

客観的自己認識の相互構築を促進する エージェントアーキテクチャの提案と実装

Proposal and Implementation of an Agent Architecture for Facilitating the Mutual Construction of Objective Self-Awareness

今井康智^{1*} 大本義正¹
Yasutomo Imai¹ Yoshimasa Ohmoto¹

¹ 静岡大学
¹ Shizuoka University

Abstract: 人間は自分自身を他者の目を通じて知覚し意識する傾向を持つとされており、こうした客観的な自己認識は他者を理解し行動する上で重要な要素となっている。この傾向が他者として認識されるエージェントにおいても適応される場合、エージェントの行動は人間からの印象、つまり他者の目を通じて変化することが望ましいと考える。そこで本研究では、この客観的な自己認識の相互構築を促進する枠組みを提案し実装する。

1 はじめに

人間と共生するコンパニオンエージェントの需要は情報技術の進歩と共に高まっている。こうしたエージェントの振る舞いはユーザの継続的な利用を促すためにも、能動的で適応的なものが良いと考える。ここで、人間の持つ能動性や適応性は、生存という根本的な欲求によって現れると考えられる。人間は生存の為、自ら危険の中に身を置き、そして環境に適応してきた。その為、人間にあるような根本的な欲求をエージェントに持たせることは、能動性や適応性の基本的な仕組みとして有用であると考えられる。しかし、人工物であるエージェントにとって生存等の欲求は意味がなく、コミュニケーション相手に違和感が生じる可能性がある。その為、欲求がエージェントにとって意味のあるもの、つまり存在意義と関連しているものが望ましいと考える。人間と情報空間を仲介することがエージェントの基本的な役割の一つであることを踏まえ、エージェントが自ら情報を集め、それを他者に共有し喜ばせることを欲求として持つことは、存在意義と行動が結びついており有用だと考えた [1]。本研究では、人間と共生するエージェントの実現という最終目標のために、情報の収集と共有という欲求をエージェントに持たせることで、人間とのインタラクションに与える影響を検討する。

欲求を持ち共生するエージェントを実現する上で考慮すべき要素は、擬人化やプライバシーの懸念など多

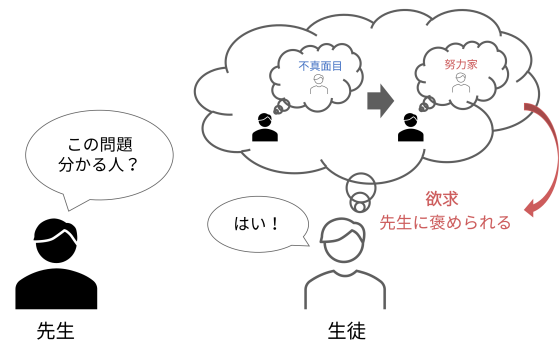


図 1: 人間同士の承認欲求の達成

岐に渡る [2][3]。その中でも本論文では、エージェントの基盤となる行動生成の枠組みを提案する。この枠組みによってエージェントは、他者を含めた環境に適応した自律的な行動をする為、情報の収集と共有という欲求を基盤として行動を生成することが可能となる。

情報の収集と共有という欲求では、共有した情報が他者に喜ばれることを一種の承認欲求が満たされることとして報酬に設定し、その為に環境内の様々な情報を自ら収集する。ここで、他者に喜ばれる情報は相手によって異なる為、エージェントは報酬が貰える行動を、相手毎に考える必要がある。その際、常に相手のモデルに合わせた行動を取っている、自身の欲求を効率的に満たすことはできない。

一方、人間同士のインタラクションでは、多くの場面で他者を介した欲求の達成が行われている。例えば、図 1 の場面では、ある学校での授業中、右側の生徒が、

*連絡先: 静岡大学
〒432-8011 静岡県浜松市中央区城北 3 丁目 5-1
E-mail: imai.yasutomo.18@shizuoka.ac.jp

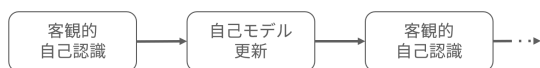


図 2: 客観的自己認識によるモデル更新

左側の先生に褒められたいという欲求を持っており、それを達成する為の行動を取っている。この場合、生徒は現状、先生の自分に対する印象があまり良くないもの（不真面目）であることを読み取っており、それを褒められるという欲求を達成する為の良い印象（努力家）へ変えるために積極的に発言している。人間はこのようにして他者を介して承認欲求を満たそうとしている。これは、他者から見た自分を想像し、その想像と、欲求が満たされた状態である自分を比較し、その差を埋める為の行動をしていると言い換えられる。こうした、自分自身を認識し、他の個体の評価の目を通じて自己について認識しようとする傾向を自己認識といい、人間同士のインタラクションでは不可欠だとされている [4]。その中でも、この場面のような他者の立場から自分を理解することを客観的自己認識と言い、そこで得た自己を踏まえ、自己モデルを更新し、また自己認識を行うというサイクルを回すことで、意図の押し付けや汲み取りを行っていると考えられている (図 2)[5]。

以上の客観的自己認識をコンセプトとしてエージェントの枠組みに取り入れる。それによりエージェントは、他者から見た自分のモデルを推定し、そのモデルを自分の欲求達成に都合の良い方向へ促す為に、自分のモデルを更新させることが可能となる。そうして、人間とエージェントが相互に他者の他者モデルを調整することで、お互いがお互いに適応し、持続的な関係が構築できるのではないかと考える。

本論文の構成として、第二章では客観的自己認識というコンセプトを達成する為、エージェントが参照可能な自己を定義する。第三章では、今回提案するモデルと、実装用のアーキテクチャを紹介する。第四章では、実装した環境、エージェントを述べ、具体的な振る舞いを内部状態の動きを含めて説明する。第五章では、本研究の展望を述べる。

2 自己の定義

エージェントが他者の中の自己を持つには、自己をモデル化する必要があり、この自己の側面として Subagdja らは以下の五つの要件の定義を想定している [5]。客観的自己認識では、他者から見た自己としてここで定義するものを参照し、自らのモデルを更新することができる。

一つ目の側面はアイデンティティであり、これはエー

ジェントを他者とは異なるユニークな個人として特徴づける情報を示している。例えば、エージェントの名前、外見、色、年齢などが挙げられる。本研究で作成するコンパニオンエージェントの見た目は架空の生物であるユニコーンにし、かつ環境内に一体のみ存在させている。その為、他者から見たエージェントはその存在自体がアイデンティティであり、エージェント自身がその個別な情報を参照することは少ない。今後、環境内に本エージェントを複数存在させる場合には、それぞれの色や名前の違いを認識させ、それを踏まえた自己認識を行えるようにデザインする。

二つ目の側面は身体性であり、これはエージェント自身の身体構造、存在する状況及び環境を特徴づける情報を示している。例えば、エージェントの大きさやポーズ、現在の位置、環境の状態を変化させる能力などを知ることが挙げられる。本研究では、客観的自己認識を行う対象から自分が見えないことを認識できるようにデザインする。また、現実的な動物の外見にしてしまうと、客観的自己認識時にそのユーザが持つイメージに引っ張られることが考えられるため、架空の生物であるユニコーンの外見を用いる。

三つ目の側面は心であり、これはエージェント自身の精神状態に基づく個人の特徴づけに関わる情報を示している。例えば、思考、感情、欲望、意図、想像を含むエージェントの心の状態に基づいて個人を特徴づける情報である。本エージェントは、冒頭で提示したように、情報の収集と共有という欲求を持つ [1]。この欲求下でエージェントは、環境内の最新の情報を収集し、そこで集めた情報を環境内に存在する他のエージェントに共有する。そして共有した情報が相手に喜ばれることがエージェントの報酬となる。この報酬を効率的に獲得する為には、他者が自分の欲求を理解していることが重要であり、その為エージェントは客観的自己認識のコンセプトを用いて他者の適応を促進させる。

四つ目の側面は関係であり、これは他の個人とどのように社会的に関係しているかに基づく個人の特徴づけを示している。主観的自己認識では、エージェントが友人、敵、所有者、好き嫌い等社会的な関係を持ち、客観的自己認識では、他者がエージェントの自己について何を考え、感じているか意識するようになる。欲求に他者の喜びを用いるエージェントにおいて他者の中の自己の印象を知ることが、関係を好意的なものへ発展させる為にも重要である。本研究でも、他者からのフィードバックを通じて、その印象を読み取り、自身の欲求達成に都合の良い方向へ操作できるようにデザインを行う。

最後の側面は記憶であり、これは自分が経験したこと、あるいはこれまでの経験に基づいて将来展開すると思われることに関する情報を示している。この側面により、エージェントは将来の見通しに関する予見を

構成することができる。本研究では、長期記憶を構築することで、将来の自分の欲求を満たす行動を自らプランニングすることが可能となる。また、プランニングの結果出力された行動により、自分の欲求が満たされたかを判定し、それを基に長期記憶をアップデートすることが可能となっている。

以上の五つの側面を定義することで、エージェントは他者から見た自己を認識し、自身の欲求を満たす為に適した自己へ操作する行動が可能となる。例えば、自身の喜びを効率的に満たす為、相手の持つエージェントモデルを操作する行動を取る。その際には、自分の場所が相手から見える位置かを確認し、また現在の相手からの印象を認識できることで、行動する場所やインタラクション内容を操作した上で行動を決定し、その結果を記憶に蓄積することでモデルを常に更新することが可能となる。

3 提案システム

以上の自己の側面を踏まえ、本研究で提案するシステムを示す。まずは概念モデルとして、エージェントがどのように客観的的自己認識を行うのかを提示する。その後、実装用アーキテクチャを提示し、モデルとの関連を説明する。

3.1 モデル

本研究で提案するモデルのイメージを図3に示す。ここでは、前述した自己モデルを、他者との間でどのように更新しているかを表している。図は大きく左右で分割され、左側では主観的的自己認識、右側では客観的的自己認識を行っている。

主観的的自己認識とは一人称視点での自己認識であり、自身の意図モデルと他者モデルを内包して他者及び自己を理解している。これにより、エージェントは自己の心の側面から、自分の実現したい状態を他者を介して持つことができる。それぞれのモデルは高次の意図と低次の行動の階層型ネットワークを持つ意図モデルとして構成されている。エージェントは、他者の行動から他者モデルを参照して意図推定を行い、その意図によって自身の意図モデルのパラメータを更新し、欲求を満たす行動を決定する。ここで、内包する他者モデルの枠組みは自身の意図モデルと同じモデルによって構成されている。その為、他者の行動や意図は、実情はどうあれエージェントの取れる行動や意図として処理されている。

主観的的自己認識モデルにより他者の行動から自身のモデルを更新することができるが、その更新は他者の実際の行動に依拠しており、自らの望む未来へ他者を

導くことが難しい。そこで、主観的的自己認識モデルによって出力される行動を、客観的的自己認識によって評価するモデルを構築した。これがモデル右側で、決定した行動をユーザの立場から評価することが可能となっている。そしてこの評価は、主観的的自己認識モデルと同じ過程を他者とエージェントの状況を入れ替えて行っている。これは自身の意図モデルと他者モデルが同じであるから可能であり、文字通り相手の立場に立って、その意図を推定し、自身の意図を変化させ、行動を出力することができる。なお、この状況として自己の側面の一つである身体性を参照し、相手から自分が見える位置かどうかは行動の生成に大きく関わっている。そして、出力された行動が、自身の行動に対するユーザの行動予測となる。

以上の二つのモデルを接続することで、エージェントは主観的的自己認識でユーザを踏まえた行動を決定し、客観的的自己認識でその行動を相手の立場から評価することができる。ここで予測される行動が、主観的的自己認識で想定した望む未来とどの程度近いかを評価値として、最も評価の期待値が高くなる行動を、主観的的自己認識の意図モデルを一時的に更新して探索する。例えば、ユーザに感謝される状態を目的として次の行動を取る際に、ユーザからの印象や位置関係を考慮し、様々な行動をシミュレートし、記憶を基に比較することができる。そして、最終的に出力された行動により実際にユーザがどう変化したかを取得し、その差を学習し、自身のモデルを更新することができる。こうしたサイクルを回すことで、図2のような客観的的自己認識による自己モデル更新が行えると考える。

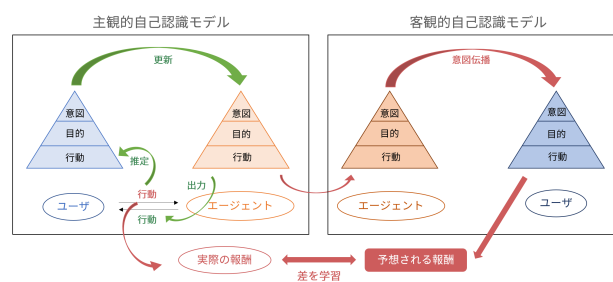


図 3: 自己認識モデルイメージ

3.2 実装用アーキテクチャ

以上のモデルに従って作成した、実装用のアーキテクチャを図4に示す。アーキテクチャ内には、前述したモデルと一致する部分として、行動決定システムと行動予測システムがあり、それぞれ主観的的自己認識と客観的的自己認識をコンセプトに作成している。以下ではアーキテクチャの振る舞いを、エージェントが意思決定を行う際の流れに沿って説明する。

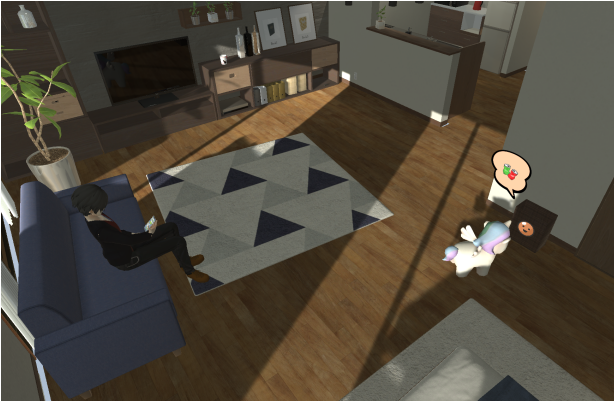


図 5: 環境の様子

物、冷蔵庫、ゴミ箱、物置、玄関に存在し、普段は隠れているが、エージェントに呼びだされることで出現し、インタラクションが開始される。エージェントによるインタラクションは、遠目からでも分かるように吹き出しを用いて行われた。なお今回の実装では、実装したエージェントがこれらの専門エージェントに対して客観的自己認識を行うことはない。

4.2 エージェント

エージェントの見た目は前述の通りユニコーンを用いている。また、言語を発話せず吹き出しによるコミュニケーションを行っている。これは、言語を発話する場合、エージェントが人間からの客観的自己認識で人間レベルの返答を期待されてしまい、実際とのギャップがインタラクション体験に影響を与えることを防ぐためである。

また本論文はモデルの枠組みを評価するものであり、実装はシミュレートの部分を学習済みモデルとして扱っている。その為、意図決定後、次の行動をプランニングする際は、その行動によって報酬がどの程度与えられるかを確率で評価している。またシミュレートとしては、自らのワーキングメモリに乗せる意図を変化させることで行動を変え、それぞれの確率を比較し最も高い、つまり最も報酬が高くなると予想される行動を取っている。また、エージェントはこのモデルの他に環境情報の評価システムを保持しており、環境内の情報変化を取得、予測し、情報の確信度が低い情報への収集行動や、情報の変化が激しくユーザが喜びそうな価値の高い情報の共有行動が行えるようになっている。

4.3 実装動作例

以下では、実装したエージェントの振る舞いを、内部状態の動きを含めて説明する。初めに、基本的な状

態の説明として情報の収集と共有という欲求を満たす行動について説明する。その後、環境内で発生するイベントの例として宅配便が来たというシナリオ上でのような振る舞いをするのかを説明する。

実装したエージェントの振る舞いは、環境の変化に適応して生成される。基本的には、自らの環境情報の評価システムを用いて、環境内の情報変化を予測し、現在持っている情報から最も変化したと予測されるものをワーキングメモリ上に格納する。それを参照しシミュレートを行い、情報収集等のプランを決定する。収集行動では、手に入れた情報を環境情報評価システムに入れ、予測との差から予測モデルを更新する。この時、予測した値と大きく異なる値を受け取った場合には、再度シミュレートを行い、確率に従ってユーザへの共有が行われる。また、それ以外の環境の変化として、ユーザからの呼びかけやインターホン等のイベントがあり、そうした環境の明示的な変化を入力とした場合も、ワーキングメモリ上に意図を格納し、プランニング、シミュレートを行っている。

ここで、エージェントが自身の欲求を満たす為に他者の他者モデルへ働きかけが行われるシナリオとして、宅配シナリオを紹介する。まず、ユーザがスマートフォンを操作しており、エージェントはその意図を推定し、ワーキングメモリに格納する。ここで、ワーキングメモリ上には、環境情報の評価システムによって格納された植物の情報収集が存在しており、スマートフォン、収集、植物という意図が浮上している。その中でプランニングした結果、[植木鉢エージェントから情報を収集する, 収集した情報をユーザに共有する]という一連のプランが選ばれ、行動を実行する。エージェントが植物の情報収集を行っている際、インターホンが鳴るといったイベントが発生する。エージェントは知覚システムを通じインターホンへの反応を行う。インターホンの原因として宅配便を推定し、エージェントのワーキングメモリ上に、インターホン、宅配便、収集、植物という意図が浮上する。それぞれが持つ行動の確率を合算し、次のプランとして [玄関エージェントから情報を収集する, 収集した情報をユーザに共有する] という一連の行動をプランニングする。それに従い、エージェントは玄関まで行き、宅配便が来たという情報を受け取る。その後、プラン通りにユーザの元へ向かうが、道中でプラン作成後の時間経過に伴うプランの再評価が行われる。この時、玄関で宅配便が来たという情報を受け取っている為、ワーキングメモリ上の意図は、宅配便、収集、植物というものに変化している。その結果、次のプランとして [印鑑を取りに行く, 印鑑をユーザに渡す, 宅配便が来たという情報をユーザに共有する] という一連の行動の予測報酬が現在の行動よりも高い為、プランを変更する。変更したプランを実行した結果、ユーザに喜ばれた為、エージェントも喜んだ。

本シナリオでは、エージェントが他者から見た自己を認識することで、印鑑を渡すという行動が、それ自体に対する感謝と、情報の収集と共有という欲求の提示によるタスク催促に繋がり、結果として予想される報酬によって評価され、選択された。こうして、他者に自身のモデル理解を促進させることが、エージェントにとっての戦略となり、この戦略が行動の端々に見られるという点が、本エージェントの特徴の一つであると考えられる。

本シナリオを含む環境上のエージェントの振る舞いを同研究室のメンバーに評価してもらったところ、エージェントが自分なりに考えて行動していることが伝わったという意見を貰った。また、振る舞いの中には、自らの欲求を優先する為環境の変化に反応しないという場面もあった為、そうした点からどういう考えなのか分からないという指摘もあった。今回は非常に短期的な振る舞いの観察だった為、こうした曖昧なモデルの提示となってしまったが、今後ある程度の時間を取って実験を行う中で、エージェントのモデルがどう構築されて行くのかを調べていきたい。

5 展望

本論文では、情報の収集と共有という欲求を持つ自律エージェントを作成する為、ユーザと客観的自己認識の相互構築を促進するアーキテクチャを作成し、実装した。実装の結果、特定のシナリオ上での振る舞いから、エージェントが自分なりに考え、ユーザを支援しようとする姿勢を取らせることができた。その為、本枠組みを基盤とし、客観的自己認識部分を自ら構築し、広範な状況下でもシミュレートが行えるように実装を進めていく。そして、リアルタイムで学習を行うエージェントに対して人間がどのような他者モデルを構築し、それがインタラクションによってどう変化していくかについて実験を通じて解明する。

参考文献

- [1] 今井康智, 大本義正: 情報の収集と共有を欲求としたインタラクションエージェントの提案, HAI シンポジウム (2022).
- [2] Roesler, E., Manzey, D. and Onnasch, L.: Embodiment matters in social hri research: Effectiveness of anthropomorphism on subjective and objective outcomes, *ACM Transactions on Human-Robot Interaction*, Vol. 12, No. 1, pp. 1–9 (2023).
- [3] Chatterjee, S., Chaudhuri, R. and Vrontis, D.: Usage intention of social robots for domestic purpose: From security, privacy, and legal perspectives, *Information Systems Frontiers*, pp. 1–16 (2021).
- [4] Rochat, P.: The ontogeny of human self-consciousness, *Current Directions in Psychological Science*, Vol. 27, No. 5, pp. 345–350 (2018).
- [5] Subagdja, B. and Tan, A.-H.: Beyond autonomy: The self and life of social agents (2019).