し指の付け根が手領域の最上部に位置する。

以上より、本システムでは、左手の人差し指を用いて、対象への視線に人差し指を重ねるようなポインティングジェスチャを行うものとする。

## 5 指示点の検出

まず、前述の方法により、床平面上の目標地点をポインティングする。このとき、手又は腕領域は、ヘッドマウントドカメラの視界の左下部分からほぼ中央に向かってのびている。次に手領域をHSVカラーを用いて2値化、ラベリングした後、面積最大の領域を手領域とし、その手領域の最上部に位置する点を指示点の画像平面上での点とする (Fig.6 参照)。

そして3章で述べたように、得られた画像上の指示点を、幾何学的不変量を用いて座標変換することによって、指示点のロボット座標系における座標値が求まる。

## 6 移動ロボットの誘導実験

提案手法の有効性を検証するため、以下のような実験を行った。

### 6.1 実験内容

床平面上の座標  $(0,0)$ ,  $(0,600)$ ,  $(600,600)$ ,  $(600,0)$  (単位 [mm]) を頂点とする正方形を一周するようにロボットを誘導する。各頂点には、指示者に分かるように  $2 \times 2$  [cm] の矩形マーカが配置されており、ロボットは初期状態で  $(0,0)$  に位置している。被験者はまず  $(0,600)$  の地点を指さし、「ここ」と入力することでロボットにその座標を知らせる (Fig.7(a))。同様に  $(600,600)$   $(600,0)$   $(0,0)$  の順に、正方形の頂点をポインティングジェスチャでロボットに入力し (Fig.7(b)~(d))、全ての点を入力した後、「走れ」と発してロボットを実際に走行させる (Fig.7(e)~(h))。

被験者が立つ位置は特に指定はせず、自分の視界にロボットとポインティング地点が入るように調節してもらった。

以上を、3人の被験者A,B,Cにつき各4回行った。なお、被験者A及びBは、この誘導実験の以前に何度かこのシステムを使用した経験をもち、被験者Cは実験で初めて本システムを使用した。

ロボットの移動軌跡は、SONY製デジタルビデオカメラ DCR-TRV20 で撮影しておき、QuickMAG System III[10] の2次元カラーシステム<sup>2</sup>を利用して、マーカ平面上での移動軌跡として記録した。記録したのは、ロボット中央部の円形マーカ (Fig.3 参照) の移動軌跡である。

### 6.2 結果と考察

実験結果の走行軌跡を Fig.8~Fig.10 に示す。

各被験者について、ロボットの走行軌跡を見てみると、被験者A,Bの軌跡は、ほぼ正

<sup>2</sup>使用した QuickMAG SystemIII のシステム構成では、2つの CCD カメラを用いて 60fps での 3次元のリアルタイム位置計測を行う能力をもつ。最高 16 の対象を色情報を用いて追跡することが可能である。今回利用したのは、2チャンネルある映像入力的一方だけを用いた平面運動の解析機能である。任意のビデオカメラについて DLT 法を用いたキャリブレーションを行って解析することが可能である。

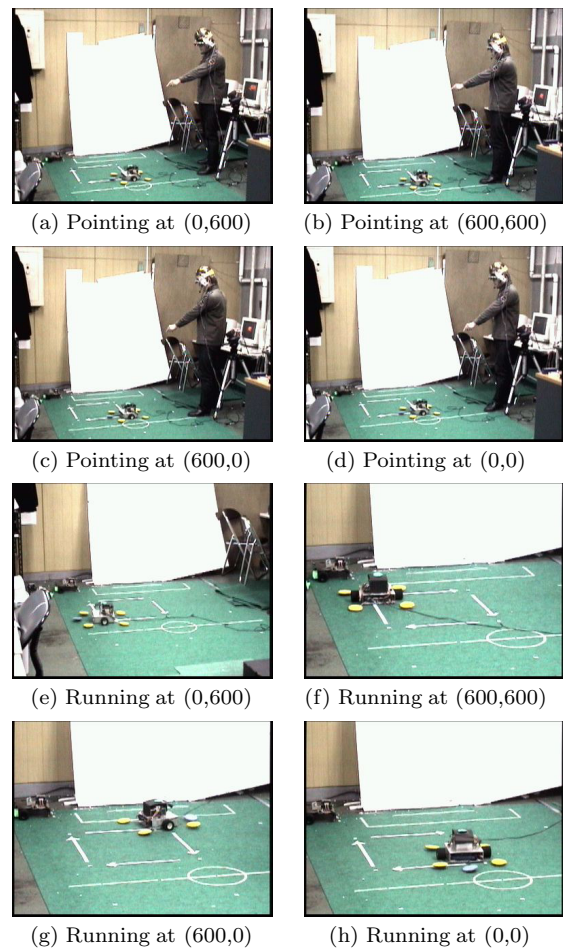


Fig.7 Overview of Experiment

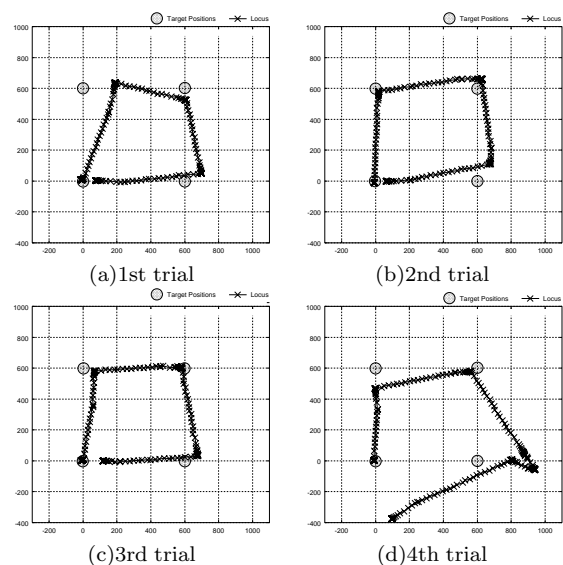


Fig.8 Experimental Results :Subject A [unit:mm]

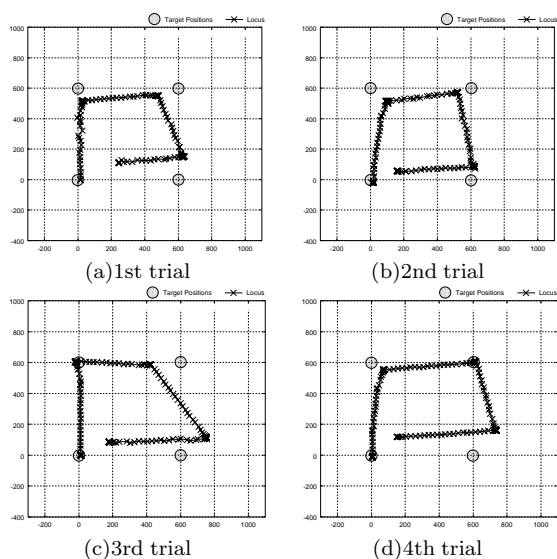


Fig.9 Experimental Results :Subject B [unit:mm]

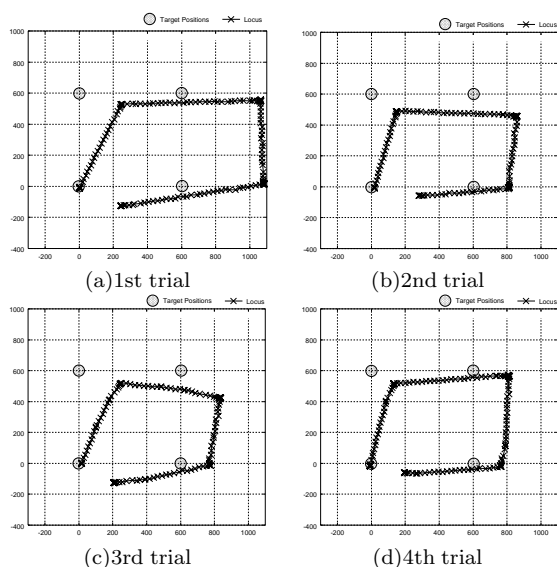


Fig.10 Experimental Results :Subject C [unit:mm]

方形を描いているのに対し、本システムに始めて触れた被験者Cの軌跡は、他の2人に比べ、目標の軌跡から大きく外れている。さらに、被験者Cの軌跡を詳しく見てみると、目標の軌跡からのy軸方向のずれに対して、x方向のずれがかなり大きい。特に、実験2の1回目の走行軌跡 (Fig.10(a)) では、それが最も顕著に現われている。この原因について考える前に、実験時の環境について以下の点を挙げておく。

- 指示者の立つ位置に関しては、特に指定はしなかったが、被験者A,B,Cともに、ほぼ全ての実験において、フィールド座標にして(0,-1200)あたりから、y軸方向正の向きに向いてポインティングジェスチャを行っていた。
- 従って頭部左側のヘッドマウントドカメラからの画像におけるx軸と、実験フィールドのx軸はほぼ平行になっているものと推測できる。

この事実を踏まえて、x軸方向の誤差がy軸方向の誤差に比べてかなり大きくなった原因について考察する。そもそも本システムでは、ポインティングジェスチャから指示点を検出する際に、ヘッドマウントドカメラの視線と人間の視線はほぼ同一のものであるという仮定をおいていた。しかし実際にはカメラは頭部左側面から、視界の軸が水平方向右向きになるように取り付けられているため、人間が視線方向に人差指を重ねるポインティングジェスチャを行った際に、カメラ中心から人差指の付け根を通る直線と床との交点、すなわちシステムが算出した目標点は、指示者が意図した点より右にかたよる。x軸方向の誤差がy軸方向の誤差に比べてかなり大きくなったのはこのためだと考えられる。

一方、被験者A,Bは同じ場所に立ってポインティングジェスチャを行ったにもかかわらず、被験者Cほどx方向の誤差が顕著に現われなかった。これは、システムを初めて使用して、全くカメラの位置について意識しなかった被験者Cに対して、被験者A,Bは、過去に数回の使用経験があるため、多少カメラの視線

を意識したポインティングジェスチャを行ったためではないかと考えられる。

## 7 まとめ

本研究では、ヘッドマウントドカメラに、人間のポインティングジェスチャに関する仮定を巧く組み合わせる事によって、1台のCCDカメラによる簡素なポインティングシステムを開発した。

また、カメラ座標とロボット座標の座標変換に平面5点の幾何学的不変量を用いることにより、カメラの位置姿勢の3次元情報を扱う手法などに比べ、より低い計算コストを実現した。

さらに、ポインティングジェスチャのトリガとして、またポインティングジェスチャ以外のロボットとのコミュニケーション手段として音声を用いる事で、ポインティングシステムが、より直感的で人に優しいインタフェースになりうるように試みた。

そして開発したシステムで、実際にポインティングジェスチャによりロボットを誘導し、有効性を検証した。

本システムの今後の課題としては

- 音声認識の認識率の向上と音声処理時間の短縮
- ポインティング部に応用するための、人間のポインティングジェスチャに関する新たな考察と仮説の提唱
- システムの小型化、モバイル化

等があげられる。

また、本システムはタスクを「ポインティングジェスチャによる移動ロボットの誘導」のみに絞ったが、発展として、例えば、ロボットに物体を把持する能力が備わっている場合を想定すると、指示者がポインティングジェスチャによって指し示した物体をとってくる等のタスクに応用することも考えられる。

## 参考文献

[1] 谷内田正彦：顔とジェスチャの認識，システム/制御/情報，Vol. 44, No. 3, pp. 97-101

(2000).

- [2] Cipolla, R. and Hollinghurst., N.: A Human-Robot Interface using Pointing with Uncalibrated Stereo Vision, in *Computer Vision for Human-Machine Interaction*, pp. 97-110, CAMBRIDGE UNIVERSITY PRESS (1998).
- [3] 渡辺博己, 本郷仁志, 安本護, 山本和彦: マルチカメラを用いた全方位ポインティングジェスチャの方向推定, *T.IEE Japan*, Vol. 121-C, No. 9, pp. 1388-1394 (2001).
- [4] *Video4Linux Kernel API Reference*,  
<http://roadrunner.swansea.uk.linux.org/v4lapi.shtml>.
- [5] 日本IBM株式会社: ViaVoice for Windows 日本語版,  
<http://www-6.ibm.com/jp/event/museum/mirai/vreco.html>.
- [6] 升谷保博, 福留正一, 宮崎文夫: 車輪型移動ロボットを用いたメカトロニクス実習, 日本機械学会ロボティクス・メカトロニクス講演会'96 講演論文集, Vol. A, pp. 401-404 (1996).
- [7] 進戸健太郎, 西川敦, 宮崎文夫: ヘッドマウントドカメラを用いた指示者の注目点と手振りによる指示を組み合わせたヒューマンインタフェース, *インタラクシオン2002 論文集*, pp. 57-58 (2002).
- [8] Bouguet, J.-Y.: *Camera Calibration Toolbox for Matlab*, MRL - Intel Corp.,  
[http://www.vision.caltech.edu/bouguetj/calib\\_doc/](http://www.vision.caltech.edu/bouguetj/calib_doc/).
- [9] 杉本晃宏: コンピュータビジョン 技術評論と将来展望, 第7章, 新技術コミュニケーションズ (1998).
- [10] 株式会社応用計測研究所: リアルタイム動作解析システム Quick MAG System III ユーザーズマニュアル.