# 主体と客体のはざまで揺らぐロボットの自己意識の 問題

谷淳

ソニーコンピュータサイエンス研究所

(141) 東京都品川区東五反田3-14-13高輪ミューズビル tani@csl.sony.co.jp, http://www.csl.sony.co.jp/person/tani.html

#### 1 はじめに

本文では、認知システムにおける主体と客体の問題について注目し、そこから自己意識の問題への取り組みについて構成論の立場から考える。知能ロボットの研究の姿勢は常に主体中心のトップダウン指向と客体中心のボトムアップ指向の両者の間をさ迷ってきたといえる。70年代から80年代前半にかけては、人工知能研究全盛の時代の中、記号による世界表現、そしてそれら記号操作に基づく推論計画といった、主体での計算を中心としてトップダウン型のロボットがひろく試みられた。しかしながら、このような研究においては、記号表現された内部世界像と、物理的実世界の間に、常にギャップが存在し、Harnad[1]がいうところの記号接地問題が深刻化した。一方80年代後半になってから、従来の人工知能研究を批判する形で現れた、Brooks[2]の提唱する behavior-based robotics が脚光をあびるようになった。この新たなロボットにおいては、内部の複雑なメンタルな計算をことごとく排除し、行動は単に現時点でのセンサー入力に対して反射的に生成していくという、客体世界からのボトムアップの流れを重視した設計がなされた。Beer[3]は、内的プロセスはセンサーモータループを通して、環境のダイナミクスに構造的に結合され、世界に無理なく"situate"すると言うが、しかしそこにはもとから環境と拮抗すべきは主体は見えない。

さて、筆者は、この両者の対立的関係性こそが重要であり、自己意識の問題はまさにその関係性から発生すると考える。つまり、主体のトップダウン的プロセスと客体からのボトムアップ的プロセスが密に相互作用しあい、その結果両者の関係性が常に変化するところに、初めて自己意識の問題がたち現れると考える。Harnad[1] は、記号システムで実現される主体プロセスと、外界からシグナルをじかに処理するアナログ的客体プロセスは、ニューラルネットなどにより実現されるカテゴライザーをもって接合できると主張する。しかし、カテゴライザーは、筆者が目指す密な相互作用を本質的に実現しえない。なぜならば、記号プロセスとアナログプロセスは、同じメトリクスの上に存在せず、そのインターフェースはつねに恣意的にならざるをえないからである。これに対する代案として筆者は力学系的アプローチに注目する。主体の予測・計画といった、従来、記号的プロセスと考えられたものを、アナログ力学系の時間発展プロセスとして実現することにより、主体的プロセスと客体的プロ

研究会報告

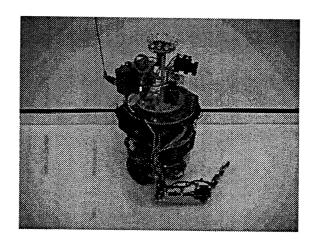


Figure 1: The vision-based mobile robot used in the experiments.

セスの両者は同じアナログのメトリクスのうえに存在でき、そこではじめて両者の本来的な 身体性をもった相互作用が実現できると考える。

以下に筆者の行ってきたロボットの行動学習実験の一部を紹介しながら、議論をより具体的に進めていく。なお実験の技術的詳細は [4,5] を参照願いたい。

## 2 実験例

筆者は図1に示す視覚つき移動ロボットを用いて、ナビゲーション学習実験を行った。ロボットはカメラを上下左右に回転させ視点を変えることができ、また左右のタイヤの速度差を制御することにより、進む方向を変える。基本的にロボットはその視覚注意を、沿って移動すべき自分の左側の壁および右側に逐次現れる色パターンのついたオブジェに交互に振り分けながら、ワークスペース内を移動していくして行く。ロボットは繰り返し同じワークスペースを移動していくうちに、経験したランドマーク(壁のコーナーおよび色模様のついたオブジェ)のシーケンスを、エピソード記憶として、逐次長期記憶に学習していく。そして、学習が進むにつれて、次にどのようなランドマークに出くわすかを、徐々に予測できるようになる。そのトップダウン的先行予測は、実際のランドマーク認識において、視覚アテンションのランドマークへの誘導、下位からのボットムアップ的パーセプションの補強などに利用される。

### 2.1 認知モデル

ランドマークの予測・学習・認識のために構成した神経回路網モデルを図 2に示す。この神経回路網モデルは、低次センサー知覚パートと、高次予測学習パートに分かれる。低次センサー知覚パートは、TE 野において、なにを見たかを画像入力から Hopfield net での associative memory で認識し、また PE 野において、その見たランドマークの前のランドマークからの相対的位置(ロボットの相対的な移動ベクター)をモータエンコーダ積分値のクラスタリングから認識する。この両者の"What and where"の情報は統合圧縮された形で、高次予測学習パートに送られる。高次予測学習パートは recurrent neural net(RNN) [6] により構成される。

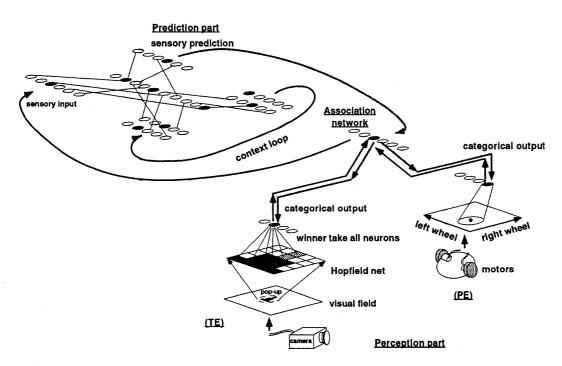


Figure 2: Proposed architecture consisting of multiple neural networks.

RNN は前に入力されたランドマーク刺激から次のランドマーク刺激を予測学習する。つまり、次にどれくらい移動したところで、どのような色形のランドマークを見るかを予測するのである。この予測は TE 野および PE 野にトップダウンの流れとして再投入され、次のランドマーク認識を補強する。RNN の学習は、ワークスペースを数回か回ったあと、経験したランドマークシーケンスを、前学習した記憶の上に追加的に学習させる。この追加学習において、過去の記憶が新しい学習により過度に歪められないように、人工的な Consolidation Learning を導入した (詳しくは [4] を参照)。TE 野および PE 野ではセンサー入力がなされると逐次回路網の結合重みが学習更新される。

さて、高次予測からの低次知覚へのトップダウン的再投入については、さらなる説明が必要である。高次からの予測は、ロボットが環境をよく学習したあとであれば、その予測は低次知覚を補強し認識時間(Hopfield net の収束時間)を早めるが、一方その学習が不完全の場合は、間違った予測が正しい知覚を阻害する可能性がある。そこで、現在の予測的中度を観測し、その値を用いてトップダウン再投入の強さを適応制御することを考えた。つまり、予測率が悪いときは、トップダウン再投入を弱め、Hopfield net の収束時間を長めにとる。一方、予測率が良いときは、トップダウン再投入を強め、Hopfield net の収束時間を短めにするという方策をとる。

#### 2.2 実験結果

実験において、ランドマークを5つ含むワークスペースを用意し、そこでロボットを約20回巡回移動する間追加学習するようすを繰り返し観測した。図3にその時の神経回路網の学習にともなうの時間発展の一例を示す。図3(a)は正規化した予測エラーを、(b)はRNNの

#### 研究会報告

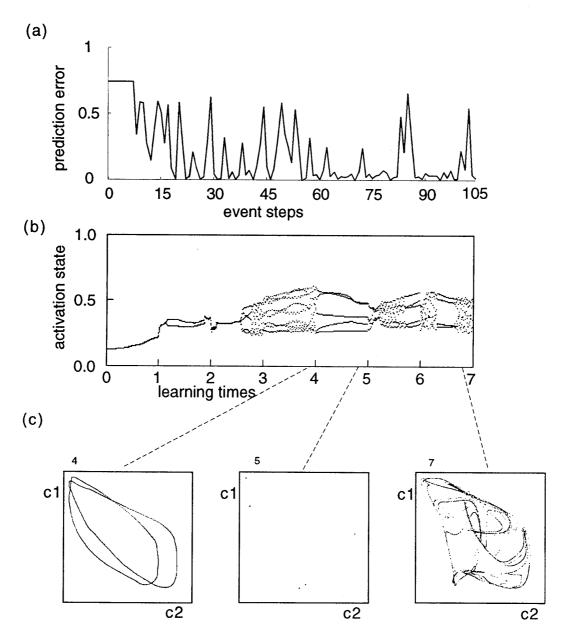


Figure 3: (a) the prediction error, (b) the bifurcation diagram of the RNN dynamics and (c) the phase plots at particular times. The times are indicated by the dashed lines.

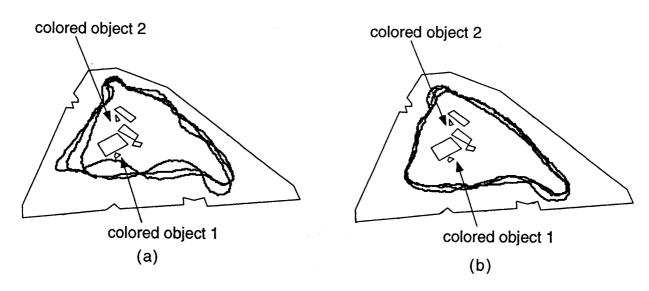


Figure 4: The robot trajectories measured (a) in the unsteady phase (60th step to 74th step) and (b) in the steady phase (75th step to 89th step).

学習に伴う内部状態の分岐ダイアグラムを、(c) は特定時点での RNN の内部状態の位相プロットを示す。予測エラー値は、初期の学習の後下がるものの、その後随時バースト的に上昇を繰り返す。分岐図の方では、それに伴い、状態の分岐が繰り返される。随所に周期5(プロットが5本並んでいる状態)が見られ、これは RNN が5つのランドマークのシーケンスを正しく学習し、周期5の周期解アトラクターにそれを埋め込んでいると考えられる。ところが、周期5は安定的に継続せず、それは時間とともに崩壊し、無周期アトラクター(分岐図での砂嵐状態の部分)が発生する。これらは(c)の位相プロットにしめすような準周期解であったりカオスであったりする。

さて、RNN の内部力学構造は、なぜこのように安定相と不安定相の間を遷移しつづける のだろうか。それを考えるため、安定相と不安定相での、ロボットの実際の移動軌跡を、比 較してみた。図 4にそれを示す。(a) に示す不安定相での移動軌跡は、(b) のそれに比べて大 きく揺らいでいることがわかる。移動軌跡が揺らぐ原因は、視覚注意のメカニズムによると ころが大きい。前に説明したように、予測エラーが大きくなると、注意はランドマーク認識 に長い時間が費やされ、沿うべき壁への注意が損なわれる。これによりロボットと壁との保 つべき距離は不安定化し、軌跡が揺らぐのである。軌跡の揺らぎは、さらなる影響を認知過 程におよぼす。まず軌跡の揺らぎはコーナーなどのランドマークの誤認識を引き起こし、1 周で5つシーケンスであるはずのランドマーク列がロボットの経験上では、4つになったり、 6つになったりする。また周期のずれは、トップダウン予測とボットムアップ知覚のあいだ にづれを生じさせ、しいては高次レベルの RNN の指し示すものと低次の知覚パターンの関 係が変わったりする。つまり、行動の揺らぎがランドマークといった認識の基底をゆさぶり、 世界の観測における文節をかえるのである。世界の見えかたの揺らぎは、そのまま学習に反 映され、RNN は観測された非決定的なランドマークのシーケンスを、世界の像として恣意 的に学習[7]する。RNNは世界を平均化すれば周期5と汎化理解したり、はたまた非決定的 だと理解しカオスを内部生成したりする [8]。追加学習が行われることにより、再度予測率 が変わり、それはロボットの動きに影響するのである。ここに見られるものは、トップダウ

研究会報告

ン、ボトムアップの相互作用を契機として発生する、身体運動、注意、認識、学習の各プロセスの相互干渉、そしてそれらを要素とする全体の動的構造変化である。全体の力学系構造は安定相でのコヒーレントな状況と、不安定相でのインコヒーレントな状況とを繰りかえしながら、常に変動していくのである。

## 3 自己意識の解釈

筆者は、実験で得られた自発的な安定相と不安定相の間で繰り返す転移に、自己意識と解釈できる現象が潜むと考える。安定相においては、客体と主体は構造的カップリング [9] の上にコヒーレントなダイナミクスを形成し、全てのプロセスがオートマテイックにスムーズにすすみ、そこでは自己を意識する機会は失われる。

一方、不安定相においては、客体と主体の構造的カップリングの上でのコヒーレンスは失 われ、そこで立ち会う予測の裏切り、矛盾らは、「私は何か?」という問い、つまり自己意識 を現前させると解釈できる。このような解釈はハイデッガーの自己に関するそれと共通する。 ハイデッガーは、職人と道具の関係の中で、作業が順調に進む中では、職人と道具はひとつ のプロセスの中に一体化するするという。しかし、一旦両者の関係に不都合が現れ、作業が つまづいたとき、自分を問う自己意識が現れ、職人と道具の独立した存在が浮かび上がる。 また、ハイデッガーは、人が日常生活の中で、日々をただの繰り返しとして無難に送ってい る状態を、タイ落という。このとき、人は固有の生き方を意識せず、一般世間におぼれ自己 を喪失しているという。しかし、人は自分の避けられない死という各々の問題に直面したと き、はじめてそれぞれの生きざまは固有であることに気がつき、自己を意識するという。よ く、末期のガン患者がホスピスなどで、残された少ない日々の中、草木の揺らぎ、鳥のさえ ずりなどに、一刻ごとの新鮮なものと感じるといった口述は、これに対応するであろう。筆 者のロボットは、もちろん己の死といった概念は持ちあわせない。しかし、ロボットは「死」 といった契機をもたづとも、世界との身体的経験は有限であることを、いやおうなしに実体 験する。この身体をもつものを必然的に拘束する経験の有限性において、ロボットの世界に ついての内部理解は汎化されず、ただ恣意的にならざるを得ないのである。実験で見られた 周期解5は安定化せず、学習の恣意性から多様な非周期解があらわれ、主体と客体の関係は 時々刻々と変化しつづけた。ロボットの経験が環境との関係において刻々と変化する固有の トラジェクトリーをすすむとき、そこにわれわれはハイデッガーがいうところの自己の本来 的存在を見いだすことが可能ではないだろうか?

## 4 構成論的方法論

意識の問題の難しさは、その定義が本来不明確であることにある。構成論者の陥りやすい間違えは、意識をあたかも実存的実態として陽に記述してしまうことにある。その「意識」は構成者の意図するがごとく作動することが実験で確かめられたとしても、その「意識」は何をもって正当化されるのだろうか?似た問題は、人工知能における記号の取り扱い、最近の人工的感情モデルの研究にも見られる。人々は、これらの起源を問うこと無しに、恣意的にその実態を与え、その上での可能なしくみ・作動についてのみ議論する。

禅において意識は「空」であると解釈され、バレーラ [9] はそれを実在に決して落とせないもの (sense of groundless) という。意識しかり記号しかり、これらは単独に実在するものではなく、ある関係性において成立するものではあるまいか。われわれが考える構成論的ア

プローチの第一歩は、まず認知における各プロセスの正しい関係性の構築にある。これは、筆者の示した例では学習、認識、行動、注意などの互いに相互干渉しあう部分からなるシステム構成することに他ならない。次に構成したシステムについて実験を行い、どのような現象が発生するか調べる。重要なのはこの実験過程は、発見の過程であるということである。発生する現象が、単純に構成者の意図を正当化するだけのものであれば、その実験は意味を持たない。最後に、発見された現象のパターン・構造を、解釈する過程に入る。この解釈の過程なくしては、実験は決められた境界条件から時間発展を数理物理として示したことにすぎない。この解釈の過程は、選られた数理の構造を、過去の認知に関するさまざまな生理、心理、哲学的議論を踏まえて、言葉の構造に置き換える作業であり、これをもって機械は本来の無機質な物理の枠組みを超えて、認知の問題の対象になりうるのである(津田[10]らは同様の議論をハイデッガーの解釈論を発展させて行っている)。こうして編みこまれた言葉のつながりを持って、初めて、本来実在しないが誰もが経験する意識・記号について言及可能になると筆者は考える。

#### References

- [1] S. Harnad: The symbol grounding problem. Physica D, 42:335-346, 1990.
- [2] R. Brooks: Intelligence without representation. Artificial Intelligence, 47:139-159, 1991.
- [3] R.D. Beer: A dynamical systems perspective on agent-environment interaction. Artificial Intelligence, 72-1:173-215, 1995.
- [4] J. Tani, J. Yamamoto, J. and H. Nishi, H: Dynamical interactions between learning, visual attention, and behavior. In P. Husbands and I. Harvey, editors, Proc. of the Fourth European Conf. of Artificial Life, Cambridge, MA: MIT press., 1997.
- [5] J. Tani: An Interpretation of the "Self" from the Dynamical Systems Perspective: A Constructivist Approach. Sony CSL TR-98-18 (submitted), 1998.
- [6] M.I. Jordan and D.E. Rumelhart: Forward models: supervised learning with a distal teacher. Cognitive Science, 16:307-354.
- [7] T. Ikegami and M. Taiji, M: Structure of possible worlds in a game of players with internal models. In Proc. of the Third Int'l Conf. on Emergence, 601-604, Helsinki, 1998.
- [8] J. Tani and N. Fukumura: Embedding a grammatical description in deterministic chaos: an experiment in recurrent neural learning. Biological Cybernetics, 72:365-370, 1995.
- [9] F.J. Varela, E. Thompson and E. Rosch: The embodied mind. Cambridge, Mass: MIT Press., 1991.
- [10] I. Tsuda and E. Koerner and H. Shimizu: Memory dynamics in asynchronous neural networks. Prog. Theor. Phys. 78:51-71, 1987.