人体姿勢アノテーション困難な映像における 類似姿勢学習の有用性

村上 達哉^{1,a)} 三輪 誠¹ 浮田 宗伯¹

概要:我々は,既存の大量映像コンテンツからの演技特徴データセット生成を目指している.こうした演 技特徴は,人の姿勢により表現可能であるため,このデータセット生成のためには,映像からの姿勢推定 が基礎技術として必要になる.しかし,現在の姿勢推定手法では,演技特徴を表現するに足るほどの高精 度な姿勢推定は実現できていない.本研究では,深層学習による姿勢推定モデルにおいて,そのモデルを 特定の動きに特化して最適化することで,姿勢推定の精度にどのような影響が与えられるかを検証する. 本稿では,特徴的な動きが色濃く表れる時代劇の殺陣を解析対象とした.時代劇の演者のほとんどは,着 物を着ており,その姿勢推定モデル学習のために姿勢アノテーションを与えることが困難である.そこで, 姿勢推定対象となる演者と同様の動きを,モーションキャプチャシステムと同期させつつ映像撮影するこ とで,各フレームへの姿勢アノテーションを自動化した.しかし,このアノテーション映像の演者は,モー ションキャプチャ計測のために着物を着用していない.このように,姿勢推定モデルが学習している情報 のうち「見え方」は大きく異なるが「人体パーツ間の相対的な位置関係」は類似している学習データを利 用することで,姿勢推定モデルの最適化がどのように影響を受けるのかを検証するのが,本稿の目標であ る.実験では,人体パーツによっては最大6%程度の性能向上が得られることを確認した.

キーワード:姿勢推定,単眼画像,深層学習,畳み込みニューラルネットワーク

1. まえがき

3DCG アニメーションの需要の高まりに伴い,キャラク ターを動かすアニメーターの負担が増大している.モー ションキャプチャシステムを用いれば,アニメーターの負 担は大きく減るが,マーカーの位置から姿勢を取得する分, 得られるデータはマーカー位置のずれによるノイズが避け られない.さらに人物の動きをより魅力的にするために, アニメーターによる調整はかかせない.このように,アニ メータにかかる負担は3DCG アニメーション制作の発展 を阻害する要因となっている.

一方で、デジタルデータである 3DCG アニメーション は、一度動きのデータセットを作成すれば制作費の大幅な 削減が可能になる.しかし、人間の演技特徴の多様性を考 慮すると、一つ一つ動きのデータを作ることは現実的では ない、そこで、既に人間の多様な演技特徴の情報を含んだ データである、既存の映像コンテンツの活用を考える. 本研究では、動きのデータセット作成のために、大量の

1 豊田工業大学 Toyota Technological Institute, Japan

^{a)} sd14072@toyota-ti.ac.jp



図 1 2次元姿勢推定の結果.RGB 画像を入力として,その画像に写 る人物のキーポイントの画像上における位置を推定している. 図中では,キーポイントを色付きの円として表現している.

映像コンテンツから演技の特徴を解析することを最終的な 目標とする.そうした演技の特徴は,人の3次元的な姿勢 で定義されるので,特徴解析には姿勢推定技術を利用する. 姿勢推定技術とは,画像に写る人物の関節や端点(例:肩, 肘,手首)といったキーポイントの位置を画像から推定す る技術(図1)である.キーポイントの位置を画像上の位 置として推定する場合は2次元姿勢推定,空間上の位置と して推定する場合は3次元姿勢推定と呼ばれる.

姿勢推定技術は,深層学習を用いた手法によって,推定 精度は大幅に向上しており,特に2次元姿勢推定において 顕著である.しかし,現状の3次元姿勢推定技術では,演 技の特徴を得るに足る,高精度な推定は実現されていない. そこで,汎化ではなく,特定の個人/動きに特化すること で3次元姿勢推定の精度を向上させる.推定精度の向上だ けでなく,その対象が頻繁にとる姿勢や,他の人には見ら れない姿勢といった,動きの特徴を解析するための特徴抽 出も可能になることが期待される.

このように,特定の個人/動きを表現するには,キーポイントの相対的な位置関係が重要である.例えば,ものまね 芸人の動きは姿勢情報がオリジナルのものと物理的に異なるにも関わらず,われわれ人間にはオリジナルと似て見える.この時,ものまね芸人はオリジナルと似たような動きをするので,両者のキーポイントの相対的な位置関係は類似している.すなわち,似ていると感じる要因の一つは,両者のキーポイントの相対的な位置関係の類似性である.したがって,姿勢推定においてもこの情報を重視することで,より特徴を捉えた姿勢推定が可能になると考えられる.

この考えに基づき,今回は修士研究に向けた予備実験として,「特定の個人/動き」のキーポイントの相対的な位置 関係を重視することによって,その「特定の個人/動き」を 対象にしてどれほど姿勢推定精度を向上できるかを検証し た.演技の特徴解析や3次元 CG への応用のためには,本 来3次元姿勢推定が必要であるが,先に述べた通り3次元 姿勢推定そのものの誤差が大きく,実験の際に誤差の原因 を特定することが困難になりかねない.そこで本課題研究 では,近年精度が急激に向上してきた2次元姿勢推定を対 象として基礎検討を行う.

2. 関連研究

スポーツや医療目的の運動解析,ゲームなどのエンター テインメント用途への応用が期待され,姿勢推定に関する 研究はコンピュータヴィジョンの分野で盛んに研究されて きた.現在の姿勢推定に関する研究は,入力となる画像の 種類によって2種類に分けることができる.RGBD画像 と,RGB画像である.

RGBD 画像では,RGB 画像に比べ,深度情報(Depth)を 利用できる分,リアルタイムの3次元姿勢推定問題を単純化 できる[1].しかし,深度情報を得るにはKinect[2],[3],[4] に代表される RGBD カメラが必要であり,構築した姿勢推 定モデルを適応する対象も RGBD 画像に限られてしまう. そうした制限を持つ RGBD 画像と比べ,RGB 画像からの 姿勢推定は,既存の大量の画像データ(Youtube, Netflix など)を利用できるので,人間の動作解析の面で期待され ている. RGB 画像からの 2 次元姿勢推定研究の初期には, グラ フィカルモデルをベースとした姿勢推定手法 [5], [6], [7] が 多い.このモデルでは,構成されるグラフが木構造となる ため,高効率な推定が可能であった.しかし,隣接してい ないパーツ間の依存関係が表現できていなかった.この問 題を解決するために,様々な手法が考案されてきた.例え ば,階層構造を取ったモデル [8] は,木構造による高効率 性という利点を維持しつつ,パーツ間の依存関係の表現力 を高めている.他には,パーツ間のエッジを増やすことに よって,表現力を高めた非木構造のものがある [9].このモ デルはパーツ間の依存関係を十分に表現できるものの,グ ラフにループが存在することで高効率な計算を行えず,そ の予測値は近似した計算結果である.

このように, グラフィカルモデルベースの手法には隣接 しないパーツ間の表現に問題があったが,現在優れた結果 を残している2次元姿勢推定手法のほとんどが, CNN ベー スの手法である.CNN を利用した姿勢推定手法の内で最も 初期の取り組み [10] では, CNN を利用して回帰的に姿勢を 推定している.CNN が利用され始めた頃にはグラフィカ ルモデルと CNN を組み合わせたものもあり,部分的なキー ポイントの画像を,隣接するキーポイント位置の推定に役 立てている[11].こうした初期の取り組みと異なる,以降 の2次元姿勢推定研究に大きな影響を与えたのが,ヒート マップと CNN を組み合わせた手法である [12], [13], [14]. 代表的なものは,多段階の手法 [15] に CNN を利用するこ とで,パーツ間の複雑な依存関係を捉えることを可能にし た [14]. このように, CNN を利用した2次元姿勢推定手法 は優れた結果を示し,その関心は「複数人に対する2次元 姿勢推定」へと広がっている[16],[17],[18].

RGB 画像からの3次元姿勢推定でも CNN を用いた手法 により、その推定精度を向上させている[19],[20],[21],[22]. 時系列情報を用いてより安定で正確な姿勢推定を実現して いるもの [23], [24] や, リアルタイム 3 次元姿勢推定 [25] の研究も行われている.一方で,3次元姿勢推定は2次元 姿勢推定と比べるといまだに難しいタスクであり続けてい る.2次元的に体の姿勢が得られたとしても,3次元的に みるとその姿勢を取りうる姿勢は幾通りも存在するからで ある.加えて,学習データが十分でないことも,3次元姿 勢推定を難しいタスクにしている.2 次元姿勢推定におけ る深層学習の成功の大きな要因の一つは大量のデータセッ ト [26], [27] が確保できたことであり,量と種類は3次元姿 勢推定用のもの [28], [29] と比べると非常に多い. CG 技術 を用いてデータセットを作る動き [30], [31] もあるが,これ は本研究の主目的ではない.本研究では演技特徴を表現で きるようなモデルを作りたいので,シーケンス情報を含ん だ実際の人物の画像データセット [32] が重要になる.

IPSJ SIG Technical Report



図 2 姿勢推定ライブラリ「OpenPose」で使われるアーキテクチャ. Branch 1(橙色)からは PCM が出力され, Branch 2(青色) からは PAF が出力される. Stage 1 において入力は画像から 取り出した特徴量(F)のみであるが, 点線で囲まれた Stage t (t ≥ 2)においては特徴量 F に加え, Stage t-1 における出 力も入力となる.黄色い四角は畳み込み層(C)を, 灰色の四 角は損失関数(L)を示している.



Stage 1Stage 3Stage 6

図 3 Part Confidence Map(PCM) と Part Affinity Field(PAF). 上段が右手首の PCM,下段が右前腕の PAF を示している. 橙色の円で囲んだ場所は,ステージを経るごとに推定できるようになっているのが分かる.赤色の円で囲んだ場所は,ステージを経るごとに間違えなくなっているのが分かる.図は[18] より引用.

3. OpenPoseの姿勢推定アルゴリズム

実験では姿勢推定ライブラリ「OpenPose」[14], [18], [33] を使用した.OpenPose は CNN ベースの精度の高い2次 元姿勢推定を可能にするライブラリであり,細部にわたり 整備され研究に際し扱いやすいため,実験に用いた.

図 2 に, OpenPose で使用される姿勢推定アルゴリズム の構造を示す. OpenPose の姿勢推定アルゴリズムは,Part Confidence Map (PCM), と Part Affinity Field (PAF) と呼ばれる二つの特徴量と,画像から直接得られる特徴 量を CNN により多段階にわたって繰り返し抽出すること で,姿勢推定の精度を高めている.ネットワークが PCM や PAF といった特徴量を Branch 1, Branch 2 から繰り返 し出力できるように,各ステージの終わりにそれぞれ異な る損失関数 f1,f2 を適用している.具体的な PCM,PAF の画像を図 3 に示す.

PCM とは,姿勢のキーポイントの位置の信頼度を示す ヒートマップのことであり,画像上のスカラー場として表 現される.図3では,上段に右手首に関する PCM の画像 を挙げている.この図を見ると,Stage1では左手首を誤っ て右手首だと判断しているケース(図中赤丸)があるが, 多段階にわたる処理を行うことで誤りを修正し,推定精度 を向上させている.このヒートマップのピーク値を取得す ることによって、姿勢の座標を取得することができるが、 OpenPose はボトムアップに PCM を作るため,画像上に 複数人が存在する場合,問題が生じる.例えば,図3の上 段には右手首の PCM を記載したのだが,最終的な Stage 6において,二人の人物の右手首が画像から検出されてお り,このままでは検出したキーポイントを結んで一人の姿 勢とすることができない.そのため,複数人の姿勢推定を 行いたい場合には PCM だけでは不十分である.人体検知 の手法を用いて,画像中の人物を一人検出したあとならば, その一人に対して CPM を用いた姿勢推定を行うことが可 能である、しかしこうしたトップダウンのアプローチでは 一人ずつキーポイント検出を行う必要があるため,時間が かかってしまう.

OpenPose の姿勢推定アルゴリズムは, PCM の他に PAF を用いることで、この問題を解決し、リアルタイムでの複 数人の姿勢推定を可能にしている.PCM を用いて特定さ れた複数人のキーポイント位置から,各人の姿勢を得るに は各キーポイントを個々人に結び付ける必要がある.そこ で使用されるのが PAF である. PAF は, ベクトル場であ り,画像における個人の四肢上においては,その四肢が結 ぶキーポイントに沿った方向を向いたベクトルになるよう に設計されている.よって, PAFを用いると PCM によっ て得られたキーポイントを結ぶことができ,各個人の姿勢 を推定することができる.さらには PAF を利用してキー ポイントを結ぶ際には貪欲法を用いた高速な処理が行われ ているため,リアルタイムでの複数人に対する姿勢推定が 可能となっている.図3の下段には右上腕におけるPAF を記載した. PCM と同様に,多段階処理からその精度を 向上させている.今回の実験では対象とする人物が一人だ けなので, PAF による姿勢推定精度の向上は限定的であ ると予想される.しかし,今後の研究にて解析対象とする 映画のような一般的な映像では, 複数人が重なり合うこと による遮蔽が頻繁に生じる.そういった状況においては, PAF が姿勢推定の精度向上に大きく寄与することが期待さ れる.

4. 類似姿勢学習の有用性の検証法

今回の実験では、「特定の個人/動き」のキーポイントの 相対的な位置関係を重視することによって、「特定の個人/ 動き」に対する姿勢推定精度をどれほど向上できるかの検 証を目標とする.

この検証のためには,実験結果における違いが際立つように,一般的な人物の動きの特徴ではなく,個人の演技特

情報処理学会研究報告

IPSJ SIG Technical Report





図 4 同じような動きのシーケンス.解析対象の動き(上段)と同じ ような動き(下段)をモーションキャプチャシステムによって 計測することで,解析対象における,キーポイントの相対的な 位置関係と似たような情報を取得できる.



図 5 モーションキャプチャシステム (ツークン研究所). 収録エリ アサイズは 10m × 7m × 2.5m で, 24 台の Vicon 社「T160」 という光学式カメラを用いて撮影できる.

徴が色濃く表れる映像が有効である.そうした映像として, 日本固有のコンテンツである時代劇,特に典型的な演技と して,殺陣を実験対象とする.理想的には,時代劇の殺陣 に対して高精度な姿勢推定モデルを作成したい場合,キー ポイントの相対的な位置関係や,見え方に対する姿勢をモ デル化するために,学習データにも時代劇の殺陣の画像を 使用するほうが良い.しかし,時代劇では,登場人物が着 物という遮蔽の大きな衣服を着用しているため,体のパー ツが画像から確認しずらく,画像に対する姿勢のマニュア ルアノテーションが困難である.そのため学習データを大 量に獲得することが難しい.そこで,解析対象とする殺陣 の動きと類似した動きをモーションキャプチャシステムを 用いて計測した(図4). これにより,見え方は大きく異な るものの,殺陣の動きにおけるキーポイントの相対的な位 置関係と類似した情報を持つ学習データを、マニュアルア ノテーションなしに大量に作成することが可能である.計 測したデータから画像上のアノテーションデータを得る具 体的な方法は 5.2 節で述べる.こうして得られた学習デー タは,姿勢推定の対象となる画像(図4上段)を解析する 際に有用な,キーポイントの相対的な位置関係に関する情 報を含んでいる.一方で,学習データ中の人物はモーショ ンキャプチャスーツしか着ておらず,見え方に関する情報 が不足している.この情報不足を補うために,作成した学 習データを用いて姿勢推定モデルを構築する際には,見え 方について多様な情報を含んだデータセット [26] から学習 したモデルをファインチューンする.作成したモデルを用 いて,純粋に類似姿勢学習の有用性を検証するには,ノイ ズや遮蔽の少ない画像が望まれる.解析対象(図4上段) はノイズや遮蔽が多いため、検証に向いていない.よって、 解析対象と似た動きを,着物を着た状態で撮影し,検証時 のテストデータとする.



図 6 計測したデータ.(a)解析対象となる映像.着物の遮蔽により キーポイントの位置が特定しにくN.(b)モーションキャプ チャシステムでマーカー計測時の映像.解析対象と見え方は 大きく異なるが,キーポイントの相対的な位置関係は類似して いる.(c)モーションキャプチャシステムから計測した,マー カーの3次元座標の可視化.橙色の点がマーカーに対応して おり,その3次元的な座標を,CGソフトウェアを介して取得 することができる.(d)テスト用の映像.解析対象の映像は, ノイズや他の人物による遮蔽が大きく,姿勢推定精度に与える 影響が強N.そのため,テスト用画像を作成するために,着物 を着た状態での類似した動きを撮影した.

5. 実験

5.1 殺陣データセットの計測

実験に必要となる殺陣のデータは,東映デジタルセン ターツークン研究所にて撮影を行った(図5).東映株式 会社の長年にわたる映像制作の実績から,ツークン研究所 では様々なジャンルのアクターを国内外から手配するこ とができる.今回は,解析したい映像における殺陣の動き (図6(a))を,プロの殺陣師にアクターとして再現しても らうことで,その殺陣の動きと類似した動きをモーション キャプチャシステムにて計測(図6(c)),それと対応した映 像(図6(b))も撮影した.また,本来であれば構築した姿 勢推定モデルの精度を検証するために,解析対象となる映 像(図6(a))でテストを行いたいが,その映像はノイズや



図 7 カチンコによる映像とモーションデータの同期. 同期のため には,カメラとモーションキャプチャ用カメラから,撮影が開 始する瞬間を確認できる必要がある.撮影開始時(図中 246 frame)に強い瞬間的な光を発する(図中青丸)ことによって, 画像上からも,モーションキャプチャシステムからも,撮影開 始のタイミングが判断できる.

他の人物による遮蔽が多く,実験の狙いに関係なく,姿勢 推定精度に影響を与えてしまう可能性がある.そのため, 今回はテスト用の映像として,着物を着た状態での類似し た動き(図 6(d))の撮影も行った.

データセットの計測に用いた機器はそれぞれ以下に述べ る.マーカーの3次元座標(図6(c))の計測に用いたモー ションキャプチャ用のカメラは Vicon 社の T-160である. T-160は1600万画素(4704×3456)の解像度を持ち,フ ル解像度で120FPSのキャプチャが可能である.一方,学 習用の映像(図6(b))と,テスト用の映像(図6(d))の撮 影に用いたカメラは,Sony社のPXW-FS7で,約1160万 画素の解像度で60FPSの撮影を行える.

計測されたデータセットは,学習用,テスト用でそれぞ れシーケンス数が5つである.本実験では,その中でも特 に類似した動きを再現できている1シーケンスのデータを 用いて,学習用データセットを作成した.

5.2 学習用データセットの作成

モーションキャプチャシステムから得られたマーカー (図 6(c))のシーケンス情報と、それに対応する映像(図 6(b)) は、フレーム間の対応がついていないため、まずはデータ 間の同期を行う必要がある、このために、映画撮影の際な どに使われるカチンコから得られる情報を利用した.モー ションキャプチャシステムにおいては,光がカチンコの役 割を果たしている (図 7). 図中では 246 フレーム上で光 が発せられたと確認できるので,学習用の画像のシーケン スには246フレーム目からの画像を使用する.撮影終了時 も同じ原理でフレームを合わせる.結果,246フレームか ら 1073 フレーム目までの計 828 枚の画像を学習データに 使用することができる.画像に対する姿勢情報のアノテー ションに関しては,計測されたキーポイントの3次元座標 を,カメラ校正により得られたカメラパラメータによって, 2次元画像座標に変換する.カメラ校正のためには,モー ションキャプチャスーツに着けるマーカーの位置情報を利 用した.まずモーションキャプチャシステムによってマー



図 8 再最適化したモデルを用いた姿勢推定結果.着物の遮蔽と形 状変形によって推定エラー(図中赤丸)が生じている.

カーの3次元座標が取得できる(図6(c)). さらにマーカー は画像上からもその2次元座標を確認することができる (図6(b)). したがって,空間上の3次元座標とそれに対応 する画像上の2次元座標のペアを取得することができる. この座標のペアを50組ほど集め,空間座標から画像座標 への透視投影行列を作成した.この作成した透視投影行列 を用いることで,モーションキャプチャシステムから計測 されたキーポイントを3次元座標から2次元座標に投影す ることができる.具体的には,「1フレームごとに3次元座 標を取得し,2次元座標へ投影する」という処理をすべて CG ソフトウェア上で行うことができるので,1つの映像 シーケンスからアノテーション付きの画像828枚を容易に かつ大量に作ることができる.

5.3 実験結果

実験では,一般的な人の画像と姿勢情報を含むデータ セット COCO[26]で学習した姿勢モデルと,この姿勢モデ ルを殺陣の画像と姿勢の教師データ(Toei+COCOと記述) で再最適化した姿勢モデルを用いて,着物を着た殺陣の画 像に対して姿勢推定を行った.COCOは,人が密集した状 態,遮蔽が大きい状態など,様々な問題ケースの画像を含 んだデータセットであり,330Kの画像に対し,250,000 個 の姿勢情報がアノテーションされている.

図 8 に推定結果の可視化を示す.赤い丸で囲んだ部分で 推定エラーが生じているが,動きの激しさから着物が複雑 に変形している画像におけるエラーが多かった.これは見 え方から生じるエラーであり,連続的な動きの特徴を表現 できていないことが分かる.

また,図9に,COCOと,Toei+COCOで学習したモ デルのPCKhによる評価結果を示す.PCKhとは姿勢推 定においてしばしば用いられる指標であり,Part Correct Keypoint (Head)の頭字語である.これは,頭の大きさ を1として正規化した距離に基づき,ある閾値を定めたと きに,キーポイントの正解座標と推測結果との間の距離が, その閾値の範囲にどれほど収まっているかを表したもので ある.図9を見ると,限定的だが微小な精度向上を確認で きる.例えば図中右図の肘に関するPCKhでは,横軸が 0.11の所において,6パーセント(16パーセントから22 パーセントに上昇)の推定精度の向上を果たしている. IPSJ SIG Technical Report



図 9 体の各部分に関する PCKh. 青色の実線が, COCO のみで学 習したモデルの推定結果,赤色の点線が, Toei+COCO で学 習したモデルの推定結果である.体全体(PCKh total)でみ るとあまり違いは見られないが,各パーツで見ると,推定精度 が向上しているところも見られる.例えば,肘に関する PCKh (PCKh elbow)では,Normalized distance(横軸)が0.11 の所において精度が6パーセント向上している(図中黒丸).

さらに,図10に,COCOのみで学習したモデルと,Toei のみで学習したモデルのPCKhによる評価結果を示す. 多様な見え方の情報を含むCOCOで学習したモデルの方 が,より優れた推定結果を出している.

これらの実験結果より,特定の動きのデータセットで姿 勢推定モデルを再最適化することで,特定の動きに対する 姿勢推定の精度を高めることができる可能性を示せたが, 大幅な精度向上は示せなかった.

5.4 考察

精度向上が限定的だった大きな理由として,姿勢を静止 画で表現し,動きを表現できていなかったことが考えられ る.着物とスーツでは見え方が大きく異なり,キーポイン トの相対的な位置関係はスーツの画像に対して学習されて いるので,着物の画像に対しては学習結果を活かした姿勢 推定を行えなかった.そのため,相対的なキーポイントの 位置関係の情報を活用した姿勢推定を行いたいのであれ ば,その見え方に大きく依存しない手法が必要である.



図 10 体の各部分に関する PCKh. 青色の実線が, COCO のみで 学習したモデルの推定結果,赤色の点線が, Toeiのみで学習 したモデルの推定結果である.

6. むすび

6.1 まとめ

一般的な人物を対象にした3次元姿勢推定ではなく, 「特定の個人/動き」に特化した3次元姿勢推定という最終 目的への取り組みとして,今回は「特定の個人/動き」の キーポイントの相対的な位置関係を重視し,2次元姿勢推定 の精度を,「特定の個人/動き」に対して向上できるかどう かを検証した.この検証では,特徴的な動きが表れる「殺 陣」を対象にして,モーションキャプチャシステムから得 られる画像上の2次元姿勢と対応する画像のデータセット から,キーポイントの相対的な位置関係をモデル化した. このモデルを用いた検証のため,モーションキャプチャ撮 影時の動きを着物を着た状態で再現し,再現画像に対して 姿勢推定を行った.その結果,姿勢推定精度の向上は限定 的であり,「特定の個人/動き」に特化した姿勢推定を行う ためには別の手法が必要だと分かった.

検証に用いた姿勢推定アルゴリズムは,画像1枚のみを 入力として画像上の人物の姿勢を推定しており,画像上の 人物の見え方に大きく依存したアルゴリズムと言える.姿 勢推定モデルの構築には,モーションキャプチャースーツ の画像が使用されているため,着物を着た状態の画像に対 しては,モデル化したキーポイントの相対的な位置関係を, その見え方の違いから上手く姿勢推定に活かせなかった

6.2 今後の発展

検証に用いた姿勢推定アルゴリズムは,画像上の人物の 見え方に大きく依存したものであった.そこで,今後の研 究においては,時系列データを用いることで,着物映像の 姿勢推定精度を向上させていく.時系列データの利用で姿 勢推定精度を向上できることが分かっており,また,個人 に特化した姿勢推定を行う観点でも,時系列データの使用 は有効である.例えば,人の歩き方からその人物を特定す る,歩容認証と呼ばれる問題では,個人の動きの特徴は時 系列データに良く表れるこが明らかになっており,動きの 時系列データは個人に固有のものである.したがって,「特 定の個人/動き」の時系列データを利用することで,「特定 の個人/動き」に特化して高精度に働く姿勢推定モデルの 実現を目指していく.

謝辞 本研究で使用したデータは東映デジタルセンター ツークン研究所から提供されている.また,本研究は京都大 学インキュベーションプログラムの支援を受けて行った.

参考文献

- Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore, R., Kipman, A. and Blake, A. Real-time human pose recognition in parts from single depth images, *CVPR* (2011).
- [2] Kinect for Xbox 360, Microsoft Corporation (2010).
- Kinect for Xbox One, Microsoft Corporation, (online), available from (http://www.xbox.com/ja-jp/xboxone/accessories/kinect/) (2013).
- Kinect SDK, Microsoft Corporation, (online), available from (https://developer.microsoft.com/jajp/windows/kinect/) (2015).
- [5] Felzenszwalb, P. F. and Huttenlocher, D. P. Pictorial structures for object recognition, *IJCV*, Vol. 61, No. 1, pp. 55–79 (2005).
- [6] Ferrari, V., Marin-Jimenez, M. and Zisserman, A. Progressive search space reduction for human pose estimation, CVPR (2008).
- [7] Andriluka, M., Roth, S. and Schiele, B. Pictorial structures revisited: People detection and articulated pose estimation, *CVPR* (2009).
- [8] Sun, M. and Savarese, S. Articulated part-based model for joint object detection and pose estimation, *ICCV* (2011).
- [9] Dantone, M., Gall, J., Leistner, C. and Van Gool, L. Human pose estimation using body parts dependent joint regressors, *CVPR* (2013).
- [10] Toshev, A. and Szegedy, C. Deeppose: Human pose estimation via deep neural networks, CVPR (2014).
- [11] Chen, X. and Yuille, A. L. Articulated pose estimation by a graphical model with image dependent pairwise relations, *NIPS* (2014).
- [12] Tompson, J., Goroshin, R., Jain, A., LeCun, Y. and Bregler, C. Efficient object localization using convolutional networks, *CVPR* (2015).
- [13] Pfister, T., Charles, J. and Zisserman, A. Flowing con-

vnets for human pose estimation in videos, ICCV(2015).

- [14] Wei, S.-E., Ramakrishna, V., Kanade, T. and Sheikh, Y. Convolutional pose machines, *CVPR* (2016).
- [15] Ramakrishna, V., Munoz, D., Hebert, M., Bagnell, J. A. and Sheikh, Y. Pose machines: Articulated pose estimation via inference machines, *ECCV* (2014).
- [16] Pishchulin, L., Insafutdinov, E., Tang, S., Andres, B., Andriluka, M., Gehler, P. V. and Schiele, B. Deepcut: Joint subset partition and labeling for multi person pose estimation, *CVPR* (2016).
- [17] Insafutdinov, E., Pishchulin, L., Andres, B., Andriluka, M. and Schiele, B. Deepercut: A deeper, stronger, and faster multi-person pose estimation model, *ECCV* (2016).
- [18] Cao, Z., Simon, T., Wei, S.-E. and Sheikh, Y. Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields, *CVPR* (2017).
- [19] Zhou, X., Zhu, M., Leonardos, S., Derpanis, K. G. and Daniilidis, K. Sparseness meets deepness: 3D human pose estimation from monocular video, *CVPR* (2016).
- [20] Brau, E. and Jiang, H. 3D Human Pose Estimation via Deep Learning from 2D Annotations, 3DV (2016).
- [21] Tome, D., Russell, C. and Agapito, L. Lifting from the deep: Convolutional 3d pose estimation from a single image, *CVPR* (2017).
- [22] Martinez, J., Hossain, R., Romero, J. and Little, J. J. A simple yet effective baseline for 3d human pose estimation, *ICCV* (2017).
- [23] Lin, M., Lin, L., Liang, X., Wang, K. and Cheng, H. Recurrent 3D Pose Sequence Machines, *CVPR* (2017).
- [24] Hossain, M. R. I. and Little, J. J. Exploiting temporal information for 3D pose estimation, arXiv (2017).
- [25] Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.-P., Xu, W., Casas, D. and Theobalt, C. VNect: Real-time 3D Human Pose Estimation with a Single RGB Camera, (online), available from (http://gvv.mpi-inf.mpg.de/projects/VNect/) (2017).
- [26] Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P. and Zitnick, C. L. Microsoft COCO: common objects in context, *ECCV* (2014).
- [27] Andriluka, M., Pishchulin, L., Gehler, P. and Schiele, B. 2D Human Pose Estimation: New Benchmark and State of the Art Analysis, *CVPR* (2014).
- [28] Sigal, L., Balan, A. O. and Black, M. J. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion, *IJCV*, Vol. 87, No. 1-2, p. 4 (2010).
- [29] Ionescu, C., Papava, D., Olaru, V. and Sminchisescu, C. Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments, *TPAMI*, Vol. 36, No. 7, pp. 1325–1339 (2014).
- [30] Chen, W., Wang, H., Li, Y., Su, H., Wang, Z., Tu, C., Lischinski, D., Cohen-Or, D. and Chen, B. Synthesizing training images for boosting human 3d pose estimation, *3DV* (2016).
- [31] Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M. J., Laptev, I. and Schmid, C. Learning from Synthetic Humans, *CVPR* (2017).
- [32] Insafutdinov, E., Andriluka, M., Pishchulin, L., Tang, S., Levinkov, E., Andres, B. and Schiele, B. ArtTrack: Articulated Multi-person Tracking in the Wild, *CVPR* (2017).
- [33] Simon, T., Joo, H., Matthews, I. and Sheikh, Y. Hand Keypoint Detection in Single Images using Multiview Bootstrapping, *CVPR* (2017).