

生物学的制約の導入による「人間らしい」振る舞いを伴う ゲーム AI の自律的獲得

藤井 叙人^{1,3,a)} 佐藤 祐一¹ 中嶋 洋輔² 若間 弘典¹ 風井 浩志^{1,b)} 片寄 晴弘^{1,c)}

概要：ゲーム AI が搭載されたコンピュータプレイヤー (NPC) の振る舞いの自律的獲得は、「人間の熟達者に勝利する」という人工知能領域の長年の目標を達成しつつある。一方で、獲得された NPC の振る舞いは、過度に最適化され機械的に感じるという課題が浮上している。人間プレイヤーの代替という視点に立てば、NPC には『強さ』だけではなく、人間を楽しませるための『人間らしさ』が構成される必要がある。本研究では、『人間の生物学的制約』を課した強化学習や経路探索により、人間らしい NPC を自律的に構成する手法について提案する。また、自動獲得された NPC の振る舞いについて、本当に人間らしいと解釈されるかどうかを主観評価実験により検証する。

Autonomously Acquiring human-like behaviors of the Game AI with Biological Constraints

FUJII NOBUTO^{1,3,a)} SATO YUICHI¹ NAKAJIMA YOSUKE² WAKAMA HIRONORI¹ KAZAI KOJI^{1,b)}
KATAYOSE HARUHIRO^{1,c)}

Abstract: While various systems that have aimed at automatically acquiring behavioral patterns have been proposed and some have successfully obtained stronger patterns than human players, those patterns have looked mechanical. When human players play video games together with NPCs as their opponents/supporters, NPCs' behavioral patterns have not only to be strong but also to be human-like. We propose the autonomous acquisition of NPCs' human-like behaviors, which emulate the behaviors of human players. In this paper, instead of implementing straightforward heuristics, the behaviors are acquired using techniques of reinforcement learning and pathfinding, where *biological constraints* are imposed. We evaluated human-like behavioral patterns through subjective assessments, and discuss the possibility of implementing the proposed system.

1. はじめに

ゲームにおけるコンピュータプレイヤー (=NPC:ノンプレイヤーキャラクター) の戦略の自律的獲得は、「人間の熟達者に勝利する」という人工知能領域の長年の目標を達成しつつ

ある。コンピュータ将棋における棋力はここ数年で飛躍的に向上しており、最強将棋プログラムが現役のトッププロ棋士に史上初の勝利を収めたのは記憶に新しい (2013 年 4 月)。コンピュータ囲碁でも 4 子局でプロ棋士に 1 勝 1 敗と迫っており (2013 年 3 月)、ボードゲームで人間がコンピュータに勝てなくなる日も遠くないと言われている。そのため、「強さを追求」した“強い NPC”の次の段階として、「人間を楽しませる」ための“人間らしい NPC”の自律的獲得に興味が集まっている。

ビデオゲームの分野に目を向けると、人間らしい NPC の実装に余念が無いことが伺える。ほとんどのビデオゲームには NPC が用意されており、その NPC のもつ戦略や

¹ 関西学院大学大学院 理工学研究科, Graduate School of Science and Technology, Kwansai Gakuin University

² 関西学院大学 理工学部, School of Science and Technology, Kwansai Gakuin University

³ 日本学術振興会特別研究員 DC2, Research Fellow of Japan Society for the Promotion of Science

a) nobuto@kwansai.ac.jp

b) kazai@kwansai.ac.jp

c) katayose@kwansai.ac.jp

振る舞いが、人間プレイヤーのプレイフィール（プレイ時の感覚や印象）に大きな影響を与える。ユーザ数の増加、ひいては、売上の増加には、人間らしいNPCの実装が欠かせないのである。（実際、2012年国内ゲーム市場規模は前年比15.3%増の9776.9億円[1]と年々躍進している）。しかし、人間らしいNPCの戦略や振る舞いのデザインは、プレイヤーのレベルにあわせた難易度の調整（レベルデザイン）も含め、ゲームプログラマによる煩多な作り込み作業により実現されている。

“人間らしいNPC”を自律的に獲得するための研究として、人間プレイヤーの戦略を記録し機械学習により模倣する手法[2],[3],[4]、強いNPCに対して恣意的にエラーを導入する手法[5]などがある。これらの研究では、開発者が意図した「強くない」NPCのデザインが可能となっており、レベルデザインの一アプローチとしては有効である。しかし、「どのような振る舞いが人間らしいか」ということ自体が、そもそも形式化されていないため、開発者の経験（ヒューリスティック）による煩多な作り込みにより実現する他ない。

本研究では、『人間の生物学的制約』の条件下での機械学習及び経路探索により、人間らしいNPCの振る舞いを自律的に獲得する手法について提案する。生物学的制約とは、人間が生得的・遺伝的にもつ特徴や性質から生じる制約を指す。人間らしいNPCの獲得にあたっては、これを、人間がゲームをするときに必ず生じる制約や欲求として扱う。人間の行動制御における制約[6],[7]や自己実現理論[8]から着想を得て、「身体的な制約：“ゆらぎ”“遅れ”“疲れ”」「生き延びるために必要な欲求：“訓練と挑戦のバランス”」の2つを生物学的制約と定義する。生物学的制約の導入は、開発者のヒューリスティックに依存せず、生理学的・心理学的知見に基づいて設定できるため、汎用的な振る舞い獲得の手法として構築が可能である。振る舞い獲得の対象としては、アクションゲームの“*Infinite Mario Bros.*”を採用し、自動獲得された振る舞いが人間らしいかどうかを主観評価実験により検証する。

以下、第2章で、関連研究を紹介し、第3章で、生物学的制約を導入する意義と、その定義を述べる。第4章で、“*Infinite Mario Bros.*”の仕様と、振る舞い獲得の方法について説明する。第5章で、獲得された振る舞いの人間らしさを主観評価実験により検証する。

2. 関連研究

2.1 強いNPCを追求した研究

振る舞いを自動的に獲得する手法として、教師あり学習による事例参照型の手法[9]と、強化学習や経路探索による非事例参照型の手法[10],[11]に分類される。

教師あり学習は、事前に与えられた大量のデータセットを教師データ（入力データに対して出力されるべきデータ

の例）とし、有用なルールを学習する手法である。教師あり学習によるアプローチの代表的な研究として、保木は、コンピュータ将棋プログラムである*Bonanza*を提案している[9]。*Bonanza*は、プロ棋士の棋譜6万局のデータを教師とし、将棋の局面における評価関数を自動学習することで、従来手法よりも良い振る舞いを得ることに成功している。この手法は*Bonanza*メソッドと呼ばれ、多くのコンピュータ将棋プログラムで採用されている画期的な手法である[12]。将棋のように、強い人間プレイヤーの膨大な棋譜データが用意できる場合には、教師あり学習による振る舞い獲得は有効である。

経路探索は、ゲーム木におけるスタートからゴールまでの、最小コストとなる経路を探索する手法である。経路探索によるアプローチの代表的な研究として、Baumgartenは、2009年のMario AI Competitionにおいて、A*アルゴリズムに基づいたNPCを構築し優勝している[10]。Mario AI Competitionとは、“*Infinite Mario Bros.*”（ランダムに生成されるステージを制限時間内に攻略する、「スーパーマリオワールド」のようなアクションゲーム。）を対象としたNPCの評価コンテストである。BaumgartenのNPCは、マリオや敵の動きを事前に解析し、A*アルゴリズムを用いた経路探索によって、ステージをほぼ最適解で攻略することが可能となっている。

強化学習は、自身の振る舞いの試行錯誤を繰り返すことで最適な振る舞いを獲得する手法である。強化学習によるアプローチの代表的な研究として、藤田らは、カードゲームの*Hearts*を題材とし、Q学習を用いて、NPCの振る舞い獲得に成功している[11]。藤田らは、巨大な状態空間となること、相手の所持するカードを観測できないこと、4人対戦のゲームであること、の3つを*Hearts*における学習の困難性と考察している。その上で、解決手法として、パーティクルフィルタによるサンプリング、相手の行動予測器、現在の戦局を評価する状態価値関数、ゲームの特徴に基づく次元圧縮を提案し、困難性の解決を図っている。実験の結果、人間の熟達者よりも優れた振る舞いを得ることに成功している。

これらの手法を用いて獲得されたNPCは、極めて最適であるが故に、人間にとっては機械的と感じる振る舞いを表出してしまう。そのため、エンタテインメント性の向上という視点に立った場合、人間プレイヤーの代替として扱うことは憚られる。ゲームAI領域では、人間プレイヤーが強いNPCに勝てなくなる日も近いと考えられており、人間らしいNPCの構築が最重要課題となりつつある。

2.2 人間らしいNPCを実装した研究

人間らしいNPCを実装した関連研究として、Schrumらは、2012年のThe 2K BotPrizeにおいて、大会史上初となる、人間よりも人間らしいと評価されるNPCの構成に

成功している [2]. The 2K BotPrize とは、FPS (一人称視点シューティングゲーム) を対象とした、NPC の人間らしさを競う評価コンテストである。人間プレイヤーの振る舞いをトレースしたデータベースを基に、人間らしいと思われる振る舞いを決定論的に定義し、ニューラルネットにおける制約として適用している。その結果、対戦相手の人間プレイヤーから「人間らしい」と評価される NPC の振る舞いが獲得できている。

池田らは、コンピュータ囲碁を対象に、既存の強い NPC に意図的に人間らしいミスをさせることで、手加減と思われない程度の「強くなさ」を実現するための初期的検討を実施している [5]。現在の局面における予測勝率と候補手の選択確率を用いた形勢の制御、楽観派や悲観派といったプレイスタイルによる獲得戦略の分析をしており、ゲームのレベルデザインにおける一アプローチを提案している。

上記の手法は、人間らしいと思われる振る舞いを、開発者が恣意的に定義したものである。そのため、振る舞い獲得における作業負荷の軽減や汎用性の確保は実現されていない。

3. 生物学的制約

3.1 生物学的制約の導入

従来手法では、人間が生得的・遺伝的にもつ特徴や性質から生じる生物学的制約を無視しているために、機械的と感じる振る舞いが表出していると考えられる。コントローラ操作の反応速度が速すぎる、コントローラのボタンの入力が正確すぎる、常に一定の行動のみを正確に繰り返すといった、人間プレイヤーでは実現不可能な振る舞いが表出するケースもある。また、レベルデザインを重視しすぎると、ゲームの途中から急に弱くなる、あからさまなコントローラ操作のミスをするといった、プレイスタイルの統一性が崩壊した振る舞いが表出するケースもある。これらの振る舞いは、「相手がいんちきをしているのでは」、「本当に自分の力で勝ったのか」という疑念を生むため、人間プレイヤーのゲームへのモチベーションを削ぐ要因となっている。

本研究で NPC に導入する『生物学的制約』は、人が生得的に持っている制約や欲求である。人間プレイヤーがゲームをするときには必ず生物学的制約が生じており、人間プレイヤーはその制約下でキャラクタを操作している。よって、生物学的制約下で操作されたキャラクタの振る舞いとは、人間プレイヤーにとっては最も一般的で見慣れた振る舞いであると言えるだろう。Zajonc は、繰り返し体験した刺激に対して好意度や印象が高まる現象として単純接触効果 [13] を提唱している。単純接触効果はあらゆる事象において生じるとされていることから、人間プレイヤーに好印象を与えるような人間らしい振る舞いの実現には、生物学的制約の導入が有効である可能性が高い。

生物学的制約としては「身体的な制約」と「生き延びるために必要な欲求」を機械学習や経路探索の制約条件として課する。生物学的制約と、戦略獲得手法との融合により、人間プレイヤーで実現不可能な振る舞いを排除し、かつ、「わざとらしさ」や「明らかな弱さ」を露呈しない、「戦略の統一性」が再現されると考えられる。身体的な制約に関する研究例として、Cabrera らは人間の指先による倒立棒の制御実験を [6]、大平らは人間の直立姿勢の制御実験を実施している [7]。人間の行動制御には「ゆらぎ」「遅れ」「疲れ」といった制約が生じるが、人間は訓練によってこれらの制約を意識的もしくは無意識的に考慮し、安全性とパフォーマンスを両立させる行動制御が獲得できると提唱している。また、生き延びるために必要な欲求について、Maslow は人間の欲求を 5 段階の階層構造で理論化した「自己実現理論」を提唱している [8]。原始的な欲求に近い階層から順に、1) 生理的欲求、2) 安全の欲求、3) 所属と愛の欲求、4) 承認 (尊重) の欲求、5) 自己実現の欲求、と人間の欲求を分類している。そして、「人間は自己実現に向かって絶えず成長する生きものである」という仮定の下、「訓練」による知識の定着や、「挑戦」による不満の解消といった行動の動機は、5) 自己実現の欲求に帰結すると考えられている。

3.2 生物学的制約の定義

本研究では、『生物学的制約』を「身体的な制約：“ゆらぎ” “遅れ” “疲れ”」、「生き延びるために必要な欲求：“訓練と挑戦のバランス”」として、以下のように定義する。

(1) センサ系、運動系における「ゆらぎ」

人間プレイヤーは、操作対象や敵オブジェクト等の位置 (座標) を正確に観測し認識することは難しく、必ず誤差 (ゆらぎ) が生じる (見間違い、操作ミスなど)。そこで、NPC が観測する操作対象の現在位置やゲームの局面情報に対し、ガウスノイズを付与することで「ゆらぎ」を再現する。

(2) 知覚から運動制御に至る「遅れ」

人間プレイヤーは、ゲームの局面を認識してから、実際に動作するまでに遅れが発生する (眼と手の協応動作における遅延など)。そこで、NPC が観測する操作対象の現在位置やゲームの局面情報を、数百ミリ秒過去の情報にすることで「遅れ」を再現する。

(3) キー操作の「疲れ」

人間プレイヤーは、ゲームのコントローラのキー操作を、極めて短時間で何度も、または、長時間連続して実施すると疲れが生じる (ボタン連打、単調な操作の繰り返しなど)。そこで、振る舞いを学習する際に、NPC にキー操作変更による負の報酬を与えることで「疲れ」を再現する。

(4) 「訓練と挑戦のバランス」

人間プレイヤーは、同じ行動を繰り返す事で「訓練」す

一方で、同じ行動の結果に飽きたり、その行動で失敗を繰り返したりすると、飽きや失敗を解消するための新奇な行動に「挑戦」する。そこで、失敗を繰り返しているゲーム局面では、新奇な行動に挑戦する傾向を高め、逆に、失敗をほとんどしないゲーム局面では、同じ行動を繰り返して訓練する傾向を高めることで「訓練と挑戦のバランス」を再現する。

4. 振る舞い獲得の手法

4.1 生物学的制約の導入

ビデオゲームにおいては、教師となるプレイデータが大量に用意できないため、非事例参照型の手法である Q 学習 [14] (強化学習) と A* アルゴリズム [15] (経路探索) を用いることにする。

Q 学習では、ゲームのある局面における最適な行動を数式 1 で、また、NPC が行動した際の Q 値の更新を数式 2 で算出する。

$$\operatorname{argmax}_{a_t} Q(s_t, a_t) \quad (1)$$

$$Q(s_t, a_t) = (1 - \alpha)Q(s_t, a_t) + \alpha((r + \gamma \max_p Q(s_{t+1}, p)) \quad (2)$$

Q 学習では、時刻 t におけるゲーム局面 s_t において、Q 値 $Q(s_t, a_t)$ が最も大きくなる行動 a_t が最適行動であると出力される。また、NPC の行動選択手法としては ϵ -greedy 法を用いる。生物学的制約の導入に関して、「ゆらぎ」と「遅れ」は、数式 2 の $Q(s_t, a_t)$ の計算の際に、数百ミリ秒過去の位置情報や局面情報にガウスノイズを付与したものを s_t とすることで実現する。「疲れ」は、数式 2 の報酬 r にキー操作変更による負の報酬を与えることで実現する (報酬 r の詳細については節 4.3 で述べる)。「訓練と挑戦のバランス」は、ランダム行動選択確率 ϵ の設定において、失敗を繰り返しているゲーム局面 s_t では大きな値を設定することで、新奇な行動に挑戦する傾向を高め、逆に、失敗をほとんどしないゲーム局面 s_t では小さな値を設定し、同じ行動を繰り返して訓練する傾向を高めることで実現する。

つぎに、有名な最短経路探索手法である A* アルゴリズムにおいて、生物学的制約の導入を試みる。節 2.1 で述べたとおり、A* アルゴリズムはアクションゲームにおいて、ほぼ最適解を獲得した実績のある手法である。A* アルゴリズムでは、以下の式によりゲーム木の経路のコストを算出する。

$$f^*(n) = g^*(n) + h^*(n) \quad (3)$$

生物学的制約の導入に関して、「ゆらぎ」と「遅れ」は、数百ミリ秒過去のキャラクターの位置情報に対してガウスノイズを付与し、その座標をスタートノードとすることで実現する。「疲れ」は、極めて短時間でのキー操作の変更を禁止することで再現する。「訓練と挑戦のバランス」は、学習フェーズを持たない A* アルゴリズムでは実現不可能であ



図 1 “Infinite Mario Bros.” のゲーム画面

るため対象外とする。

4.2 “Infinite Mario Bros.” の仕様

NPC の振る舞いを獲得するにあたり、対象とするゲームとしては、1) 同じ局面を何度も再現できる、2) ゲームの明確な目標が設定できる、かつ、3) ビデオゲームを代表する有名なゲームである必要がある。

本研究では、上記条件を満たし、かつ、ゲームの仕様やゲーム環境パラメータが公開されている、“Infinite Mario Bros.” を対象とし、振る舞いの獲得と、その比較検証、主観評価を実施する。“Infinite Mario Bros.” は、世界的に有名なゲームである“スーパーマリオワールド”を模したアクションゲームであり、そのゲーム画面を図 1 に示す。また、“Infinite Mario Bros.” における仕様は以下のとおりである。

- ステージの自動生成

事前に与えた疑似乱数のシード値に従って無限にステージが生成される。

- NPC の操作キャラクタ (マリオ)

NPC はマリオ (図 1 中央) を操作する。NPC によるマリオの操作はコントローラのキー入力 (LEFT, RIGHT, DOWN, SPEED, JUMP) により行う。毎フレームのキーの押下状態により、マリオは対応した行動を行う (毎秒 24 フレームで動作)。

- 敵キャラクタ

ステージには数種類の敵が登場し、敵はそれぞれ独自の動作をしている。NPC は、これらの敵を避けて進むか、倒して進むかを決定しなければならない。

- スコアの獲得

マリオが死亡する、または、設定された制限時間に達すると攻略は終了し、スコアを獲得する。スコアは Mario AI Competition で規定されている評価関数で計算され、ステージを攻略した距離に応じてスコアが上昇する。

- NPC の観測情報

NPC は、マリオの座標、マリオの状態、画面内の敵の種類および座標、ステージの地形座標を観測すること

ができる。NPC の観測する地形座標は、ステージに配置されているブロックのうち、画面内にある 22×22 のブロックの配置座標となる。NPC は毎フレーム観測情報を受け取り、マリオの行動制御を行うためのキー入力を返す必要がある。

4.3 “Infinite Mario Bros.”での振る舞い獲得

Q 学習での “Infinite Mario Bros.” の扱い方として、まず、現実的な時間で学習が収束し NPC の振る舞いを獲得できるように、ゲーム局面 s の次元を圧縮する方法を述べる。ゲームの攻略にあたって重要となる情報を削減すると、学習が正常に動作しなくなってしまうことを考慮し、NPC の観測できるゲーム局面 s を以下のとおりに圧縮する。

- **マリオを中心に 7×7 ブロックの地形と敵配置**

NPC が観測可能な地形座標と敵座標は画面を 22×22 ブロックに分割したものである。しかし、1 フレームあたりのマリオの移動距離は小さく、画面内全ての地形座標や敵の配置がマリオの行動に影響することはない。そこで、学習に使用する地形情報と敵の配置は、マリオを中心とした 7×7 ブロックとする。これにより、ゲーム局面 s の次元数は大幅に削減される。

- **マリオの進行方向**

敵や地形との関係性を把握するための重要な要素であるため、NPC は 8 方向+停止の 9 次元としてマリオの進行方向を把握しておく必要がある。

- **「でかマリオ」か「ちびマリオ」か**

「でかマリオ」でダメージを受けた場合は「ちびマリオ」に変化するだけで攻略を続行できるが、「ちびマリオ」でダメージを受けた場合は死亡となる。より長く攻略を進めるうえで重要な要素であるため、NPC はマリオの状態を把握しておく必要がある。

- **マリオが地上にいるか**

マリオは、地上にいる場合はダッシュやジャンプができるが、空中にいる場合はできない仕様である。マリオが地上にいるかどうかは、行動選択にあたって重要な要素であるため、NPC はマリオが地上にいるかどうかを把握しておく必要がある。

次に、Q 学習における選択可能な行動 a の設定方法について述べる。マリオの行動は、コントローラのキー入力によって決定される。そこで、マリオの行動制御に影響があるキー入力の組み合わせ（例：RIGHT+SPEED+JUMP=右に走りジャンプ）である全 11 パターンを行動 a として定義する。

続いて、Q 学習における報酬 r の設定方法について述べる。敵を可能な限り避け、ステージをより早く、より遠くまで攻略するためには、ステージを早く攻略することに対して正の報酬を与え、逆にダメージを受ける、死亡するといった、攻略を阻害する要因に対して負の報酬を与えるこ

とが望ましい。また、キー操作変更による疲れを実現するため、キー操作を変更した場合は負の報酬を与える必要がある。そこで、報酬 r を以下のとおりに設定する。

$$r = distance + damaged + death + keyPress \quad (4)$$

数式 4 において、 $distance$ は行動によって進んだ距離であり、そのまま正の報酬とする。 $damaged$ は行動によってダメージを受けた場合に与える負の報酬、 $death$ は行動によって死亡した場合に与える負の報酬である。また、 $keyPress$ は前フレームから行動を変更した場合に与える負の報酬である。予備実験の結果、本研究における $distance$ は進んだ距離 $\times 2.0$ 、 $damaged$ は -50.0 、 $death$ は -100.0 、 $keyPress$ は -5.0 とした。

最後に、A* アルゴリズムにおけるゲーム木の作成方法と、経路のコスト算出について述べる。節 2.1 で述べた A* アルゴリズムに基づく NPC[10] を参考にする。スタートノードを現在のマリオの位置座標（ただし「ゆらぎ」や「遅れ」が付与された座標）、ゴールノードを画面の右端とし、マリオが取り得る行動によってゲーム木を作成する。 $g^*(n)$ としては、スタートノードから現在ノードまでの時間を、 $h^*(n)$ としては、現在ノードから画面の右端に到達するまでの推定時間を算出している（詳細は [10] を参照）。

5. シミュレーションと主観評価実験

5.1 計算機シミュレーション

Q 学習によって生成された NPC について、正常に学習が進んでいるかどうかを確認するため、「ゆらぎ」と「遅れ」のパラメータを変更し、獲得スコアの推移を調べる。「疲れ」は報酬の一部であり、学習性能に大きな影響を与えないため、今回は考慮外とした。用意した 3 つの Q 学習エージェントを表 1 に示す。毎試行ランダム生成されるステージを対象として学習試行を行い、学習試行回数は 10 万ゲーム、200 ゲームごとの獲得スコアの平均をとる。獲得スコアの比較対象として、2009 年の Mario AI Competition において優勝を収めた、Baumgarten の A* エージェント [10] を用いた。学習シミュレーションの結果を図 2 に示す。図 2 から、どのパラメータセットにおいても正常に学習が進んでいることが示された。

表 1 学習性能を比較する Q 学習エージェント

Table 1 Three video game agents for evaluating performance.

Q 学習エージェント	Q 学習 1	Q 学習 2	Q 学習 3
生物学的制約の導入	無し	あり	あり
ゆらぎ (block)	0	0.5	1.0
遅れ (frame)(秒)	0(0)	6(0.25)	12(0.5)
学習率 α	0.2	0.2	0.2
割引率 γ	0.9	0.9	0.9
ランダム選択確率 ϵ	0.05	0.05	0.05

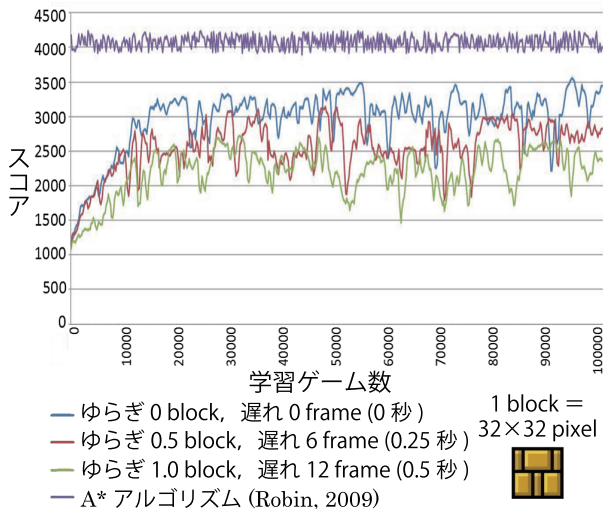


図 2 異なるパラメータセットで学習した際の獲得スコアの推移
Fig. 2 Scores with different biological constraints.

5.2 獲得された振る舞い

本研究の振る舞い獲得手法を用いて獲得された NPC の振る舞いを図 3 に示す。生物学的制約の導入の有無によって、表出した振る舞いの特徴に以下のような差異があった。

触れることができない敵を回避する場面 (図 3 上段)

- 導入なし (左) : 最小限のジャンプ, かつ, ノンストップで攻略
- 導入あり (右) : 大きくジャンプし, 途中で一瞬止まるような行動をしつつ攻略

5 体の敵が段差の上に存在する場面 (図 3 中段)

- 導入なし (左) : 正確な行動制御で敵が大量に存在する区間を攻略
- 導入あり (右) : 区間の手前で待機し, 安全に進める状態に変化してから攻略

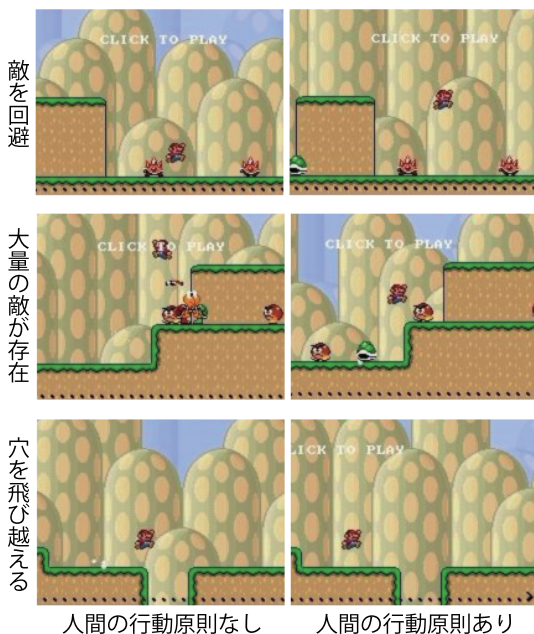


図 3 導入なし (左列) と導入あり (右列) での振る舞いの比較

穴を飛び越える場面 (図 3 下段)

- 導入なし (左) : 穴に落ちる寸前のところから最小限のジャンプで攻略
- 導入あり (右) : 穴の少し手前から大きくジャンプし余裕を持って攻略

これらの振る舞いの特徴は, Q 学習によって生成された NPC (以降, Q 学習エージェント), A* アルゴリズムによって生成された NPC (以降, A* エージェント) の双方において共通であった。以上の結果から, 生物学的制約の導入なしでは, パフォーマンスのみを重視しているが, 導入有りでは, 安全性も考慮した振る舞いが獲得できているといえる。

5.3 主観評価の実験計画

生物学的制約を導入したエージェントによって獲得された振る舞いが, 本当に人間らしいかどうかを検証するため, 20~24 歳の男女 20 名 (男性 13 名, 女性 7 名) を対象に主観評価実験を実施した。被験者 20 名における, 横スクロール型マリオのプレイ時間の累計は平均 $\mu = 34$ 時間, 標準偏差 $\sigma = 29$ 時間であった。そこで, 本実験においては, 横スクロール型マリオの熟練度を 3 つのグループに分類した。横スクロール型マリオのプレイ時間が 5 時間 ($\mu - \sigma$) 未満の被験者 4 名を「初級者」(うち, 3 名はプレイ時間が 0 時間の初心者), 63 時間 ($\mu + \sigma$) 以上の被験者 2 名を「上級者」, 5 時間以上 63 時間未満の被験者 14 名を「中級者」と定義した。

実験手続きは以下の通りである。まず, 被験者に「ブロック, アイテム, コイン等は無視して, ステージの先に進め」と教示し, “Infinite Mario Bros.” を 10 回プレイ (1 プレイ 25 秒) させた。次に, プレイ動画を 2 つずつ比較させ「どちらのマリオが人間らしいプレイか」を 7 段階で評価させた。最後に, プレイ動画を 1 つずつ見せ「どのような振る舞いが人間らしい (人間らしくない) と感じたか」を自由記述で回答させた。

実験に使用したプレイ動画を表 2 に示す。本実験では, Q 学習エージェントによるプレイ動画を 3 つ, A* エージェントによるプレイ動画を 2 つ, 人間が操作したプレイ動画を上記熟練度を考慮して 3 つ用意した。Q 学習エージェントに関しては, 生物学的制約の導入ありと導入無しの 2 つに加えて, 訓練をせず失敗に対する挑戦のみを実施するエージェントも用意した。この Q 学習エージェントにおけるランダム選択確率 ϵ は 0, 失敗を繰り返しているゲーム局面での ϵ は 0.2 と設定した。人間の操作者に関しては, 初級者動画は横スクロール型マリオのプレイ時間が 5 時間の人間プレイヤー, 中級者動画は 50 時間の人間プレイヤー, 上級者動画は 200 時間の人間プレイヤーとした。また, 敵, 土管, 穴といった障害物の有無や, マリオが敵に接触しダメージをうけるシーンが, 人間らしさの評価に大きく影響

表 2 プレイ動画のラベルと内容

ラベル	操作者	生物学的制約	再生時間	スコア
[強化, 無し]	強化学習 (NPC)	導入なし	10.62 秒	5448
[強化, 導入]	強化学習 (NPC)	導入あり	14.25 秒	4069
[強化, 導入, 挑戦のみ]	強化学習 (NPC)	導入あり (挑戦のみ)	15.57 秒	3458
[探索, 無し]	経路探索 (NPC)	導入なし	7.29 秒	7926
[探索, 導入]	経路探索 (NPC)	導入あり	9.34 秒	3118
[中級者]	中級者 (人間)	-	10.08 秒	6031
[初級者]	初級者 (人間)	-	14.25 秒	3644
[上級者]	上級者 (人間)	-	7.68 秒	7371

を与えると考えられる。そこで、全ての動画でプレイ区間を統一し、マリオが敵に接触しダメージを受けたプレイ区間は不採用とした。これ以降、プレイ動画を表 2 のラベル名で表記する。

5.4 分析手法と結果

本実験では、ランダムな順序で呈示される 2 つのプレイ動画を比較し、人間らしさについて 7 段階で評価する。統計的分析手法としてシェッフェの一対比較法 [16] (中屋の変法) を使用し、分散分析で主効果の有無を確認する。その後、ヤードスティック法によりプレイ動画の嗜好度を一本の直線上にプロットし、動画同士の相対的な関係性と、信頼区間について検討する。本実験では、NPC における生物学的制約の導入の有無による比較、NPC と人間プレイヤーとの比較に焦点を当てるため、Q 学習エージェントと A* エージェントを分けて分析することとした。

図 4 は、人間らしさに関する相対的嗜好度をプロットしたものである。上の直線は Q 学習エージェントと人間プレイヤーの比較、下の直線は A* エージェントと人間プレイヤーの比較である。まず、Q 学習エージェント同士の比較結果を述べる。生物学的制約を導入した [強化, 導入] (相対的嗜好度: 0.66) は、生物学的制約を導入していない [強化, 無し] (相対的嗜好度: 0.29) と比較して、人間らしいという結果が得られた。しかしながら、相対的嗜好度の差 ($0.66 - 0.29 = 0.37$) が 95% 信頼区間である 0.48 より小さ

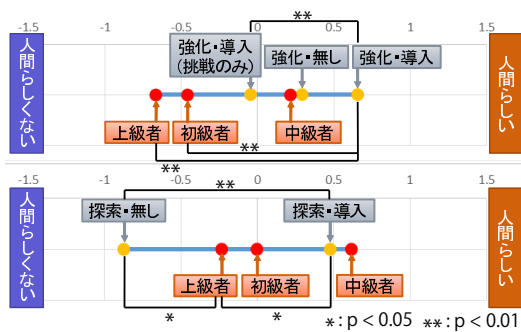


図 4 人間らしさに関する相対的嗜好度

いため、5%水準の有意差は認められなかった。この結果を、以降 (差: $0.37 < 95\%$ 信頼区間: 0.48) と表記する。次に、A* エージェント同士の比較結果を述べる。生物学的制約を導入した [探索, 導入] は、生物学的制約を導入していない [探索, 無し] と比較して、1%水準で有意に人間らしいという結果が得られた (差: $1.35 > 99\%$ 信頼区間: 0.72)。最後に、NPC と人間プレイヤーの比較結果を述べる。生物学的制約を導入した Q 学習エージェント [強化, 導入] は、人間プレイヤーの [初級者][中級者][上級者] より人間らしいという結果が得られた。また、生物学的制約を導入した A* エージェント [探索, 導入] は、人間プレイヤーの [初級者][上級者] より人間らしいという結果も得られた。ただし、有意差が認められたのは、[強化, 導入] と [初級者] (差: $1.12 > 99\%$ 信頼区間: 0.58), [強化, 導入] と [上級者] (差: $1.33 > 99\%$ 信頼区間: 0.58), [探索, 導入] と [上級者] (差: $0.71 > 95\%$ 信頼区間: 0.59) のみであった。

6. 考察

主観評価実験の結果から、生物学的制約を導入することで、『人間らしい』と解される NPC を自律的に構成できることが示された。では、「どのような振る舞いが人間らしいのか」について、主観評価実験の結果 (図 4 と自由記述質問の回答) から考察していく。

[強化, 導入] は全動画中で最も人間らしいと評価されている。また、[探索, 導入] は [探索, 無し] と比較すると人間らしい (1%の有意水準で有意差あり) という評価である。自由記述質問では、人間らしいと感じる理由として「敵や穴を飛び越える時に一瞬後ろを向く」、「敵や穴を大きく飛び越える」、「ときどき不必要な行動をとる」という回答があった。この結果から、「ためらい」や「余裕」、「熟慮 (試行錯誤)」を感じさせる要素として、『生物学的制約』を導入することの妥当性が示された。

[上級者] は人間プレイヤーの操作であるにもかかわらず、人間らしくないと評価されている。また、ほぼ最適解である [探索, 無し] は [上級者] よりもさらに人間らしくない (5%の有意水準で有意差あり) という評価である。自由記述質問では、人間らしくないと感じる理由として「敵や穴をギリギリまで避けない」、「無駄な行動が一切ない」、「動きが一定である」という回答があった。この結果は、「過度に最適化された振る舞いは人間らしくない」ことを意味する。節 2.1 で述べた、強い NPC を人間プレイヤーの代替として扱うことができない根拠が示された。また、このことから、[強化, 導入] と [強化, 無し] で有意差が認められていない理由も説明できる。[強化, 無し] は [上級者] や [探索, 無し] のスコアに遠く及んでおらず (表 2)、最適化された振る舞いの獲得に至っていないためと考えられる。Q 学習エージェントの改良には、節 4.3 で述べたゲーム局面 s の観測情報を拡張する必要がある。

[強化, 導入, 挑戦のみ] は, 生物学的制約を導入しているにもかかわらず, 比較的人間らしくないという評価である. 自由記述質問では「段差や土管にぶつかってからジャンプする振る舞いが人間らしくない」という回答があった. この動画は, スコアがかなり低く, その「たどたどしい」振る舞いは, コントローラ操作やゲームルールに慣れていない, あたかもゲーム初心者の操作のようであった. この結果は, 「初心者相当の下手すぎる振る舞いは人間らしくない」ことを意味する. また, 「訓練と挑戦のバランス」を変化させることで, 人間の熟達過程を再現できる可能性が示された.

人間は誰しも, 身体的な制約を生得的に持っており, また, 生きていくためには訓練や挑戦といった自己実現の欲求が必要不可欠である. 人間は, これらの制約や欲求が考慮された振る舞いからは, 「ためらい」や「余裕」, 「熟慮 (試行錯誤)」といった感情を想起し, その結果, 人間らしい振る舞いであると解釈していると言える. 逆に, それらの制約や欲求を無視した「過度に最適化された振る舞い」からは, 情緒的な反応が惹き起こされることがなく, 人間らしさも感じないのであろう. もちろん, 「人間らしさの評価基準は被験者間で異なるのではないか」という疑問もある. 被験者の横スクロール型マリオの熟練度や, ゲームに対するプレイスタイルにより, 被験者を群分けすることで, 被験者間の評価基準の差異を検証する必要がある. その結果, 人間が人間らしさを解釈するための評価基準の策定が可能であると考えられる.

7. おわりに

ビデオゲームにおけるプレイフィールドの向上には, 人間らしい NPC の実装が必要不可欠であり, その自律的獲得には, 従来, ゲームジャンルやゲームタイトルに合った人間らしさの解析が必要であった. 本研究では, 生物学的制約を導入することで, 人間らしい NPC の振る舞いを自動獲得できることが示された. 生物学的制約を課した強化学習や経路探索により, 「人間プレイヤーがゲームをしている」かのような振る舞いが表出され, また, 主観評価実験により, それらの振る舞いが人間プレイヤーよりも人間らしいことを示した. 生物学的制約は, 開発者のヒューリスティックや, 人間らしさの解析に依拠しない要素である. そのため, あるゲーム状況を入力とし, そのゲーム状況で最適な行動を出力する必要があるゲームであれば, ゲームジャンルや振る舞い獲得の手法を問わず, 人間らしい NPC の振る舞いを獲得できると考えられる

本研究の振る舞い獲得手法を使用することで, 人間らしい NPC を実装したいゲームプログラマにとって, 1) ヒューリスティックの導入に係る煩雑な作業負荷 (開発コスト) を削減できる, 2) 人間が持つ生物学的制約であるため生理学的・心理学的知見に基づいて設定できる, 3) 様々なゲー

ムジャンル, 様々な機械学習手法に対しても, 汎用的に導入できる, という3つのメリットがある. また, 人間らしい NPC が実現されることで, 人間プレイヤーの満足感の確保やエンタテインメント性の持続といった, ユーザエクスペリエンスの向上につながると考えられる. 今後の展望としては, 被験者を群分けし人間らしさの評価基準を特定する, アクションゲーム以外のジャンルにも生物学的制約の導入を試みる.

参考文献

- [1] エンターブレイングローバルマーケティング局: ファミ通ゲーム白書 2013, エンターブレイン (2013).
- [2] Schrum, J., Karpov, I. V. and Miikkulainen, R.: Human-like Behavior via Neuroevolution of Combat Behavior and Replay of Human Traces, *2011 IEEE Conference CIG' 11*, pp. 329-336 (2011).
- [3] Soni, B. and Hingston, P.: Bots Trained to Play Like a Human are More Fun, *2008 IEEE International Joint Conference on Neural Networks*, pp. 363-369 (2008).
- [4] Ortega, J., Shaker, N., Togelius, J. and Yannakakis, G. N.: Imitating human playing styles in Super Mario Bros, Vol. 4, pp. 93-104 (2013).
- [5] 池田心, Viennot, S.: モンテカルロ基における多様な戦略の演出と形勢の制御〜接待基 AI に向けて〜, *GPW2012*, pp. 47-54 (2012).
- [6] J.L.Cabrera and J.G.Milton: On-Off Intermittency in a Human Balancing Task, *Physical Review Letters*, Vol. 89, No. 15 (2002).
- [7] 大平徹, 保坂忠明: 不安定な状況でのノイズと遅れの役割と制御への考察, 交通流のシミュレーションシンポジウム, pp. 19-22 (2004).
- [8] Maslow, A. H.: A Theory of Human Motivation, *Psychological Review*, Vol. 50, pp. 370-396 (1943).
- [9] 保木邦仁: 局面評価の学習を目指した探索結果の最適制御, *GPW2006*, pp. 78-83 (2006).
- [10] Togelius, J., Karakovskiy, S. and Baumgarten, R.: The 2009 Mario AI Competition, *Evolutionary Computation (CEC) 2010 IEEE*, pp. 1-8 (2010).
- [11] Fujita, H. and Ishii, S.: Model-based reinforcement learning for partially observable games with sampling-based state estimation, *Neural Computation*, Vol. 19, pp. 3051-3087 (2007).
- [12] Hoki, K. and Kaneko, T.: The Global Landscape of Objective Functions for the Optimization of Shogi Piece Values with a Game-Tree Search, *Advances in Computer Games 2012, Lecture Notes in Computer Science*, Vol. 7168, pp. 184-195 (2012).
- [13] Zajonc, R. B.: Attitudinal Effects Of Mere Exposure, *Journal of Personality and Social Psychology*, Vol. 9, pp. 1-27 (1968).
- [14] Watkins, C.: Learning from Delayed Rewards, *PhD thesis, Cambridge University, Cambridge, England*. (1989).
- [15] Hart, P. E., Nilsson, N. J. and Raphael, B.: A Formal Basis for the Heuristic Determination of Minimum Cost Paths, *IEEE Transactions on Systems Science and Cybernetics*, Vol. 2, pp. 100-107 (1968).
- [16] Scheffe, H.: An Analysis of Variance for Paired Comparisons, *Journal of the American Statistical Association*, Vol. 47, No. 259, pp. 381-400 (1952).