DA3:マルチエージェント深層強化学習における 協調行動の解釈性確立と対ノイズ性能の検証

元川 善就^{1,a)} 菅原 俊治^{1,b)}

概要:本稿では、マルチエージェント深層強化学習におけるノイズの影響度を学習し、自律的に抑制する distributed attentional actor architecture model for multi-agent system (DA3)を提案する. ノイズは多 くの分野に存在し、長年の間その特性の解明や抑制手法が考案されてきた. 特に、周囲の限られた情報を もとに他のエージェントとの協調・調整,さらに相互協力を前提とした互恵行動を獲得する必要のあるマ ルチエージェントシステムにおいても、ノイズ抑制や排除は重要な研究テーマである. 本提案である DA3 は attention mechanism を内包しているため、エージェントは観測情報を取捨選択しながら環境に適応す ることが可能である. 実験結果から、観測情報内のノイズや協調行動に無関係な第三者的エージェントに 惑わされることなく、他エージェントとの協調行動を獲得し、attention mechanism を持たないベースラ イン (ここでは DQN) よりも高い学習性能を達成することを示す. また、DA3 内部に存在する attentional weights を解析することで、これまでブラックボックスとされてきたエージェントの行動決定過程における ノイズの影響度などを把握できることを示す.

キーワード:マルチエージェント深層強化学習,分散型自律学習,注意機構,協調行動解釈性,ノイズ抑制

1. まえがき

ノイズ抑制とは、様々な領域において取り組むべき重要 な研究テーマである.長年、多くのノイズ特性の研究が続 けられており、その特性に合わせたノイズ抑制手法が考案 されてきた.たとえば、ロボティクスにおける微振動、電 気回路における熱雑音、物理学における不確実性など研究 領域ごとにノイズの定義や性質は異なるが、実世界では 無視できないという点では共通であり、これは(深層)強 化学習でも同様である. Wang ら [23] によると、強化学習 におけるノイズは複数のカテゴリーに分類できる.例え ば, inherent noise (内部雑音) はエージェントが情報を観 測する際に,光線や熱雑音などの外的要因によって発生す るノイズを指す. 同様に, 限られた局所情報から協調性を 創発・学習するマルチエージェントシステム(multi-agent systems, MAS) あるいはマルチエージェント深層強化学習 (multi-agent deep reinforcemnet learning, MADRL) にお いても、システム全体の信頼性・堅固性を向上させるため にノイズ抑制は非常に重要な課題となる.しかし,発生す るノイズ特性は想定する応用環境によって多種多様である

^{b)} sugawara@waseda.jp

ため,既存研究のほとんどはノイズのない理想的な環境を 仮定して学習エージェントのモデル化や協調性の獲得を検 証することが多かった.

シングルエージェント深層強化学習の研究ではノイズを 想定しているものもあるが([4,12,19–21,25]など),多く の場合 surrogate loss function を利用するため,既存の強 化学習アルゴリズムへの統合は容易ではない.MADRLに おいてはノイズを想定した研究は少なく,存在しても基本 的にはエージェント間で共有する集中型ネットワークモデ ルでの学習[5,7,17]を前提としている.また実世界で発 生するほとんどのノイズの性質や発生場所は未知であるた め,エージェントは学習前からノイズ特性が既知と想定す る先行研究では実用性に欠ける.また,エージェントの行 動場所やその処理能力・目的などによりノイズの影響度も 異なるため,エージェント毎にノイズの特性や影響を自律 的に把握し,それに惑わされず分散的に学習できる仕組み が必要である.

そこで本研究では、各エージェントが学習に不要なノイズの存在を同定し、そのノイズ抑制を考慮したマルチエージェント深層強化学習手法である distributed attentional actor architecture for a multi-agent system (DA3)を提案する. DA3 は MAT-DQN [10] と呼ばれる、attention mechanism を活用してエージェントの行動解釈性の確立

¹ 早稲田大学 基幹理工学研究科 知識ソフトウェア研究室, 東京

^{a)} y.motokawa@isl.cs.waseda.ac.jp

を目指した学習モデルの拡張であり、これをノイズの同定 に適用拡大している. DA3 の特徴として, saliency vector と呼ぶパラメータの集合をネットワーク内部に導入し、行 動決定過程における観測情報を超次元空間で表現する.本 手法の重要な利点として、DA3を内部に保有するエージェ ント(以下 DA3 エージェント)へ観測情報のどの領域が ノイズであるか、或いはどのエージェントが協調関係であ るかという情報を事前に与える必要がない点である. 例え ば、あるエージェントの故障、異なる目的を持つエージェ ントとの共存, エージェントが持つセンサーの不良や, 不 良ではなくても距離に応じてセンサのデータにノイズが含 まれ信頼性成果などの不確実性が発生する. このようなノ イズを含んだ観測情報は直接 DA3 エージェントに与えら れ, DA3 エージェントは観測情報のどの部分が(非)信頼 できるかを自ら学習する. この観測情報は saliency vector として表現され、その表現に基づいて DA3 エージェント は行動を決定し学習する.

本研究ではノイズ環境下における DA3 エージェントの 学習性能を検証・評価するために, object collection games を使用し、もっとも単純な vanilla deep Q-network (以下, DQN)との学習性能を比較する.実験結果から、DA3エー ジェントはノイズを同定するとともに、そのノイズの混在 比率に応じて信頼度を決定し、ノイズに惑わされることを 最小限に抑え、その影響を下げることが分かった.そして、 比較手法よりも高い学習性能を達成できることが分かった. 実際に, DA3 ネットワーク内の attention mechanism が学 習した観測情報に対する attentional weights を分析した結 果,エージェントの行動決定に強く関連する観測情報には より大きな重みを与え、無関係なエージェントの情報やノ イズを多く含む情報にはより小さい重みを与えるように学 習していることを確認した. この解析結果は, DA3 エー ジェントはノイズに混乱することなく行動を決定できたか を可視化することも可能としている.

2. 関連研究

強化学習の分野ではノイズを考慮した研究は少ない [9,11,15,18].たとえば Raeisy ら [15] と Sabzevari ら [18] は active noise control (ANC)問題を強化学習の枠組みで 検証しているが、この ANC とは電気信号中のノイズを アクチュエータの集合で緩和する技術であり [6],単音や 複数音の周期的な正弦波音響ノイズにのみ着目している. Natarajan ら [11] は不偏推定と重み付き損失関数を用いて ノイズ環境下におけるエージェントの耐性を検証した.他 にも、不偏な surrogate loss function を定義することで、 破損データやラベルノイズでの学習に関する研究もいくつ かある [4,12,19–21,25].しかし MADRL の分野、特に各 エージェントに与えられた分散型のネットワークではノイ ズを考慮した研究は筆者らの知る限り存在しない.また、 我々の目標は損失関数の変更やノイズの事前知識を一切与 えずにエージェントがノイズ環境下で学習可能にする手法 を提案することであるため,上記の先行研究とも目的が大 きく異なる.

ノイズの存在は必ずしも有害とは限らず,有用なことも あることを主張した先行研究もある [2,13,24].たとえば, Fortunatoら [2] はガウスノイズをネットワークの重みに 意図的に追加することで,エージェントの環境探索を促進 する NoisyNet を提案している.さらに Han ら [3] はエー ジェントの方策関数に確率性を導入した NoisyNet に拡張 することで,エージェントの学習安定を実現した.しかし これらの手法はエージェントの観測が常に正確でノイズを 含まないことを前提としており,内在的な inherent noise などの観測情報のノイズがある場合は検証されていない. また Wang ら [23] はノイズによって変動のある報酬設計を 提案し,強化学習モデルの堅固性向上を確認したものの, マルチエージェント環境での適用は想定されていない.

MADRL におけるノイズに関する研究は少ないものの, 全エージェントで共有する集中型ネットワークを用いた学 習モデルを想定したものが若干存在する [5,7,8,16,17]. し かし前述のとおり,エージェントの個体差や役割に応じて 観測するノイズは異なり,またそれに対する影響や耐性も エージェントごとに異なる.一方,本稿で提案する DA3 は個々のエージェントが行動を学習する際に,観測情報の どの部分がノイズかノイズでないかを自律的に識別し,学 習過程におけるノイズ発生による行動決定の混乱を低減す ることを検証する.また,エージェントが学習過程でノイ ズの多いデータや矛盾・不合理な行動をとる他のエージェ ントの位置を推定するために DA3 が役立つことを示すこ とで,DA3 エージェントの学習性能向上の説明を試みる.

3. 問題設定のモデル化

3.1 Decentralized POMDP

Markov decision process (MDP) とは、次ステップの状態 が現在の状態のみに依存して遷移する動的システムの確率 モデルのことを指す [14]. この時, n 体のエージェントによ る decentralized partially observable Markov decision process (dec-POMDP) を $\langle \mathcal{I}, \mathcal{S}, \{\mathcal{A}_i\}, p_T, \{r_i\}, \{\Omega_i\}, \mathcal{O}, H \rangle$ で 定義する. ここで要素は以下を表す.

- (1) *I* = {1,...,*n*}: エージェントの有限集合
- (2) S: 状態の有限集合
- (3) A_i : エージェント *i* のとりうる行動の有限集合 (*i* \in *I*)
- (4) $p_T(s'|s, a)$: 遷移関数 $(s, s' \in S, a \in A)$
- (5) $r_i(s,a)$: エージェント $i \in \mathcal{I}$ の報酬関数 ($s \in S, a \in \mathcal{A}$)
- (6) Ω_i : エージェント $i \in \mathcal{I}$ の観測の有限集合 ($i \in \mathcal{I}$)
- (7) $\mathcal{O}(o|s,a)$: 条件付き観測確率関数 $(o \in \Omega, s \in S, a \in A)$
- (8) H: 確率過程の最大範囲 (≥ 0)





このとき,エージェント*i*は割引累積報酬 $R_i = \sum_{t=0}^{H} \gamma^t r_i(s,a)$ を最大化するような方策 π_i を学習する. ただし, γ を割引率 ($0 \le \gamma < 1$)とする.本実験は決定的 で離散的なパトロールタスクを想定しているため,離散的 な時間ステップをt (≥ 0)と定義し,全ての確率関数は0か1のどちらかを取り得るとする.

3.2 問題設定と観測情報

本研究では図1に示すような $G_x \times G_y$ ($G_x, G_y \in \mathbb{Z}^+$, ここで Z⁺ は正の整数の集合) のグリッド状の環境を想定 し、数体のエージェント(図 1a では数字の振られた青い マス)が移動しながら、環境で発生する回収対象物(星型) を回収する object collection game を考える. 各エージェ ントは,自身を中心とした $R_x \times R_y$ $(R_x, R_y \in \mathcal{Z}^+)$ の範 囲内を観測できる. 各時間ステップにおいてエージェント $i \in \mathcal{I}$ は観測情報 $o_i \in \Omega_i$ を取得し、これを $N_C \in \mathcal{Z}^+$ チャ ネルの観測行列 (R_x, R_y) で表現し、内部のネットワーク ヘ与える. 観測行列は視野内の他エージェントの相対的な 位置,回収対象物(タスク),壁や障害物などの情報を含 む. エージェント $i \in I$ は得られた観測情報をもとに行動 $a_i \in \mathcal{A}_i = \{up, down, right, left\}$ を決定し、実行する.な お A_iの各行動は隣接する上・下・右・左のいずれかへの 移動を表す.移動先の事象によって以下の報酬を得るもの とする.

- (1) r_e > 0: 対象物を回収(同じ座標上に移動)して得られる報酬
- $(2) r_c < 0$: 壁やその他エージェントに衝突した場合に与



図 2: DA3の構造

えられる負の報酬

(3) r = 0: その他

各エージェントiは割引累積報酬 R_i を最大化するために 衝突回数を減らし対象物を多く回収するような方策 π_i を 自律的に学習する.なおこのゲームは単純な問題だが、本 提案である DA3 の効果を確認するために他の要因を極力 減らすために選択した.

3.3 自己注意機構 (Self-Attention)

Self-attention mechanism は attention mechanism の一 種であり,観測行列内の異なるベクトル間の依存関係・類 似性を示す [22]. Self-attention mechanism では,入力さ れる系列から query, key, value の3つのベクトルが計算 される. ここで query と key の各要素毎の類似性は,それ ぞれの内積によって求められる.その後 softmax 関数に よって類似性を [0,1]間に正規化して表し,これを attentional weights と呼ぶ.最後に value を掛けることで,以下 の attention mechanism の計算が完了する.

Attention
$$(Q, K, V) = \operatorname{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) \cdot V,$$
 (1)

ただし, Q, K, V, および d_k はそれぞれ query, key, value, そして value の次元数を示している. 更なる詳細は [22] を 参照していただきたい.

さらに, attention mechanism をh個の attentional heads へ拡張し, multi-head attention mechanism を導入する:

$$MultiHead(Q, K, V) = Concat(head_1, \dots, head_h)W^O$$
$$head_l = Attention(Q \cdot W_l^Q, K \cdot W_l^K, V \cdot W_l^V),$$
(2)

ただし、 W_l^Q , W_l^K , W_l^V , W_l^O は $head_l$ ($1 \le l \le h$) とする 際の学習対象パラメータ行列である. 各 attentional head が異なる入力状態を特徴づけるため、h の値を大きくする ことで複雑な環境を効率よく学習することが可能になる.

4. 提案手法

提案手法である DA3 は全てのエージェントが独立した ネットワークを所有することを想定している. 各エージェ ントがもつネットワークの構造を図2に示す.この図が 示すように, DA3 ネットワークには multi-head attention mechanism を含む transformer encoder が内在している. まずエージェントの観測データは state encoder で N_C チャ ネルの観測行列 (R_x, R_y) に変換(エンコード化)される. このとき各チャネルでは視野内にいるエージェント、回収 対象物 (タスク), また壁や障害物などの情報が2値({0,1} あるいは {0,-1}) で表現される. 図 1b は図 1a における エージェント 2 の観測情報 (図 1b の最左) と $N_C = 3$ で のエンコード化した情報を示す.ここで,壁は黒色,壁向 こうの観測できない領域は灰色で表現されている. 左から 2番目と3番目の行列は、エージェント3および対象物の 相対的な位置が1で示され、4番目の行列ではエージェン ト2の視野が壁などで観測不能な領域を-1で埋めている.

カーネルサイズ $P \in Z^+$ の畳み込みニューラルネット ワークにより、エンコード化された N_C チャネルの観測行 列は state embedder によって C チャネルの $\lfloor \frac{R_v}{P} \rfloor \times \lfloor \frac{R_v}{P} \rfloor$ 行列に展開される. P を patched size と呼ぶ. その後, 展 開された行列はさらに $\lfloor \frac{R_v}{P} \rfloor \cdot \lfloor \frac{R_v}{P} \rfloor + 1 \times C$ に平坦化され る. このとき、長さ C の saliency vector を追加するため、 1 が連結されている. Saliency vector は BERT [22] におけ る class token に近い役割を担っており、vision transfomer (ViT) [1] における classification token と同様の手法で実装 される. しかし class token や classification token は入力 情報を要約した形でベクトル化される一方、saliency vector は更に入力情報の顕著・重要な情報(協調エージェント、 回収対象物など)を表現している. その後、ViT と同様に position embeddings が行列に付与される.

State embedder で得られた $\lfloor \frac{R_x}{P} \rfloor \cdot \lfloor \frac{R_y}{P} \rfloor + 1 \times C$ 行列を 図 2a に示す通り transformer encoder へ入力する. このと き, $\lfloor \frac{R_x}{P} \rfloor \cdot \lfloor \frac{R_y}{P} \rfloor + 1 \times \lfloor \frac{C}{h} \rfloor$ の query, key, value がそれぞ れ h 個計算される. 観測情報の各要素と saliency vector の類似性は query と key の内積で計算され, その他の計 算は先行研究 MAT-DQN [10] と同様である. Transformer encoder での処理は $L \in \mathbb{Z}^+$ 回行われるが, この値は環境 の学習難易度によって決定するハイパーパラメータである.

最後に図 2b のように, transformer encoder から出力さ れた行列から saliency vector のみを抽出して *DRL head* へ 入力する. DRL head は学習対象に応じて変更することが 可能であり (たとえば Double DQN, Dueling DQN など), エージェントの環境適応力を促進する狙いがある. 図 2b の通り saliency vector のみが DRL head へ伝播するため, saliency vector と観測情報の要素との類似性を解析するこ とで, エージェントが行動決定過程において着目・無視し

表 1: 使用されたハイパーパラメータ

概要	パラメータ	値
エージェント数	n	6
対象物回収時の報酬	r_e	1
衝突時の報酬	r_c	-1
エピソード長	H	200
エポック数		10,000
対象物数		25



(a) R = 7 (b) R = 9 (c) R = 7 (d) R = 7図 3: ノイズ発生時の同一観測情報の例. (a) ノイズなしの 観測情報 (R = 7),ただし中心の緑が観測するエージェン ト,黄が対象物,白が壁である. (b) Large marginal noise 発生時の観測情報 (R = 9). (c) Small marginal noise 発 生時の観測情報 (R = 7). (d) Small full noise 発生時の観 測情報 (R = 7). いずれもオレンジ枠の外側でノイズが発 生する.

ている情報を可視化することが可能になる.

DA3 エージェントは内部に transformer encoder を持つ ことで2つ利点がある.第1に,エージェントが観測情報 の「取捨選択」が可能になる.つまり,エージェントが学 習を経て有益な領域のみ着目(あるいは不要な部分を無視) するように観測情報を制限するため,効率的な環境学習に 繋がる.第2に,エージェントの attentional weights を可 視化することでこれまでブラックボックスだった行動決定 過程に解釈性を持たせることが可能になる.前述の通り, saliency vector は観測情報を表現したベクトルであること から,エージェントの行動決定過程において影響のある情 報が表現される.この点は先行研究 [10] でも認められた が,他にもノイズ環境下での attentional weights を解析す ることで学習に対するノイズの影響度なども多角的に確認 できる.そこからの発見はシステムの更なる堅固性,効率 化へと拡張でき,有意義と考えている.

5. 実験と結果

5.1 実験設定

DA3 エージェントのノイズ環境下での学習性能を検証 するため、図 1a の環境にて object collection game を実行 し、比較手法との差異を調査した.本実験の目的は、エー ジェントが環境観測時に inherent noise が発生した場合の 学習性能の評価とノイズの影響を特定することである.こ の他に、ランダムに行動するエージェントを無視し、協調 すべきエージェントを同定し、それらと協調行動を学習で

きるかを検証する実験も行ったが,ここではページ数の関 係で省略する.実験に使用されたパラメータを表1に示 す.以下に示す学習データは20試行の平均値である.

本実験では、各 DA3 エージェントは 4 つの attentional heads (h = 4) を所有し transformer encoder は 1 回のみ (L = 1) ループする.比較手法である DQN の構造は図 2 とほぼ等しいが、transformer encoder は含まれていない ものとする.DA3 の DRL head には全結合層を採用し、 一方、比較手法では 2 層の畳み込みネットワーク、maxpooling 層、そして全結合層とする.Patch size は P = 1と設定し、観測情報の全ての要素が行動決定過程に関わる ように学習させる.また、saliency vector の長さを C = 64と設定した.エージェントの観測領域は一辺の長さが R $(= R_x = R_y)$ の正方形と想定した.

実験は図 1a に示す 3 部屋のグリッド環境(25×25)で 実施し、ノイズ環境でのエージェントの振る舞いを検証した.環境内には 6 体のエージェントが存在し、それぞれが 同じアルゴリズム (DA3 または DQN)で学習する.報酬 設計として $r_c = -1$, $r_e = 1$ を設定した.観測情報は図 1b の通り $N_C = 3$ でエンコード化され、前述の通り各チャネ ルでエージェント、対象物、壁あるいは観測できない領域 が表現される.観測情報は直接エージェントのネットワー クへと入力される.

各エピソードの開始時に, エージェントは環境中の固定番 号位置に配置され(図1), 25 個のオブジェクトがベージュ 色の領域にランダムに発生する.エピソード長を H = 200 ステップとし,それまで,エージェントは環境を探索,学 習する.エージェントが対象物を回収すると直後に別の対 象物がベージュ色領域の重ならないランダムな位置に発生 する.

5.2 観測情報におけるノイズ

本実験では、エージェントの観測情報をある確率で反転 させるノイズを導入する. このノイズは [11] に基づくもの であり、観測したエージェント、対象物、壁などが確率的 に反転(行列に応じて0と1あるいは0と-1)する.つま り、エージェントは観測している対象物の存在の有無に不 確実性を持つことになる.図3にて観測情報におけるノイ ズの内容を示す.図 3a はあるエージェントのノイズのな い観測情報(R=7)である.図の中央の緑は観測してい るエージェント自身、青は観測領域内のその他のエージェ ント, 黄は回収対象物, 白は壁を表す. 図 3b (R = 9)と 図 3c (R = 7) は、センサーによる観測範囲のうち、観測 可能限界付近の境界でノイズが高確率で発生する特徴を表 した. これを marginal noise と呼ぶ (各図でオレンジ線の 外側が確率に従って反転している). また, 観測範囲の境界 のみならず領域内全体にノイズは存在するが、距離に応じ て反転確率が徐々に小さくなる full noise を想定した場合

表 2: 各環境における学習性能の内訳

(a) DA3 による学習性能

発生ノイズ名	対象物回収量	エージェント間の衝突回数	壁との衝突回数
ノイズなし	271.45 ± 12.55	0.20 ± 0.12	2.98 ± 0.61
Large marginal	271.24 ± 11.88	0.24 ± 0.09	3.15 ± 1.11
Small marginal	208.59 ± 11.94	0.21 ± 0.10	2.47 ± 0.43
Small full	92.84 ± 6.96	2.56 ± 0.45	4.53 ± 0.52

(b) DQN(比較手法)による学習性能

発生ノイズ名	対象物回収量	エージェント間の衝突回数	壁との衝突回数
ノイズなし	239.01 ± 14.78	0.32 ± 0.12	3.95 ± 0.70
Large marginal	233.31 ± 12.85	0.56 ± 0.17	4.46 ± 0.86
Small marginal	172.22 ± 11.20	0.29 ± 0.11	3.61 ± 0.52
Small full	29.41 ± 4.63	6.33 ± 1.39	13.49 ± 1.12

でも検証した (図 3d). この full noise の場合, 観測者自身 に近い距離の情報は低い確率, 遠い位置の情報は高い確率 でノイズが発生するとする. なおデータの反転確率は, 以 下の実験結果の節で説明する. また, すべてのエージェン トは同じ性能のセンサーを装備しており, ノイズの発生確 率は同等と仮定した.

5.3 学習性能比較

DA3 エージェントが 3 種類の観測ノイズが発生する環 境で(図 3b-d)に適応して学習できることを、ノイズなし 環境での学習(図 3a)と比較することで検証した.図 3b では、観測領域を通常より一回り大きく設定し(R = 9)、 代わりに縁の情報が 0.5 の確率で反転するノイズ(以下、 large marginal noise)を想定した.これは観測領域の縁に 情報量は無いため、実質的にはR = 7 のノイズなし観測情報と等価となる.しかし、エージェントはこのノイズに関 する事実は与えられていない.図 3c では = 7 とし、同様 に縁部分の情報のみ 0.2 の確率で反転させる(以下、small marginal noise と呼ぶ).つまり縁の情報は、信頼できる情 報と出来ない情報がランダムに混在していることを意味す る.最後に図 3d の small full noise 設定では、縁では 0.2、 その内側では 0.1、さらにその内側(観測エージェントの 周囲)では 0.05 の確率で情報を反転させる.

図4にエージェントの学習推移に応じて、1 エピソード ごとに獲得した報酬の推移と、エージェント同士・壁との 衝突回数を示す.また、表2は実験の最後の100 エピソー ド間における対象物の平均回収量、衝突の内訳を集計した ものである.一般にノイズが増えるとその影響によって エージェントの学習性能はいずれの強化学習アルゴリズム に関わらず低下することが予想されるが、本実験でも同様 の現象を確認した.ノイズなしの観測情報を得たエージェ ントが最も高い学習性能を達成し、特に small full noise の ようなノイズが観測領域全体に発生する環境では学習性能 が最も低いことが分かる.

情報処理学会研究報告



図 4: 観測情報にノイズが発生する際の学習性能比較

図 4 から、DA3 エージェントが比較手法より多くの対象 物をより少ない衝突数で回収することができており、DA3 エージェントがどの環境下でも高い学習性能を達成したこ とが分かる.表2を参照すると、各環境におけるエージェ ントの対象物回収量の差はそれぞれ 32.44 (ノイズなし、 13.57%の性能向上) 37.93 (large marginal noise, 16.26% の性能向上), 36.37 (small marginal noise, 21.12%の性 能向上), そして 63.44 (small full noise, 215.72%の性能 向上)となっている.また,DA3 エージェントの衝突回数 が比較手法よりも大きく抑制できたこと(特に small full noise)も確認できる.

さらに図 4 から, DA3 エージェントはノイズなしの場合 と large marginal noise の場合でほぼ同程度の学習性能を 達成している. このことは表 2a からも確認でき,対象物 回収量(差分は 0.21),そして衝突回数の差(0.04)は僅か である. 実際に large marginal noise の場合は *R* = 9 であ るが外周は事実上情報量がなく,ノイズ無しと同等の情報 量であり,提案手法ではこの事実を反映してほぼ同程度の 学習性能を示した.一方,比較手法の DQN は,これらの 間で性能劣化が確認された.その他のノイズの場合も DA3 エージェントが比較手法よりも高い学習性能を発揮してお り,ノイズの頻度(確率)に応じたノイズ抑制ができてい ると言える.

5.4 Attention の解析・考察

DA3の transformer encoder 内部に存在する attentional weights から attention heatmap を作成することで, 観測 情報のノイズが如何にエージェントの行動決定過程に影響 を与えるかを解析した. 図 5 はノイズなし環境における (a) エージェントの観測情報, (b) 4 つの attentional heads (h = 4) からの平均 attention heatmap, (c) 各 attentional head の attention heatmap を表す.

ノイズなし環境での観測情報の例 (図 5a) からエージェ ントは左手に壁,前方に2つの対象物を観測している.こ







図 6: 観測情報と attention heatmap (large marginal noise)

のとき attention heatmap (図 5b) と照合すると,エージェ ントは自身の次位置の可能性のある四方と対象物に対して 比較的強く着目していることが分かる.このエージェント は自身の上方に最も高い attentional weight を与えており, 上へ移動して即時報酬を得ることを期待していると考えら れる.また,図 5cより,各 attentional head が異なる観測









図 8: 観測情報と attention heatmap (small full noise)

領域に着目している.したがって h をある程度増やし,よ り効率的に環境を捉えた学習が可能になると思われる.こ の現象は MAT-DQN [10] でも確認されている.

次に図 6 は large marginal noise 環境での観測情報と attention heatmap である.先と同様にエージェントは自 身の四方,そして対象物に対して着目していることが図 6a と図 6b から確認できる.一方ノイズの発生する観測領域 の縁では,attentional weights がほぼ 0.000 となっている ことも併せて確認できる.これは非常に顕著であり,エー ジェントが観測領域縁における情報量の皆無性を認め,意図 的にその領域を無視するように学習できたと言える.DA3 エージェントは large marginal noise が発生する環境でも ノイズなし環境での学習性能と遜色ない結果を達成してい るが,この「ノイズを無視する」現象が要因と考えられる.

Small marginal noise 環境での結果である図 7 から, DA3 エージェントは、同様に自身の四方(0.263, 0.039, 0.032, 0.095)と近くの対象物に比較的高く着目するよう学習し ている.また、観測領域の縁における attentional weights の値は [0.005, 0.021] で推移しており、先の環境(図 5b, 図 6b) での attentional weights の値と比較すると若干高 い. これは small marginal noise に影響を受けている縁の 情報を全く無視するわけではなく,むしろある程度考慮し て次の行動を決定していると思われる.

最後に, small full noise 環境でのエージェントの観測情 報と attention ヒートマップを図 8 に示す.図 8b からエー ジェントは自身の周辺8つを重点的に着目して行動決定す ることが示唆される. 観測領域の縁に対する attentional weights はおおよそ 0.007 程度であり、一回り中心に近い 領域は 0.021, エージェントの周辺には 0.062 程度である. これはエージェントは遠くの対象物にはあまり着目せず, 代わりに比較的信頼のできる観測領域(自身の周辺8マス) に着目し、ノイズの発生率が高くなるほど低い attetntional weights を与えるように学習しているからである. この特 徴は先の環境で見られたエージェントの隣接する位置に対 して強く着目する状況とは大きく異なる. これは近視眼的 な行動となり、学習の性能をやや下げることになる. しか しその場合でも、DA3エージェントはエピソード間に100 以上の対象物を回収できた一方で、比較手法は 50 未満と なり、その差は大きい. また興味深いことに、図 8c から各 attentional head がそれぞれエージェントの上下左右の一 方向を担当していることも確認できる.

5.5 考察

我々の実験結果から DA3 エージェントは, DA3 の部分 で観測情報のどの領域が行動決定に重要かを判断しながら 学習することで,ノイズの度合い応じた処理をしていると 思われる.実験では,DA3 エージェントは比較手法より も優れた学習性能を示し,ノイズを含まないときも何を着 目するべきかを学習し,観測情報にノイズが含まれる場合 もその度合いに応じて重要さ (attentional weights)を学習 し,ノイズに対する高い堅固性を示した.このような学習 結果の説明は未知のことが多かったが,本研究によりエー ジェントの協調や競合回避行動の解釈性に役立つと考えら れる.

6. むすび

本稿では、distributed attentional actor architecture model for multi-agent system (DA3)を提案した.提案 手法によるノイズ抑制を検証するため、エージェントの観 測情報にノイズを様々な確率で発生させた環境での学習 性能をベースラインである (vanilla) DQN と比較した.結 果、DA3 エージェントは行動決定に重要な情報に着目する ように学習し、ノイズがある程度含まれる場合も、それに 応じて重みを学習し、比較手法よりも優れた性能を達成で きた.DA3 エージェントが保有する transformer encoder から attentional weights を可視化・解析したところ、観測 情報に発生するノイズへの適応も分かるようになった.今 後の展開としてはノイズの特性の幅を広げ,現実に近い環 境での検証へと拡張する予定である.

謝辞 本研究は JSPS 科研費 17KT0044 と 20H04245 の 助成を受けたものです.

参考文献

- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J. and Houlsby, N.: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, *International Conference on Learning Representations* (2021).
- [2] Fortunato, M., Azar, M. G., Piot, B., Menick, J., Hessel, M., Osband, I., Graves, A., Mnih, V., Munos, R., Hassabis, D., Pietquin, O., Blundell, C. and Legg, S.: Noisy Networks For Exploration, *International Conference on Learning Representations* (2018).
- [3] Han, S., Zhou, W., Liu, J. and Lü, S.: NROWAN-DQN: A Stable Noisy Network with Noise Reduction and Online Weight Adjustment for Exploration, *CoRR* (2020).
- [4] Hendrycks, D., Mazeika, M., Wilson, D. and Gimpel, K.: Using Trusted Data to Train Deep Networks on Labels Corrupted by Severe Noise, Advances in Neural Information Processing Systems (Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N. and Garnett, R., eds.), Vol. 31, Curran Associates, Inc. (2018).
- [5] Kilinc, O. and Montana, G.: Multi-agent Deep Reinforcement Learning with Extremely Noisy Observations (2018).
- [6] Kuo, S. M., Kuo, K. and Gan, W. S.: Active noise control: Open problems and challenges, *The 2010 International Conference on Green Circuits and Systems*, pp. 164–169 (online), DOI: 10.1109/ICGCS.2010.5543076 (2010).
- [7] Li, S., Wu, Y., Cui, X., Dong, H., Fang, F. and Russell, S.: Robust Multi-Agent Reinforcement Learning via Minimax Deep Deterministic Policy Gradient, *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33, No. 01, pp. 4213–4220 (online), DOI: 10.1609/aaai.v33i01.33014213 (2019).
- [8] Lin, W., Zhixin, L. and Lei, G.: Robust Consensus of Multi-agent Systems with Noise, 2007 Chinese Control Conference, pp. 737–741 (online), DOI: 10.1109/CHICC.2006.4347503 (2007).
- [9] Minutti, C., Gomez, S. and Ramos, G.: A machinelearning approach for noise reduction in parameter estimation inverse problems, applied to characterization of oil reservoirs, *Journal of Physics: Conference Series*, Vol. 1047, p. 012010 (online), DOI: 10.1088/1742-6596/1047/1/012010 (2018).
- [10] Motokawa, Y. and Sugawara, T.: MAT-DQN: Toward Interpretable Multi-agent Deep Reinforcement Learning for Coordinated Activities, *Artificial Neural Networks* and Machine Learning – ICANN 2021 (Farkaš, I., Masulli, P., Otte, S. and Wermter, S., eds.), Cham, Springer International Publishing, pp. 556–567 (2021).
- [11] Natarajan, N., Dhillon, I. S., Ravikumar, P. K. and Tewari, A.: Learning with Noisy Labels, Advances in Neural Information Processing Systems (Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z. and Weinberger, K. Q., eds.), Vol. 26, Curran Associates, Inc. (2013).
- [12] Patrini, G., Rozza, A., Menon, A., Nock, R. and

Qu, L.: Making deep neural networks robust to label noise: a loss correction approac, *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, pp. 2233–2241 (online), DOI: 10.1109/CVPR.2017.240 (2017).

- [13] Plappert, M., Houthooft, R., Dhariwal, P., Sidor, S., Chen, R. Y., Chen, X., Asfour, T., Abbeel, P. and Andrychowicz, M.: Parameter Space Noise for Exploration, *International Conference on Learning Representations* (2018).
- [14] Puterman, M. L.: Markov Decision Processes: Discrete Stochastic Dynamic Programming, John Wiley & Sons, Inc., USA, 1st edition (1994).
- [15] Raeisy, B. and Golbahar Haghighi, S.: Active Noise Controller with reinforcement learning, *The 16th CSI International Symposium on Artificial Intelligence and Signal Processing (AISP 2012)*, pp. 074–079 (online), DOI: 10.1109/AISP.2012.6313721 (2012).
- [16] Raeisy, B., Golbahar Haghighi, S. and Safavi, A. A.: Active Noise Control System via Multi-Agent Credit Assignment, J. Intell. Fuzzy Syst., Vol. 26, No. 2, p. 1051–1063 (2014).
- [17] Ryu, H., Shin, H. and Park, J.: Multi-Agent Actor-Critic with Hierarchical Graph Attention Network, Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, No. 05, pp. 7236–7243 (online), DOI: 10.1609/aaai.v34i05.6214 (2020).
- [18] Sabzevari, S. A. H. and Moavenian, M.: Application of reinforcement learning for active noise control, *Turkish J. Electr. Eng. Comput. Sci.*, Vol. 25, pp. 2606–2613 (2017).
- [19] Scott, C.: A Rate of Convergence for Mixture Proportion Estimation, with Application to Learning from Noisy Labels, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics* (Lebanon, G. and Vishwanathan, S. V. N., eds.), Proceedings of Machine Learning Research, Vol. 38, San Diego, California, USA, PMLR, pp. 838–846 (2015).
- [20] Scott, C., Blanchard, G. and Handy, G.: Classification with Asymmetric Label Noise: Consistency and Maximal Denoising, *Proceedings of the 26th Annual Conference on Learning Theory* (Shalev-Shwartz, S. and Steinwart, I., eds.), Proceedings of Machine Learning Research, Vol. 30, Princeton, NJ, USA, PMLR, pp. 489–511 (2013).
- [21] van Rooyen, B. and Williamson, R. C.: Learning in the Presence of Corruption (2015).
- [22] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u. and Polosukhin, I.: Attention is All you Need, *Advances in Neural Information Processing Systems* (Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S. and Garnett, R., eds.), Vol. 30, Curran Associates, Inc., (2017).
- [23] Wang, J., Liu, Y. P. and Li, B.: Reinforcement Learning with Perturbed Rewards, AAAI (2020).
- [24] Xu, H., Zhang, C., Wang, J., Ouyang, D., Zheng, Y. and Shao, J.: Exploring Parameter Space with Structured Noise for Meta-Reinforcement Learning, *IJCAI* (2020).
- [25] Zhang, Z. and Sabuncu, M.: Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels, Advances in Neural Information Processing Systems (Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N. and Garnett, R., eds.), Vol. 31, Curran Associates, Inc. (2018).