

## Web 動画ショットの動作分類のための 時空間特徴抽出手法の提案

野口 顕嗣<sup>†1</sup> 柳井 啓司<sup>†1</sup>

Web には多種多様な動画ショットが存在している。しかしこれらは人手でつけられたタグによって分類されており、動画の内容では分類されていない。そこで本研究では Web 動画分類のための新しい時空間特徴抽出手法について提案し、実際にその特徴を利用することで Web 動画ショット分類を行う。また Multiple Kernel Learning (MKL) に基づく特徴統合による精度の向上についても検討した。Web 動画分類には教師信号ありの動作ランキング付け (あるショットがその動作を含んでいる度合いをランキング付けする) に加えて教師信号なしのクラスタリングによる自動シーン分類を行った。結果として教師信号ありランキング付けの実験では、時空間特徴のみでランキング上位 200 位までの適合率が 57.2%、MKL で特徴を統合することで 68.8% という結果が得られた。

### Extracting Spatio-Temporal Local Features for Classifying Web Video Shots

AKITSUGU NOGUCHI<sup>†1</sup> and KEIJI YANAI<sup>†1</sup>

Nowadays, there is a large amount of video data on the Web. Most of videos on the Web have keyword tags for text-based search. However, tags do not always reflect the contents of videos, so that content-based video search is needed. Then, in this paper, we propose a new spatio-temporal feature and a feature fusion method based on Multiple Kernel Learning (MKL) for classifying Web video shots. We made experiments on supervised shot ranking and unsupervised shot clustering. As results, we obtained 57.2% as the precision rate regarding top-ranked 200 shots using only the spatio-temporal feature. By utilizing Multiple Kernel Learning, we obtain the 68.8% precision.

<sup>†1</sup> 電気通信大学  
The University of Electro-Communication

## 1. はじめに

近年 Web 上の動画の数は爆発的に増えてきている。しかし一方でユーザが見たいシーンを探すことが非常に困難な問題となってしまう。現状の動画検索システムはタグなどによるテキストベースな手法が用いられているが、これはユーザの主観によってつけられるもので、タグのみで動画を特定することは非常に難しい。そのため動画の内容を解析するコンテンツベースな検索手法が今後求められてくると考えられる。そのなかでも、動画の内容を解析し、分類を行うことは非常に意義のあることである。しかし現状このような Web 動画の分類を行った研究は少ない。

そこで本研究では、新しい時空間特徴の提案を行い、その特徴を利用することで、Web から大量に収集した動画を分類する。分類には教師信号ありの分類に加え、教師信号なしクラスタリングの二種類の分類を行う。また複数特徴を MKL によって統合することによる有効性について検証を行う。

教師信号あり分類 walking という単語で収集されたショットには walking を含むショットも存在するが、含まないショットも大量に存在する。そこで本実験では、そのショットが walking を含むかどうかランキング付けを行い、ランキングの上位には walking を含むショットが、ランキングの下位には walking とは関連のない動作を順位付けることを目的としている。

教師信号なしクラスタリング soccer という単語で収集された動画には様々なシーンがある。そこで本実験では教師なしクラスタリングを行うことで、soccer を「試合のシーン」、インタビューシーンなどのシーンに自動で分類することを目的としている。

## 2. 関連研究

近年、動画の解析のために時空間特徴が注目を集めている。時空間特徴とは、動画から抽出される特徴の一つで動き情報と視覚情報を同時に表現することが可能な特徴である。

まず時空間特徴の抽出の考え方として、画像認識の三次元拡張の手法が挙げられる。Kobayashi らは自己相関特徴を三次元に拡張し、サーベイランスの分野に特化した cubic higher-order local auto-correlation (CHLAC) を提案した<sup>1)</sup>。また Laptev らはハリス点検出を三次元に拡張した検出手法を提案している<sup>2)</sup>。

次に主要な手法として、cuboid と呼ばれる立方体を抽出し、それを特徴化する手法がある。Dollar らは空間軸にガウシアンフィルター、時間軸にガボールフィルタを適応することでこの cuboid を抽出する手法を提案した<sup>3)</sup>。また cuboid の記述子として、Laptev や Dollar らは Histogram of Orient Gradient (HoG) や Histogram of Orient Flow (HoF) で表現することを提案している。一方で Kläser らは、記述子として三次元的な HoG を利用することを提案している<sup>4)</sup>。

しかしこのような cuboid 主体の手法は計算コストが高くなる傾向があり、本論文で行うような大量データを扱うのに向いていない。更に cuboid のサイズを求めることは非常に困難な問題である。そこで本研究では、特徴的な点と、その点の微小時間の動きを特徴化する新たな時空間特徴を提案する。この特徴は計算時間が少なく、かつ正確な分類が期待できる。

Web 動画に対する動作認識を行った研究は少ない。Cinbis らは Web 上から動作を自動学習する手法を提案し、Youtube データを動作認識に利用している<sup>5)</sup>。この研究では、まず query word をもとにして画像を収集し、その画像から特徴を抽出することで動作モデルを構築し、実際の映像において認識を行っている。しかしこの手法の学習データは Web 上から収集された静止画像であり、動作の記述も全て動画像ではなく静止画像ベースで行われている。本研究では画像の視覚特徴のみではなく、動き特徴も考慮して動作を分類することを考える。

最も本研究の手法と類似している手法として、Liu らの研究が挙げられる<sup>6)</sup>。この手法は特徴量として 3) で提案された時空間特徴を、視覚特徴として SIFT 記述子<sup>7)</sup> を利用し、Adaboost に基づき統合する手法を提案している。またこの研究では、Page Rank に基づく重要な特徴の選択を行っている。本研究では、動作を認識するために新たな時空間特徴を提案し、利用するが、この手法では既存の特徴をどのように扱うかに重点が置かれている。

### 3. 提案手法概要

本節では、時空間特徴抽出手法、特徴統合による Web 動画分類手法の概要について述べる。

#### 3.1 時空間特徴抽出手法

時空間特徴には 8) を拡張した手法を利用する。まずフレーム画像から、SURF<sup>9)</sup> に基づく視覚特徴を抽出する。SURF とは、回転、スケール変化に頑健な局所特徴である。次に、抽出された SURF の座標に対して、Lucas-Kanade アルゴリズム<sup>10)</sup> に基づいて動きを計算する。ここで動きのなかった点は時空間特徴として適さないのので、これらの点は排除し、動きのあった点のみを取り扱うこととする。その後、決定された時空間特徴点に関して Delaunay 三角分割法を適用する。以降特徴は三角形の頂点を構成する三点のペアで表現される。更に時空間特徴点の微小区間の動きを特徴化することにより動き特徴を抽出する。最後に抽出された視覚特徴と動き特徴を重みをつけて結合することで時空間特徴を抽出する。特徴抽出の手順は、手順 1 にて述べられている。

ここで抽出された特徴は Dollar らの手法<sup>3)</sup> と同様に bag of spatio-temporal features (BoSTF) 表現によって認識に利用される。

#### 3.2 Web 動画分類

図 1 は Web 動画のショット分類の概要を示している。まず特定のキーワードで Youtube から動画を収集する。次に集められた動画を色特徴に基づきショット分割する。その後、そ

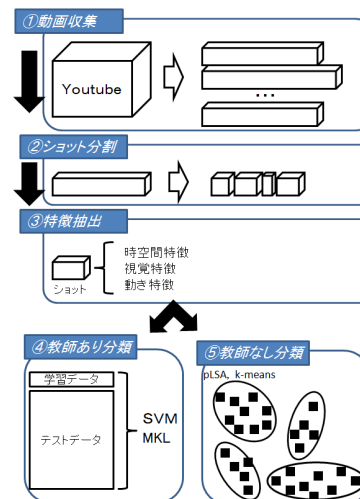


図 1 Web 動画分類手法の概要



図 2 教師信号ありの Web 動画分類

それぞれのショットから時空間特徴、動き特徴、視覚特徴を抽出する。最後に分類を行うが、本研究では SVM や MKL による教師あり分類の他に、pLSA や k-means による教師なしクラスタリングも試みる。

教師信号ありのショットランキング付けでは、図 2 に示すように、予め用意されている running 学習セットで SVM のような分類器を学習する。その分類器によって収集した Web 動画ショットの分類を行う。その結果、ランキングの上位には running を含むショットが、ランキング下位には running を含まないショットがランキングされることが期待される。本研究では提案した特徴単体での分類と、MKL に基づく特徴統合による分類を行う。統合に利用した特徴は、時空間特徴、動き特徴、視覚特徴の三種類の特徴である。

次に教師なしの分類では、クラスタリングによるシーンの自動分類を行う。これは例えば soccer という単語で集められたショットには様々なサッカーのシーンが含まれると考えられる。そこで pLSA や k-means のようなクラスタリング手法に基づき、クラスタリングすることで「ドリブルシーン」、「シュートシーン」、「インタビューシーン」のようなシーンに分類されることが期待される。

## 4. 提案手詳細

### 4.1 時空間特徴抽出

本研究では、Web 動画分類に適した新しい時空間特徴抽出手法を提案する。多くの動作認識の研究は、カメラモーションに対する対応がない。しかし Web 動画を分類するためには「カメラモーション」をどのように扱うかは非常に重要な問題となってくる。そこで本研究では、特徴を抽出する前にカメラモーションの検出を行う。

また既存の抽出手法は計算コストが高く、本研究のような大量のデータから抽出することに向いていない。本研究では、その問題を解決するために、点の周辺パターンと、その点の動きで動画を特徴化する手法を提案する。この手法は既存手法に比べ計算コストが低く、大量のデータから特徴を抽出することが可能である。

手順 1 に提案する特徴抽出手法の流れを示す。

手順 1 時空間特徴抽出手法の流れ	
step1	: カメラモーション検出
step2	: 時空間特徴における視覚特徴抽出
▷step2-1	: SURF 抽出
▷step2-2	: 時空間特徴点の決定
▷step2-3	: Delaunay 三角分割
step3	: 時空間特徴における動き特徴抽出
▷step3-1	: Lucas-Kanade 法による動き抽出
▷step3-2	: SURF の dominant rotation による方向の正規化
step4	: 視覚特徴ベクトルと動き特徴ベクトルの結合

本手法は大きく 4 つに分けることができる。まず、カメラモーション検出部 (step 1)、二つ目に視覚特徴抽出部 (step 2)、三つ目が動き特徴抽出部 (step 3)、最後に特徴統合部 (step 4) である。以下ではそれぞれのステップ毎に詳しく述べていく。

(Step 1) カメラモーション検出部 動画は撮影者の意図によって、ズームやパンなどのカメラモーションを含むことがある。しかし Web 動画においてのカメラモーションは手振れなどの、撮影者の意図しないカメラモーションが多く含まれている。また Web 動画は解像度が低い。これらのことから収集した動画の正確なカメラモーションを求めることは非常に難しい問題である。

Liu らの手法ではカメラモーションを検出した場合そのフレームを破棄することで対応をしている<sup>6)</sup>。本研究でも、これを参考にし、カメラモーションを検出した場合、その特徴を破棄する。しかしカメラモーションを全て破棄しては、「動画全てにカメラモーションを含んでいるような動画から特徴が抽出されない」、「カメラモーション中に存在する重要な特徴を抽出できない」、などの問題点もある。

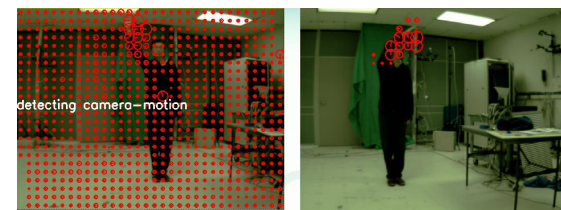


図 3 カメラモーション検出例 (左)、カメラモーションが検出されない例 (右)



図 4 視覚特徴抽出の様子

検出手法として、図 3 のようにグリッド上に Lucas-Kanade アルゴリズムに基づいて、動き情報を計算する。その中で動きのあった領域が一定以上だった場合カメラモーションの検出とする。

(Step 2) 時空間特徴における視覚特徴抽出部 図 4 は視覚特徴の抽出の様子を示している。各ステップ毎に見ていくと、(1) フレーム画像から SURF を抽出する。本研究では動画としての特徴を抽出したいので、動きがない点は時空間特徴として適していない。(2) よってそれぞれの特徴点の動きを計算し、動きがなかった点を削除する (時空間特徴の決定)。(3) 最後に残った点について Delaunay 三角分割を行う。この時空間特徴を三角分割することは本手法独自のものである。これを行うことで、その点のみではなく隣接する特徴も考慮した特徴を構築することが可能となる。

視覚特徴には三角形の頂点を構成する三点の SURF 記述子を使用する。それぞれの三角形の頂点の点を SURF 記述子のスケール毎に整列させる。ただし同一スケールの特徴が存在した場合は、 $x$  座標の最も小さい特徴点から見て右周りに特徴を配置する。SURF の次元数は 64 次元なので、結果として視覚特徴は  $64 \times 3 = 192$  次元で表現される。

(Step 3) 時空間特徴における動き特徴抽出部 動き情報として、本手法では三角形の頂点を構成する、それぞれの点の動きと、三角形の面積の変動を特徴化する。点の動きは視覚特徴の時と同様で SURF のスケールによって整列させる。

図 5 は動き特徴の抽出の概要を示している。まず図 5(左) が時空間特徴点の決定の様子を示している。最初に SURF を抽出したフレームから  $N$  フレーム先までのフレームをと

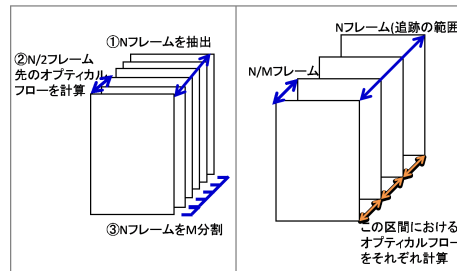


図 5 全体の動き特徴抽出概要 (左), 局所追跡の概要 (右)

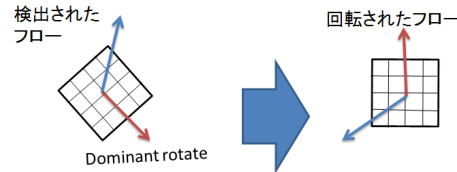


図 6 オプティカルフローの回転

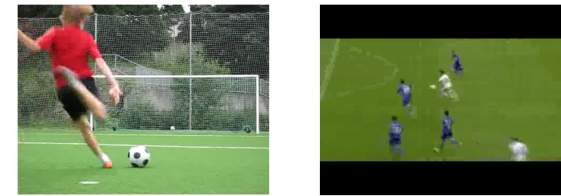


図 7 動き特徴が有効な例 (左), 視覚特徴が有効な例 (右)

り, 以降この区間から動き特徴を抽出する. 次に, 抽出された SURF の点に対して,  $N/2$  フレームの動き情報を計算する. この動き情報に基づいて時空間特徴は決定される (手順 1 の step 2-2).

次に動き特徴の抽出に関して説明する, 図 5(右) がその様子を示している. 選択された  $N$  フレームを,  $M$  分割する. そして分割したそれぞれの区間で Lucas-Kanade アルゴリズムによって, 特徴点のオプティカルフローを計算する. ただし, インターバル区間  $i$  の特徴点の座標  $L_i$  は, 前の区間  $i-1$  で推定された, 動き情報によって求められる. この分割数  $M$  が  $N$  の値に近いほど詳細な追跡が可能になり, 逆に  $M$  が 1 に近くなると簡略な動き特徴が抽出される.

各区間の動き情報は  $x^+, x^-, y^+, y^-$ , 及び動きなし, の 5 次元で表現される. ただし  $x^+$  はオプティカルフローの  $x$  成分の正の方向を,  $x^-$  は  $x$  成分の負の方向を示している. 実験において特徴を抽出する区間  $N=5$ , 分割数  $M=5$  に設定しているので, それぞれの点における次元数は  $(M-1) \times 5$  で 20 次元で, 三角形の面積変化は 5 次元で表現される. よって動きの次元数は  $20 \times 3 + 5 = 65$  次元となる.

ただしこのままの計算では動き特徴は回転に関して敏感な特徴になってしまう. 同じ「歩く」という動作においても, 動作の方向によって異なる特徴が抽出されてしまう. そこで本研究では視覚特徴で得られる dominant rotation を利用し, オプティカルフローを回転させる. 特徴は三点で一組になっているが, ここではそれぞれの dominant rotation を利用して, 動き情報の回転を行う. その様子を図 6 に記載する.

特徴点の座標を  $(x_1, y_1)$ , オプティカルフローが検出された座標を  $(x_2, y_2)$  とした, SURF の dominant rotation を  $\theta$  とした場合, 回転されたフローの座標  $(x, y)$  は式 1 で定義されるものとなる.

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta & x_2 \\ \sin\theta & \cos\theta & y_2 \end{bmatrix} \begin{bmatrix} x_1 - x_2 \\ y_1 - y_2 \\ 1 \end{bmatrix} \quad (1)$$

最後に二つの特徴を単純に結合することによって時空間特徴を構築する.

#### 4.2 特徴統合

ショットによって重要となる特徴は異なる場合がある. 例えば図 7(左) のような場合は視覚特徴より動き特徴が有効であると考えられる. しかし図 7(右) のような試合の中におけるシュートシーンの場合, カメラモーションが含まれる場合が多く, 動き特徴の信頼性は低い.

よって本研究では Multiple Kernel Learning(MKL) によって特徴の重みを自動で推定する手法を利用することで, これらのショットの分類を行った. 統合には視覚, 動き, 時空間特徴の 3 種類の特徴を利用した.

##### 4.2.1 Multiple Kernel Learning

本研究では動作認識を行うために Multiple Kernel Learning(MKL) を利用する. これは複数のサブカーネルを線形結合することで, 新たな最適なカーネルを求める手法で, 統合カーネルは以下の式で求められる.

$$K_{combined}(x, x') = \sum_{j=1}^K \beta_j k_j(x, x') \quad (2)$$

with  $\beta_j \geq 0, \sum_{j=1}^K \beta_j = 1.$

各サブカーネルの重み  $\beta_j$  の設定によって, 統合された新たなカーネルの出力は変化する. そのためこの最適な重みをどのように求めるかが問題となり, この問題は MKL 問題と呼ばれている<sup>11)</sup>. この MKL 問題は全ての  $\beta_j$  の組み合わせを実際に試すことで解くことが出来る. しかし, 特徴やカーネルの数が増えるにつれ,  $\beta_j$  の組み合わせも爆発的に増えてしまう. 実時間内にこの問題を解決するために, 近年 MKL 問題を凸面最適化問題として解く手法が

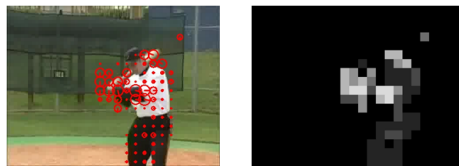


図 8 抽出された動き特徴

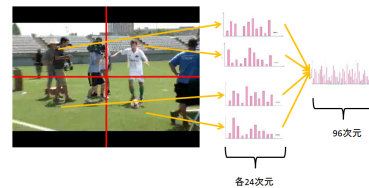


図 9 視覚特徴の抽出

表 1 特徴表現手法

特徴	次元数	表現手法	コードブックサイズ
時空間特徴	257	BoSTF	5000
視覚特徴	96	BoFr	3000
動き特徴	56	BoFr	3000

提案されている<sup>12)</sup>。Sonnenburg らは単一カーネルの SVM 学習を反復することによって最適なカーネルの重み  $\beta_j$  を求める手法を提案している<sup>12)</sup>。

#### 4.2.2 特徴抽出

以下では統合に利用する動き特徴，視覚特徴，及び特徴の表現手法について述べる。

**動き特徴** 時空間特徴は，動き情報も内包しているが，それだけでは局所的な点の動きしか表現できない。よって本研究ではフレーム全体から動き特徴を抽出する。この動き特徴は，全体的な動きを表現出来るので，時空間特徴点の局所点の動き情報とは，異なる識別能力が期待できる。

動き特徴として，本研究ではグリッド点におけるオプティカルフローを利用する。それぞれの検出されたフローは，8 方向，7 段階の大きさからなるヒストグラムに投票されていく。図 8 は実際に抽出された動き特徴を示している。

**視覚特徴** 視覚特徴には 6 方向，4 周期のガボール特徴を使用する。その際に図 9 に示す通り，画像を  $2 \times 2$  に分割し各領域から 24 次元ベクトルを抽出する。領域ごとのベクトルを結合することによって一つの特徴とする。よって特徴の次元数は各領域のベクトル  $24 \text{ 次元} \times 4$  で表される。

この二つの特徴を時空間特徴と統合することで動作認識を行う。しかしフレーム一枚で認識しようとした場合，選ばれたフレームによって結果が大きく変わってしまうことがある。特に動き特徴の場合この傾向が顕著に現れる。そして，最も重要なキーフレームを選択することは非常に難しい問題である。よって本研究では bag of frames 表現を導入することで，この問題を解決する。

**特徴表現** 表 1 に特徴の表現手法についてまとめた。抽出された 257 次元の時空間特徴は BoSTF で表現される。この時のコードブックサイズは 5000 に設定した。

表 2 大量 Web 動画ショット分類のためのデータ

動作	動画数	ショット数	平均時間 [秒]	合計時間 [時間]	学習データ	
					positive	negative
batting	174	8,980	5.9	14.6	31	75
running	170	7,342	6.6	13.4	28	66
walking	174	6,567	7.4	13.4	23	63
shoot	164	7,718	5.3	11.3	14	75
eating	142	3,442	7.7	7.3	22	64
jumping	160	3,130	6.6	5.8	27	40
dancing	185	8,235	6.1	13.9	-	-
soccer	178	10,430	4.5	12.9	-	-
合計	1,247	55,844	6.3	92.6	145	383

視覚特徴，動き特徴を表現するために，従来研究ではキーフレームから特徴抽出を行うことが多いが，本研究では bag of frames(BoFr) 表現で視覚，動き特徴を表現する。これはフレームから抽出された特徴をベクトル量子化して，その特徴の出現頻度で動画を表現する手法である。ただしこの時のコードブックサイズは 3000 に設定した。

カメラモーションが検出された場合，時空間特徴と動き特徴では特徴が破棄されるが，視覚特徴は抽出が行われる。ショット全体がカメラモーションを含んでいるショットでは，時空間特徴や動き特徴が検出されない場合がある。この場合，時空間特徴，動き特徴のベクトルは 0 ベクトルで表される。

## 5. 実験

### 5.1 データセット

教師信号ありランキングのデータセットとして，独自に収集した batting, running, walking, shoot, jumping, eating の 6 つの動作を利用した。表 2(上) に詳細なデータを記した。合計 984 の動画をショット分割して，得られた 37,179 のショットを本実験では使用した。現時点で Web 動作分類で用いたデータセットは最も多い Liu らの Wild Youtube データセットでも 1500 本程度のショット数であったが，本データセットはその 20 倍以上の 37,179 本ものショットからなっている。そのデータの中から教師あり学習では，表 2 右のショット数を学習データとして使用している。

教師なし Web 動画クラスタリングでは Youtube から soccer, dancing のキーワードで収集されたものを利用する。データの詳細については表 2(下) に記載する。

### 5.2 実験結果

**KTH データセットにおける動作分類** まず特徴の分類精度の評価を行うために KTH データセットを用いて，動作分類を行った。結果を表 3 に記載する。ただし表 3 の VMR というのは本論文で提案した時空間特徴を，Point は 8) で提案した時空間特徴をそれぞれ

表 3 KTH データセットにおける分類率の比較

	Point <sup>8)</sup>	VMR	MKL	Liu <sup>6)</sup>	Lin <sup>13)</sup>	Gilbert <sup>14)</sup>
walking	0.94	0.98	0.99	0.96	0.93	0.94
jogging	0.76	0.84	0.94	0.83	0.91	0.91
running	0.81	0.83	0.85	0.84	0.85	0.89
boxing	0.91	0.92	0.96	0.99	0.96	1.0
waving	0.90	0.98	0.98	0.95	0.96	0.99
clapping	0.86	0.91	0.93	0.94	0.99	0.94
average	0.86	0.91	0.94	0.92	0.93	0.95

表 4 特徴抽出速度の比較 (単位は秒)

ours	CHLAC	Kläser
1.38	18.62	125

示している．比較に用いた手法は Liu らの研究<sup>6)</sup>, Lin らの研究<sup>13)</sup>, Gilbert らの研究<sup>14)</sup> の 3 つである．本提案手法が 94.0%, Liu らの手法で 91.8%, Lin らは 93.3%, Gilbert らは 94.5% の分類率であった．これらのことから, 提案した分類手法は非常に有用であることが分かる．これらの結果の詳細については 15) に述べる

次に表 4 は特徴抽出手法毎の抽出速度の比較を行ったものである．ただし抽出に利用したビデオはフレームサイズ  $80 \times 60$ , フレーム数 312 の映像である．比較に利用した特徴は, CHLAC<sup>1)</sup> と Kläser らによる勾配ベースの手法<sup>4)</sup> である．ただし実行環境は AMD Phenom II X4 3.0GHz, 8G メモリで行った．CHLAC では 125 秒, 勾配ベースの手法で 18.62 秒だったのに対し, 本提案手法では 1.38 秒で抽出を完了した．これは他の手法に比べ, 高速に特徴を抽出が出来ていることを示している．

教師信号ありの分類 ここではランダムに選択された動画をアノテーションし, 学習データとする．その後テストデータをランキング付けをする．図 10 はその結果を示している．

この図の横軸はランキング N 位を, 縦軸は適合率を示している．例えば, batting ではランキングの上位 20 位までの適合率が 1.0, 上位 40 位まで見ると適合率が 0.975 であることを示している．

実験として, 提案した時空間特徴, 視覚特徴, 動き特徴, ランダムサンプリング及び MKL による特徴統合の 5 つの手法の比較を行っている．ランダムサンプリングは 200 のショットをランダムに選択し, その適合率を計算したものである．

結果は eating 以外では提案した時空間特徴が最も良い結果で上位 20 位までの平均は 0.91, 上位 200 位までの適合率の平均は 0.57 であった．

eating で動きが良かった理由としては, eating 特有の上下への動きが抽出されたためと考えられる．全体的にランキングが下るにつれてランダムに近づく傾向があるが, 視覚特徴で running の分類を行った際, ランダムとほぼ同等の値となっている．これは running が

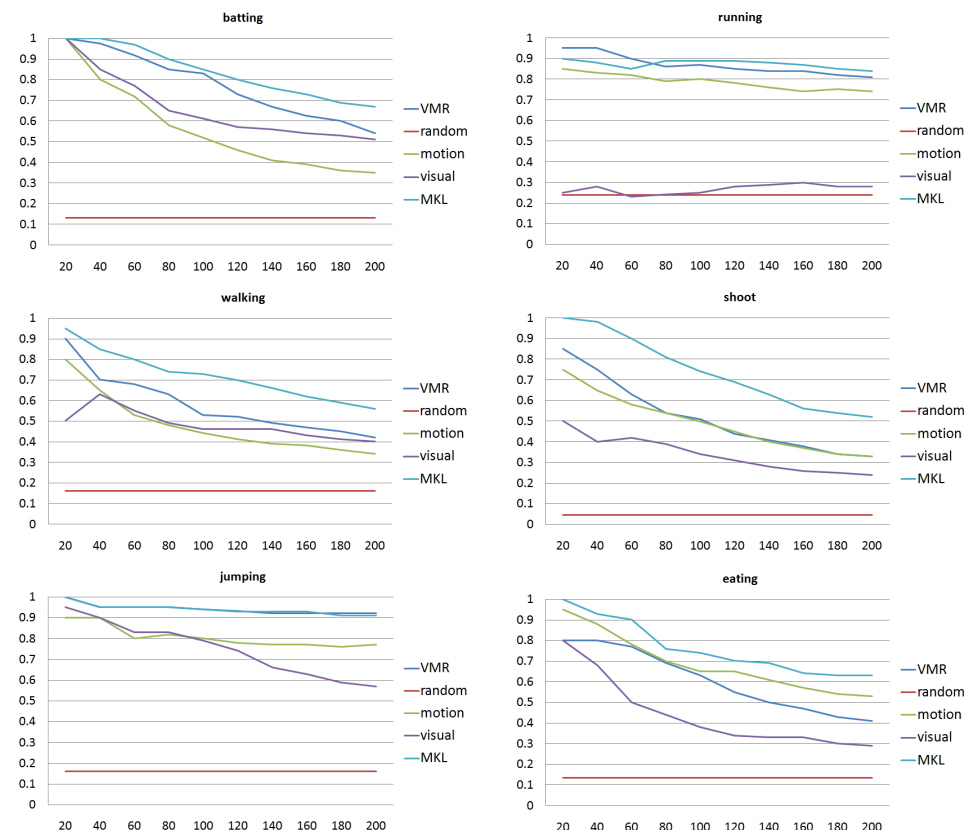


図 10 Web 動画ランキング付けの結果

行われる場所の多様性によるものである．

図 11 は各動作のランキング上位 100 位までの適合率を示している．すべての動作において MKL による特徴統合の値が最も良くなっている．特に walking や shoot のように単一の特徴で上手く検出が出来なかった動作においては精度の向上が顕著に現れている．一方で running や jumping のような単一特徴での適合率が高いものでは大きな変化はなかった．また図 12 はそのときの各特徴の重みを示している．

教師信号なしクラスタリング まず, pLSA や k-means クラスタリングでクラスタリン

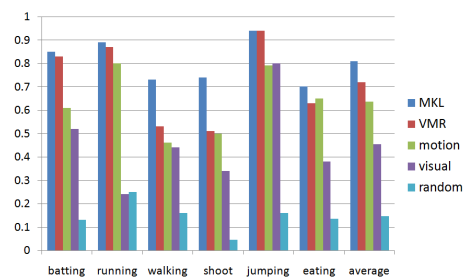


図 11 ランキング上位 100 位までの適合率

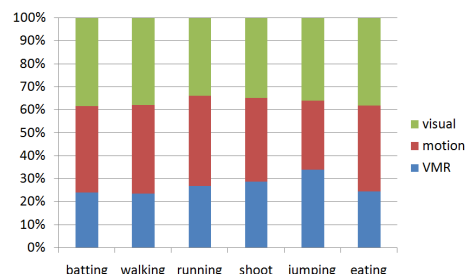


図 12 ランキング特徴毎の重み



図 13 各クラスに付与されたラベル:上が dancing に付与されたラベル, 下が soccer に付与されたラベル

グした後, 最も要素の多かった 10 のクラスにラベルをつける. つけられたラベルに関しては図 13 に示す通り, dancing では「踊っている (dancing)」、「歌っている (sing)」、「映画やドラマの映像 (acting)」、「その他 (others)」に分ける. また soccer では「遠いアングルで撮影された試合動画 (play-far)」、「近いアングルで撮影された試合動画 (play-near)」、「インタビューなど人の会話シーン (talking)」、「その他 (others)」のラベルを各クラスに付与した.

図 14 は各クラスタリング結果の要素数の大きかった 1 位から 10 位までのクラスのショットに, このラベルを付与したときの各ラベルの割合を示している. 例えば k-means で dancing をクラスタリングした結果, 最も要素数の大きかったクラス 1 は, dancing ショットを 65.0%含み, sing ショットを 6.0%含んでいることを意味している. ただし k-means, pLSA 共にクラスタ数は 200 に設定した.

同一のクラスのショットは類似していることが望まれるので, 本実験では, このラベルの割合が偏っているほど優れた結果であるといえる.

それぞれの動作毎に見ていくと, pLSA ではどのクラスにおいても 4 つのラベルが多く混在してしまっている. 一方で, k-means クラスタリングでは, クラス 5, 8, 10 を除

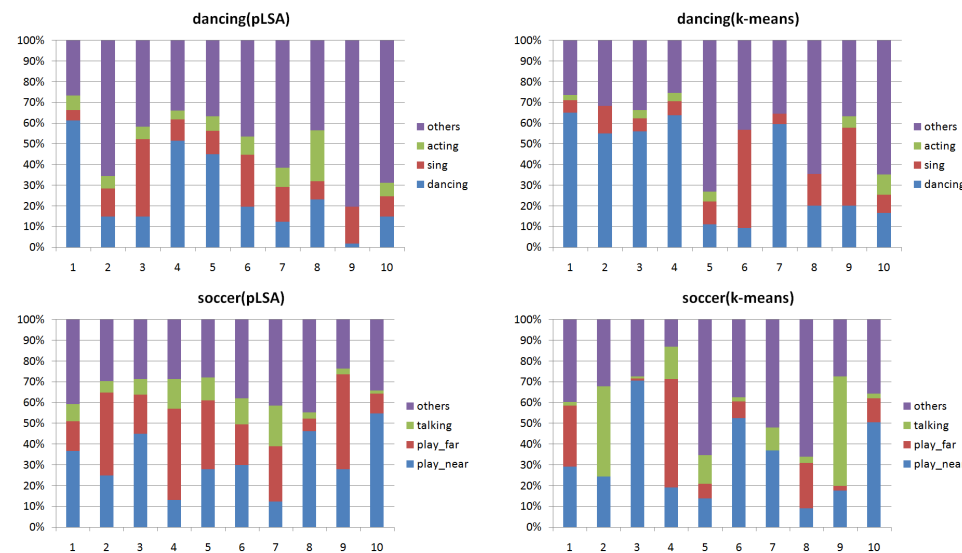


図 14 教師なしクラスタリングの結果

いて, dancing や sing のような特定のラベルの割合が高い場合, 他のラベルの割合は下がる傾向にあり, より正しい分類が出来ている. ここで k-means クラスタリングで, 5, 8, 10 のクラスの結果が悪かった理由として, k-means はその性質上, ノイズで構成されるクラスが構築されることがある. 実際これらのクラスでは others が他のクラスより高くなっている. よってこれらのクラスはノイズが多く集められたクラスであったため, 正しい分類が出来なかったと考えられる.

次に soccer のクラスタリング結果であるが, pLSA では play-far と play-near を正しく分類することが出来ず, 全てのクラスでこの二つのラベルが混在している. また, talking ラベルがどのクラスにも均一に存在してしまっている. 一方で k-means クラスタリングでは pLSA 分類に比べ, 各ラベルの偏りが大きくなっており, 正しく分類が出来ていることが分かる.

今回使用したデータは大量のノイズを含んでおり, いずれの分類手法にしても, その影響を受けてしまっている. 特に pLSA ではそれが顕著に現れてしまい, 精度が大幅に下がってしまった.

図 15 は k-means クラスタリングによる dancing のクラスタリング結果の一例を示している.

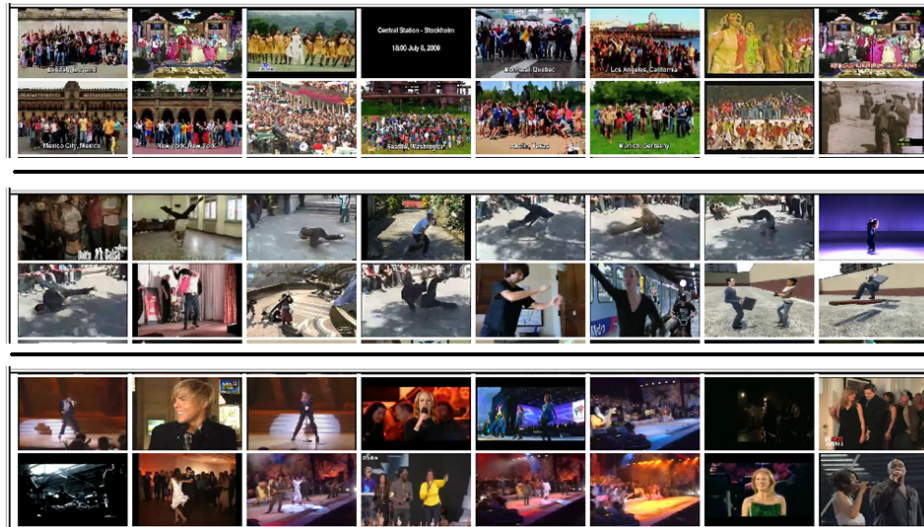


図 15 k-means クラスタリングによる dancing の分類結果の例

## 6. おわりに

本研究では Web 動画分類のための時空間特徴抽出手法について提案し、その特徴を利用した Web 動画分類を行った。提案した時空間特徴は視覚的な局所特徴と、その点の動きを特徴化することで抽出を行った。

実験として Web から大量に収集したショット分類を行った。分類には教師信号ありと、教師信号なしの 2 種類の分類を試みた。教師信号あり分類実験では、ランキングの上位 20 位までについて 91.7%，上位 200 位までについて約 57.2% の適合率であった。また教師信号なしのクラスタリングでは、pLSA と k-means クラスタリングによる分類を行った。k-means クラスタリングは pLSA より正確に分類できた。しかしデータセットのノイズの影響が大きく受けてしまい、正確な分類は非常に困難であった。

時空間特徴の改良点として、現状ではカメラモーションを検出した場合、その特徴を排除することで対応してきた。しかし、それでは有用な特徴を抽出出来ない場合がある。よってカメラモーションをより精巧に取り扱うため、動き補正を入れることでモーションを検出した場合も正しい特徴を抽出するシステムを構築していくことが必要である。

また提案手法では一つのショットから大量の特徴が抽出されてしまっている。より良い分

類を行うために、有用な特徴を選択して抽出するようしなければならない。

Youtube のような動画には複数の人間が動作を行っている「歩いている人」の隣で「走っている人」がいるときもある。しかし現在のシステムでは、このように同時に起こる、異なる動作を検出することができない。このような動作を検出するために、人検出とトラッキングに基づく、抽出手法を構築していくことも、精度を向上させるために有効なことである。また動画像からの特徴だけでなく、ほとんどの Web 動画に付けられているタグなどのメタデータとの融合も教師なし分類の精度向上ためには有効な方法である。

## 参 考 文 献

- 1) T.Kobayashi and N.Otsu. A three-way auto-correlation based approach to human identification by gait. In *Proc. of IEEE Workshop on Visual Surveillance*, pp. 185–192, 2006.
- 2) I.Laptev and T.Lindeberg. Local descriptors for spatio-temporal recognition. In *Proc. of IEEE International Conference on Computer Vision*, 2003.
- 3) P.Dollar, G.Cottrell, and S.Belongie. Behavior recognition via sparse spatio-temporal features. In *Proc. of Surveillance and Performance Evaluation of Tracking and Surveillance*, pp. 65–72, 2005.
- 4) A. Kläser, M. Marszałek, and C. Schmid. A spatio-temporal descriptor based on 3d-gradients. In *British Machine Vision Conference*, pp. 995–1004, sep 2008.
- 5) R.I. Cinbins, R.Cinbins, and S.Sclaroff. Learning action from the web. In *Proc. of IEEE International Conference on Computer Vision*, pp. 995–1002, 2009.
- 6) J.Liu, J.Luo, and M.Shah. Recognizing realistic action from videos. In *Proc. of IEEE Computer Vision and Pattern Recognition*, pp. 1–8, 2009.
- 7) D.Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, pp. 91–110, 2004.
- 8) A.Noguchi and K.Yanai. Extracting spatio-temporal local features considering consecutiveness of motions. In *Proc. of Asian Conference on Computer Vision (ACCV)*, 2009.
- 9) B.Herbert, E.Andreas, T.Tinne, and G.Luc. Surf: Speeded up robust features. *Computer Vision and Image Understanding*, pp. 346–359, 2008.
- 10) B.Lucas and T.Kanade. An iterative image registration technique with an application to stereo vision. In *Proc. of International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.
- 11) G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan. Learning the kernel matrix with semidefinite programming. *Journal of Machine Learning Research*, Vol.5, pp. 27–72, 2004.
- 12) S.Sonnenburg, G.Rätsch, C.Schäfer, and B.Schölkopf. Large scale multiple kernel learning. *Journal of Machine Learning Research*, Vol.7, pp. 1531–1565, 2006.
- 13) Z.Lin, Z.Jiang, and L.S.davis. Recognizing action by shape-motion prototype trees. In *Proc. of IEEE International Conference on Computer Vision*, pp. 444–451, 2009.
- 14) A. Gilbert, J. Illingworth, and R. Bowden. Fast realistic multi-action recognition using mined dense spatio-temporal features. In *Proc. of IEEE International Conference on Computer Vision*, pp. 925–931, 2009.
- 15) 野口顕嗣, 柳井啓司. 多種類特徴統合による動作認識手法の提案. 電子情報通信学会研究会報告: パターン認識・メディア理解研究会, 2010.