

# 他者の内部状態推定を実現する自己観測原理<sup>\*†</sup> —コミュニケーションの新しい計算理論—

牧野 貴樹

合原 一幸

東京大学大学院 新領域創成科学研究科 複雑理工学専攻

mak@sat.t.u-tokyo.ac.jp

aihara@sat.t.u-tokyo.ac.jp

2003年10月24日

## 概要

人間のような高度なコミュニケーションの基礎と考えられる、自己と同等な他者の内部状態推定に関する計算理論を提案する。他者の内部状態推定においては、推定者のパラメータ次元の制約、主観的情報と客観的情報の相互変換という2つの本質的な困難が存在する。我々が提案する自己観測原理は、自己観測での学習に基づいて、この2つの問題を同時に解決する。すなわち、獲得が容易な自己のダイナミクスを学習することで他者のダイナミクスに関する先行知識を得られるので、パラメータ次元の制約を緩和できる。また、自己の主観的状态と自己を客観的に観測した情報の関連付けを学習することで、他者を客観的に観測した情報と他者にとっての主観的情報の相互変換を可能にする。本論文では、力学系観測問題の枠組でコミュニケーション問題を定式化し、そのうえで2つの困難

とその解決法について論じる。また、進化心理学・神経科学との関連についても考察する。

**Keywords:** communication, computational theory, model estimation, mirror neuron, self-observation

## 1 はじめに

コミュニケーションは、人間にとってはたやすいけれども、コンピュータにとっては未だに困難な、未解決問題のひとつである。人間だけではなく、イヌやサルといった動物も、ごく自然にコミュニケーションをとっている。そして、高度なコミュニケーションの場では、駆け引きや、協調、文化といった新しい動きが創発する。しかし、コンピュータは、人間とのコミュニケーションはもちろん、コンピュータ同士のコミュニケーションでも、プログラムされた範囲のことしかできないのが現状である。

今までは、コンピュータの「人間らしさ」の欠如という、技術的な限界がコミュニケーションの障壁になっていると考えられてきた。しかし、そうした限界は次々と克服されてきている。CGによる表情の合成や音声合成や音声認識といった技術も格段の進歩を遂げた。最近では、人型ロボット [IOIK03, AHP<sup>+</sup>00] など、人間とのコミュニケーションを目指した意欲的な取り組みも現れてきている。しかし、新しい動きの創発が

\* 本論文は Makino, T. and Aihara, K. *Self-observation Principle for Estimating the Other's Internal State*. Mathematical Engineering Technical Reports METR 2003-36, the University of Tokyo, October 2003. [MA03] の非公式和訳版です。

† この研究は、文部科学省 科学研究費補助金 特定領域研究 (2) 脳科学の先端的研究 (No. 15016023)、文部科学省 21 世紀型革新的ライフサイエンス技術開発プロジェクト (動的インタラクションによるコミュニケーション創発機構の構成と解明)、日本学術振興会特別研究員制度の補助を受けています。

起きるような高度なコミュニケーションは、いまだに実現の糸口すら見出せていない。

それでは、人間などの生物のなかで、コミュニケーションを実現している情報処理機構が解明されるのを待つべきだろうか。しかし、手掛かりとなる仮説がない現状ではそれも困難である。たとえ神経細胞の結合がすべて分かったとしても、それがどのようにコミュニケーションを実現しているかは分からない。認知科学や神経科学の実験では、すでにある仮説の証拠を探すことはできても、全く未知の仕組みを解明することはできないのである。

いま必要とされているのは、コミュニケーションの根幹をなす計算理論である。計算理論とは、コミュニケーションの目的と必要条件や、そのための情報処理の入出力を規定する、脳の情報処理機構の研究にとって基礎となる枠組である [Mar82, 川人 96]。そのような計算理論が定義されれば、我々はその要件を満たすアルゴリズムをいくつも提案することができる。そして、こうしたアルゴリズムは、人工知能としてコンピュータに実装するためにも使えるし、脳内処理の仮説として認知科学や神経科学の研究の手掛かりとすることもできるだろう。

これまでのコミュニケーションの計算理論としては、他者の内部状態の推定を通じたインタラクションである、というものがある [川人 01]。しかし、その理論を満たす有効なアルゴリズムはいまだ見つかっていない。これが示唆していることは、この計算理論が不十分である、つまり実際のアルゴリズムを実現するために鍵となる計算要素が足りない、という可能性である。もしそうならば、その空白を埋めた、もう1段詳細な計算理論がなければ、研究を進めることができない。

本稿では、コミュニケーションについての従来の計算理論をダイナミクス推定の枠組で検討し、実現を困難にしている2つの要素として、他者の内部ダイナミクスモデルの獲得におけるパ

ラメータ次元の制約と、主観的情報と客観的情報の変換があることを指摘する。そして、その困難を解決する方法として、自己を客観的に観察し学習したダイナミクスモデルを他者の内部状態の推定に利用するという、新しい計算原理を提案する。また、提案した方式を心理学・神経科学的視点から考察し、自己の行動を観察し学習するモデルは人間の自己意識やミラーニューロンの役割と対応することを指摘する。

## 2 ダイナミクス推定の枠組からみるコミュニケーション

本節では、コミュニケーションの計算理論を示し、それをダイナミクス推定問題の形で定式化する。

### 2.1 従来のコミュニケーションの計算理論

コミュニケーションという言葉はさまざまな意味で使われているので、混乱を招くおそれがある。ここでは、本論文におけるコミュニケーションの意味を定義し、そこで必要となる他者の内部状態の推定を計算理論の形で定義する。

本論文が対象とするコミュニケーションは、生物が自然淘汰を勝ち残る手段として発達させた、情報処理に基づく同種の仲間とのインタラクションである。ここでは、単純な観察、模倣といった行為などで、生存率・繁殖率の向上などの自己の利益に反映されないものは、コミュニケーションからは外す。目的が不明だと、構成するアルゴリズムの検証ができなくなるからである。また、アリのフェロモンのように、本能に基づく自動的なものもここでは除く。脳における情報処理のメカニズムに研究の焦点を絞るためである。そして、コミュニケーションの対象は同種の仲間、つまり自分と同等の構成や複雑性を持つ相手と限定する。

ここで定義するようなコミュニケーションでは、他者の内部状態の推定という情報処理が行われていることが、これまでの研究で指摘されている [川人 01]。ここで、内部状態というのは、

外側から直接観察することができないが、相手の行動を決定しているような要素のことで、感覚や、感情、欲求、意図、知識といったものを含む。

生き残るためには、自分だけで頑張るより、仲間に自分を助けてもらうほうがいい。これを促進するには、自分の利益になる行動を仲間がとってくれるよう働きかけることが必要であり、これがコミュニケーションの目的となる。そのためには、仲間がどのような場合にどう行動するかが予測できなければ、うまく働きかけることができない。そのためには、働きかけに先立って、行動の予測がコミュニケーションのための情報処理では必要になる。さらに、観察できることだけから仲間の行動を予測することは困難であるため、行動の予測に先立って、その仲間の内部状態を推定することが必要になってくる。

たとえば、外側から観察できる状態(痩せた男がよろよろ歩いている)では、その相手が将来どのような行動をとるかは予測しがたい。しかし、そこから内部状態(腹が減っていて、食べ物を探している)を推定できれば、その相手の行動を予測する役に立つ(食べ物を見つけたら取って食べるだろう)。そうした予測は、自分の生存に役立つ(残り少ない食料を隠せる)し、相手の内部状態を変更するようなインタラクションにもつなげられる(食べ物をあげれば、空腹が癒されて喜ぶだろう)。

さらに高度になると「『相手が自分の内部状態をどう推定しているか』ということも、相手の内部状態の一部として推定する」再帰的な推定が行われる。たとえば、あるチンパンジーが、いたずらを怒りにきた別のチンパンジーの気をそらすために、まるで迫る危険に気づいたかのように遠くを見るしぐさをしたという、欺き行動[PW78]の事例が知られている。こうした行為の背後にある高度な状態推定は、人間のコミュニケーションにおいても重要な役割を担っていると考えられる。

このように重要であるにもかかわらず、人間の脳のどの部分がどうやって他者の内部状態を推定しているかは分かっていない。この理由として、計算論的研究が進展していないことが考えられる。人間の脳が行うどんな情報処理でも、どのように情報を表現し、どのようなアルゴリズムを処理することで実現できるのかが分かっていないのでは、脳内の対応物を探すことは困難である。計算論的研究では、その情報処理を実現できる表現やアルゴリズムを構成することで、脳内における処理の仮説を提示することができる。

そのときに、それらのアルゴリズムの共通の基盤として、計算の目的や入出力の構造に対する要求を計算理論として定義する必要がある。本研究では、次のようにコミュニケーションの計算理論を定義する。

1. 行為者の生存率・繁殖率が向上するような行動を選択すること(他者の行動が行為者の生存率・繁殖率の向上につながるような行動・インタラクションを行為者が選択することを含む)。
2. 1. を実現するために、自分と同等の存在である他者の行動を推定すること。
3. 2. を実現するために、自分と同等の存在である他者の内部状態や、内部状態と環境・行動との関係を推定すること。

本論文の残りでは、この計算理論を満たすアルゴリズムの原理を探求していく。

## 2.2 ダイナミクス推定問題

直接観測できないパラメータの推定には、さまざまな要素の推定が必要になる。ここでは、ダイナミクス推定問題の一般的な枠組みを示すことで、推定が必要な要素を明確にする。

ある時刻  $t$  におけるある推定対象の状態を、ベクトル  $x(t) \in S$  で表す。その対象のダイナミクス、つまり周囲の環境に応じて状態が刻々と変化するさまは、各時刻の入力  $a(t)$  と状態遷移写

像  $f$  を用いて次のように表現できる。

$$\mathbf{x}(t+1) = f(\mathbf{x}(t), \mathbf{a}(t)) \quad (1)$$

しかし、我々には現在の状態  $\mathbf{x}(t)$  を直接観測することはできない。観測できるのは、次のように  $\mathbf{x}(t)$  から出力写像  $g$  によって決まる出力  $\mathbf{y}(t)$  だけである。

$$\mathbf{y}(t) = g(\mathbf{x}(t)) \quad (2)$$

ここで、入力  $\mathbf{a}(1), \mathbf{a}(2), \dots$  と出力  $\mathbf{y}(1), \mathbf{y}(2), \dots, \mathbf{y}(t)$  だけから、将来の出力  $\mathbf{y}(t+1), \mathbf{y}(t+2), \dots$  を予測することを考える。精度のいい予測を行うためには、予測者は、現在の状態の推定  $\hat{\mathbf{x}} \in \hat{S}$  を求めるだけでなく、対象のダイナミクスモデル、すなわち状態遷移写像の推定  $\hat{f}$  と出力写像の推定  $\hat{g}$  を求める必要がある。これがダイナミクス推定問題である (図 1)。

もし、対象のダイナミクス、すなわち状態遷移写像  $f$  や観測写像  $g$  が既知であれば、問題は現在の状態の推定に集約され、多くの場合は比較的容易に解決できる。 $f$  と  $g$  のどちらか一方しか分からない場合は、もう片方を推定する必要があり、問題としてはより難しい。

しかし、 $f$  も  $g$  も未知であるような場合、問題は飛躍的に難しくなる。元の状態空間  $S$  に対する情報がまったくないため、状態の推定  $\hat{\mathbf{x}}$  が動く空間  $\hat{S}$  を推定者が再構成しなければいけないからである。一般的には、入力が一定の決定論的な系において、 $\hat{S}$  の次元は、元の状態空間  $S$  の次元の 2 倍もしくはアトラクタの容量次元の 2 倍よりも大きくしなければ、再構成が保証されないことが知られている (埋め込み定理 [合原 00])。状態空間  $\hat{S}$  が高次元になればなるほど、推定写像  $\hat{f}$ 、 $\hat{g}$  を構成するパラメータ数も増大するため、推定に必要なデータ量は飛躍的に増大し、推定はより困難になる。逆に、 $f$  や  $g$  に関して、少しでも情報があれば、 $\hat{S}$  の次元を落とすことで、必要となるデータ量を減らす (推定を容易にする) ために使える可能性がある。

## 2.3 ダイナミクス推定問題としてのコミュニケーション

コミュニケーションにおける他者の内部状態推定も、このダイナミクス推定問題を非定常系に拡張したものの一種であると考えることができる。コミュニケーションする個体は互いに等価であるが、ここではその一方に焦点をあて、ダイナミクス推定者としての振る舞いを考える。この論文では、焦点をあてる側を自己、自己が推定する対象を他者と呼ぶことにする。

図 2 にその枠組を示す。ここで、自己は、他者の内部ダイナミクスの推定者であると同時に、他者から観測されるダイナミカルシステムでもある。これは、他者のダイナミクスモデル  $\hat{f}_2$ 、 $\hat{g}_2$  や他者の内部状態の推定  $\hat{\mathbf{x}}_2$  は自己の内部状態  $\mathbf{x}_1$  の一部であり、状態遷移写像  $f_1$  の中にあるモデル構成器  $M$  の働きによってモデル構成が行われる、という形で表現することができる。他者も自己に関して同様にモデル構成を行っているので、全体の枠組みは対称的になる。また、再帰的推定が行われる、すなわち他者の内部状態  $\mathbf{x}_2$  に含まれている、他者が構成している自己のダイナミクスモデル ( $\hat{f}_1$  など) も推定に含まれる場合、フラクタル構造が現れることになる。

ここで、他者に対する入力  $\mathbf{a}_2$  や、出力  $\mathbf{y}_2$  について、直接自己が知ることはできないことに注意する。相手が何を受け取り、何を出力しているかは、あくまで自己の入力  $\mathbf{a}_1$  を通して、間接的にしか知ることができない。視点の問題は 4 節で詳細に取り扱うが、ここでは暫定的に、プライムをつけた記号  $\mathbf{a}'_2$ 、 $\mathbf{y}'_2$  によって、間接的、部分的情報であることを表現している。

## 3 パラメータ次元に関する困難

本節では、他者の内部状態推定における第一の困難、パラメータ次元に関する困難について述べる。

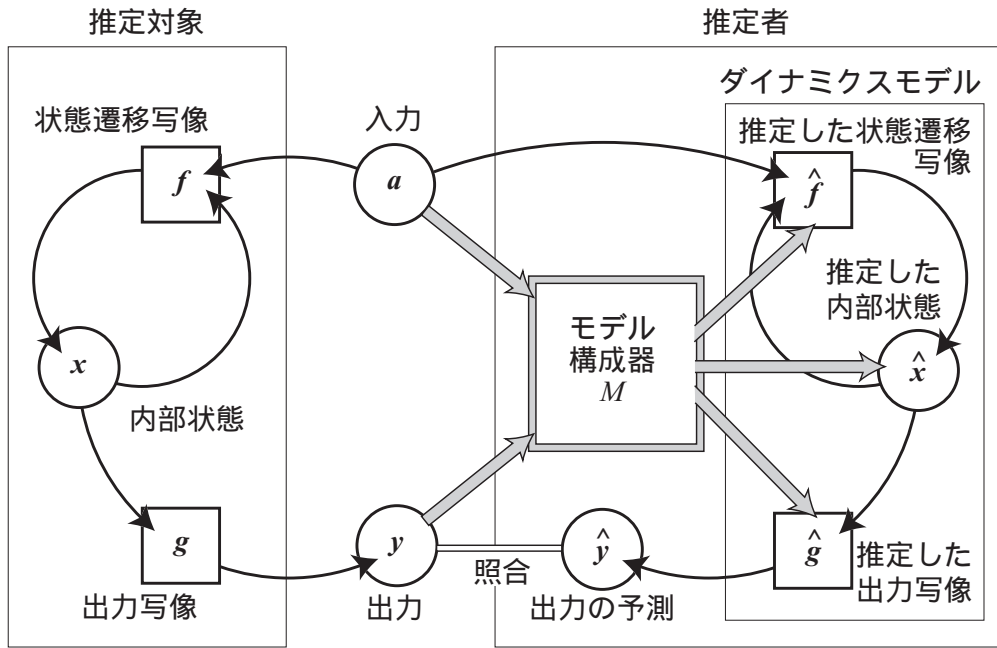


図1 ダイナミクス推定問題の枠組み

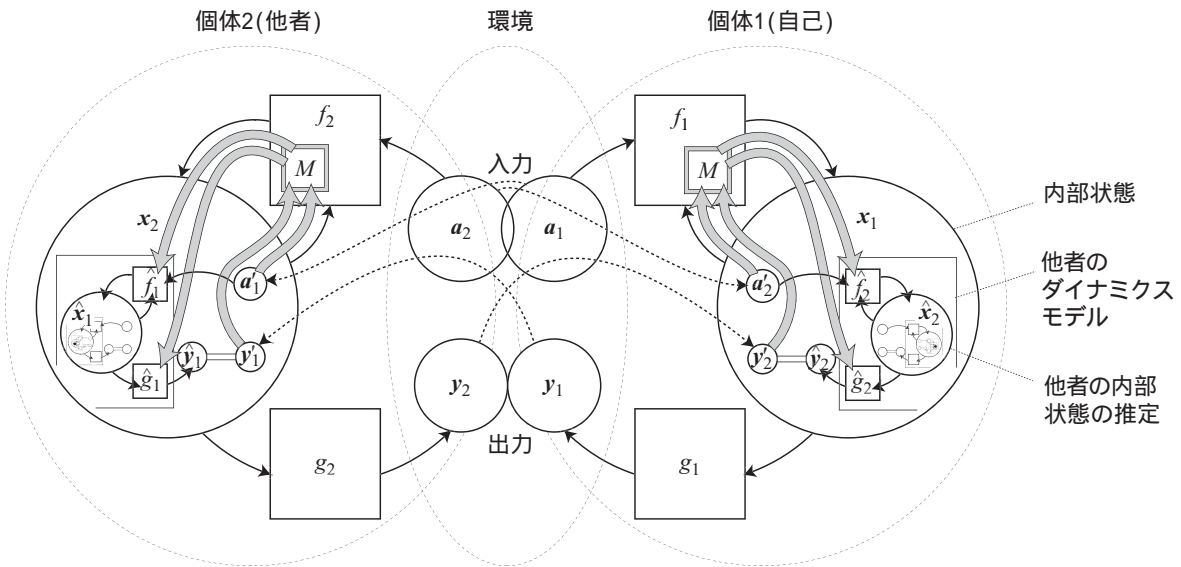


図2 ダイナミクス推定問題としてのコミュニケーションモデル

### 3.1 パラメータ次元の限界

前節では、コミュニケーションにおける他者の内部推定を、ダイナミクス推定問題として捉えられることを示した。

しかし、一般的なダイナミクス推定は、コミュニケーションとは無関係に、生物の脳が古くから行なってきたことである。運動のフィードバッ

ク制御、環境の地図の構成など、どれも自分の外界に関する状態変化の予測の例である。直接観測できない内部状態とその状態遷移・観測写像を推定することで、入力と出力の複雑な関係を解決し、最適な行動を実現する手がかりにしてきた。実際に、そのような脳の処理を再現しようとする計算論的研究もある [WK98]。

それでは、こうした能力をより発展させ、複雑にすれば、他者の予測モデルというコミュニケーションの基礎が築けるのだろうか。これまでのコミュニケーションに関する計算論的研究の分野では、それが可能である、という暗黙の了解が存在しているように思える。

しかし、筆者はそうは思わない。なぜなら、自分と同等な他者は複雑すぎて推定が困難であるからである。人間のような個体は、内部状態の次元が非常に大きいと考えられる。また、それらのダイナミクスも単純ではないうえに、状態空間の再構成による次元の増加が加わる。観察するだけでは、他者の状態遷移写像  $\hat{f}_2$ 、出力写像  $\hat{g}_2$ 、現在の内部状態  $\hat{x}_2$  をすべて推定することは、非常に困難である。

理論的には、大量の観測データを用意することで推定を行うことは不可能ではない。しかし、推定する側のパラメータの次元は、推定される側の内部状態の次元よりもはるかに大きくなる。言い換えると、推定する側の次元が決まっているとすると、それよりもかなり単純な、低次元で記述できる相手しか推定することができない。これを超えるような複雑な相手は、大雑把に近似する程度のことしかできない。

前述のサルの子の欺き行動のように、再帰を含むような高度な推定を可能にするためには、自分と同等の他者、つまり予測者に近いオーダーのパラメータ数をもつ他者の推定が必要になる。これは、観察だけから推定できる次元を大幅に超えており、また大雑把な近似での推定も難しいと思われ、実現が事実上不可能ではないかと思われる。

実際、人間であっても、隠れ変数のダイナミクス推定を含む予測問題の解決能力は決して高くない。他の人間のような複雑な対象を推定することは難しく、これ以上の高度な予測能力がコミュニケーションのために必要であったとは考えにくい。

### 3.2 ダイナミクス推定の手がかり

前述の困難は、観測のみからダイナミクスをすべて推定する、という制約があったために生じたものであった。2.2節で論じたように、ダイナミクスの手がかり（内部状態の遷移写像や出力写像についての情報）があれば、この困難を回避できる可能性がある。

もっとも単純には、これらの手がかりが先天的に与えられていると考えることができる。つまり、人間の脳には、相手がどのような状態のときにどのように見えるか、という基本情報処理構造が生得的に構築されている、という考え方である。これは十分に納得できることであり、表情や声色などのプリミティブな要素に関して先天的な知識や反応があってもおかしくはない。

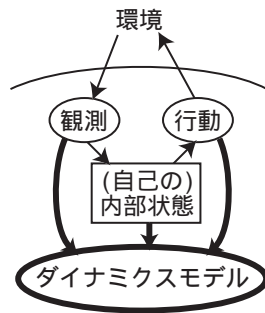
しかし、これだけでは不十分である可能性が高い。人間は、自分自身に関する知識であっても多くを後天的な学習によって獲得する。他者に関する知識、さらに再帰的な推論といった複雑な推定が、すべて先天的に持っていると考えことは難しい。実際、再帰的な「心の理論」は、人間の場合でも、年齢を積むにしたがって発達することが知られている [鯨岡 97]。何らかの学習によって後天的にダイナミクスを学習していると考えほうが自然である。

しかし、いったいどうしたら学習できるのだろうか。他者の観察では不足なのだから、どこかに学習対象を探さなければいけない。

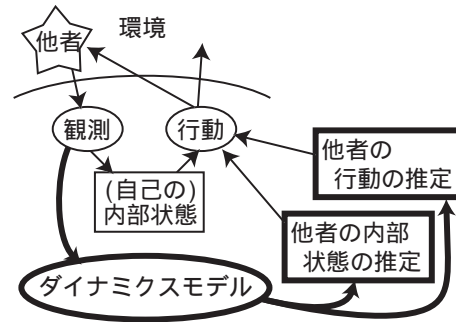
### 3.3 自分自身の利用による解決

この問題を解決する筆者らの提案は、「他者のダイナミクスを自分自身から学習する」というものである (図 3)。すなわち、他者は、基本的には自分自身と同じダイナミクスに従うと仮定するのである。この仮定が正しければ、他者の内部状態について大まかな推定をすることができる。細かい差異はその後で学習すればよい。

この提案のポイントは、次の点にまとめられる。



(a) 自分自身のダイナミクスからダイナミクスモデルを獲得



(b) 他者を観測するときには、獲得したダイナミクスモデルを利用し、相手の内部状態・行動を推定

図3 自己を利用する解決方法

- 「自分と同等である他者」の内部推定という、コミュニケーションの計算理論に適合した方法である。ダイナミクスが同等なのだから、上記の仮定は満たされる。というより、「自分と同等の他者」という制約をもっとも素直に表現した形態であるといえる。
- 他者よりも自分自身の方が、観測可能な情報の量はずっと多い。感情や欲求などの内部状態や、痛みや触覚といった他者では直接観測できない入力であっても、自分自身であれば直接観測できる。こうした豊富な情報から構築したダイナミクスモデルを使うことで、相手の痛みや感情などを推定することはずっと容易になる。
- 内部状態空間の再構成による次元の増加が抑えられる。自分自身について分かっている内部状態や状態遷移写像を利用すれば、内部状態空間をゼロから再構成する必要がない。そして、推定対象の内部状態空間の構造が分かっているなら、観測写像の学習もはるかに容易になる。

この構成が働く様子を、具体的な人間のコミュニケーションの場面で例示すると、次のようになる。まず、ダイナミクスモデルは、自分自身の

体験から、入力と状態変化の関係（手をぶつけると痛む）や、観察と状態の関係（血が出ていると痛い）を学習する。そして、その学習結果を利用し、他者や他者の環境の観察から他者の状態を推測することができる（手をぶつけて出血しているから痛そうだ）。また、行動決定の方法によっては、他者の状態に働きかけ（薬を塗る）、変えていく（痛みをおさえる）ことができる。

ここで注意すべき点は、このダイナミクスモデルは、自分自身の完全な反映である必要がない、ということである。自分自身を観測し学習する枠組みには、自分自身の複雑性を自分の中で表現するという本質的な困難があることが知られている [Rös98]。しかし、本研究で必要となるのは、あくまで他者の内部状態の推定の“種”として使えるダイナミクスモデルであり、そのモデルが自分自身を適切なレベルで近似していれば十分である。

#### 4 主観的-客観的情報変換に関する困難

本節では、もうひとつの困難、主観的情報と客観的情報の変換について述べる。自分自身のダイナミクスを単に利用するだけでは、他者の内部状態を推定することできない。そのモデルの入力である、他者にとっての主観的入力を、自分

が観察した情報から再構成しなければならないからである。しかし、主観的入力のリ構成は、内部状態の推定と同レベルの困難であり、簡単には解決できない。ここでは、ナイーブな手法では主観的情報の必要性から逃れられないことと、それを解決できる方法としての自己の客観的観測について述べる。

#### 4.1 自己適用原理

前節の方法をもっともナイーブに実現するには、自分自身の行動規則をそのまま適用して、他者の行動の予測に使う、という方法が考えられる。ここでは、この方式を自己適用原理と呼ぶ(図4)。自分自身の状態遷移写像  $f_1$  と出力写像  $g_1$  を、他者予測のためのダイナミクスモデル内の対応する写像、 $\hat{f}$  と  $\hat{g}$  にそれぞれ投影(コピー)する。そして、そのダイナミクスモデルを用いて、相手の内部状態  $\hat{x}$  を推定し、出力  $\hat{y}$  を予測する。

この原理で問題となるのは、自己と他者の立場の違いである。他者に関して自己が観測できるのは、あくまで客観的に観察した情報でしかない。ここでは、客観的に観察した情報には \* を付して  $a_2^*$  の形で示す。しかし、これは自分自身の状態遷移写像をコピーした写像  $\hat{f}$  の入力として適合しない。必要な入力  $a_2'$  は、自分自身がその他者の立場であったときに受け取るであろう入力、すなわち主観的な入力の情報である。この違いを埋めるには、 $a_2^*$  から  $a_2'$  に情報を変換する過程  $V_a$  が必要になる。

また、出力に関しても、ダイナミクスモデルの出力  $\hat{y}_2$  は主観的な出力である。行動決定に利用するには、その出力を自分の立場から見たもの  $y_2^*$  が必要になる。この変換のためには、 $\hat{y}_2$  を  $y_2^*$  に変換する過程  $V_y$  が必要になる。

この過程、すなわち主観的情報と客観的情報の変換機構を用意することが、コミュニケーションの計算理論における第2の困難である。我々の直感に反して、この変換は自明なプロセスではない。2種類の情報表現は全く異なる様式で

なされているからである。視覚の場合でも、自分の視点から見えることが他者の視点でどう見えるかは、複雑な視点変換の計算を経る必要がある。触覚・痛覚などは、ほとんどを視覚や聴覚などのほかの感覚から再構成しなければいけない。

たとえば、相手が誤って手を机にぶつけたのを見たとする。自分をその立場において考えると、手を机にぶつくと驚くだろうということがわかる。その規則をそのまま適用し、相手も驚くだろう、と推論するのが自己適用方式である。しかし、そのためには、自分の見ているものから、相手からの主観的な手と机の見え方への変換ができなければいけない。また、その視界から痛みの感覚を再構成しなければならない。このように複雑な変換や再構成は、人間が先天的に持っているとも考えにくいし、他者の観察からだけでは獲得することも難しい。

別の言い方をすれば、主観的入力の推定は、内部状態の推定と同等の困難であるといえる。他者の内部状態は私秘的 (private) であり、他者の主観的入力も私秘的である。 $\hat{f}$  や  $\hat{g}$  などの私秘的な処理の枠組みが分かっても、観察者に見えることから私秘的な情報への橋渡しはできない。筆者らの主張は、自己適用原理は、内部状態推定の問題を主観性-客観性変換問題に置き換えるだけであり、解決になっていない、ということである。

#### 4.2 自己観測原理

本研究では、この問題を解決する方式として、自己観測原理を提案する。これは、自分自身の行動を予測するダイナミクスモデルを構築して、他者の行動の予測に使う、という方法である。

この方式の概要を図5に示す。自分自身の状態  $x$  の中に、他者に適用するダイナミクスモデル  $\hat{f}^*$ 、 $\hat{g}^*$  が表現されている。自己適用原理と異なるのは、ダイナミクスモデルが、客観的に観測した情報をそのまま扱うとすることである。

まずモデル学習の段階(図5(a))では、モデル



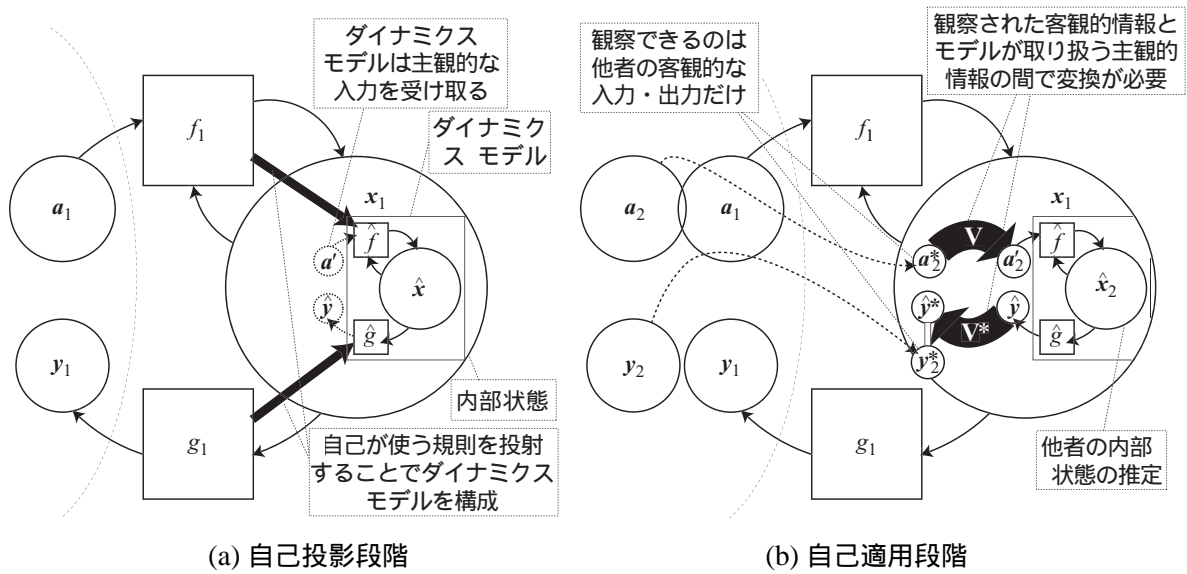


図4 自己適用原理

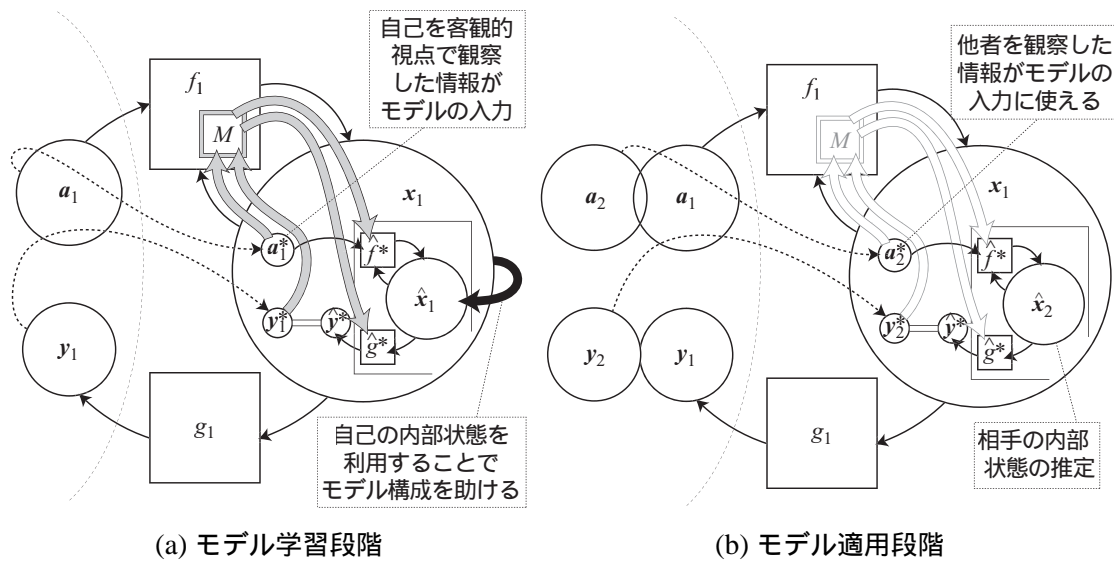


図5 自己観測原理

学習器  $M$  が、自分自身の入力や出力を改めて客観的視点で観測した観測結果  $a_1^*$ 、 $y_1^*$  についてのダイナミクスを学習する。ここで、自分自身の客観的観測というのは、外界を経由した自己観測で、他者に対しても同様の観測が可能なものを指す。例えば、手の動きの見え方、ぶつかる音、客観的空間内での机と手の位置関係などで

ある。この段階では、モデルの内部状態  $\hat{x}$  は、実際の自分の内部状態  $x_1$  を教師データとして利用できるため、学習は容易である。

こうして、自分自身の行動を学習したダイナミクスモデル  $\hat{f}^*$ 、 $\hat{g}^*$  ができると、それは他者に対する観察に適用することができる(図5(b))。他者の入力を客観的に観測した情報  $a_2^*$  を、ダイ

ナミクスモデルに通すことで、他者の内部状態の推定  $\hat{x}$ 、他者の行動を客観的に観測したものの予測  $y_2^*$  が得られる。

上記と同様の例で考えると、このダイナミクスモデルは、最初に「手をぶつくと痛い」ということを自分の経験から学習する。このモデルでは、自分自身に対する客観的観測である、自分から見える手の動きやぶつかる音、客観的空間上での机と手の位置関係などと、内部状態における痛みの間の相関も学習する。その後で、相手が手を机にぶつくところを観察したときには、自分から見える相手の手の動きや、ぶつかる音、客観的空間内での机と手の位置関係などから、その相手は痛みを感じるだろうと推定できる。

こうした相関は、自分自身の行動を決める上では直接には役に立たないことに注意されたい。自分自身が次からは手をぶつけないようにするためには、主観的な机との位置関係、腕を動かす筋肉信号と、痛みという結果の間の相関を学習すれば十分である。しかし、そうした主観的情報だけでは、相手の内部状態や行動の予測に役立たない。客観的情報と獲得し、主観的情報・内部状態との相関を学習することではじめて他者の内部状態の推定が可能になる。

#### 4.3 考察

本研究で提案する自己観測原理は、一見複雑に見えるが、実際は自己適用原理よりも単純である。その理由は、まず主観-客観の変換  $V_a$ 、 $V_y$  が不要になるため、計算が必要な要素が減るからである。

そのうえ、自己適用原理の考え方の元になっている「自分のダイナミクスを相手に適用する」というアイデアは、暗黙に自己観測原理に含まれている。自分のダイナミクスをダイナミクスモデルに転写する機構が、単純なコピーではなく、モデル学習器  $M$  によって実現されていると考えればよい。その意味で、自己観測原理は、ナイーブな自己適用原理において主観的情報の

再構成問題を解決した原理であると考えられることもできる。

もうひとつ、自己観測原理が優れている点は、他者に対しても同じモデル学習器  $M$  を使うことで、実際の他者を自分のもつモデルに反映できることである。自己適用原理では、ある入力について他者に対する予測が外れたときに、視点変換  $V_a$  とダイナミクスモデル  $\hat{f}$  のどちらを修正すべきか、あいまい性があるが、自己観測原理では、その2者は一体化した  $\hat{f}^*$  になっており、あいまい性が発生しない。

#### 4.4 新しい計算理論として

提案した自己観測原理は、推定者のパラメータ次元の限界、主観的情報と客観的情報の変換という、他者の内部状態の推定における2つの理論的困難をクリアする方式である。しかし、ここで示した原理は一般的なものであり、その実現方法は多数存在すると考えられる。また、実装に際してはさまざまな問題が表面化することもあるだろう。

こうした研究を推進するためには、提案した原理を組み込んだ新しい計算理論を考えることが有効である。2.1 節で述べたような従来の計算理論では、これらの理論的困難を打ち破るアルゴリズムを見つけることは難しかったが、新しい計算理論では、その問題とアプローチが明確にされており、研究として取り組みやすい形になっていると考えられる。こうして、さまざまなアルゴリズムが提案されれば、人工知能の開発などにもつながるだろうし、計算論的脳研究における仮説としてさまざまな脳研究の基盤ともなるだろう。

ここでは、コミュニケーションの計算理論として、提案した原理を統合した新しい計算理論を提案する。具体的には、2.1 節における計算理論に、次のような1行を付加したものである。

- 
4. 3. を実現する (= 他者の内部状態を推定する) ために、自分自身を客観的に観測す

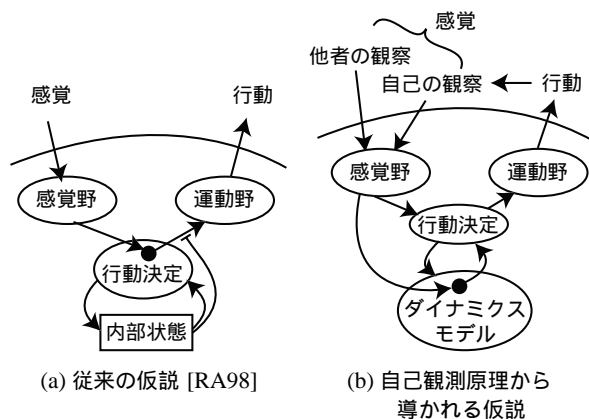


図6 ミラーニューロンの役割に関する仮説の比較。黒丸がミラーニューロンを表す。

ることでダイナミクスモデルを構成し、それを他者にあてはめること。

もちろん、この計算理論は唯一絶対のものではない。パラメータ次元の制約を越える有効な方法が別に存在する可能性もあるし、そもそも本論文が基盤としたコミュニケーションの定義や元の計算理論に問題がある可能性もある。そうした議論も含めて、本論文がコミュニケーションの計算理論に関する研究のきっかけになればいいと思っている。

## 5 関連研究

自己観測原理は、非常にシンプルであるが、様々な分野の研究と関係する。本章では、われわれの提案とそれらの研究との関連について述べる。

### 5.1 心理学

「他者の内部を推定する道具としての自己観測」という考え方は、1970年代後半に、進化心理学者 Humphrey によって提案されている [Hum78, Hum84]。しかし、彼の理論は純粋に理論的な仮説、すなわち進化の過程でなぜ意識が生まれたのか、という問いに対する答えのひとつでしかなかった。自己のどの部分を観測すべきかという彼の指定にはあいまいさが残っていた。

本研究は、コミュニケーションプロセスの定式化を通して、コミュニケーションにおける困難を解決するには自己観測原理が有効であること、そしてその際には自己の主観的状态と客観的観測とを関連付けることが必要であることを示したものであり、Humphrey の理論を精緻化したものである。また、この原理を計算理論として定式化したことは、Humphrey の仮説を新たな学問領域に広げたものであるといえる。

Baron-Cohen は、自閉症児の研究に基づく心の理論の提案 [BC95] の中で、Humphrey の仮説は自己適用原理に相当すると主張し、そのような仕組みは自閉症患者における症状の解離と矛盾することを指摘した。そして、内観による方式(自己観測原理に相当すると思われる)を考えれば自閉症患者の症状との矛盾がないことを短く示唆している。筆者らは、Humphrey の仮説は自己観測原理も含む大きなものであると考えているが、自閉症研究によっても自己投影原理より自己観測原理が適しているという結果は興味深い。

### 5.2 計算論的脳研究

計算論的神経科学におけるコミュニケーション研究のほとんどは、自己観測原理と合致しない。川人ら [川人01] は、ダイナミクスとのインタラクションモデルの拡張としてのコミュニケーションモデルを議論したが、このモデルでは主観的情報と客観的情報の表現が等価であると暗黙に仮定されている。

我々の知る限り、Humphrey の仮説や、それに類した枠組みに基づいてモデルを構成した研究は存在しない。これは、他者の内部状態の推定を成立させる原理が、これまで計算理論の形で定式化されてこなかったことが原因であると考えている。本研究は、今後の構成論的アプローチを推進する原動力となるであろう。

コミュニケーションと離れた文脈では、谷は、意識や自己の存在を追及するための構成論的研究を提案している [Tan98]。しかし、この研究で

は、環境内に単独に配された個体によって行われ、個体間のインタラクションや他者の内部状態の推定については考慮されていない。

### 5.3 実験的神経科学

本研究で提案する計算理論に基づくと、ミラーニューロン [DFP<sup>+</sup>92, GFFR96] の役割について、新しい仮説を考えることができる。ミラーニューロンとは、自分がある行為（物をつかむなど）をしているときと、他者がその行為をしているのを観察しているときに活動を示すニューロンである。サルでは腹側運動前野にあり、人間では同様のニューロンが Broca 野にあるとされており [RFGF96]、コミュニケーションや言語の能力と深い関連があると考えられているが、具体的にどのような働きをしているのかは分かっていない。

Rizzolatti らは、このニューロンは行動を指令する役割を担っていたものが、コミュニケーション能力の発達に従って他者の同様の行動にも反応するようになったニューロンであるという仮説 [RA98] を唱えた (図 6 (a))。しかし、どのように他者の行動を関連付けるのかについては不明である。また、観察した行動にミラーニューロンが反応しても、それだけで行動が誘発されないように抑制がかかっているとされているが、この抑制の仕組みについても説得力のある議論はない。

それに対して、本研究のモデルからミラーニューロンの役割を考えると図 6 (b) のようになる。3.3 節で述べた仕組みがミラーニューロンの発生を非常にうまく説明できる。自己の行動を観察するダイナミクスモデルの中で、動作の要素を学習したニューロンが、他者の行動に対して適用される結果、ミラーニューロンとして活動するようになる。また、このようなニューロンであれば、直接行動には影響しないため、抑制する仕組みは不要である。

Oztop らは、自己の行動における誤差フィードバック回路のニューロンが、他者の行動にも反

応するようになり、ミラーニューロンになったという説を唱えている [OA02]。たとえば、コップをつかむときの手の運動の誤差制御をするために、自分の手の運動とコップとの位置関係を客観的に捉えるニューロンが発達し、それが外適応 (exaptation) の結果、他者のコップをつかむ運動を見たときにも反応するようになった、というものである。この説は、本研究のモデルとも合致することから、誤差フィードバック機構が人間のコミュニケーションにおける自己観測機構の進化的原点であったと考えることができる。

### 5.4 人工知能

チェスや将棋などのゲームの思考アルゴリズムでは、他者の行動の推定をいかに行うかが重要である。これまで提案されている思考アルゴリズム ( $\alpha$ - $\beta$  法 [石畑 89] など) は、他者は自分と同じ価値関数に基づいて行動を決定するという仮定に基づいて設計されている。これは、自己適用原理を利用していると捉えることができる。こうしたゲームでは、規則に対称性が成立しているため、視点変換メカニズムを提供することが容易であり、自己適用原理が適していると考えられる。しかし、自己観測原理を利用しても思考アルゴリズムが構成することも可能性である。こうした思考アルゴリズムによるプログラムならば、対戦の繰り返しから対戦相手の特徴を学習して、戦略を修正することができるかもしれない。

## 6 むすび

本論文では、コミュニケーションについての従来の計算理論である、他者の内部状態の推定を検討し、その実現を妨げる 2 つの困難、推定者のパラメータ次元の制約と、主観的-客観的情報の変換があることを指摘した。そして、それを解決する方法として、自己の行動を学習したダイナミクスモデルを他者の内部状態の推定に利用するという自己学習原理を提案した。この原理は、自己を客観的に観察した情報と、自己の

内部状態を関連付けることで、他者を客観的に観察した情報から他所の内部状態の推定を得られるというものである。本研究では、従来あいまいだった学習の対象と目的を明確にし、自己観測原理を計算理論に統合することで、今後のコミュニケーションの構成的研究への扉を開いたといえる。さらに、この計算理論と、自己意識やミラーニューロンなど、さまざまな研究分野との関連について考察した。

## 参考文献

- [AHP<sup>+</sup>00] C. G. Atkeson, J. Hale, F. Pollick, M. Riley, S. Kotosaka, S. Schaal, T. Shibata, G. Tevatia, S. Vijayakumar, A. Ude, and M. Kawato. Using humanoid robots to study human behavior. *IEEE Intelligent Systems: Special Issue on Humanoid Robotics*, Vol. 15, pp. 46–56, 2000.
- [合原 00] 合原一幸, 池口徹, 山田泰司, 小室元政. カオス時系列解析の基礎と応用. 産業図書, 2000.
- [BC95] S. Baron-Cohen. *Mindblindness: An Essay on Autism and Theory of Mind*. MIT Press, 1995.
- [DFP<sup>+</sup>92] G. Dipellegrino, L. Fadiga, L. Fogassi, V. Gallese, and G. Rizzolatti. Understanding motor events - a neurophysiological study. *Experimental Brain Research*, Vol. 91, pp. 176–180, 1992.
- [GFFR96] V. Gallese, L. Fadiga, L. Fogassi, and G. Rizzolatti. Action recognition in the premotor cortex. *Brain*, Vol. 119, pp. 593–609, 1996.
- [Hum78] Nicholas Humphrey. Nature’s psychologists. *New Scientist*, pp. 900–904, June 1978.
- [Hum84] Nicholas Humphrey. *The inner eye: Social intelligence in evolution*. Faber and Faber, 1984.
- [IOIK03] H. Ishiguro, T. Ono, M. Imai, and T. Kanda. Development of an interactive humanoid robot “Robovie” – an interdisciplinary approach. In R. A. Jarvis and A. Zelinsky, editors, *Robotics Research*, pp. 179–191. Springer, 2003.
- [石畑 89] 石畑清. アルゴリズムとデータ構造: 岩波講座 ソフトウェア科学 3. 岩波書店, 1989.
- [川人 96] 川人光男. 脳の計算理論. 産業図書, 1996.
- [川人 01] 川人光男, 銅谷賢治, 春野雅彦. モザイクの拡張とコミュニケーション – ヒト知性の計算神経科学 (第 5 回) –. 科学, Vol. 71, pp. 197–204, 839–843, 2001.
- [鯨岡 97] 鯨岡峻. 原初的コミュニケーションの諸相. ミネルヴァ書房, 1997.
- [MA03] Takaki Makino and Kazuyuki Aihara. Self-observation principle for estimating the other’s internal state. Mathematical Engineering Technical Reports METR 2003–36, Department of Mathematical Informatics, Graduate School of Information Science and Technology, the University of Tokyo, October 2003.
- [Mar82] David Marr. *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. W. H. Freeman, San Francisco, 1982.
- [OA02] E. Oztop and M. A. Arbib. Schema design and implementation of the grasp-related mirror neuron system. *Biological Cybernetics*, Vol. 87, pp. 116–140, 2002.
- [PW78] D. G. Premack and G. Woodruff. Does the chimpanzee have a theory of mind?

- Behavioral and Brain Sciences*, Vol. 1, pp. 515–526, 1978.
- [RA98] Giacomo Rizzolatti and Michael A. Arbib. Language within our grasp. *Trends in Neuroscience*, Vol. 21, pp. 188–194, 1998.
- [RFGF96] G. Rizzolatti, L. Fadiga, V. Gallese, and L. Fogassi. Premotor cortex and the recognition of motor actions. *Cognitive Brain Research*, Vol. 3, pp. 131–141, 1996.
- [Rös98] Otto E. Rössler. *Endophysics: the world as an interface*. World Scientific Publishing, 1998.
- [Tan98] Jun Tani. An interpretation of the ‘self’ from the dynamical systems perspective: a constructivist approach. *Journal of Consciousness Studies*, Vol. 5, No. 5–6, pp. 516–542, 1998.
- [WK98] D. M. Wolpert and M. Kawato. Multiple paired forward and inverse models for motor control. *Neural Networks*, Vol. 11, pp. 1317–1329, 1998.