

映像中の人物の手と物体のインタラクション検出 Detecting Hand-Object Interaction Based on Movements around Hand

小西 陸斗¹⁾ 阿部 亨¹⁾²⁾ 中村 隆喜¹⁾²⁾ 菅沼 拓夫¹⁾²⁾
Rikuto Konishi Toru Abe Takaki Nakamura Takuo Suganuma

1 はじめに

映像からの人物と物体のインタラクション検出は、監視・防犯、人物の行動分析など多くのアプリケーションでの利用が期待されている [1, 2]. 人物と物体のインタラクションには様々なものがあるが、その中でも、人物が物体を手で動かす動作は、インタラクションの中核を担う動作であり [3], 人物と物体の接触を示す確実な手がかりとしても有益である。

映像からのインタラクション検出のアプローチには、1 段階目で、映像内の全ての人物と物体の抽出と対応する組み合わせを決定し、2 段階目で、決定した組み合わせからインタラクションの有無を求める two-stage method と、映像から抽出した特徴量に基づきインタラクションの有無を直接求める one-stage method の 2 つがある。two-stage method は、独立した 2 段階の処理が必要であり、全体の性能が第 1 段階の精度に制限されるため、one-stage method がより効果的と考えられる [2].

One-stage method における特徴量の設計には、サンプル画像から機械学習で特徴量を決定する learning-based アプローチと、検出対象に関する事前知識に基づき決定する handcrafted アプローチの 2 つがある。learning-based アプローチは、多様な検出対象への対応が可能になるが、映像から得られる人物の骨格情報や動き情報など複数の情報を利用する場合に、各々の情報を異なるストリームで学習し統合するフレームワーク (multi-stream framework) が用いられており、学習のためのネットワークが大規模かつ複雑になることから、処理に必要なリソースや処理時間が増加していた [4, 5]. これに対し、handcrafted アプローチは、事前知識の利用が可能な特定の動作に検出対象が限定されるものの、事前知識に基づき複数の情報を効果的に統合した特徴量が設計でき、単一のストリームで学習を行うネットワークにより効率的にインタラクションを検出できる可能性がある。

筆者らは、手が動かす物体は前腕と類似した動きを示すという事前知識に基づき、映像中の人物の骨格情報と動き情報を統合した特徴量を設計し、前腕と同様の動きが手周辺に生じているかによりインタラクションの検出を効率的に行う手法を提案している。本稿では、この提案手法に関し、その効果を検証した結果を示す。

2 関連研究

2.1 learning-based アプローチ

learning-based アプローチによる人物の行動認識のための特徴量の決定手法を 2 つ取り上げる。1 つ目は、映

1) 東北大学大学院情報科学研究科

Graduate School of Information Sciences, Tohoku University

2) 東北大学サイバーサイエンスセンター

Cyberscience Center, Tohoku University

像から得られる動き情報を利用する Simonyan ら [6] の手法である。この手法では、図 1 のように、映像の各フレームの元画像とオプティカルフローをそれぞれ別の CNN に入力して特徴を抽出し、2 つのストリームから得た出力を統合することで行動の認識を行う。2 つ目は、映像中の人物の骨格情報を利用する Haroon ら [7] の手法である。この手法では、図 2 のように、映像の連続するフレームの元画像と人物の骨格 (キーポイント) を別々の LSTM で処理し、出力を統合することで認識を行う。これらの手法では、多様な人物行動の認識が可能であるが、複数の情報を異なるストリームで処理するため、ネットワークが大規模かつ複雑になり、学習や認識を行う際の処理量が増加してしまう。

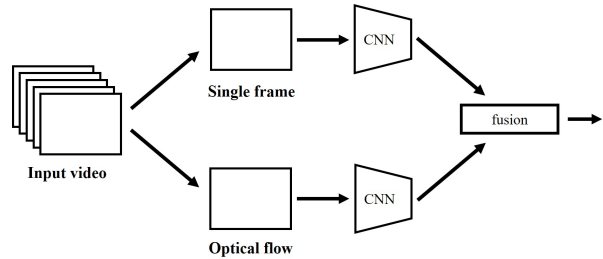


図 1 動き情報を利用する手法

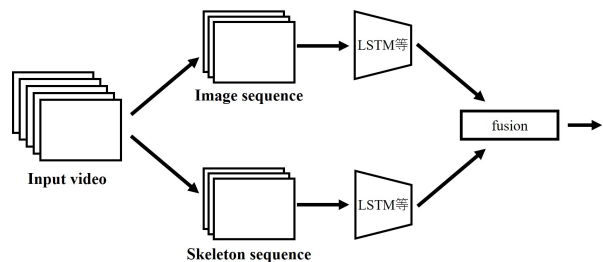


図 2 骨格情報を利用する手法

2.2 handcrafted アプローチ

手で物体を動かす動作を検出するための特徴量を handcrafted アプローチにより決定するものとして Tsukamoto ら [8] の手法がある。この手法では、手で物体を動かしている場合は前腕と同様の動きが手の周辺に生じているという事前知識に基づいて、骨格情報と動き情報を統合した特徴量を設計し、物体領域を抽出することなく、インタラクションの判定を行う。これにより、映像から得られる情報を効率的に利用可能となる。しかし、前腕の一部の動き (画面に対して直交する動き) が特徴量に反映されないことや、骨格情報と動き情報を統合した特徴量が手周辺での動きの状態の詳細を表現できていないという問題がある。

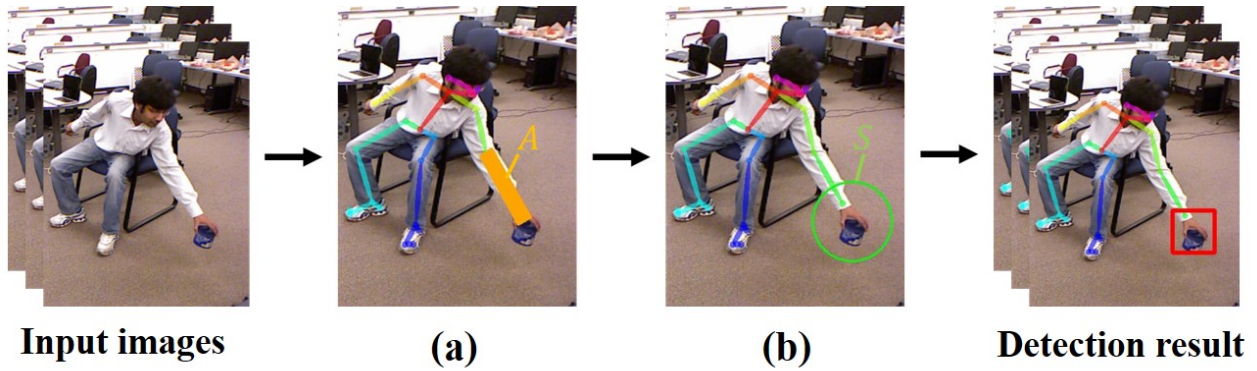


図 3 提案手法の処理概要

3 提案手法

提案手法では、Tsukamoto らの従来手法 [8] に着目し、映像より得られる動き情報と骨格情報を統合した特徴量から、機械学習を用いて手が物体を動かしているかを判定する。提案手法の処理の概要を図 3 に示す。まず、入力画像から人物の骨格を抽出し、図 3 (a) のように、前腕領域 A を設定する。前腕領域 A 内の動きを評価し、前腕が動いていると判定された場合、図 3 (b) のように、手の周辺領域 S を設定する。ここで、 S 内に前腕と同様の動きを示す領域が生じていれば、手が物体を動かしていると判定する。提案手法における各処理の詳細を、従来手法と比較しながら以下に説明する。

3.1 前腕領域 A の設定

はじめに、映像中の人物の各前腕領域を設定するために、入力画像から画像内の人物の骨格（キーポイント）を抽出する。抽出した肘 P_E 、手首 P_W のキーポイントの位置に基づき、 $P_E - P_W$ を前腕の先端方向に $\Delta L = \alpha \times L$ 延長した点を手の中心 O として設定する。ここで、 L は $P_E - P_W$ の距離を表す。続いて、 $P_E - O$ に沿った $(L + \Delta L) \times (w \times L)$ の矩形領域として前腕領域 A を設定する。

3.2 前腕の動きの判定

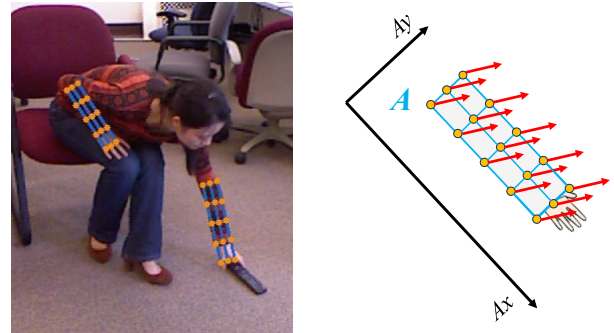
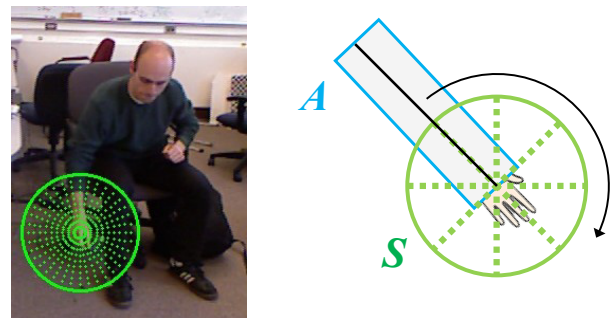
設定した各前腕領域 A 内のオプティカルフローの状態から、前腕が動いているか判定を行う。

従来手法 従来手法では、 A 内の全画素のオプティカルフローの平均値が閾値 T_A 以上であれば前腕が動いていると判定する。

提案手法 提案手法では、図 4 に示すように、 A を等間隔にサンプリングし、各箇所オプティカルフローを前腕領域の長軸・短軸方向に分解した成分で特徴ベクトルを構成する。その特徴ベクトルを入力として、機械学習により前腕の動きの有無を判定する。

3.3 周辺領域 S の動きの分析

前腕が動いていると判定された場合、図 5 に示すように、手の中心 O を中心とした半径 R の円を手の周辺領域 S として設定する。続いて、 S 内に手で動かしている物体が存在するかを判定するために、前腕の動きと S 内の動きの類似性を評価する必要がある。そこで、まず、

図 4 前腕領域 A のサンプリング例図 5 周辺領域 S のサンプリング例

前腕領域 A の各画素 $p_n = (x_n, y_n)$ の動き $va(p_n)$ をモデル化する。

従来手法 従来手法では、前腕が画面に対して平行に角速度 ω で回転、速度 $T = (t_x, t_y)$ で並進すると仮定し、 p_n での動き $va(p_n)$ を式 (1) でモデル化する。

$$va(p_n) = (-\omega y_n + t_x, \omega x_n + t_y) \quad (1)$$

提案手法 提案手法では、アフィン変換を用いて p_n での動き $va(p_n)$ を式 (2) によりモデル化する。

$$va(p_n) = (c_1 x_n + c_2 y_n + c_3, c_4 x_n + c_5 y_n + c_6) \quad (2)$$

上記のパラメータ（従来手法では式 (1) の ω, t_x, t_y 、提案手法では式 (2) の c_1, c_2, \dots, c_6 ）は、 A の各 p_n で実際に観測されたオプティカルフロー $vo(p_n)$ との差の二乗和が最小になるように設定する。提案手法では、アフィン変換により動きをモデル化しているため、従来手法とは異なり、画面に対して平行に回転、並進する前腕の動き

だけでなく、画面に直交する方向の動きも近似的にはあるが表現できるようになる。

続いて、 S 内の各画素 q_n で、手で動かす物体に生じると予想される動き $ve(q_n)$ と実際に観測されたオプティカルフロー $vo(q_n)$ との差 ndv を式 (3) により求める。

$$ndv(q_n) = \frac{\|vo(q_n) - ve(q_n)\|}{\|ve(q_n)\|} \quad (3)$$

物体は、手でしっかり把持されおり、手と一体となって動く場合には前腕領域 A と同様の動き $va(q_n)$ を、手からぶら下がった状態で動く場合には手の中心 $O = (x_O, y_O)$ と同様の動き $va(x_O, y_O)$ を示すと予想されるが、実際は、二つの状態が混合した動きを示す場合が多いと考えられる。そこで、二種類の動きの線形和 $ve(q_n) = \alpha \times va(q_n) + (1 - \alpha) \times va(x_O, y_O)$ で予測値を表現し、 $vo(q_n)$ との差が最小になるように重み α を設定する。このようにして得られた $ndv(q_n)$ の値が小さければ、 S 内の q_n は、手で動かされた場合と同様の動きを示しており、手で動かされる物体に対応する画素である可能性が高いと考えられる。

手の周辺領域 S 内で得られた ndv の状態により、手が物体を動かしているか（手で動かされた物体が S 内に存在するか）を判定する。

従来手法 従来手法では、 ndv の値が閾値 T_d 以下となる画素の総数と S 内の画素の総数の比が閾値 T_S 以上であれば手が物体を動かしていると判定する。

提案手法 提案手法では、図 5 に示すように、 S 内をサンプリングした箇所の ndv の値から特徴ベクトルを構成し、機械学習で手が物体を動かしているかを判定する。

提案手法では、従来手法とは異なり、 S 内の状態をベクトルで表現しているため、 ndv の分布等のより詳細な状態に基づいた判定が可能となっている。

4 評価実験

提案手法の有効性を検証するために、以下 2 つの実験を行った。

実験 1 従来手法 [8] との検出精度の比較

実験 2 学習データ数と検出精度の関係の検証

4.1 対象動画

実験には、人物行動の認識実験用に一般公開されているデータセット CAD-120 [9] の全 10 カテゴリーの動画計 124 本（1 カテゴリーにつき被験者 4 人、各被験者同様の動作を 3 回（making_cereal のみ同様の動作を 4 回））を用いた。各動画（RGB 映像、640×480 画素、30fps）では、1 人の人物が物体を手で動かしており、動作はカテゴリー毎に異なる。

各動画の各フレームに対し、OpenPose [10] を適用して人体のキーポイントの検出を行い、信頼度が 0.5 以上となる肘と手首のキーポイントの位置から抽出された前腕を実験の対象とした。抽出された各前腕について、物体を動かしているか否かを人手でラベル付けし、検出結果を評価する際の正解とした。実験に用いた動画、フレーム、抽出された前腕、物体を動かしている手の総数

表 1 実験に使用した対象動画

Category	# of videos (frames)	# of extracted forearms (hands moving objects)
picking_objects	12 (2501)	4818 (669)
arranging_objects	12 (3781)	6629 (1413)
unstacking_objects	12 (5586)	10986 (2751)
taking_food	12 (5614)	8616 (2129)
stacking_objects	12 (5813)	11472 (2972)
microwaving_food	12 (6350)	9946 (2870)
taking_medicine	12 (6394)	12887 (3461)
cleaning_objects	12 (7406)	11539 (3799)
having_meal	12 (9829)	19066 (4933)
making_cereal	16 (11647)	22500 (771)
Total	124 (64921)	118459 (32708)

を表 1 に示す。

4.2 評価方法

実験では、動画中で抽出された前腕を対象に、物体を動かしている手の検出を行い、検出結果を正解と比較することで、True Positives の数 TP 、False Positives の数 FP 、False Negatives の数 FN を求め、式 (4)~(6) で得られる適合率 P 、再現率 R 、F 値 F により検出精度を評価した。

$$P = \frac{TP}{TP + FP} \quad (4)$$

$$R = \frac{TP}{TP + FN} \quad (5)$$

$$F = \frac{2 \times P \times R}{P + R} \quad (6)$$

4.3 実験の設定

従来手法、提案手法ともに、OpenPose で検出されたキーポイントに基き前腕の骨格を決定し、前腕の長さ L に対し、前腕に沿って手首から $\Delta L = 0.35 \times L$ の箇所に手の中心 O を設定し、前腕領域 A の長さは $L + \Delta L$ 、幅は $0.25 \times L$ 、手の周辺領域 S は O を中心とした半径 $R = 1.10 \times L$ の円とした。

従来手法では、前腕の動きの判定と手と物体のインタラクションの判定を閾値処理により行う。各閾値は予備実験により、 $T_A = 0.02$ 、 $T_d = 0.55$ 、 $T_S = 0.15$ として設定した。

提案手法は、前腕が動いているかの判定のため、 A を 5×3 箇所等間隔にサンプリングし、各箇所で見測されたオプティカルフローの成分で構成される $5 \times 3 \times 2$ 次元の特徴ベクトル VA に基づき、SVM [11] により行う。手が物体を動かしているかの判定は、 S を円周方向に 36 箇所、半径方向に 10 箇所等間隔にサンプリングし、各箇所で見測された ndv の値で構成される 36×10 次元の特徴ベクトル VS に基づき SVM により行う。

4.4 実験 1 従来手法との検出精度の比較

実験 1 では、データセットの多様な動作に対して提案手法の特徴量の有効性を確認するため、従来手法との検出精度の比較を行った。

表 2 従来手法との検出精度の比較結果

Category	Method	w. affine t.	w. SVM	TP	FP	FN	P	R	F
picking_objects	proposed	○	○	591	828	78	0.42	0.88	0.57
		×	○	599	930	70	0.39	0.90	0.55
		○	×	461	570	208	0.45	0.69	0.54
	existing	×	×	425	543	244	0.44	0.64	0.52
arranging_objects	proposed	○	○	1183	581	230	0.67	0.84	0.74
		×	○	1107	672	306	0.62	0.78	0.69
		○	×	1056	487	357	0.68	0.75	0.71
	existing	×	×	963	398	450	0.71	0.68	0.69
unstacking_objects	proposed	○	○	2225	383	526	0.85	0.81	0.83
		×	○	1982	518	769	0.79	0.72	0.75
		○	×	1579	950	1172	0.62	0.57	0.60
	existing	×	×	1195	829	1556	0.59	0.43	0.50
taking_food	proposed	○	○	835	208	1294	0.80	0.39	0.53
		×	○	711	276	1418	0.72	0.33	0.46
		○	×	1302	461	1827	0.40	0.14	0.21
	existing	×	×	229	404	1900	0.36	0.11	0.17
stacking_objects	proposed	○	○	2512	488	460	0.84	0.85	0.84
		×	○	2235	619	737	0.78	0.75	0.77
		○	×	1891	1033	1081	0.65	0.64	0.64
	existing	×	×	1538	923	1434	0.62	0.52	0.57
microwaving_food	proposed	○	○	1609	329	1261	0.83	0.56	0.67
		×	○	1399	490	1471	0.74	0.49	0.59
		○	×	1111	732	1759	0.60	0.39	0.47
	existing	×	×	845	611	2025	0.58	0.29	0.39
taking_medicine	proposed	○	○	1923	626	1538	0.75	0.56	0.64
		×	○	1742	840	1719	0.67	0.50	0.58
		○	×	1561	785	1900	0.67	0.45	0.54
	existing	×	×	1252	659	2209	0.66	0.36	0.47
cleaning_objects	proposed	○	○	1506	271	2293	0.85	0.40	0.54
		×	○	1266	330	2533	0.79	0.33	0.47
		○	×	588	446	3211	0.57	0.15	0.24
	existing	×	×	410	345	3389	0.54	0.11	0.18
having_meal	proposed	○	○	3408	386	1525	0.90	0.69	0.78
		×	○	3629	438	1304	0.89	0.74	0.81
		○	×	2608	447	2325	0.85	0.53	0.65
	existing	×	×	2618	358	2315	0.88	0.53	0.66
making_cereal	proposed	○	○	5265	1537	2446	0.77	0.68	0.73
		×	○	5015	1762	2696	0.74	0.65	0.69
		○	×	4542	1833	3169	0.71	0.59	0.64
	existing	×	×	4035	1716	3676	0.7	0.52	0.60
Total	proposed	○	○	21057	5637	11651	0.79	0.64	0.71
		×	○	19685	6875	13023	0.74	0.60	0.66
		○	×	15699	7744	17009	0.67	0.48	0.56
	existing	×	×	13510	6786	19198	0.67	0.41	0.51

4.4.1 実験方法

提案手法については、前腕の動きを、従来手法と同様、並進・回転のみで表現する場合と、前腕が動いているかの判定、手が物体を動かしているかの判定を、従来手法と同様、閾値処理で行う場合についても実験を行った。

提案手法の SVM の学習は対象とする動画毎に、対象以外の 123 本の動画から得られた VA, VS を用いて行った。

4.4.2 実験結果

従来手法との検出精度の比較の結果を表 2 に示す。また、提案手法による検出結果の例を図 6 に示す。

Total の検出結果では、アフィン変換を用いることで、FP も増加するが、TP がより大きく増加している。これは、オプティカルフローの推定誤差等も前腕の動きとして過度に近似される場合があり FP が増加するものの、画面に直交する前腕の動きがより良く表現された効果の方が高く TP の増加が大きくなったものと考えられる。

一方、SVM による判定を行うことで、FP がわずかに増加したものの、TP が大きく増加している。これは、前腕領域や周辺領域での動きの状態の詳細が判定に利用されたことと、SVM による機械学習で判定を行ったことの効果であると考えられる。

カテゴリごとの検出結果では、提案手法、従来手法ともに、taking_food と cleaning_objects では、R が低いいため F が低くなっている。これらのカテゴリでは、前腕の動きが小さい場合や遅い場合が多く、前腕が動いていないと判定されてしまうため、インタラクションの検出ができなかったためと考えられる。また、picking_objects では、FP が多く P が低いいため、F も低くなっている。picking_objects は、他カテゴリと比べて、フレーム数が少なくなっており、SVM を学習する際に、カテゴリ間でデータ数に大きな偏りが生じ、検出結果に影響を与えた可能性がある。そこで、実験 2 では、学習に用いるデータの数をカテゴリ間で揃えた場合の検出精度を検証した。

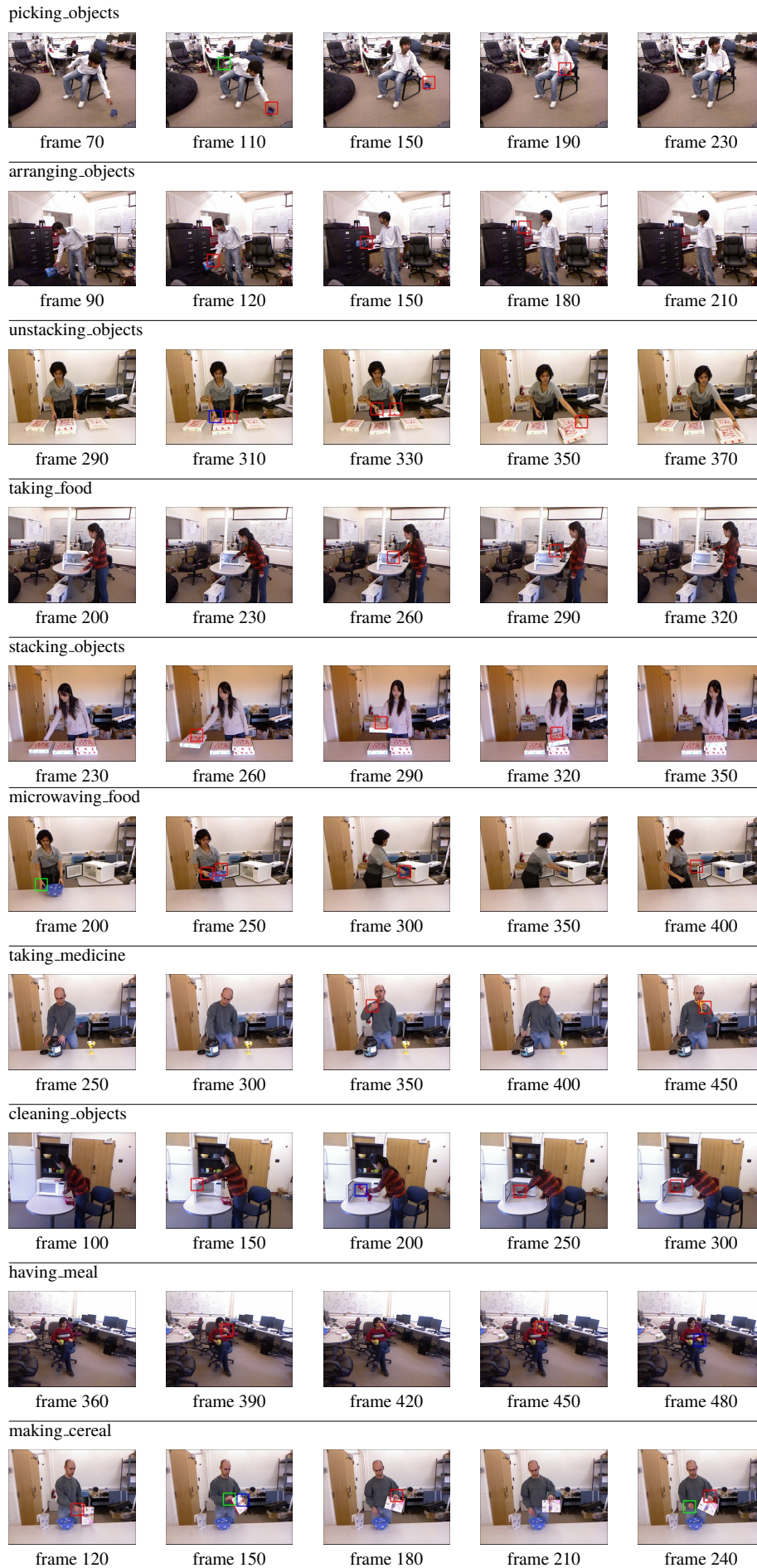


図 6 提案手法による検出結果の例 (赤の矩形 True Positive, 青の矩形 False Negative, 緑の矩形 False Positive)

4.5 実験 2 学習データ数と検出精度の関係の検証

実験 2 では、実験 1 の結果に基づき、学習に用いるデータの数が提案手法の検出精度へ及ぼす影響について検証した。

4.5.1 実験方法

実験 2 では、データセットを被験者ごとに 4 つのグループに分割し、各グループを検出対象とする際に、他 3 グループのデータを *SVM* の学習用に用いた。また、カテゴリ間での学習データ数を揃えるために、動画のフレームをサンプリングし、学習用の特徴ベクトルと正解ラベルを選択した。今回は、1 カテゴリ当り 850 フレーム、1700 フレームとなるようサンプリングした場合について実験を行った。

4.5.2 実験結果

実験結果を表 3 (1 カテゴリ当り 850 フレーム)、表 4 (1 カテゴリ当り 1700 フレーム) に示す。

表 2 の結果と比較して、全体的な傾向としては、学習データの総数が少なくなると、*TP* だけでなく *FP* も減少するため *P* の変化は小さいものの、*TP* の減少がより大きく、その結果、*R* が大きく低下し *F* も低下している。また、1 カテゴリ当り 850 フレーム分のデータで学習した場合と 1700 フレーム分のデータで学習した場合では、*P*, *R*, *F* に大きな違いは生じなかった。

また、カテゴリ間の学習データ数を揃えても、*picking_objects* に対する *P*, *F* は低くなっており、モデル化した前腕の動きでは対応できない動きが検出対象である動作に含まれている可能性がある。その原因については、今後さらに検証を行う必要がある。

5 おわりに

本稿では、人物が手で物体を動かす動作に着目し、映像中の人物の骨格情報と動き情報を統合した特徴量から機械学習を用いて動作の検出を行う手法を提案した。実験では多様な動作に対して、従来手法と提案手法での検出精度の比較を行い、提案手法の特徴量の有効性を示した。またカテゴリ間での学習データ数を揃えることで、学習データ数と検出精度の関係を検証した。

今後は、*SVM* 以外の機械学習手法での判定方法の検討や、映像中に多数の人物や物体が存在する場合やより複雑な動作で物体を動かす場合などの多様な状況を対象とした実験を行う予定である。また複数の情報を各々別のストリームで処理し、インタラクションの検出を行う *multi-stream framework* に基づく手法との、検出精度、学習効率、計算コストなどの比較検証を行う予定である。

参考文献

- [1] Khaire, P. and Kumar, P.: Deep learning and RGB-D based human action, human-human and human-object interaction recognition: A survey, *Journal of Visual Communication and Image Representation*, Vol. 86, p. 103531 (2022).
- [2] Antoun, M. and Asmar, D.: Human object interaction detection: Design and survey, *Image and Vision Computing*, Vol. 130, p. 104617 (2023).
- [3] Fan, H., Zhuo, T., Yu, X., Yang, Y. and Kankanhalli, M.: Understanding Atomic Hand-Object Interaction With Human Intention, *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 32, No. 1, pp. 275–285 (2022).

表 3 検出結果 (1 カテゴリ当り 850 フレーム分のデータで学習した場合)

Category	<i>TP</i>	<i>FP</i>	<i>FN</i>	<i>P</i>	<i>R</i>	<i>F</i>
<i>picking_objects</i>	537	692	132	0.44	0.80	0.57
<i>arranging_objects</i>	1091	461	322	0.70	0.77	0.74
<i>unstacking_objects</i>	2023	332	728	0.86	0.74	0.79
<i>taking_food</i>	494	153	1635	0.76	0.23	0.36
<i>stacking_objects</i>	2341	494	631	0.83	0.79	0.81
<i>microwaving_food</i>	1305	269	1565	0.83	0.45	0.59
<i>taking_medicine</i>	1608	481	1853	0.77	0.46	0.58
<i>cleaning_objects</i>	877	223	2922	0.80	0.23	0.36
<i>having_meal</i>	2872	255	2061	0.92	0.58	0.71
<i>making_cereal</i>	4664	1273	3047	0.79	0.60	0.68
Total	17812	4633	14896	0.79	0.54	0.65

表 4 検出結果 (1 カテゴリ当り 1700 フレーム分のデータで学習した場合)

Category	<i>TP</i>	<i>FP</i>	<i>FN</i>	<i>P</i>	<i>R</i>	<i>F</i>
<i>picking_objects</i>	543	697	126	0.44	0.81	0.57
<i>arranging_objects</i>	1101	474	312	0.70	0.78	0.74
<i>unstacking_objects</i>	2060	314	691	0.87	0.75	0.80
<i>taking_food</i>	555	174	1574	0.76	0.26	0.39
<i>stacking_objects</i>	2387	457	585	0.84	0.80	0.82
<i>microwaving_food</i>	1340	288	1530	0.82	0.47	0.60
<i>taking_medicine</i>	1673	509	1788	0.77	0.48	0.59
<i>cleaning_objects</i>	963	241	2836	0.80	0.25	0.38
<i>having_meal</i>	2920	261	2013	0.92	0.59	0.72
<i>making_cereal</i>	4779	1302	2932	0.79	0.62	0.69
Total	18321	4717	14387	0.80	0.56	0.66

- [4] Wang, H., Yu, B., Li, J., Zhang, L. and Chen, D.: Multi-Stream Interaction Networks for Human Action Recognition, *IEEE Transactions on Circuits and Systems for Video Technology*, Vol. 32, No. 5, pp. 3050–3060 (2022).
- [5] Wang, C. and Yan, J.: A Comprehensive Survey of RGB-Based and Skeleton-Based Human Action Recognition, *IEEE Access*, Vol. 11, pp. 53880–53898 (2023).
- [6] Simonyan, K. and Zisserman, A.: Two-Stream Convolutional Networks for Action Recognition in Videos, *CoRR*, Vol. abs/1406.2199, (2014).
- [7] Haroon, U., Ullah, A., Hussain, T., Ullah, W., Sajjad, M., Muhammad, K., Lee, M. Y. and Baik, S. W.: A Multi-Stream Sequence Learning Framework for Human Interaction Recognition, *IEEE Transactions on Human-Machine Systems*, Vol. 52, No. 3, pp. 435–444 (2022).
- [8] Tsukamoto, T., Abe, T. and Suganuma, T.: A method for detecting human-object interaction based on motion distribution around hand, in *Proc. 15th Int Joint Conf. Comput. Vision, Imaging Comput. Graphics Theory Appl.*, pp. 462–469 (2020).
- [9] Koppula, H. S., Gupta, R. and Saxena, A.: Learning human activities and object affordances from RGB-D videos, *Int. J. Rob. Res.*, Vol. 32, No. 8, pp. 951–970 (2013).
- [10] Cao, Z., Hidalgo, G., Simon, T., Wei, S.-E. and Sheikh, Y.: OpenPose: Realtime multi-person 2D pose estimation using part affinity fields, *IEEE Trans. Pattern Anal. Mach. Intell.*, Vol. 43, No. 1, pp. 172–186 (2021).
- [11] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E.: Scikit-learn: Machine learning in Python, *J. Mach. Learn. Res.*, Vol. 12, No. 85, pp. 2825–2830 (2011).