

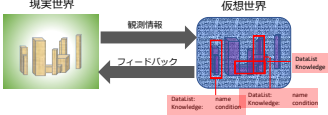
脳のマルチモーダル情報処理に着想を得た 物体推定手法の提案と評価

関 良我† 小南 大智† 下西 英之†† 村田 正幸†
藤若 雅也†† 野上 耕介††
†大阪大学 大学院情報科学研究科
††NEC システムプラットフォーム研究所

1

研究背景

- ◆デジタルツインへの期待
 - 現実世界のあらゆる物体をリアルタイムに仮想空間上にマッピングすることで環境を再現
 - 仮想空間でのシミュレートによって得られる様々な情報を現実世界へフィードバック
- ◆想定されるデジタルツインの活用例
 - 自動運転：人や自動車、障害物の動きを収集、自動車の運転制御
 - 公共空間：空港など公共空間内の人の動きを把握、不審者や不審人物を検出
 - 荷物管理：倉庫などの荷物の動きを把握、ARゴーグルによる情報提示や自動搬送



2

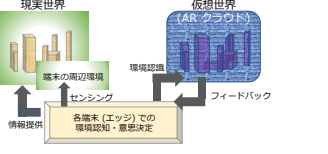
デジタルツインに必要な物体推定手法

- ◆高速なセンシングデータ処理
 - 画像分析などによる物体やその状態の識別
 - センサー機器自体の自己位置推定と物体位置推定
 - リアルタイムに複数の物体を個別に追跡して識別/推定
- ◆従来の物体推定手法
 - CNN などの物体推定手法では高精度な識別が可能であるがデータをクラウド上に集約して処理
 - 実世界のすべてのセンサー情報をクラウドに送信して処理することは困難
 - 逐一クラウドに送信せず各端末で物体推定処理を行う

3

AR クラウドを利用したエッジ - クラウド型デジタルツイン構想

- ◆それぞれの端末で集めた現実世界の情報を統合して仮想世界を構築
 - 各端末で環境認識を行い、局所的な仮想世界を構築
 - AR クラウド上で各端末の情報を統合し、巨大な仮想世界を構築
 - AR クラウドの情報をもとに各端末は行動を決定
- ◆より小さな時間粒度で処理するため各端末で物体推定
- ◆刻一刻と変化する現実世界の環境に対応するため柔軟性のある設計
 - 軽量・ロバストな端末での処理が必要



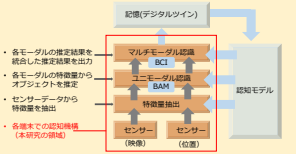
4

研究の目的・アプローチ

目的 端末単位で処理可能な軽量性と現実世界の環境変化に適応可能なロバスト性を持った物体推定手法の提案 / 評価

アプローチ


- ◆脳の優れた認知機構に着想を得た物体推定手法の提案
 - 目や耳、三半規管などから得られた不確実な情報のマルチモーダル統合処理
 - 階層構造の軽量かつロバストな意思決定
- ◆脳の認知機能の数理モデルを利用
 - 脳や生体の「ゆらぎ学習」によるユニモーダル処理
 - Bayesian Attractor Model (BAM)
 - 脳の知覚による因果推論に基づいたマルチモーダル処理
 - Bayesian Causal Inference (BCI)



5

BAM の物体推定手法への利用

- ◆映像モダリティと位置モダリティをもとに意思決定
 - 特徴量抽出
 - ◆ SiameseRPN [1]) によりテンプレート画像をもとにバウンディングボックスとして位置を出力
 - ◆ 抽出された特徴量を 128 次元に圧縮 → 映像モダリティ入力
 - ◆ カメラ方向ベクトルとデータセットの深度情報から算出した 3 次元ワールド座標系 → 位置モダリティ入力
 - 特徴量変換
 - ◆ 平均・分散を一定にする標準化処理
 - ◆ 映像特徴量の領域が定まっておらず、BAM で処理しやすい形に整形



[1] B. Li, J. Tan, W. Wu, Z. Zhu, and K. Yu, "High performance visual tracking with siamese region proposal network," in Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8972-8980, 2018.

6

Bayesian Attractor Model^[2]

- 人が観測した情報をもとに意思決定する脳の振る舞いを数理モデル化
 - 情報を観測し続け、事前に学習した選択肢のどれに近いかを判断
 - 状態更新**: 観測情報 x_t に基づいて内部状態 z_t をベイズ推定によって更新
 - ベイズ推定であるため z_t は点ではなく脳の不確実性を反映した確率分布 $P(z_t)$
 - 意思決定**: 事前に学習した選択肢 (アトラクター) のいずれかを採用
 - 状態空間上にアトラクター ϕ_1, \dots, ϕ_n を事前に用意
 - 状態値 z_t が選択肢 i である事後確率 (確信度 $P(z_t = \phi_i | x_{0:t})$) を算出
 - 確信度が閾値 λ を超えたものを採用

[2] S. Wiles, I. Brunberg, and S. I. Khalil, "A Bayesian attractor model for perceptual decision making," *PLoS Comput. Biol.*, vol. 11, no. 8, p. e1004442, 2015.

7

Bayesian Causal Inference^[3]

- 人がマルチモーダルな知覚刺激をもとに行う認知を数理モデル化
 - 2つの入力刺激が同じ刺激源から出たものであるかを確率的に識別 (**Causal Inference**)
 - その確率をもとにそれぞれの入力刺激を統合する (**Model Average**)

- 2つの入力刺激 x_1, x_2 でそれぞれ個別に認知結果 S_1, S_2 を計算
 - 本研究ではBAMによる認知を行う
- 両モーダルの入力刺激を統合した認知結果 S_{12} を計算
- 両モーダルで同じものを認知しているか ($C=1$) 別のものを認知しているか ($C=0$) を確率的に計算する (**Causal Inference**)
- $C=1$ であれば S_{12} を出力、 $0 < C < 1$ であればその割合に応じた認知結果を出力 (**Model Average**)

[3] K. P. Körding, U. Tenenbaum, W. J. Ma, S. Quartz, J. B. Tenenbaum, and I. Shams, "Causal inference in multi-sensory perception," *PLoS one*, vol. 2, no. 9, p. e493, 2007.

8

Causal Inference

- 観測対象が同じであるか計算
 - 従来は刺激源の位置 (連続値) を対象としたが今回はオブジェクトの識別 (離散値) が対象
- 離散値への拡張

<p>従来の Causal Inference (連続値)</p> $p(X, Y C = 1) = \int p(X, Y s) p(s) ds$ <p>別オブジェクトを観測</p> $p(X, Y C = 0) = \int p(X s) p(s) ds \int p(Y s) p(s) ds$	<p>拡張 Causal Inference (離散値)</p> $p(X, Y C = 1) = \sum_{k=1}^K p(X, Y O_k) p(O_k)$ $p(X, Y C = 0) = \sum_{k=1}^K p(X O_k) p(O_k) \sum_{l=1}^K p(Y O_l) p(O_l)$
<ul style="list-style-type: none"> s: オブジェクトの位置 $p(s)$: 位置の確率分布 $p(X s)$: オブジェクトが s にあるときに X として観測される確率 	<ul style="list-style-type: none"> O_k: オブジェクト k (識別ラベル) $p(O_k)$: オブジェクト O_k が観測される確率 $p(X O_k)$: オブジェクト O_k を観測しているときに X として観測される確率

[3] K. P. Körding, U. Tenenbaum, W. J. Ma, S. Quartz, J. B. Tenenbaum, and I. Shams, "Causal inference in multi-sensory perception," *PLoS one*, vol. 2, no. 9, p. e493, 2007.

9

Model Average

- $0 < C < 1$ の場合の重み付け処理
 - BCI ではどちらのモダリティが優先されるべきかは定義しない
 - 映像 / 位置モダリティを優先した場合の2通りの結果を出力

$$Cost_x(O_x) = p(C=1) \sum_{k=1}^K |O_x - O_k|^2 p(O_k | X, Y) + p(C=0) \sum_{k=1}^K |O_x - O_k|^2 p(O_k | X)$$

$$Cost_y(O_y) = p(C=1) \sum_{k=1}^K |O_y - O_k|^2 p(O_k | X, Y) + p(C=0) \sum_{k=1}^K |O_y - O_k|^2 p(O_k | Y)$$

[3] K. P. Körding, U. Tenenbaum, W. J. Ma, S. Quartz, J. B. Tenenbaum, and I. Shams, "Causal inference in multi-sensory perception," *PLoS one*, vol. 2, no. 9, p. e493, 2007.

10

シミュレーション評価方式

- 公開データセットの映像データをもとに映像モーダル・位置モーダルの特徴量を抽出
 - 1111 枚の映像フレームと 4 つのオブジェクト
 - 4444 枚をオブジェクトごとに直列で入力
 - アトラクターには各オブジェクトのはじめの1枚を記憶
- ユニモーダルによる物体推定
 - BAMのみの推定精度を評価
- マルチモーダルによる物体推定
 - BCI によって拡張された BAM の推定精度を評価

[4] オブジェクトの検知と公開データセット (Pascal3D+ Berkeley (YCB) Object and Model set)

11

ユニモーダルによる物体推定

- 1 モダリティのみを BAM に入力したときの確信度出力
 - 映像モダリティ
 - 位置モダリティ

正答率: 79.41 %
 オブジェクト 2 を上手く認知できず

正答率: 81.66 %
 オブジェクト 3 を上手く認知できず

12

