機械学習とのハイブリッド逆投影法による 鍵盤楽器演奏時の高精度マルチカメラハンドトラッキング

山本 和彦†

†ヤマハ株式会社



図1 鍵盤楽器演奏時の複数視点のカメラ画像と音が与えられたとき,両手の手関節の3次元位置とカメラ姿勢,鍵盤平 面を非接触で同時推定する(左).これは逆投影問題として定式化することによって効率的に解くことができる.複数視点 カメラ画像からのバンドル調整と機械学習とのハイブリッドな手法によって,従来手法では困難な自己遮蔽の大きい状況 でも正確に手の3次元形状を絶対スケールとともにキャプチャすることができる(右).

概要

モーションキャプチャにおいて,外部物体との複雑なイン タラクションを伴う手指動作の詳細なトラッキングをする ことは最も難易度の高いタスクである.とくに,楽器演奏時 には,手にグローブやセンサ類を装着する従来手法では自然 な動きを阻害してしまう,楽器や手そのものによる遮蔽が発 生してしまう,という致命的な問題がある.そこで本論では, 複数の RGB カメラと音の情報のみを使い,手に一切のセン サ類を装着することなく鍵盤演奏時の両手の詳細な3次元 モーションキャプチャをおこなう手法を提案する.これを実 現するため,バンドル調整で一般的な再投影誤差の最小化で はなく逆投影問題として定式化し,機械学習とのハイブリッ ドな手法でこれを効率的に解く.検証の結果,演奏動作を妨 げることのない高精度の非接触ハンドトラッキングが実現 できたことを示す.

1 はじめに

楽器演奏時の手指動作の詳細なキャプチャをすることは 高い需要があるにもかかわらず,非常に難易度が高く,挑戦 的なタスクとなっている.従来,手にマーカーやグローブな どのセンサを装着する手法が一般的であるが,手の自然な動 きを阻害したり,手の自己遮蔽によるセンサ認識ミスなどの 致命的な問題が発生する.手に一切のセンサ類を装着せず画 像情報 (深度含む) からの動作認識 [1] も深層学習の発展に 伴い性能が向上しているが, それぞれの手の関節と両手の絶 対位置関係を同時に推定することが困難であったり, 手首な どを基準とした相対位置しか計算できない, 一部が遮蔽され た手関節の推定精度が著しく低下するといった問題がある.

そこで本論では鍵盤楽器演奏時の両手の手指の関節の詳 細な動きを複数の RGB カメラと音情報を使って非接触で トラッキングすることのできる手法を提案する.提案手法 では複数視点のカメラ画像と音の情報から手関節の3次元 位置とカメラ姿勢, 鍵盤平面を同時推定する. これを実現す るために、カメラ画像から2次元の手関節位置推定をまずお こない、世界座標系にある3次元手関節点からのレイがそれ ぞれのカメラへ到達し、推定された2次元関節点へと投影さ れるプロセスを逆に辿ることでつじつまを合わせる, 逆投影 問題として定式化をおこなう.このとき、補助的に2次元手 関節変位を入力として3次元化をおこなう事前学習済みの 深層学習モデルを使う. これによって, 楽器や手自体によっ て遮蔽された関節においても安定して良好な3次元位置推 定をおこなうことができる.また,音の情報を利用すること によって従来手法では決定することができなかった絶対ス ケールを定めることができ、かつ鍵盤表面の平面も求めるこ とができるようになる.

提案手法の評価のため,手と手を除いた全身に小型の光学 式のマーカーを貼り付け,提案手法と同時に鍵盤楽器演奏の 収録をおこなって比較した.その結果,完全非接触で演奏を 妨げることのない高精度のハンドトラッキングが実現でき たことを示す.

2 関連研究

複数視点のカメラ画像が与えられたときにカメラ姿勢と オブジェクト位置を同時推定する手法は,バンドル調整 [2] として知られており、SLAM (Simultaneous Localization and Mapping) や3次元再構成をおこなう COLMAP [3] な どのための基礎技術となっている. 一般的なバンドル調整で は、まず複数視点の画像が与えられ、それぞれの画像からの 特徴点抽出をおこなう. 抽出された特徴点のうち, 異なる画 像間で対応関係にある可能性が高いものを選出し. それらが 同じ3次元位置から投影されたものとみなして辻褄が合う ように再投影誤差の最小化をおこなう.しかし、カメラ配置 が非常に疎であり視点が大きく離れている場合には画像か らの特徴点マッチングが機能しなくなる,最適化が困難にな る、という問題がある.特に本論で対象としている楽器演奏 時には前景となる手が激しく動き、かつ背景となる鍵盤は繰 り返しを多く含んでいるため特徴点マッチングが機能しな い. また, 楽器周辺に密にカメラを配置することは現実的で はない. このような状況に従来法では対応が困難である.

3 提案手法

3.1 **逆投影問題**

複数視点のカメラ画像において 2 次元の手関節位置 $f_i^j \in \mathbb{R}^2$ が事前学習済みの機械学習モデルによって得られたとする. これを特徴点として用いると, 関節それぞれが一意に定まるため, 画像からの特徴点抽出で問題となる画像間でのマッチングエラーは回避できる. 本論では片手につき 21 点の関節点を使う. ここから従来のバンドル調整では, 世界座標系での未知の 3 次元手関節点 $\mathbf{P}_i \in \mathbb{R}^3 = [x_i, y_i, z_i]^T$ をそれぞれのカメラへ投影 $\mathcal{P}()$ したときの画像上の点と各特徴点 f_i^j との再投影誤差を最小化することが一般的である.

$$\underset{P,W,K}{\operatorname{arg min}} \sum_{i \in \mathbb{C}} \sum_{j \in \mathbb{F}_i} \| K \left(\mathcal{P}(\mathbf{W}_i P^j) \right) - f_i^j \|_2 \tag{1}$$

ここで、 \mathbb{C} , \mathbb{F} , \mathbf{W}_i , $\mathbf{K}()$ はそれぞれ手を認識したカメラ集 合, カメラ画像 *i* での特徴点集合. カメラ *i* に対する原点か らの剛体変換, カメラ内部パラメータに依存する関数であ る. しかし, この最適化はカメラ視点が極端にスパースで あったり, 特徴点が局所的 (本論の場合は手) である場合に 解くことが困難となり, 実際に筆者らの環境で実験したとこ ろ正常な値へ収束することはなかった. この問題を解決する ため本論では, 以下のように光線軌跡をカメラ側から世界座 標系へ逆に辿る逆投影問題として定式化する.

$$\underset{z,\mathbf{W},\theta^{k}}{\operatorname{arg min}} \sum_{\substack{i\neq j\\i,j\in\mathbb{C}}} \sum_{k\in\mathbb{F}} \|\mathcal{P}^{-1}\left(K_{\theta^{k}}^{-1}(f_{j}^{k}), z_{j}^{k}\right) - \mathbf{W}_{ij}\mathcal{P}^{-1}\left(K_{\theta^{k}}^{-1}(f_{i}^{k}), z_{i}^{k}\right)\|_{2}$$
(2)

ここで $\mathbf{W}_{ij} \in \mathbb{R}^{4 \times 4}$ はカメラ *i* から *j* への回転 $\mathbf{R}_{ij} \in \mathbb{R}^{3 \times 3}$ と並進 $\mathbf{T}_{ij} \in \mathbb{R}^3$ を表す行列,

$$\mathbf{W}_{ij} = \begin{bmatrix} \mathbf{R}_{ij} & \mathbf{T}_{ij} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$
(3)

 $\mathcal{P}^{-1}(u_i^k, v_i^k, z_i^k) = [u_i^k z_i^k, v_i^k z_i^k, z_i^k, 1]^T$ は逆投影関数であ る.また,関節の深度 z_i^k については直接最適化するのでは なく, $[z_i^1, ..., z_i^{21}] = z_i^{scale} \cdot LiftNet(f_i^1, ..., f_i^{21}) + z_i^{wrist}$ と モデル化する. LiftNet()は片手すべての関節の 2 次元座 標を入力とし,手首を原点としたそれぞれの関節の深度を出 力する事前学習済みニューラルネットワークである.これに よって本来の 3 × 2 × 21 次元の関節推定問題は 2 × 2 × $|\mathbb{C}|$ 次元の最適化問題にまで簡略化される.我々はこの問題を Levenberg-Marquardt 法 (LM 法) [4] で解く.ここで,カメ ラ内部逆関数 $K_{\theta^k}^{-1}$ ()についてはレンズ中心からの距離に対 する多項式関数としてモデル化した.

3.2 拘束

深度の時間的な一貫性を保つために以下の拘束を導入し て最小化する. *t* はフレーム番号である.

$$\alpha \cdot \|z_i^t - 0.5 \cdot (z_i^{t-1} + z_i^{t+1})\|_2 \tag{4}$$

さらに, 関節の長さが大きく変わらないようにするために As-Rigid-As-Possible 拘束 [5] を導入する.

$$\beta \cdot \sum_{lk \in \mathbb{B}} \left\| \left(\mathcal{P}^{-1}(K^{-1}(f_i^l)) - \mathcal{P}^{-1}(K^{-1}(f_i^k)) \right)^2 - \left(\mathcal{P}^{-1}(K^{-1}(f_j^l)) - \mathcal{P}^{-1}(K^{-1}(f_j^k)) \right)^2 \right\|_2$$
(5)

B はボーン集合である.また,隣り合う指の関節は同じ方向 にしか稼働できず,これは速度が大きいほど顕著である.

$$\gamma \cdot \sum_{lk \in \mathbb{B}^{f}} \left(\overline{(\vec{P}_{i}^{l,t+1} - \vec{P}_{i}^{l,t})} \cdot \overline{(\vec{P}_{i}^{k,t+1} - \vec{P}_{i}^{k,t})} - 1 \right)^{2} \\ \cdot \left| (\vec{P}_{i}^{l,t+1} - \vec{P}_{i}^{l,t}) \right| \left| (\vec{P}_{i}^{k,t+1} - \vec{P}_{i}^{k,t}) \right|$$
(6)

ここで上線はベクトルの正規化, B^f は同じ方向にしか動か ないボーンの集合である. 最後に当然ではあるが, 提案手法 ではカメラ位置がほぼ固定のため以下のようにカメラ姿勢 変換を一周すると元の姿勢に帰還する.

$$\mathbf{R}^{\mathbf{T}}{}_{1C} \prod_{C-1 \ge i, j=i+1 \ge 1} \mathbf{R}_{ij} = \mathbf{I},$$
$$\mathbf{T}_{C1} + \sum_{i=1, j=i+1}^{C} \mathbf{T}_{ij} = \mathbf{0}$$
(7)

最後に、スケールを一意に決定するため、 $|\mathbf{T}_{1C}| = 1$ とする.

3.3 ファインチューニング

上記の逆投影問題を LM 法で解き, カメラ姿勢とそれぞ れのカメラからの手関節の深度を求めたのち, 最後にすべて の手関節 3 次元座標の微修正をおこなう. それぞれの時間 フレームにおいて最も逆投影誤差の小さいカメラペア ij を選び,それぞれのカメラから逆投影した 3 次元関節位置を $P'_i \approx P_i + \delta_i, \ \delta_i \in \mathbb{R}^3$ とおいて以下の最小二乗問題を解く.

$$\arg \min_{\delta_{i},\delta_{j}} \left\| P_{j}' - P_{i}' \right\|_{2} + \alpha \cdot (|\delta_{i}| + |\delta_{j}|) + \beta \cdot \sum_{lk \in \mathbb{B}} \left\| (P_{i}'^{l} - P_{i}'^{k})^{2} - (P_{j}'^{l} - P_{j}'^{k})^{2} \right\|_{2}$$
(8)

α, β は重みである. このとき δ 以外のパラメータは固定す る. また, δ の大きさは高々隣り合う関節の長さの半分を超 えないようにバウンドする. 上式の第二項は正則化項, 第三 項は As-Rigid-As-Possible 拘束である. 最終的に得られた *P*^{*i*} が手関節の 3 次元座標となる.

3.4 絶対スケールと鍵盤平面の決定

RGB 画像のみによるバンドル調整には絶対スケールが不 定となる問題がある. これを決定するために鍵盤楽器から出 力される MIDI ノート (音符情報) *o* ∈ [0,...,127] を利用す る. MIDI の出力できない楽器の場合は音を録音し, MIDI ノート化する (ここでは多分にエラーを許容する). ここで 前提として, カメラ画像上での鍵盤の位置は未知であり, ま た鍵盤はオクターブごとに繰り返されており, さらに楽器に よって何オクターブの鍵盤があるのか不明であるため特徴 点マッチングは困難である. 現在どのキーが押されているか は得られているが, どの指がどのキーを押さえているかも未 知である.

まず, それぞれの鍵盤の数分だけカメラ画像と同じ解像 度のグレースケール画像 $I_k \in \mathbb{R}^{H \times W}, k = [0, ..., 127]$ を用 意し、これを鍵盤分布画像と呼ぶ.フレーム t-1 から t 間 において新しく押された鍵盤の集合を K^t とし, t において 手の平と直交する内側方向へ動いたすべての指先の位置の ピクセル位置を平均 μ として適当な分散 σ をもつガウス分 市 $\mathcal{N}(\mu, \sigma)$ を $I_k^t \leftarrow I_k^{t-1} + v \cdot a \cdot \mathcal{N}(\mu, \sigma), k \in \mathbb{K}^t$ のように それぞれの鍵盤分布画像に足す. ここで v, a は指先の速度 と音の強さである. このとき $t \to \infty$ において、それぞれの 鍵盤分布画像の極値の位置 pk を求め, 隣り合う鍵盤の鍵盤 分布画像の値を引いたもの $I_k(p_k) - I_{k-1}(p_k) - I_{k+1}(p_k)$ が 正の値となる視点が2つ以上存在しているとき、この点を鍵 盤kの表面特徴点候補とする.1分間鍵盤を適当に演奏した 後に得られたこの候補点を鍵盤番号と色相を対応させて可 視化したのが図2上段である.繰り返しの影響の無い特徴点 が得られていることがわかる.

最終的に残った表面特徴点候補群から, 鍵盤表面の 3 次元 位置を前段階までで得られたカメラ姿勢とパラメータを用 いて逆投影問題を最小化して求める.このとき押下された鍵 盤位置の集合 (表面特徴点) は平面となることを拘束条件に 含める.最終的な鍵盤表面位置は有効な特徴点マッチングを RANSAC [6] で選択することによって決定する.

こうして求まった鍵盤平面に対して各特徴点を射影し, 主成分分析をおこなう.この主軸にさらに各点を射影した とき,最も距離の大きい2点間の距離 D に対応している音



図 2 上段: 1分間演奏後の鍵盤特徴量の収束結果. 音高 が色相に対応している. 下段: 最適化された鍵盤表面特徴 点 (緑点) とその主軸と副軸 (青線). 赤点はカメラ位置, 赤線はカメラ角度, 青点とオレンジ点は手を示す.



図3 実験のセットアップ.手と全身に光学マーカーを貼 り付けて提案手法と同時収録をおこなった.カメラは鍵盤 真上と左右の3台を設置した.

程 H から, ピアノ鍵盤の幅の規格が1オクターブ165mm であることを使い, このバンドル調整の絶対スケールが 165H/12Dと算出できる.図2下段は赤点がカメラ位置, 赤 線がカメラ角度, 青線が推定された鍵盤平面の主軸と副軸, 緑点が最終的に選択された鍵盤表面特徴点を示す.なおここ では鍵盤真上のカメラ位置を原点としている.

4 実験

4.1 実装

本システムは、マルチカメラでの撮影システム、撮影動画 からの2次元手関節推定システム、提案手法による3次元ハ ンドトラッキングシステムの3つから構成される. 画像から の2次元関節位置推定には MediaPipe [1], 3次元化をおこ なう *LiftNet()*の構造には Graph U-Net [7] を採用し、事 前学習のためには複数視点の手関節データセット [8, 9] を使 用した. 実装は C++ でおこない、実行は Apple M1 チップ 上最大 8 並列でおこなった. また、カメラには Sony RX01 II (フル HD, 60fps) を 3 台. ピアノには YAMAHA C3 を 使用した.

4.2 比較

比較のため、手の提案手法と同じ関節位置、および手 以外の全身に光学式のマーカーを貼り付け、OptiTrack (https://www.optitrack.jp)のシステムでトラッキングを おこない提案手法と同時収録した(図3).演奏時には手の マーカーの遮蔽が多分に発生し、関節のいくつかが認識でき ないことが多発したため、手動での修正をおこなった.

まず OptiTrack で得られた手関節の結果と提案手法で得 られた関節出力のプロクラステス距離を比較した結果、ショ パンのエチュードやスケールなどの6分間の演奏で関節位 置の最大誤差 20mm 以下となった (図 4). 再投影誤差の最 小化をはじめ, 従来手法はすべて最適化が失敗したため比較 が困難となったが非接触での手関節認識のひとつのベース ラインを得たと言える.この誤差の原因は単純な推定誤差だ けではなく,貼り付けたマーカー位置で2次元関節推定座標 がそのまま得られるわけではないことにもある,絶対位置/ スケールの評価として、OptiTrack で測定された現実スケー ルのデータと重ね合わせたものを添付動画に示す. 動画では 2手法の結果が高い精度で一致していることがわかる.この とき左右の手が最も離れたときの手首位置の距離の誤差が 約 30mm となった. 深度情報を一切利用しないバンドル調 整でこの精度を達成することは従来法では困難である. 最後 に、逆投影問題の最適化時間は1フレームあたり約0.2秒と なった.計算が正常な値へ収束しなかったが,再投影誤差の 最小化問題として実行した場合は1フレームあたり3秒程 度かかったことと比較すると大幅に計算が効率化されてい ることがわかる.

5 まとめ

難易度の高い楽器演奏時の手指の詳細なトラッキングを 非接触で実現するために, 画像特徴量ではなく2次元手関節 推定の結果を使ってその逆投影問題として定式化し, 機械学 習とのハイブリッドな最適化によってこれを効率的に解け ることを示した.また, 音情報を併用することによってバン ドル調整だけでは得られない絶対スケールや鍵盤平面を決 定できることを示した.



提案手法にはまだいくつかの問題が残されている.まず, 演奏するフレーズによっては片手が完全にもう片方の手に 遮蔽されてしまう.これに対応するためにはカメラを下側か らボトムアップアングルで撮影する必要がある.しかしこれ を実際に試してみたところ,従来の画像からの2次元手認識 手法では左右の手を混同してしまうなどの問題が発生した ためかえって精度を悪化させた.これを解決するためには楽 器下側からの手の画像を収集しての手認識モデルの学習を おこなう必要があるが今後の課題である.

謝辞

本論での実験収録については, Seoul National University と KAIST 研究グループからの支援を受けた.

参考文献

- F. Zhang, V. Bazarevsky, A. Vakunov, A. Tkachenka, G. Sung, C.-L. Chang, and M. Grundmannl, "Mediapipe hands: On-device real-time hand tracking," in *CVPR*, 2020.
- [2] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon, "Bundle adjustment — a modern synthesis," in *ICCV*, 1999.
- [3] J. L. Schönberger and J.-M. Frahm, "Structure-frommotion revisited," in CVPR, 2016.
- [4] L. Kenneth, "A method for the solution of certain nonlinear problems in least squares," in *Quarterly of Applied Mathematics*, 1944.
- [5] T. Igarashi, T. Moscovich, and J. F. Hughes, "As-rigid-aspossible shape manipulation," in *SIGGRAPH*, 2005.
- [6] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," in *Communications of the ACM*, 1981.
- [7] B. Doosti, S. Naha, M. Mirbagheri, and D. J. Crandall, "Hope-net: A graph-based model for hand-object pose estimation," in *CVPR*, 2020.
- [8] F. Gomez-Donoso, S. Orts-Escolano, and M. Cazorla, "Large-scale multiview 3d hand pose dataset," in arXiv:1707.03742, 2017.
- [9] Z. Yu, L. Yang, S. Chen, and A. Yao, "Local and global point cloud reconstruction for 3d hand pose estimation," in *BMVC*, 2021.