

脳型ロボット研究に基づく意識及び自由意志の統合的な理解

谷淳 沖縄科学技術大学院大学

I. 緒言

私が研究している問題は、人間の心はいかにして、感覚運動の経験の蓄積を通じて「世界内存在¹⁾」の理解を築くのか、ということです。この問題について、ロボット実験をとおした構成論的理解というアプローチを用いて、25年ほど研究をしてきました。最近、その内容を、オックスフォードから出版しました私の本[1]、**Exploring Robotic Minds** の中にまとめましたので、もしご関心ありましたら、一読を頂ければ幸いです。さて、「世界内存在」については、みなさんはすでによくご存じだと思いますので、説明の必要はないでしょう。私は、対象としての物体、他者、そして自己が世界内存在することを理解する神経認知のメカニズムに興味を持っています。このメカニズムの基本的原理理解することは、意識と自由意志の基本的原理を理解することと同時に成されると考えています。

私は自由意志の起源が決定論的なのか、それとも確率論的なのかということも問題にします。そして、リベットの問い[2]にも関心があります。つまり、なぜ意識は遅延するのかということです。これは昨日議論されていた 新しい可能性や非決定性にも関係します。すべてに答えることはできませんが、以上の問題を論じてみたいと思います。

私がやっているのは学際的アプローチと呼ばれるものです。認知神経科学で調べられているような、皮質野の様々な領域の相互作用から生じるマクロレベルの脳のメカニズムにも関心を持っています。また遺伝と環境 (**nature and nurture**) を問題とする発達心理学、それに関連して、自閉症スペクトラムや統合失調症などの発達障害のメカニズムにも関心があります。さらに複雑系、すなわち大規模かつ非線形な動的システムにおける創発の科学は、非常に重要だと私は考えています。それから私は素人ながら、現象学とりわけハイデガー、フッサール、メルロ＝ポンティなどの思想に興味を持っています。ベルクソンに関しては、本文を書き終えるころ、編者の平井先生と詳しく話をする機会があり、私どもの研究と深い関連性があることを理解したしで、その点について本文の最後に触れたいと思います。

私が試みているのは、計算シミュレーションされた脳・身体が環境とどのように相互作用するかを、ロボットプラットフォームを利用して吟味するという、構成論的脳型ロボット研究と、私が名付けたものです。なぜこのような研究をしているのかというと、これはある意味で私自身の「拡張された現象学」なのだと思最近気が付きました。多分、脳の中での非線形現象として立ち上がる心の作動は、我々が現象学的に頭の中で考えられることよりも、数段複雑なものだと思います。だから、そのような複雑な現象を理解するには、計算機、ロボットなどを利用して、現象を再構成しながら理解していくことが必要だと思います。またフ

¹⁾ [編者注] (簡単な説明)

ランシスコ・ヴァレラ[3]は、90年代の後半から神経現象学という立場を表明していて、私の研究はその影響を受けています。

II. 主観と対象世界のトップダウンとボトムアップを通しての相互作用

これから私がお話するのは、人工的な身体、つまりロボットと、それと相互作用する神経回路システムについてで、私はこれから、抽象的な心の構成論的モデルを構築しようと試みています。図1に示されている、抽象的な神経回路モデルは、階層的な構成を持ち、高次のレベルは遅く、低次のレベルは速く作動するという、時間的階層性の制約を受けています。上位層からは、対象世界に予測的に働きかけようとする、主体のもつ意図のトップダウンの志向性があると考えられます。そして下位層からは、その予測される知覚に対して、対象世界から得られる、実際の知覚との誤差がボトムアップの志向性が立ち上がり、その両者が相互作用するという様相を仮定します。

この図において、上から入っている線が主観的な心、下からの線は客観的世界を表していますが、両者が相互作用する場面において、両者は不可分になるということが、重要です。このような考えはメルロ＝ポンティが示すところの身体性の現象学に準じたものだと思います。そしてこの様相が出発点となります。さて、もし心を構成する部分が、非計量的(non-metric)で形而上学的な何か、あるいは記号であったら、両者の密な相互作用は構成されないと思われま

す。大事なことは、どちらかが片方を押したら、押された方は弾力的に少し引っ込む、でもその引っ込んだ方は、その後また相手を押し返すといった、押し合いへし合いの力の相互作用が、主観的な心と、対象世界との間に構成されることです。なぜならば、両者は多くの場合に折り合いが悪く、矛盾は力の相互作用を経て、どこかに動的に妥協点を見つけ出して解決しないといけないからです。そのためには、主観的な心は、物理的な環境と同様に、距離が定義された計量空間に構成されないといけません。そのような共有できる計量空間を持つことにより、両者は密に相互作用できるようになります。そのために、主観的な心は、その作動が計量空間に定義できるように、神経力学系モデルで構成することが必要となります。これが私の議論の一つの要点です。

さて、ここからは私どものグループが提案する、予測符号化原理に基づく、神経力学モデルについて詳しく述べていきます。予測符号化原理がいわゆる再帰型ニューラルネット

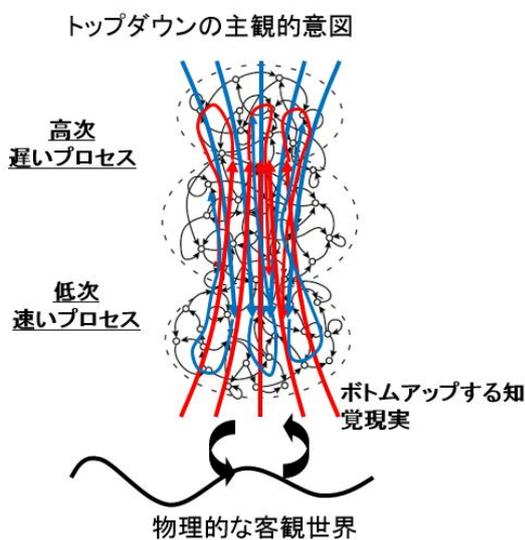


図1主観的意図によるトップダウン志向性と対象世界からのボトムアップ志向性の相互作用のようす。

ワーク(Recurrent Neural Network: RNN)に実装されることで、潜在変数であるコンテクスチュアルなダイナミクスを持ち始めるようになり、隠された次元が現れます。そして隠された次元は学習を介して自己組織化が可能になります。この点が非常に重要です。これなしでは、認知的な心について議論を進めていくことができません。

こうした再帰型ニューラルネットワークは注目を集め始めており、特に最近はやりのディープラーニングでの主役になりつつあります。これはエルマンが30年前に考案したのですが、私自身も25年ほど使っており、特に多重時間スケールのような階層を導入した型のものについて研究をしてきました。私たちはこのようなモデルを、多重時間スケールの再帰型ニューラルネットワーク(Multiple Timescale Recurrent Neural Network: MTRNN)[4]と呼んでいます。

図2の上の部分が高次のレベルで、低速で作動しています。

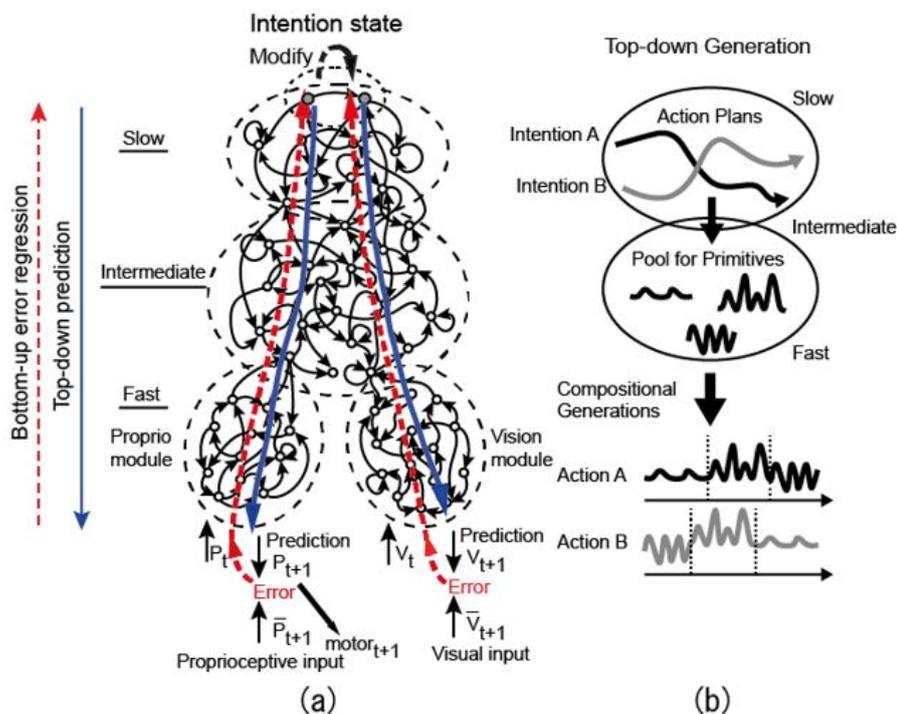


図2. (a) MTRNNモデル、(b)運動基本パターンの順序的組み合わせのようす。

中程が中レベルで、下の方が低次のレベルです。階層はもっと増やすこともできます。時間スケールが大きくなると作動は遅くなり、時間スケールが小さくなるとシナプス入力のインパルスにより早く反応するようになります。高次の階層から、意図がトップダウン的に低次の階層に伝わり、その結果、その意図に基づいた、知覚の予測が行われます。本図では、知覚には視覚と固有受容感覚の2チャンネルがあります。ちなみに固有受容感覚を予測するということは、自己の体の動きを予測することになります。

学習プロセスについては説明する時間はありませんが、高次のレベルが何らかの意図を持った場合、トップダウンの意図のダイナミクスが始まります。これは予測の動きでもある、低速のダイナミクスです。意図ダイナミクスがある初期状態から出発する場合、低速で変化する神経活動として、高次レベルを回帰的な結合を伝わりながら、下層の低次レベルに伝わっていきます。低次レベルには、学習で獲得されている、再利用可能なルーチンとしての基本的パターンが分散的に埋め込まれており、そこに低速の意図のダイナミクスが、これらの基本パターンを、意図した順番で励起していくこととなります。これを物理の言葉で説明すると、高次レベルのダイナミクスが、低速で変化する分岐パラメータとなり、低次レベルの速いダイナミクスに、パラメータ分岐を順次引き起こしていくこととなります。始まりのトップダウンの意図の初期値が異なれば、異なった低速のダイナミクスが高次に現れ、それは異なった基本パターンの順序組み合わせを合成します。

私が論じたいのは、このような認知過程にける組み合わせ合成可能性は、あからさまな記号的操作のプログラムを用いなくても、神経回路のもつアナログのダイナミクスで再構成可能だということです。私たちは現象学的に、あたかも記号があつてそれを操作する過程を、認知的な過程であると、考えがちですが、実際には上述のような非線形ダイナミクスの自己組織化的な過程を通して実現可能であり、そのことは昔の現象学者が頭の中でいくら考えても想像できるものではないかもしれません。このようなことは、計算機とロボットを組み合わせた構成論的な実験を行うことで、はじめて発見できることだと思います。ただし、ウィリアム・ジェームスなどの書いたものには、そのようなダイナミクスに直感的に気が付いていたような節が、読み取れます。

上述のトップダウンの予測のダイナミクスは、低次レベルに降りてきて、現実の知覚と向き合うことになり、そこでは大なり小なりの予測エラーが発生します。その予測エラーのシグナルは高次レベルに逆伝搬してゆき、その誤差が最小になる方向に、高次の意図の状態を変化させます。その誤差が最小になった時に、世界が再認され、それに対応するように意図が変更されます。例えば、白い大福もちが机の上にあると思って、掴んだら、温かい毛触りの触感を感じて、思わず手を離れたのは白いネズミだったというように。つまり、大福もちに手を伸ばすという意図から、ネズミを避けるという意図に、現在の知覚の整合性の合う方に意図が変化し、それに伴い行為の組み合わせ方も動的に変化するわけです。ここでは、将来に対する予測(prediction)と過去に対するポストディクション(postdiction)が同時進行する形になります。そして「今」は、その間に存在することになります。

以上の説明は簡略化したものですが、すべてはダイナミックな現象であり、それは予測エラーを最小化するというひとつの原理で作動しています。学習も、同様に、予測誤差最小の方向に、神経結合重みを変化させていきます。そこに複数の意図に関するダイナミックな構造が自己組織化されるわけです。またこのネットワークモデルは自身が出力する知覚予測を知覚入力のチャンネルに再入力することにより、学習した意図に関する長期の時間ステップに渡る、心的なシミュレーションが可能になります。

III. 構成論的ロボット実験

脳型ロボット実験の例をいくつかお見せします。図3は脳型ロボットの実験のようすを示しております、10年ほど前に、上述のモデルの原型となるモデルを使って、ソニーと共同研究したものです [5]。私たちはこのロボットに二種類のボール遊びを教えました。ひとつはボールを左右に転がす遊び、もうひとつはボールを持ちあげて下に落とす遊びです。ロボットの腕をとってこのボール遊びの

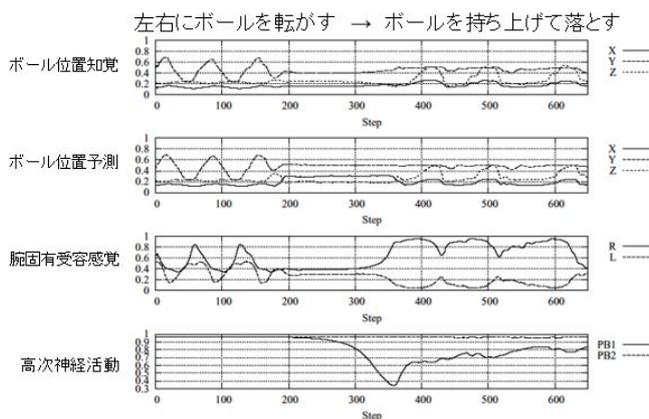
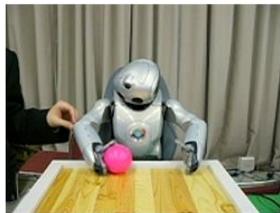


図3. (a)ソニー製ヒューマノイドロボットを用いたボール遊び学習実験のようす、(b)遊びパターンの自然発生的シフトのようす。

やり方を教えます。繰り返し教示したうえで、ロボットを自律的行動生成のテストにかけます。(以下のビデオを参照 https://www.youtube.com/watch?v=a_auIoksGN0)

まず注目すべき点は、ロボットは自分の腕の運動の予測に合わせて、常にボールの動きを予測しているということです。予測がうまくいけば、プロセスは自動的に行われます。しかし時としてボールが予測よりも大きく跳ね返り、予測エラーが生じることがあります。そのような時、このエラーは上層に向かい、高次の神経活動の状態を別のものに变化させ、その結果遊びのパターンが移り変わります。長い時間観察していると、同じボール遊びパターンが5分以上続くこともあれば、数十秒でもう一つの遊びパターンに切り替えることもあります。いったい誰がロボットの遊びのパターンを切り替えて変えているのでしょうか？ロボット自身が変えているのか、それともこれはただのノイズなのでしょう。このような議論は自由意志の起源にも迫る問題だと思われま

ここで起きているのは、連続的な知覚の流れを意識的に操作可能な対象へと分節化する過程です。もともと、知覚の流れは連続ですが、認識されるものは、遊びパターン A から遊びパターン B に移ったというような分節化された記述です。その流れの分節化は、生成される予測誤差が高次の意図を書き換えることによって起きる、低次の運動ダイナミクスの相転移により可能となります。また大事なことは、その予測誤差が上がり、それを最小化

しようとする努力が、意識を生むということです。この説明は、ハイデガーの著作で言われている大工と釘の話に触れることにより、理解が進むと思われます。大工がハンマーで釘をリズムカルにうまくいっている時は、大工は自分のことも、自分が使っている道具のこともまったく考えません。しかしたまに釘を打ちそこなうことがあると、予測エラーを伴い、「私に何が起きたのだろうか？釘はどうした？」と考え、そこに自分、そして釘といったものが分離した対象として浮かび上がるわけです。

流れ自身は操作ができませんが、それが意識的な過程を通して分節化されるようになった時に、初めて意図によって操作可能な対象となります。この意図による操作は階層方向に進むので、これはフッサールの言え、縦の志向性にあたると思います。そこには、さらに時間方向の横の志向性が絡んできます。意図の書き換えは、常に後付け的になされるために、過去への志向性を持つこととなります。それが前述したポストディクションです。一方で、その過去の後付け的な認識を基盤として、将来の予測、志向性が成り立ちます。縦と横の志向性が交互に織り込まれていく過程を通して、自己の行為は認識され、また生成されていくのだと考えられます。

次の問題は、脳そして機械はどのように確率というものをとらえるのか考えていきたいと思えます。これはこれまでの例で注目してきたのは、決定論的構造の学習だったのに対して、ここからは確率論的構造のそれについて論じるということです。実はこの問題は、自由意志の基は確率的なものか、決定論的なものなのかという問題をも射程に入れることとなります。まずここで紹介するのは決定論的で作動する前述の MTRNN を用いて、教示される行為の組み合わせについての確率的な構造を学習することに関する実験です。ロボットに物体が中央にある場合、50%の確率で物体を左に移動、50%の確率で右に移動させるような運動パターンをロボットの腕を掴んで直接教示します。さらに、物体を右に動かした時は、次の運動の確率は 50%で左、50%で中央、ということになります。このような、確率的な運動の切り替えの繰り返しに基づく、視覚と固有受容感覚刺激の長い時系列パターンを予測学習できるように教示していきます。

学習が終了した後に、ロボットが行為を生成する様子を以下のビデオで見ることができます。<https://www.youtube.com/watch?v=y-XAXaYaSGM> ロボットは確率論的に物体の動かし方を切り替えている様子がうかがえます。詳しく観測してみると、その切り替えの確率は、教えた 50%切り替えにかなり近いものになっていることが判明しました。いったいどのような内的メカニズムが発達したのでしょうか。私たちが発見したのは図 4 に示すようなものです。一番下が低次のレベルで、横方向が時間を表しています。運動パターンが、右 (R)、左 (L)、中央 (C) に物体を移動と切り替えられていくようすが分かります。真ん中の部分が中間層の神経活動で、100 個の活動状況を示しています。この階層は、基本運動パターンをエンコードしていることが、解析からわかりました。しかし、その上に表示されている高次のレベルでの神経活動からはそのような対応はわかりません。高次のレベルの

活動は一見するとランダムなように思えますが、実はカオスダイナミクスの生成がここには確認されています。ごくわずかな初期状態の違いが、物体の動かし方の違った時系列組み合わせを多様に生成するという

ようになっていきます。つまり MTRNN は観測される確率的な構造をカオスに埋め込む形で学習したといえます。脳科学的に考えると、カオスの力学構造をもつ神経活動が高次のレベル——おそらく前頭葉前部——に生成されて、それが頭頂葉に伝わり、それが一方では一次視覚野で視覚的予測を生成させ、もう

一方では自己受容性感覚に信号を送って姿勢予測を生み出しているということです。この自然発生的にあたかもランダムに行為が選ばれてくるところに、自由意志のメカニズムを考えることが可能かもしれません。

しかし私たちはリベットの実験[2]を考慮に入れる必要があります。非常に有名な実験なので説明はしませんが、この実験の結果で面白いのは、まず第一に、自由な行動の選択は脳の高次レベルから発生していることが分かったことです。それ以上に面白いのは、第二の発見で、その行為の選択が意識されるのは、その高次レベルでの神経活動が開始されたかなり後になるということです。つまり自由意志は無意識のうちに生まれ、行動を開始する直前になってはじめて、それが意識されるという構図になります。我々の自由な行為の選択は意識的に決定されないということになります。これをどう説明すればいいのでしょうか。私はこの意識されることの遅延を、身体性に基づくポストディクションで説明できるのではと考えています（図5参照）。まず

高次レベル脳、例えば前頭葉における神経活動がカオスによって自発的に摂動し、ある意図を生み出されると考えます。その過程は無意識的です。その突然生み出された意図が、低次レベルに伝搬し、ある運動パターンを励起しようとすると考えます。しかしながら、実世界、そして身体とより密接につながっている、下位層での神経活

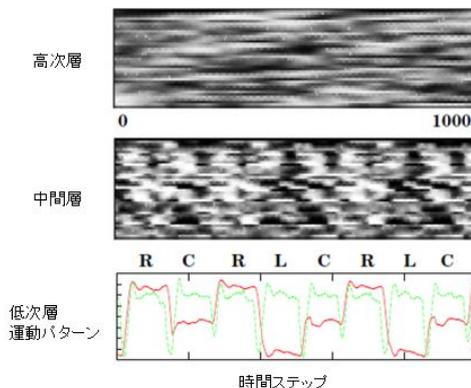


図4. 確率的な運動パターン切り替えに伴う、各層での神経活動。

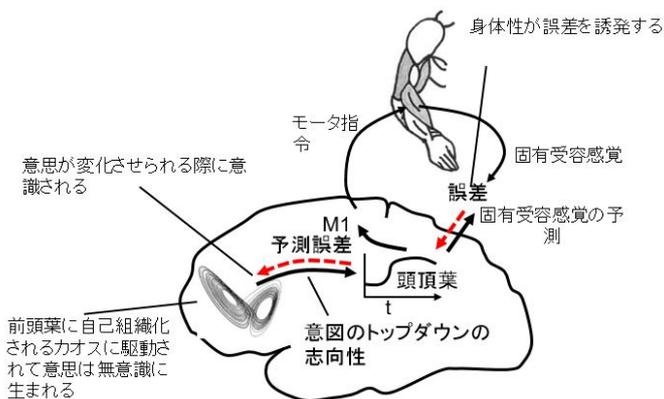


図5. 自由意志の無意識での自然発生的な生成と、それが後付け的に意識される過程。

に、新たな運動パターンが生成される契機となることがあります（図7に概念図を示します）。神経回路に複数の記憶を植え付けると、強い非線形効果で、えくぼのように、偶発的に為記憶が生成されることはよく知られています。多くの場合、この新しく創り出された行動パターンは無意味で、特別興味深いものではありません。しかし10回も20回

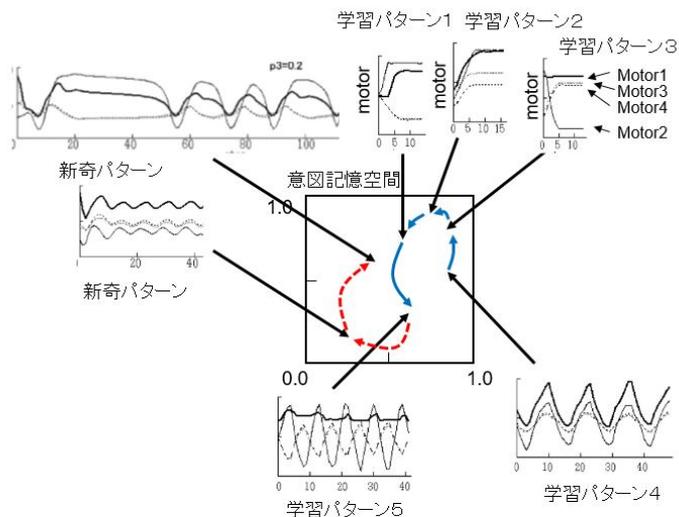


図7. 記憶空間の探索概念図。探索の過程に学習されたパターン以外に、為記憶としての新奇パターンの発見のようす。

も繰り返していると、たまに面白いものが出てきます。そういうものが見つかった時、教える側がその行動をロボットの記憶に強く学習させることにより、新たな行動パターンのレパートリーを増やすことができます。

この **Human in the robot loop** の実験は、新たな現象学的な体験の場を、実験参加者に提供します。私自身の体験では、ロボットと私の相互作用のありようは動的に変化します。ある時は、両者はどちらかが他方に従うというように、矛盾無く一体となって、同じような運動パターンを繰り返します。その時は、自分自身は川の流れに任せるように、ゆったりとした気持ちになります。でもそういった状態は長くは続きません。どこからか、それは崩壊していきます。ロボット側に小さな予測誤差が生じて、それが徐々に拡大して崩壊していく場合もあるし、私自身が同じことを繰り返すのに飽きてきて、ちょっと違ったことを試みようとする気持ちが芽生えることによることもあります。その両者にギャップが発生したときに、自分の手を自分が意図していない方向に押すロボットの手の力に、そのロボットの自由な意思を強く感じとる場合があります。

さらに興味深い点は、自分とロボットが無意識の中に一体化したと考えられる相と、それが崩壊して、矛盾の中にロボットと自分が独立した主体だと強く意識される相が、間欠的に繰り返すということです。前出のウィリアム・ジェイムズは、意識の流れにおいては、鳥がある時は木に長く留まっていたり、かと思うと、急に飛んで行ったりするように、ある意識的な内容が立ち現れては消えてを繰り返すと言っています。この言葉は、意識の持つ自律性という、その重要な特徴に言及していると思います。私どもの”**Human in the robot loop**”の実験でも、それに類似した、ロボットと人間の協働の中に立ち現れる、意識と無意識の相の間欠的な繰り返しが体験されることが良くあります。なぜそのようなことが起きるのか？それは、循環的な因果性が、ロボットと私の身体と「心」から構成される全体の系に、形成されたことによると考えられます。私の新たな運動パターンは、ロボットのトップダウン

ン予測に誤差を生成させ、それによりロボットの意図を変更され、その結果としてロボットの新たな運動パターンが誘発されます。そのロボットが生成した新たな運動パターンは、今度は私の意図を変更させ、そしてそれに対応する運動を誘発します。そのような複雑な循環性の因果の中で、系は疑似安定な解を見つけてそこにしばらく停滞したり、またそこから離れて行ったりするのではと考えられます。このような循環的な因果はよく見られる現象で、私自身は90年代に行った、ロボットと環境の相互作用的学習[6]においてはじめて観測しました。ここで紹介しました、**Human in the robot loop**の実験は、未だ体系的に研究を進めることができていません。それは、実験系に人間の主観が入るために、客観的な評価が難しいことにもよります。今後このような実験をどう進めたら良いかについてはいろいろと議論があるかと思います。ただ、個々の実験参加者が自己の主観を探索し見つめ直す、新たな機会を与える装置になると考えられ、今後も研究を続けていきたいと考えています。

V. ベルグソンとの接合

緒言でも述べましたように、最近平井先生からベルグソンの思想について、深くご教示いただく機会があり、上述の私どもの研究は、ベルグソンの考えと多くの点で接続可能だと思いましたので、それらについて最後に簡単に述べたいと思います。

平井先生からご教示いただいた範囲で私が感じたことですが、ベルグソンは、過去、記憶、そして特に時間について深い考察をしていると思いました。しかもそれらは、私自身はベルグソンに触れたことは過去にありませんでしたが、私自身の時間のとらえ方にかかなり近いということが分かりました。ベルグソンは人間と特質として、遠い過去の経験が記憶となり、それが遅れをもって現在の行為の生成に強く影響することがあることを述べているようです。このことは、私らが多時間スケールモデルで提案している、高次層の意図に関する遅い神経活動が、低次層での運動パターンの生成に、文脈依存的に強い影響を与えているという考えにかかなり近いと思われます。

またベルグソンの考えている、意識の創発のしかた、つまり意識とは長い時間スケールでの応答の欠損から創発するという考えは、私どもが予測誤差を逆伝搬させて、それが、高次層の遅い時間スケールで作動するところの意図を変化させる刹那において、意識が生じるというモデルでの説明にかかなり近いと考えられます。ベルグソンの言うところの欠損は、私が考える、予測誤差です。予測誤差が、高次層の時間スケールが遅い部分で増幅される場合、それはより強く意識に登ることになります。

最後に、ベルグソンが唱えている、記憶からの創造性の発現は、正にリカレントネットワークに多数の異なり時系列パターンを学習させていくと、それらの般化が起きると同時に、為記憶も多数生成され、それらから創造的なパターンが生まれる可能性があるという私どもの考えに対応すると考えることは可能です。ただ私は、ベルグソンが未来に働きかける意図のようなものを、どのように捉えていたのかが不明であり、その部分を起点とした、私の考えとの相違もあるかもしれないと思っています。以上は私の解釈で、間違っている

ころもあるかもしれません。それにしても、計算機もロボットもない時代に、このような思考をすることができたベルグソンは、素晴らしいと思います。しかしまた同時に、現代に生きる我々は、最先端の知識と科学技術を活用して、ベルグソンを超えるような思考を積み上げていくことが必要だとも思います。私自身、このような気持ちで、研究をさらに続けていきたいと思っています。

参考文献

- [1] Tani, J. (2016). *Exploring Robotic Minds: Actions, Symbols, and Consciousness as Self-Organizing Dynamic Phenomena*. New York: Oxford University Press.
- [2] Libet, B. (1985). Unconscious cerebral initiative and the role of conscious will in voluntary action. *Behavioral and Brain Sciences*, 8, 529-539.
- [3] Varela, F.J. (1996). Neurophenomenology: a methodological remedy to the hard problem. *Journal of Consciousness Studies*, 3, 330-350.
- [4] Yamashita, Y., & Tani, J. (2008). Emergence of functional hierarchy in a multiple timescale neural network model: a humanoid robot experiment. *PLoS Computational Biology*, Vol.4, Issue.11, e1000220.
- [5] Ito, M., Noda, K., Hoshino, Y., & Tani, J. (2006). Dynamic and interactive generation of object handling behaviors by a small humanoid robot using a dynamic neural network model. *Neural Networks*, 19, 323-337.
- [6] Tani, J. (1998). An interpretation of the 'Self' from the dynamical systems perspective: a constructivist approach. *Journal of Consciousness Studies*, 5(5/6), 516-542.