

小特集 「言語獲得」

# ロボットによる言語獲得

—言語処理の新しいパラダイムを目指して—

Language Aquisition by Robots

— Towards New Paradigm of Language Processing —

岩橋 直人  
Naoto Iwahashi

(株) ソニーコンピュータサイエンス研究所  
Sony Computer Science Labs, Inc.  
iwahashi@csl.sony.co.jp, <http://www.csl.sony.co.jp/person/iwahashi/>

**Keywords:** mutual belief, interaction, experience, coupling, robot, language, acquisition.

## 1. はじめに

日常のコミュニケーションは、コミュニケーションに参加するものが互いに共有する信念（相互信念）に基づいて成立する。言語はそうした信念の一部であり<sup>\*1</sup>、ほかの相互信念との関連に基づき、意味を伝達するために使用される（e.g. [Sperber 95]）。信念は、主体と環境や他者との認知的な相互作用によって形成されるので、発話の意味は、対話参加者が共有するそうした相互作用の中に埋め込まれる。言語は、ほかの認知プロセスや信念との関わりを含めた全体であり、表現における記号時系列規則、経験の中に埋め込まれた意味、および状況に応じた使用法といった側面を有する。

ところで、ある信念が相互に共有されているということに対話者間で論理的に確認しようとする、無限の深さの入れ子の信念が成立していなければならない。実際には、それは何らかの手掛りによって確信されるだけであるから、各人が想定している相互信念の一致は、完全に保証されるものではない（e.g. [Clark 96]）。発話生成、理解といった環境や他者との相互作用は、その時点での相互信念に基づいて生じ、一方で、相互信念はそうした相互作用の系列によって形成される。このように相互信念は自らの構成要素をつくり出しながら自律的に発展する。また、聞き手が発話を解釈しようとするプロセス<sup>\*2</sup>は、聞き手が想定する相互信念との関連に基づいて実行されるので、このプロセスを通して、聞き手は相互信念を形成するための情報を受け取ることができ、さらに、話し手も聞き手の反応により即座に同様の情報を受け取ることができる。すなわち、発話は相互信念を形成する

ための情報の送受信を同時に行おうとする行為であり、各主体の相互信念はこの行為を通して他者の相互信念とカップリングしていると言える。

言語獲得は相互信念の形成である。幼児は、発話を解釈しようとするプロセスを通して、言語が相互信念を形成するという機能に基づいて、相互信念として言語自体を形成していく。文法など、しだいに相互作用の影響を受けにくくなり固定化されていく信念もあるが、経験に埋め込まれたことばの意味を支える信念の全体は、相互作用の影響を受けダイナミックに変化し続ける。

さて、このように共有する経験を反映する日常の言語コミュニケーションを人とロボットの間で実現したいのである。従来の言語処理手法には何が足りないのだろうか。近年の音声対話処理技術（e.g. [中野 02]）の向上に貢献している大規模テキストコーパスを用いた統計的言語処理手法は、共同体の中で共有された静的な言語信念を高い精度で抽出、運用することを可能にするが、経験を通してダイナミックに変化する意味を扱うメカニズムを提供するものではない。Schank は、固定された構造によってのみ解釈され他者の理解を参照しない「有意味」と、互いのすべての動きや動機がわかる親しい者の「完全な感情移入」のそれぞれを端とする理解のスペクトルを示した [Schank 84]。ディスプレイ中の積木の世界について人とキーボードを介して対話するシステム SHRDLU [Winograd 72] は、言語、知覚、行動の情報を統合し、入力文を状況に応じて適切に理解をすることができたことから、有意味のレベルを実現したといえる。有意味レベルの理解の次に目指すべきものは、経験に基づいて学習したり変化する、現在の経験を過去の経験に知的に関連付ける、などができる認知的理解のレベルであると Schank は述べている。しかし、SHRDLU とは対極的な特徴をもつ対話システム ELIZA [Weizenbaum 66] は、非常に簡単な規則のみを用いて応答文を生成し、対話者の感情移入をもたらすことができってしまう。このことは、言語理解システムの研究を、有意味→認知的理解→完全

\*1 例えばクワインの言語観 [丹治 96] では知識と信念とは区別されない。

\*2 例えば [Schank 84] や [Reddy 93] は言語とは意味をつくり出す手続きであると主張する。



図1 人とロボットシステムのインタラクションのようす

なる感情移入のレベル，という方向で進めていくことが必ずしも正しくはないことを示唆しているように思われる．実際，幼児が言語を獲得する過程はこれとは逆向きに進むように見える．つまり，まず母子間の共感を基盤としたインタラクションが生じ，これから限定された経験に埋め込まれた言語が形成され，その後この言語が社会的な相互作用を通して汎化されていくのである．

では，共感することを最も根本的な目的とする言語処理はどのように実現できるのであろうか．共感するには，まず互いにわかり合う，つまり相互信念が形成されなければならない．すでに述べたように，各主体の相互信念は他者の相互信念とカップリングしながら，自律的に発展するものである．したがって，言語処理は，このようなダイナミクスに基づき，ロボットの相互信念が，人とロボットの共有する経験を通して形成されることを可能としなければならない．また，相互信念の形成の基盤として，環境が共有されることが望ましく，そのために言語と知覚と行動が認知的に適切な方法で統合的に処理されなければならない．

これらのことを実現するために，著者らはロボットが対話者との共同知覚経験や共同行為を基盤にして，言語信念を含む相互信念を獲得する方法の研究を進めてきた[羽岡 00, 羽岡 01, Iwahashi, 金 00, 金 01, 宮田 01]．これは相互信念の自律的発展のメカニズムを言語形成の初期段階で実現しようとするものである．主眼を，所望の機能を表面的に実現するようなシステムをつくるのではなく，情報論的な原理の追及とした．本稿では次章以降，まず開発された言語獲得手法について述べ，これについて複数の視点から考察する．次に，人工知能研究の流れの中での位置付けや関連研究の紹介などを通して今後の展開について述べる．

## 2. ロボットによる言語獲得

### 2.1 タ ス ク

ロボットによる言語獲得のタスクを次のようなものとした．図1のようにロボットが配置され，人とロボットが互いにテーブルの上のオブジェクト（ぬいぐるみや箱など）を用いながら言語によるコミュニケーションを行

う．ロボットは，初期状態ではオブジェクトやその動かし方についての概念と概念に対応する単語および文法などの言語信念をもたない．ロボットはこうした言語信念を，まず発話と行動を介した人からの教示により受動的に学習し，さらに人とのインタラクションを通して能動的に学習する．ここでのインタラクションは，人とロボットのそれぞれが他方に対して，テーブルの上の一つのオブジェクトを動かすように指示する発話を行い，聞き手がその発話を理解して行動する，といったものとする．このとき，人は接話マイクに向かって単語間に短いポーズを入れてゆっくりと発話することとする．ロボットが学習を通して人の発話を状況に応じて適切に理解し，また自然な発話ができるようになることが目的である．

### 2.2 アルゴリズム概要

ロボットは，アームの横に配置されたステレオカメラユニットの前 90 cm 以内にあるオブジェクトを検出し認識の対象とする．人が指差したオブジェクトがどれであるかを判別する．オブジェクト画像の特徴量は，色（3次元），大きさ（1次元），形（2次元）で表される．オブジェクトの動きを検出し，動きはじめてから静止するまでの連続した動きの軌道を抽出する．軌道はオブジェクトの一定時間ごとの位置情報の時系列によって表される．人がオブジェクトを指差したり動かしながら発話すると，そのときの画像情報と音声情報を関連付ける．ロボットが人の発話を理解しその結果に従って行動したときに，人に手を軽く叩かれると，ロボットはその行動が間違っていたと判断する．

言語獲得に関しては，音韻と語彙，関係の概念，文法，および語用のそれぞれの相互信念を四つのアルゴリズムが別々に学習する．音韻と語彙，関係の概念，および文法の学習では，人がロボットに対してオブジェクトを提示したり動かして見せたりすることによる共同知覚経験を基盤にして，対応付けられた音声情報と画像情報の結合確率密度関数を推定することを基本原理とする．語用論的相互信念の学習では，互いに相手の発話に従って行動する共同行為を基盤にして，人の発話をロボットが正しく理解する確率を最大化することと，ロボットが生成した発話を人が正しく理解する確率を推定することを基本原理とする．

なお，アルゴリズムは人が協力的に振る舞うことを前提にしている．アルゴリズムの基本原理の追求を目的としているため，各相互信念はかなり単純なものとなっている．すべてのアルゴリズムを通して学習規準の一貫性があるべく保たれるように考慮されているが，四つのアルゴリズムは個別に評価されており，全体の統合は行われていない．次章以降で各アルゴリズムを説明していくが，特に共同行為を通じた語用論的相互信念の学習に焦点を当てることにする．

### 2.3 音韻と語彙の学習

幼児の語彙獲得は2.7節で述べるように、とても複雑で興味深い特徴をもっている認知現象であるが、ここでは学習問題を次のように単純化した。人がオブジェクトを机の上に置いたり指差したりしながら、そのオブジェクトについての単語を発話するものとし、ロボットはこのときの音声とオブジェクトを関連付ける。これを繰り返すことで得られる音声の特徴量  $s$  とオブジェクト画像の特徴量  $o$  の対の集合を学習データとする。語彙  $L$  は、各語彙項目に対応した音声の確率密度関数 (pdf) とオブジェクト画像の pdf の対の集合、 $p(s|c_i)$ ,  $p(o|c_i)$  ( $i=1, \dots, M$ )、で表されるものとする。ここで、 $M$  は語彙項目の数であり、 $c_1, c_2, \dots, c_M$  は語彙項目を表すインデックスである。このとき、語彙項目数  $M$ 、および語彙を構成するすべての pdf  $p(s|c_i)$ ,  $p(o|c_i)$  ( $i=1, \dots, M$ ) を表すパラメータを学習することが目的である。この問題の特徴は、二つの連続特徴量空間におけるクラスメンバーシップ関数の対の集合を、対の数が未知という条件で教師なし学習により求めることである。

学習は次のように行う。各語彙項目に対して単語の音韻列が決められていても音声は発声ごとに変動するが、通常、各発声におけるその変動は、その発話が示しているオブジェクトの特徴を反映しないので、

$$p(s, o|c_i) = p(s|c_i)p(o|c_i)$$

と置くことができ、全体での音声とオブジェクト画像の結合 pdf は、

$$p(s, o) = \sum_{i=1}^M p(s|c_i)p(o|c_i)P(c_i)$$

と書ける。そして上記語彙学習問題を、このように表現された  $p(s, o)$  に対して最適なモデルを選択し確率分布パラメータの値を推定する統計的学習問題とみなす。語彙は、正確な情報伝達が行えて、かつなるべく少ない語彙項目数で構成されていることが望ましいであろうという考えに基づき、語彙項目数  $M$  を音声とオブジェクト画像の相互情報量を規準にして選択したところ、色、形、大きさ、ぬいぐるみの名前を意味する十数単語程度を学習する実験で良好な結果が得られた [金 00, 金 01]。また、音韻の pdf を表す隠れマルコフモデル (HMM) の結合により単語音声の pdf を表すことで、音韻 pdf の集合も同時に学習できることと、動かされたオブジェクトの軌道を画像特徴量として使用できることも示されている。

### 2.4 関係の概念の学習

言語の意味構造は、モノと二つ以上のモノの関係に分けることができる。2.3節においてモノの概念は語彙項目が与えられたときのオブジェクト画像の条件付き pdf で表された。関係には、最も際立つもの (トラジェクタという) とトラジェクタの基準点として働くもの (ランドマークという) が関与する [Langacker 91]。例えば、図2のようにぬいぐるみが動かされているとき、この動

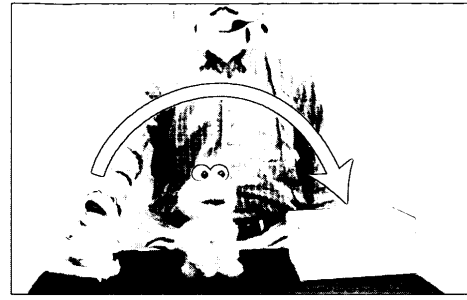


図2 提示動作の例

きは、動かされているぬいぐるみをトラジェクタとし、真中に置かれているぬいぐるみをランドマークとみなせば、『飛び越える』という意味に解釈でき、右側の箱をランドマークとみなせば、『のる』という意味に解釈できる。このような情景の集合を学習データとして用いて、オブジェクトの動かし方に関する概念を、トラジェクタとランドマークの位置関係の変化のプロセスとして学習する [羽岡 00]。動きの概念は、語彙項目  $c$ 、トラジェクタオブジェクト  $t$  の初期位置  $o_{t,p}$ 、およびランドマークオブジェクト  $l$  の位置  $o_{l,p}$  が与えられたときの動きの軌道の条件付き pdf  $p(u|o_{t,p}, o_{l,p}, c)$  で表される。アルゴリズムは、情景の中でどれがランドマークとなるオブジェクトであるかという非観測情報を推定しながら、動きの概念の条件付き pdf を表す HMM を学習する。同時に、動きの軌道を適切に記述する座標系の選択も行われる。例えば、『のる』の軌道は、ランドマークを原点、垂直と水平方向に軸とする座標系、『離れる』の軌道は、ランドマークを原点、トラジェクタの初期位置とランドマークを結ぶ線の一つの軸とする座標で記述される。

### 2.5 文法の学習

発話中の単語が表す概念の間の関係を表すための単語の並びの規則である文法の学習および運用において、前節で述べた関係の概念が重要な役割を果たす。ロボットに文法を学習させるとき、人がオブジェクトを動かしながら、その動作を表す発話をするものとする。これを繰り返すことで得られる、動作前の情景情報  $O$  と音声  $s$  と動作  $a=(t, u)$  の組  $(s, a, O)$  の集合が学習データとして用いられる。ここで  $O$  は情景の中のすべてのオブジェクトの位置と画像特徴量の集合で表される。 $t$  は、各情景の中の各オブジェクトに対して一意に与えられるインデックスのうち、トラジェクタオブジェクトを示すものである。 $u$  はトラジェクタの軌道である。発話は機能語を含まないかなり単純なものとした。例えば図2の動作を表す発話は、『大きい カーミット 茶色 箱 のせて』である。情景  $O$  と動作  $a$  は発話の意味構造を推測するために用いられる [羽岡 00, 羽岡 01]。意味構造  $z$  はトラジェクタとランドマークと軌道を構成要素とし、各要素への発話中の単語の対応付けによって表される。例えば上記発話の意味構造は次のようになる。

$$\left[ \begin{array}{l} \text{トラジェクタ：大きいカーミット} \\ \text{ランドマーク：茶色箱} \\ \text{軌道：のせて} \end{array} \right],$$

文法  $G$  は発話におけるこれらの構成要素の出現順序の生起確率分布によって表され、音声  $s$  と動作  $a$  と情景  $O$  の結合 pdf  $p(s, a, O; L, G)$  の尤度を最大にするように学習される。対数結合 pdf は語彙  $L$  と文法  $G$  のパラメータを用いて次のように表される。

$$\begin{aligned} \log p(s, a, O; L, G) & \\ & \approx \max_z (\log p(s|z; L, G) + \log p(a|z, O; L) \\ & \quad + \log p(z, O)) \\ & \approx \alpha + \max_{z, l} (\log p(s|z; L, G) \quad \text{[音声]} \\ & \quad + \log p(u|o_{l,p}, o_{l,p}, W_M; L) \quad \text{[動き]} \\ & \quad + \log p(o_{l,f}|W_T; L) + \log p(o_{l,f}|W_L; L)) \\ & \quad \text{[オブジェクト]} \end{aligned}$$

ここで、 $W_M$  と  $W_T$  と  $W_L$  は、それぞれ意味構造  $z$  の中の軌道とトラジェクタとランドマークに対応する単語(列)、 $o_{i,f}$  はオブジェクト  $i$  の特徴量を表す。 $\alpha$  は正規化項である。

## 2.6 語用論的相互信念の学習

語彙  $L$  と文法  $G$  を学習すれば、ロボットは結合 pdf  $p(s, a, O; L, G)$  の最大化を規準にしてある程度の発話理解ができるようになるが、より状況に依存した発話の理解と生成を可能とするために、共同行為を通して語用論的相互信念をオンラインで漸増的に学習する。ここで相互信念を用いた発話の生成と理解は例えば次のようなものである。図2において、人が直前の動作でカーミット(左のぬいぐるみ)をテーブルの上に置いたとしよう。そして、人がロボットにカーミットを箱の上に置いてもらいたいとき、「カーミット箱のせて」と限定的に言ってもよいが、もし、直前の動作で動かされたオブジェクトが次の動作の対象になりやすいという信念をロボットがもっていると人が想定していれば、「箱のせて」と断片的に言うかもしれないし、さらに、箱には何かかのせられやすいという信念を想定していれば、「のせて」と言うだけかもしれない。ロボットがこのような断片的な発話を理解するためには、ロボットも同じような信念をもっていて、それらを人と共有しているのだと想定していなければならない。ロボットが発話を生成する場合も同様である。

### §1 相互信念の表現

アルゴリズムにおいて相互信念は以下の二つの部分によって表される。

- (1) 発話と動作の対応の適切さを表す決定関数  $\Psi$ 。重み付けられた信念の集合によって表される。重みは各信念が対話者とロボットに共有されているこ

とに対するロボットの確信度を表す。

- (2) 決定関数  $\Psi$  に対するロボットの確信度を表す全体確信度関数  $f$ 。ロボットの発話を対話者が正しく理解する確率の推定値を出力する。

アルゴリズムはさまざまな信念を扱うことが可能であるが、前節までで述べた音声、オブジェクト、動きのおおのに関する信念(これらは語彙と文法によって表される)に加えて、次の二つの非言語的信念を例として扱った[宮田01]。

- 行動コンテキスト効果  $B_1(i, q; H)$   
行動コンテキスト  $q$  のもとで、オブジェクト  $i$  が発話による指示対象になるという信念。 $q$  は、各オブジェクトが直前の動作においてトラジェクタまたはランドマークとして関わったかどうか、およびジェスチャによって注意が向けられているかどうかなどについての情報で表される。この信念は二つのパラメータ  $H = \{h_c, h_g\}$  で表され、 $q$  に応じて対応するどちらかのパラメータの値または0を出力する。

- 動き-オブジェクト関係  $B_2(o_{i,f}, o_{l,f}, W_M; R)$   
オブジェクトの特徴量  $o_{i,f}$  と  $o_{l,f}$  が、それぞれ動き概念  $W_M$  におけるトラジェクタとランドマークの特徴量として典型的なものであるという信念。条件付き結合 pdf  $p(o_{i,f}, o_{l,f} | W_M; R)$  によって表される。この共起 pdf はガウス分布で表現され  $R$  はそのパラメータ集合を表す。

これらの信念モデルの出力の重み付け和として決定関数は次式で表される。

$$\begin{aligned} \Psi(s, a, O, q, L, G, R, H, \Gamma) & \\ & = \max_{l, z} (\gamma_1 \log p(s|z; L, G) \quad \text{[音声]} \\ & \quad + \gamma_2 \log p(u|o_{l,p}, o_{l,p}, W_M; L) \quad \text{[動き]} \\ & \quad + \gamma_2 \log p(o_{l,f}|W_T; L) + \log p(o_{l,f}|W_L; L) \\ & \quad \quad \quad \text{[オブジェクト]} \\ & \quad + \gamma_3 \log p(o_{l,f}, o_{i,f} | W_M; R) \\ & \quad \quad \quad \text{[動き-オブジェクト関係]} \\ & \quad + \gamma_4 (B_1(t, q; H) + B_1(l, q; H))) \\ & \quad \quad \quad \text{[行動コンテキスト]} \end{aligned}$$

ここで、 $\Gamma = \{\gamma_1, \dots, \gamma_4\}$  は各信念に対する重みパラメータの集合である。発話  $s$  の解釈としての行動  $a$  は  $\Psi$  の値を最大化するものとして決定される。

次に全体確信度関数  $f$  について述べる。まず、情景  $O$  と行動コンテキスト  $q$  のもとで動作  $a$  を表す発話  $s$  の生成を決定する際の決定関数の値のマージンを次式のように定義する。

$$\begin{aligned} d(s, a, O, q, L, G, R, H, \Gamma) & \\ & = \min_{A=a} (\Psi(s, a, O, q, L, G, R, H, \Gamma) \\ & \quad - \Psi(s, A, O, q, L, G, R, H, \Gamma)) \end{aligned}$$

そして全体確信度関数  $f$  は、このマージンの値を入力として発話が正しく理解される確率を出力するもので次式で表される。

$$f(x) = \frac{1}{\pi} \arctan\left(\frac{x - \lambda_1}{\lambda_2}\right) + 0.5$$

$\lambda_1$  と  $\lambda_2$  は関数  $f$  を表すパラメータである。マージンが大きいと正しく理解される確率が高くなるであろうということが前提とされている。もしマージンが少なくても高い確率で正しく理解されるならば、ロボットが想定する相互信念が、人のそれとよく一致していることを意味する。

## § 2 学習アルゴリズム

決定関数  $\Psi$  と全体確信度関数  $f$  は、発話理解と生成のそれぞれの過程で別々に学習される。まず、決定関数は次のようなエピソードの繰返しを通してオンラインで漸増的に学習される [宮田 01]。

- (1) 人がロボットに一つのオブジェクトを動かすことを表す発話を行う。
- (2) ロボットが発話を理解し、そのとおりに行動する。
- (3) もしロボットが発話内容に対して正しい動作を行ったならば終了。そうでないなら人がロボットの手を軽くたたく。
- (4) ロボットが1回目とは異なる動作を行う。
- (5) もしロボットが間違った動作を行ったならば、人がロボットの手を軽くたたく。終了。

ロボットは各エピソードにおいて正しい動作が行えたとき、そのエピソードで最初に与えられる情景  $O$ 、コンテキスト  $q$ 、発話  $s$ 、動作  $a$  を対応付け、これを教師データとする。教師データが得られる度に、動作-オブジェクト関係の信念  $R$ 、行動コンテキスト効果の信念  $H$ 、および重みパラメータ  $\Gamma$  が逐次的に更新される。 $i$  番目の対応付けサンプル  $(s_i, a_i, O_i, q_i)$  が得られた後の、重みパラメータの集合  $\Gamma_i$  は、理解誤り率を最小化 [Juang 92] するように、次の規準で最適化される。

$$\sum_{j=i-K}^i w_{i-j} g(d(s_j, a_j, O_j, q_j, L, G, R_i, H_i, \Gamma_i)) \rightarrow \min$$

ここで

$$g(x) = \begin{cases} -x, & \text{if } x < 0 \\ 0, & \text{otherwise} \end{cases}$$

$K$  と  $w_{i-j}$  はそれぞれ、用いられる過去の学習サンプル数と各サンプルに対する重みである。

一方、全体確信度関数  $f$  は、次のようなエピソードの繰返しを通してオンラインで漸増的に学習される [Iwahashi 03]。

- (1) ロボットが人に一つのオブジェクトを動かすことを表す発話を行う。
- (2) 人が発話を理解し、そのとおりに行動する。

- (3) ロボットが人の動作が正しいかどうか判断する。終了。

各エピソードにおいて、ロボットは、人に行わせようとする動作に対して、全体確信度関数の値があらかじめ決められた値とにできるだけ近くなるように発話を決定する。ロボットは、単語の多い発話をすることによって正しく理解される確率を高めることができるし、ある程度の確率で理解されればよいならば必要以上の単語を使う必要がなくなる。ここで重要なことは、単語数を節約できるということではなく、断片的な発話が生成、理解されることで相互信念の形成が促進されるということである。各エピソードにおいて、発話生成の際に生じたマージン  $d$  の値に対して、その発話が対話者に正しく理解されたかどうかの情報が対応付けられ、これを学習データとする。 $i$  番目のエピソードが終了したときの  $f$  のパラメータ  $\lambda_{1,i}$ 、 $\lambda_{2,i}$  は次のように更新される。

$$[\lambda_{1,i}, \lambda_{2,i}] \leftarrow (1 - \delta)[\lambda_{1,i-1}, \lambda_{2,i-1}] + \delta[\tilde{\lambda}_{1,i}, \tilde{\lambda}_{2,i}]$$

ここで

$$\begin{aligned} & (\tilde{\lambda}_{1,i}, \tilde{\lambda}_{2,i}) \\ & = \arg \min_{\lambda_1, \lambda_2} \sum_{j=i-K}^i w_{i-j} (f(d_j; \lambda_1, \lambda_2) - e_j)^2 \end{aligned}$$

$e_i$  は発話理解が正しいなら 1、誤りなら 0 をとる。 $\delta$  は学習速度を決定する値である。

## § 3 実験

まず決定関数の学習実験について述べる。エピソード系列のはじめの部分では、比較的完全な発話を与え（例えば、「みどり カーミット あか 箱 のせて」）、その後、徐々に断片的なもの（例えば、「のせて」）を与えていった。このような与え方をしたことによって、比較的完全な発話が正しく理解され続けている間に信念パラメータ  $R$  の推定が良好に行われ、後に断片的な発話を与えたとき、この推定値に基づき重み  $\Gamma$  の学習が効果的に行われた。発話のあいまい性が増すに従って重みの値が変化し

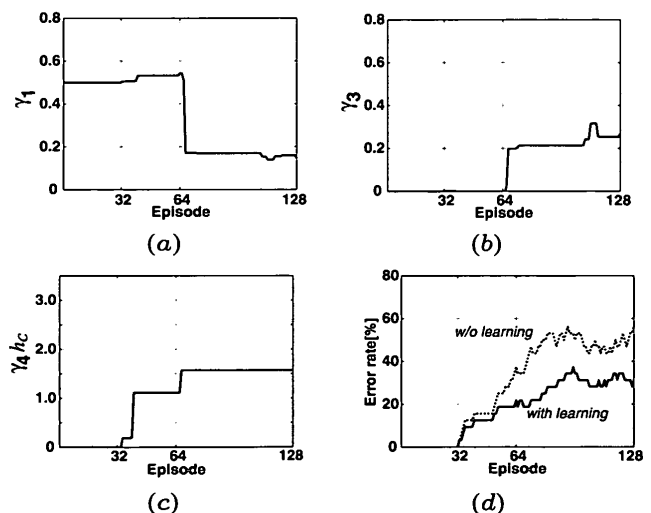


図3 重み値の変化 (a)~(c) と理解誤り率の変化 (d)

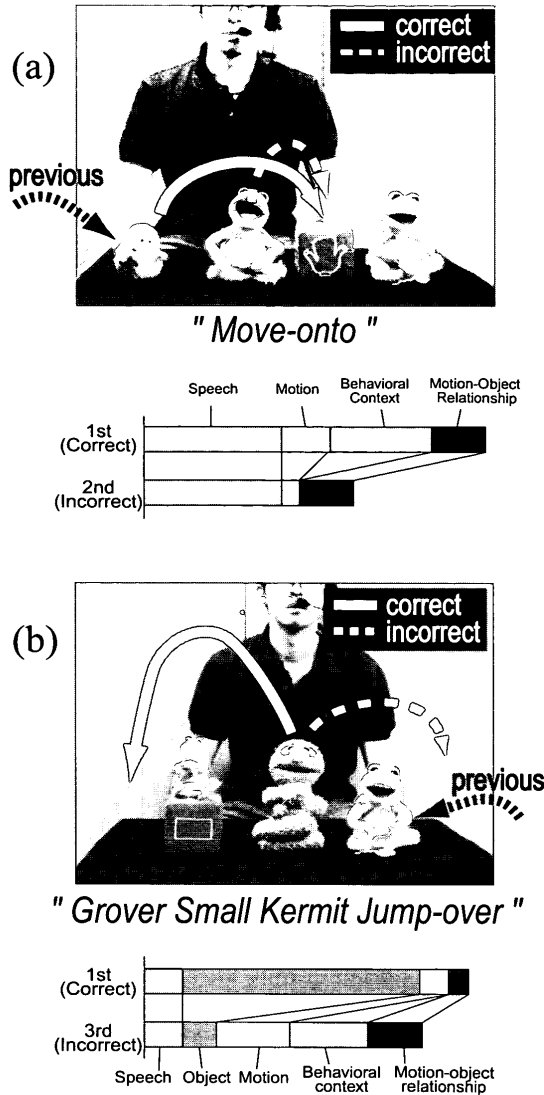


図4 断片的な発話が正しく理解されたときの行動の例

ていくようすを図3(a)~(c)に示す。また、理解誤り率の変化を、学習をしない場合の変化と合わせて図3(d)に示す。また、発話を正しく理解したときの行動例を、各信念モジュールからの重み付けられた出力値とともに図4に示す。比較のために低い順位として得られた理解候補も示す。図4(a)では行動コンテキスト効果の信念が、また、図4(b)ではオブジェクトに関する信念が、それぞれ正しい理解を導くために効果的に使われたことがわかる。

次に、全体確信度関数  $f$  のシミュレーションによる学習実験について述べる。 $f$  の初期形状は、発話が理解されるために大きなマージンを必要とするような状態、すなわち相互信念の全体的な確信度が低い状態、を表す形状に設定した。発話生成に使われる  $\xi$  の値は 0.75 に固定した。 $\xi$  を固定しても実際に得られる  $f$  の出力値は  $\xi$  の周りでばらついて、かつ、発話が正しく理解されたりされなかったりしたので、 $f$  が  $f^{-1}(\xi)$  のまわりの比較的広い範囲で良好に推定できた。 $f$  の変化と、動作に関するすべてのオブジェクトを記述するために使用された

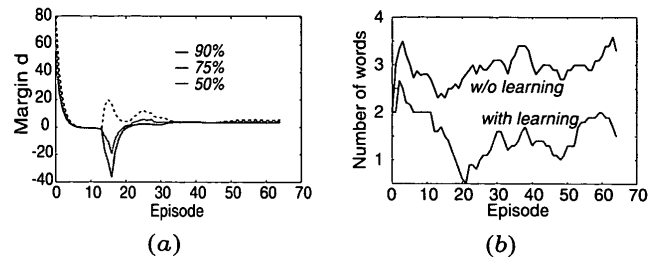


図5 学習過程における全体確信度関数の変化 (a) と各発話でオブジェクトを記述するために用いられた単語の数 (b)

単語数の変化のようすを図5に示す。 $f$  の形状の変化がわかりやすいように、 $f^{-1}(0.9)$ ,  $f^{-1}(0.75)$ ,  $f^{-1}(0.5)$  の三つの値をプロットした。学習開始からすぐにこれらの値が急速に0に近づいていき、使用する単語の数が少なくなった。これは相互信念に対する確信が強まったことを意味する。その後、15エピソード付近で単語数が少なくなり過ぎて間違えて理解されることが多くなったため、全体確信度関数の傾きが小さくなっており、いったんは相互信念に対する確信を弱めている。

2.7 言語処理の観点からの考察

幼児は、ある物体に対して新規のことばを与えられたとき、その物体に関する非常に多くの可能なことばの意味の中から、ただ一つだけの意味を選択することができる。このようなことばの意味付けは即時マッピングと呼ばれ、意味の決定においてさまざまな原理\*3が働いていることが心理実験により確かめられている [今井 97]。これらの原理は、すでに幼児が知っている語彙との関連や状況に応じて適応的に作用する [針生 00]。即時マッピングは、メルロ=ポンティの「対象に名前をつける働きは、対象をそれだと知る作用の後にやってくるのではなく、認知作用そのものである」という言語観 [長谷川 86] を支持し、経験に基づいた言語信念の再帰的な発展性が顕著に観測できる現象であるように思われる。一方で、即時マッピングは工学的実用性の観点からも重要である。あらかじめ動作環境が定められていないロボットは、環境中の事物に対応したことばを人とのインタラクションを通して効率良く学べたほうがよいからである。上記アルゴリズムでは即時マッピングを扱っていないが、現在、そこで用いられた学習原理を発展させて、即時マッピングの制約が言語学習経験によって形成されることの情報論的なモデルを構築し実験を行っているところである。

発話の生成、理解では、トラジェクタ、ランドマーク、および軌道を構成要素とする概念構造が使われた。この

\*3 例えば、子どもは未知の名詞が物体全体の名前であり、その部分性や属性を指示することばでない想定するバイアス [Markman 90] など。

概念構造を用いて学習された動きの概念モデルは、与えられた空間においてロボットがカメラを通して観測した軌道のデータを統計的に表現するものであり、ロボットの身体性を直接反映しているものである。このような物理的な領域の概念モデルを、抽象的領域の概念の理解に利用することは、これから取り組まなければならない研究テーマであろう。発話を理解するイメージスキーマ—人間が身体を介して環境と相互作用する経験で繰り返し生じる比較的単純な概念構造—は、世界の中での経験と思考を結びつける重要な媒介であり、抽象的な概念の理解を支えていると考える認知言語学 (e.g. [Lakoff 87, 山梨 00]) との関係は極めて密接である。

次に語用論的信念の形成のアルゴリズムにおける行為の誤りと修復の意味について考察する。ロボットの発話理解過程における学習では、1回目では誤った動作を行って、かつ、2回目で正しい動作が行えたエピソードで、相互信念のパラメータが比較的大きく更新される。また、ロボットの発話生成による学習では、 $\xi$  を 0.75 にした場合の実験結果を示したが、これを 0.95 にした別の実験では、ほとんどすべての発話が正しく理解されたことから  $f$  の推定を適切に行うことができなかつた。発話理解と発話生成の両方のアルゴリズムにおいて、発話がときどき間違えて理解されることが相互信念の形成を促進していることがわかる。相互信念を形成するには、発話が意味を正しく伝達するだけでは不十分で、そこに誤解されるリスクが付与されていなければならないのである。そのようなリスクを人とロボットで共有することが、発話が相互信念の情報を同時に送受信するという機能を支えていると言える。アルゴリズムではこのようなメカニズムに基づいて、ロボットの相互信念の構造的な全体が自律的に発展し、人の相互信念とカップリングすることを単純な形でではあるが実現している。また、相互作用に基づく学習における誤りと修復の重要性は、機械学習の分野でも、強化学習の主要なテーマの一つ、探査か知識利用かのジレンマ [Dayan 96]、として研究されてきている。強化学習を行うエージェントは、現在の知識を利用して報酬を獲得しながら、将来的に行動選択を改善するためには探査も行わなければならない。探査か知識利用かのジレンマとは、探査も知識利用も与えられた作業の失敗なしには独自に遂行されることはないということである。

従来から行われている相互信念の形成に関する理論研究 (e.g. [Clark 96]) では、具体的なインタラクションの分析を詳細に行うための手段を提供してくれていて大変興味深い。このような理論研究が相互信念の形成を手続きや規則で表現しようとしているのに対して、上記アルゴリズムは、相互信念の形成を認知プロセスを反映した言語のダイナミクスとして数理的に表現しようとするものであり、両者のアプローチの違いは大きい。

相互信念の形成を実用的な言語処理システムの設計の

観点で考察する。一般に、言語処理システムは、その開発者自身が使うと大変良く動作するが、開発者以外の事前知識がない人が使うと予期しない事態 (例えば、ユーザがシステムの知らない単語を使用するなど) が頻繁に起きてなかなかうまく動作しないものである。その原因の一つは、システムとユーザが互いに想定する相互信念が一致していない、すなわち相互理解が欠如していることである。もし、ユーザがもっている言語信念を、極めて簡単な手段でユーザがシステムに伝えることができれば、誰が使っても開発者自身が使うときのように期待どおりに心地良く動いてくれるシステムとなるであろう。信念の伝達手段として言語を用いるならば、言語による相互信念の伝達特性を考慮する必要が生じるのである。

## 2・8 機械学習の観点からの考察

いわゆる No Free Lunch Theory [Wolpert 95] によれば、問題に関する事前情報がまったくない場合、ある学習アルゴリズムが別のものよりも優れているという根拠がない。つまり、あらゆる問題において効率的に学習できるアルゴリズムなど存在しないのである。このことは、認知処理システムにおける学習アルゴリズムの構築において汎用性を考慮するときの重要な指針となる。汎用性よりも領域固有性に注目してよいのである。

環境との相互作用を重視したロボティクスのアプローチでは、サブシステムを並列に配置するアーキテクチャ [Brooks 91] の妥当性が主張されている。このアーキテクチャによれば、物理世界と接触するシステムが直面する二つの問題、情報の部分性と処理の実時間性に対して柔軟に対応できる [松原 89]。また、部分的な情報を用いた決定は Bayesian Networks において理論的に研究されている [Pearl 88]。そこでの決定はネットワーク上での Belief Propagation と呼ばれる局所的な確率情報の連鎖的な受渡しに基づいており、サブシステムの配置の並列性は特に必要とされない。発話生成と理解のために使われた決定関数  $\Psi$  は、観測情報の結合 pdf を表す Bayesian Networks に対して、その主要なノードの集合に重み  $\Gamma$  を付与して誤り率最小化規準により最適化するように、修正を加えて構成されたものなのであるが、2・6節で示した  $\Psi$  の式で表されているように主要なノードだけに注目すればサブシステムが並列に結合されたアーキテクチャであるとも見られる。

幼児による言語獲得では、さまざまなモダリティからの情報を効果的に用いることで、少ない言語経験からの学習を可能にしていると考えられる。一方、機械学習理論によれば一般に、観測データの次元数が増すに従い、学習に必要な学習データ数が爆発的に増えてしまう。機械学習がこの問題を克服して、幼児のように環境の中に存在する大量の情報を扱うためには、環境に存在する利用可能な情報構造の抽出、多数のサブシステム間の協調のメカニズム、認知的に妥当な制約の利用、などが適

切に実現されていなければならない。例えば、前記アルゴリズムにおいて音声情報に視覚情報を組み合わせて効率的な学習を可能としている要因をあげてみる。語彙学習では、まず2・3節で述べたように、各単語において音声信号とそれが指し示す具体的なオブジェクトの特徴とが独立であるという情報構造の存在が利用されている。これに加えて学習の最適化過程においては、人間は、一つのオブジェクトに対して複数の概念を対応させることはできるが、一つの音声セグメントに対して一つの音韻カテゴリ列しか知覚できない、という認知的に妥当な制約が利用されている。また文法学習では、トラジェクタとランドマークの空間的関係の時間変化のスキーマに基づいて、行動を解釈するという情景理解における概念的制約と、発話中の個々のオブジェクトに対応する単語の間に別のオブジェクトに対応する単語が入らないという音声理解における構文的制約とを協調させている。

本来、語彙、文法、および語用論的信念のそれぞれの学習は、学習フェーズの変化を伴いながら互いに協調して行われるべきである。そのような協調構造も、環境との相互作用によってある程度適応されることが望ましいと思われるが、人とロボットの間での少ないインタラクションで相互信念を効率的に形成することを可能とするには、適切に事前知識を設定することが重要であると考えられる。

### 3. 展 開

ここで、機械が日常的な経験を通して相互信念を形成することの困難さについて基本的な議論に立ち戻ってみよう。まず、ハイデガーおよび後期ウィトゲンシュタインの思想に依拠した Dreyfus の重要な指摘[Dreyfus 88]を取り上げよう。Dreyfus によれば、常識的背景は、技能、習慣、識別、等々の組合せであって、それらは、志向的な状態ではなく、したがってなおさらのこと、要素と規則によって説明され得るようないかなる表象的内容をもっていない。これは物理記号システム仮説——物理記号システムは一般的な知的活動のための必要かつ十分な手段を有しているという仮説——に基づいた計算主義的な AI への批判である。さらに、コネクショニズムに対しても次のような指摘をしている。人間の知能の大半は、コンテクストに対して適切な仕方で一般化することで成り立っているが、このような能力をネットがもつためには、ネットは我々がもっている出力の適切さの感覚を共有していなければならない。それが意味するものは、人間がもっているような身体を、ネットがもっていないなければならないということである。また、SHRDLUを開発した Winograd は、機械による言語理解の困難さに関して、こうしたことに加えて、Maturana が言うところの、個人が集まった集団の中で存在し、言語的活動と、その活動から生成される構造的カップリングを通じ

て絶えず再生される言語 [Maturana 78] を、表象を用いる計算主義的なアプローチに基づいたシステムに与えることはできない、と指摘した [Winograd 86]。これらの批判が依拠している知能に関する思想を積極的に取り入れているのが、主体と環境の相互作用を知能の本質としたロボティクスのアプローチ [Brooks 91] や、認知の複雑にカップリングした振舞いを力学系理論を道具にして理解しようとするアプローチ [Gelder 95] などである。これらのアプローチは表象を操作することは認知にはまったく不要だと主張し、伝統的計算主義的アプローチへの代替案として期待されている。しかしながら一方で、複雑ではあるが表象を必要としない振舞いに関する研究事例における有効性は示されているが、表象が不可欠な種類の問題領域に注意がまったく払われていない [Clark 94]、と指摘されている。Clark によれば、表象が不可欠な問題領域とは、次の二つの条件のうち一つもしくは両方が当てはまる領域である。① 現前しない、または存在しない対象、あるいは反事実的な事態についての推論が必要である。② 周囲の物理的環境において複雑で扱いにくい形で現れるパラメータに対して選択的に反応する能力が主体に要求される。この領域は、知覚的認識や知覚に導かれたさまざまな事例を含んでいる。言語活動は言うまでもなくこの問題領域にあり、しかもこれまでに述べてきたように環境や他者とのカップリングが本質的なのである。これは本稿で相互信念の形成を取り上げて議論していることである。本稿で示したアプローチは、言語という知能に対して、言語固有の特質よりはむしろ、環境との相互作用に基づく認知の基本的かつ一般的な性質に注目しようというものである。そこでは、相互信念間のカップリングを基盤にして、相互信念の自律的發展、概念構造の形成、および言語の意味付けがダイナミックに行われるシステムの実現をねらいとしている。これは、1章で述べたように、機械による言語理解における伝統的なパラダイムとは、アプローチの仕方が逆という点で根本的に異なるのである。

次に、関連した興味深い研究のいくつかを紹介しよう。[稲邑 01] では、自律的に行動するロボットが、行動を指示するユーザとの対話を通してその行動を学習する方法が示されている。行動および対話の戦略はセンサ情報とユーザの指示の共起関係を表すベイジアンネットワークによって制御される。[Singh 00] などでは、強化学習を用いて最適な対話戦略を学習する方法が示されている。対話のパフォーマンスを、タスクの達成、対話の短さなどで評価する。[Gorin 91] では、対話によりシステムがユーザの発話が指示するシステムの反応を学習する方法が示されている。単語をノードとする確率ネットワークを学習することで単語やフレーズの意味を獲得する。小野らは、ロボットの発話と動作の関係が、人間の理解にどのような影響を与えるかを研究している[小野 01]。人がロボットの動作に引き込まれるように同調的



な動作をしたとき、円滑なコミュニケーションが行われることが示されている。小嶋らは、ロボットを用いて人間の認知発達過程の研究をしている [Kozima 01]。視線・指差し・表情などの前言語的コミュニケーション能力を与えられたロボットが、人と共同注意を形成する過程について調べられている。

これらの関連研究と本稿で示したアプローチはすべて、対話者を含めた環境との相互作用を通して、相互理解を形成することを重視するものであり、システムが対話者と協調して振る舞うことを可能とする。さてここでもう一度 ELIZA を思い出してみよう。ELIZA (正確にはそのうちのひとつ) はある心理療法の診察をシミュレートするパターンで動かされており、その意味で人間の心のモデルを基本にしたと言える。そして実際に対話者に対してかなり協調的に反応することができた。ところが、多くの対話者は ELIZA と少しの間会話を楽しんだ後、対話があまりにも空虚であることに気付き、すぐに相互理解という会話の目的を放棄してしまうのである。ここで、たとえ新しいシステムがシステム自身の身体をもって人と協調的に行動できて、かつ環境や状況の変化に応じて適応的に行動を変化させたとしても、ELIZA の場合と同じようなことが起きてしまうのではないかという疑問が生じる。こうしたことはシステムの知能のレベルや複雑さを高めていけば自然に解決されていくものではないだろう。人が知能的な機械とのコミュニケーションにおいて何を期待するかについて真剣に考えるべきであろう。

#### 4. おわりに

機械との言語コミュニケーションによる相互理解を実現するためには、まだ果てしなく長い道りがあるだろう。しかし、非常に困難であると考えられている問題でも、パラダイムを変えるとあまり深刻な問題でなかったことに気づくことがある。いま、広範な関連領域において多くの刺激的な研究が進められており、それらから学ぶべきことが多い。

#### 謝 辞

研究の初期段階から熱心に議論していただいた飯田 仁氏 (東京工科大学)、言語発達心理学の視点から多くの貴重なコメントをいただいた今井むつみ氏 (慶應義塾大学)、認知ロボティクスの魅力を教えていただいた谷淳氏 (理化学研究所)、草稿の段階で詳細なコメントをいただいた富浦洋一氏 (九州大学) に感謝いたします。

#### ◇ 参 考 文 献 ◇

- [Brooks 91] Brooks, R.: Intelligence without Representation, *Artificial Intelligence*, Vol.47, pp. 139-159 (1991)
- [Clark 94] Clark, A. and Tribio, J.: Doing without representation?, *Syntheses*, Vol.101, No.3, pp. 401-431 (1994); 金杉武司 訳: 表象なしでやれるのか?, 門脇俊介, 信原幸弘 編: ハイデガーと認知科学, pp. 205-251, 産業図書 (2002)
- [Clark 96] Clark, H.: *Using Language*, Cambridge University Press (1996)
- [Dayan 96] Dayan, P. and Sejnowski, T. J.: Exploration Bonuses and Dual Control, *Machine Learning*, Vol. 25, pp. 5-22 (1996)
- [Dreyfus 88] Dreyfus, H. L. and Dreyfus, S. E.: Making a Mind Versus Modeling the Brain — Artificial Intelligence Back at a BranchPoint, *Daedalus*: 心をつくるか, それとも, 脳のモデルをつくるか. 分岐点に戻る人工知能, 門脇俊介, 信原幸弘 編: ハイデガーと認知科学, pp. 19-66, 産業図書 (2002)
- [Gelder 95] Gelder, van T.: What might cognition be, if not computation?, *The Journal of Philosophy*, No.7, pp. 345-381 (1995); 中村雅之 訳: 認知は計算でないならば, 何だろうか, 門脇俊介, 信原幸弘 編: ハイデガーと認知科学, pp. 151-203, 産業図書 (2002)
- [Gorin 91] Gorin, A.L. and Levinson, S.E.: Adaptive acquisition of language, *Computer Speech and Language*, Vol.5, pp. 101-132 (1991)
- [羽岡 00] 羽岡哲郎, 岩橋直人: 言語獲得のための参照点に依存した空間的移動の概念の学習, 電子情報通信学会技術研究報告 PRMU2000-105 (2000)
- [羽岡 01] 羽岡哲郎, 岩橋直人: 認知言語知識に基づく音声理解, 日本音響学会 2001 年春季発表会講演論文集, pp. 159-160 (2001)
- [針生 00] 針生悦子, 今井むつみ: 語意学習メカニズムにおける制約の役割とその生得性, 今井むつみ 編, 心の生得性, pp. 131-171, 共立出版 (2000)
- [長谷川 86] 長谷川宏: 言語の現象学, 世界書房 (1986)
- [今井 97] 今井むつみ: ことばの学習のパラドックス, 共立出版 (1997)
- [稲邑 01] 稲邑哲也, 稲葉雅幸, 井上博充: PEXIS: 統計的経験表現に基づくパーソナルロボットとの適応的インタラクションシステム, 電子情報通信学会論文誌, Vol. J84-D-I, No.6, pp. 867-877 (2001)
- [Iwahashi] Iwahashi, N.: Language Acquisition through a human-robot interface by combining speech, visual, and behavioral information, *Information Sciences*, to appear.
- [Iwahashi 03] Iwahashi, N. and Sundberg, P.: 投稿予定
- [Juang 92] Juang, B.-H. and Katagiri, S.: Discriminative Learning for Minimum Error Classification, *IEEE Trans. on Signal Processing*, Vol. 40, No. 12, pp. 3043-3054 (1992)
- [金 00] 金 影柱, 岩橋直人: 知覚情報の統合に基づく言語音声単位の獲得アルゴリズム, 電子情報通信学会技術研究報告, TI2000-21 (2000)
- [金 01] 金 影柱, 岩橋直人: 知覚情報の統合に基づく階層構造を有する音声単位の獲得, 日本音響学会 2001 年春季発表会講演論文集, pp. 99-100 (2001)
- [Kozima 01] Kozima, H. and Yano, H.: A Robot that Learns to Communicate with Human Caregivers, *Int. Workshop on Epigenetic Robotics* (2001)
- [Lakoff 87] Lakoff, G.: *Women, fire, and dangerous things: What categories reveal about the mind*, University of Chicago Press (1987); 池上嘉彦, 河上哲作 訳: 認知意味論—言語から見た人間の心, 紀伊國屋書店 (1993)
- [Langacker 91] Langacker, R.: *Foundation of cognitive grammar*, Stanford, University Press, CA (1991)
- [Markman 90] Markman, E. M.: Constraints children place on word meanings, *Cognitive Science*, Vol. 14, No. 1, pp. 57-77 (1990)
- [松原 89] 松原 仁, 橋田浩一: 情報の部分性とフレーム問題の解決不能性, 人工知能学会誌, Vol. 4, No. 6, pp. 695-703 (1989)
- [Maturana 78] Maturana, H. R.: Biology of language — The epistemology of reality, in Miller, G.A. and Lenneberg, E. (eds.), *Psychology and Biology of Language and Thought — Essay in Honor of Eric Lenneberg*, pp. 27-64 (1978)
- [宮田 01] 宮田篤人, 岩橋直人, 樽松 明: ロボットによる発話理解過程に基づく相互信念の形成, 電子情報通信学会技術研究報告, SP2001-98 (2001)

- [中野 02] 中野幹生, 堂坂浩二: 音声対話システムの言語・対話処理, 人工知能学会誌, Vol. 17, No. 3, pp. 271-278 (2002)
- [小野 01] 小野哲夫, 今井倫太, 石黒 浩, 中津良平: 身体表現を用いた人とロボットの共創対話, 情報処理学会論文誌, Vol. 42, No. 6, pp. 288-294 (2001)
- [Pearl 88] Pearl, J.: *Probabilistic reasoning in intelligent systems: Networks of Plausible Inference*, Morgan Kaufmann (1988)
- [Reddy 93] Reddy, M.J.: The conduit metaphor: A case of frame conflict in our language about language, in Ortony, A. (ed.), *Metaphor and Thought*, Cambridge University Press (1993)
- [Schank 84] Schank, R.: *The cognitive computer — on language, learning, and artificial intelligence*, Addison-Wesley (1984); 淵 一博 監訳, 石崎 俊 訳: 考えるコンピューター人間の脳に近づく機械, ダイヤモンド社 (1985)
- [Singh 00] Singh, S., Kearns, M., Litman, D.J. and Malker, M.A.: Empirical Evaluation of a Reinforce Learning Spoken Dialogue System, *Proc. AAAI*, pp. 645-651 (2000)
- [Sperber 95] Sperber, D. and Wilson, D.: *Relevance (2nd Edition)*, Blackwell (1995)
- [丹治 96] 丹治信春: 言語と認識のダイナミズム—ウイトゲンシュタインからクワインへ, 勁草書房 (1996)
- [Weizenbaum 66] Weizenbaum, J.: ELIZA — A computer program for the study of natural language communication between man and machine, *Communications of the Association for Computing Machinery*, Vol. 9, No. 1, pp. 34-45 (1966)
- [Winograd 72] Winograd, T.: *Understanding Natural Language*, Academic Press, New York (1972)
- [Winograd 86] Winograd, T. and Flores, F.: *Understanding Computers and Cognition — New Foundation for Design*, Alex, Norwood, N. J. (1986); 平賀 譲 訳: コンピュータと認知を理解する人工知能の限界と新しい設計理念, 産業図書 (1989)
- [Wolpert 95] Wolpert, D. H.: The relationship between PAC, the statistical physics framework, the Bayesian framework, and the VC framework, in Wolpert, D. H. (ed.), *The mathematics of Generalization*, Addison-Wesley, Reading, MA (1995)
- [山梨 00] 山梨正明: 認知言語学原理, くろしお出版 (2000)

2002年11月18日 受理

---

### 著者紹介

---

岩橋 直人 (正会員)



1985年慶応義塾大学理工学部計測工学科卒業。同年ソニー(株)入社。1990～93年ATR自動翻訳電話研究所。1998年(株)ソニーコンピュータサイエンス研究所入社。工学博士。人-ロボット言語コミュニケーションの研究に従事。電子情報通信学会, 日本音響学会, 日本認知科学会, 日本認知言語学会各会員。