

研究報告

質点の物理モデルを用いた動画像からのソニフィケーション SONIFICATION FROM THE VIDEO USING PHYSICAL MODEL OF A MATERIAL POINT

原 一歩

Kazuho HARA

東京電機大学

Tokyo Denki University

小坂 直敏

Naotoshi OSAKA

東京電機大学

Tokyo Denki University

概要

筆者らは、背景と異なる色を持つ物体の動きのある動画を対象とし、これと同期する音響を付与する、動画像からのソニフィケーションシステムを構築している。このシステムは、まず、同画像内の物体の移動を質点運動ととらえ、フレーム間差分により速度と加速度を算出する。ここで、フレーム画像から追跡対象の物体以外の情報を除去し、重心点を算出し、これを質点の座標とした。次に、加速度から質点の衝突判定を行い、衝突判定時にリアルタイムにFM合成する。この方法を球と人の動きに適用して、その効果の有効性を確認した。

1. はじめに

一般に、映像作品は映像に同期して音響を付与することにより、視聴者により効果的な印象を与える。しかし、映像内のある物体の動作に合わせて音響を付与する作業の負担は大きい。森下らは動画像内の物体を追跡し、入力した衝突音を、物体の移動に合わせて半自動的に付与するシステムを構築した [1]。しかし、このシステムは、物体の移動に対応した、衝突音の自動的な音量の変化が行われており、効果音の音色の変更は行わない。また、システムは、入力された動画像にリアルタイムに音響を付与しない。そこで、本稿では、ソニフィケーションの対象の動画像を、背景と異なる色を持つ物体が1つだけ移動しており、30fpsであるものに限定し、以下の2点を満たすシステムを構築する。1) 入力された動画内の物体の動きを質点運動と捉え、加速度から衝突判定を行う。2) 衝突判定時に、加速度の大きさを音合成のパラメータとして入力し、リアルタイムにFM合成を行う。本稿は構築したシステムの評価実験を行い、物体追跡の精度を測定する。ま

た、システムを様々な動画に応用し、その実用性を考察する。

2. 物体追跡と衝突判定

本稿では、フレーム間差分と HSV 色空間による色抽出を行い、動画像内を移動する特定の色の物体を各フレーム画像から抽出する。次に、抽出後のフレーム画像の重心を計算し、質点の座標とする。

2.1. 色重心を用いた動画像の物体追跡

まず、動画像の各フレーム画像と、そのフレーム画像を基準とした時の、過去の2フレームの画像をグレースケールに変換し、フレーム間の差分画像を求め、2値化する。次に、2回のオープニング処理を行うことにより雑音を除去する。これにより、各フレーム画像から移動している物体を抽出する。

移動している物体が抽出された2値化画像と同フレーム番号におけるRGB画像で論理積を取り、移動している物体のRGB画像を抽出し、各画素値をHSV色空間に置き換える。最後にHSV色空間において閾値の判定を行い、色に取まっているかどうかで2値化することにより、フレーム画像から移動する特定の色の物体を抽出する。画像処理後の各フレームの2値化画像において、縦方向と横方向でそれぞれヒストグラムを算出し、平均を求めることにより重心の座標を求める。次に、各座標値を0~100に正規化し、質点の画面座標とする。

2.2. 質点の加速度からの衝突判定

各フレーム画像における質点の座標と、そのフレームを基準として、過去の2枚のフレーム画像の質点の

座標を、前後のフレーム同士で差を取り、質点の速度を2つ算出する。さらに、速度ベクトルの差を取り、各フレーム画像における加速度とする。衝突時の力に比例する量として、二次元加速度ベクトルの絶対値を用いて、絶対値が指定された閾値を越えた時のフレーム番号において、物体の衝突が発生したと判定する。

3. 音合成

衝突判定時に、加速度を任意に定数倍した値を、周波数変調 (FM: Frequency Modulation) のキャリア周波数と変調周波数および変調指数に代入し、式 (1) により音波形を合成する。また、FM による合成音波形を減衰させるためにエンベロープをかけ、これを衝突音とする。

$$y(t) = A \sin(2\pi f_c t + I \sin(2\pi f_m t)) \quad (1)$$

ここで、 t : 時刻 [sec], f_c : キャリア周波数 [Hz], f_m : 変調周波数 [Hz], I : 変調指数, A : 振幅, $y_n(t)$: 合成波形である。

4. システムの構成

ソニフィケーションシステムの実装に使用したツールを表1に示す。Pure Data[2]はMiller Pucketteにより開発された、マルチメディアの表現に特化したビジュアルプログラミング言語である。OpenCV[3]はIntelにより開発された、画像処理のライブラリ群である。OpenSound Control[4](以下、OSCと略す)は、M. Writeが開発した、音響関係のデータの送受信に特化した通信プロトコルである。本稿では、物体追跡の処理と、動画再生に同期した音合成の処理にシステムを分割し、同一のマシン内で実装を行った。各処理部のシステム構成については、4.1章および4.2章で説明する。

表1. 使用ツールの一覧

ツール名	用途
C++	質点の画面座標からの加速度の算出
OpenCV	画像処理による物体追跡
Pure Data	加速度に対応した音響信号の合成
OSC	プログラム間の加速度情報の送受信

4.1. 物体追跡部のシステム構成

物体追跡部のシステム構成を図1に示す。システムには動画画像と、HSV色空間における抽出したい色の範囲設定のデータを入力するようにした。各画像処理の

実装は、OpenCVのライブラリを用いて行った。画像処理後の各フレームのそれぞれの重心点から、フレーム毎の加速度を計算し、リストファイルとして外部出力する。

4.2. 動画再生及び音合成のシステム構成

動画再生と音響信号の合成のシステム構成を図2に示す。動画再生部はC++により実装し、音響合成部はPure Dataを用いて実装した。また、OSCを用いて、C++プログラムによる動画再生と、Pure Dataプログラムによる音響信号の合成を同期させた。図2のシステムには、動画画像と、それに対応するフレーム毎の加速度のリストファイルを入力として与える。

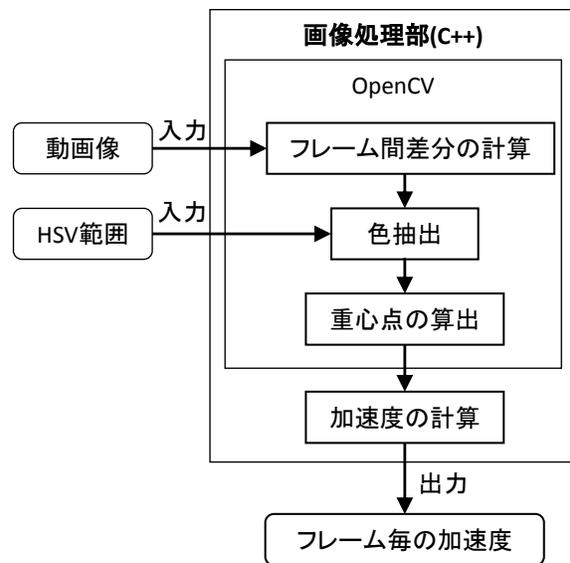


図1. 物体追跡部のシステム構成

5. 評価実験

本システムの物体追跡と衝突判定の精度を、背景やカメラの環境毎に調査するため、評価実験を行った。

5.1. 実験手法

まず、表2に示す要因と水準に基づいた、それぞれの環境下で色のある球体を落下させ、床に衝突して跳ね返る動画を撮影し、実験データとする。実験に使用した背景画像 [5] を図3に示す。次に、実験データの動画画像から目視できる衝突のフレーム番号の一覧を作成する。また、システムに動画画像を入力して得られる加速度が、閾値を超えた時のフレーム番号の一覧を用意する。2つの一覧を照合し、前後1フレームの誤差

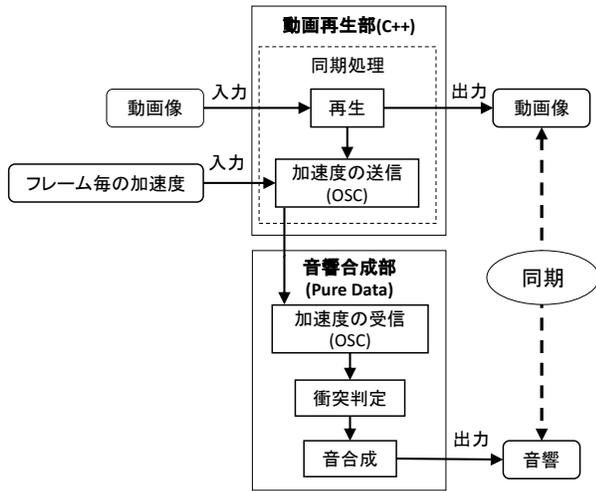


図2. 動画再生及び音合成のシステム構成

内であれば、システムが正しく判定した衝突フレームとする。最後に、適合率と再現率およびF尺度の各評価を、式(2)、(3)から算出し、実験結果とする。今回の実験では、動的特徴量の閾値を1に設定した。

$$p = \frac{R}{N}, r = \frac{R}{C} \quad (2)$$

$$F = \frac{2pr}{p+r} \quad (3)$$

ここで、 p : 適合率, r : 再現率, F : F尺度, R : システムが正しく判定した衝突のフレームの枚数, N : システムが判定した衝突フレームの枚数, C : 目視できる衝突フレームの枚数である。

表2. 実験データの要因と水準

要因	水準			
	背景画像	有	無	
カメラの高さ	0m	1m		
カメラの俯瞰角度	0°	45°		
ボールの色	赤	青	黄	白

5.2. 実験結果と考察

実験により得られた動画画像毎の評価値を、背景画像の有無毎に集計したものを表3に示す。背景画像がない場合のF尺度は0.667であり、この性能は厳密に質点の動作と対応させる目的では満足できないが、全体のグローバルな動きに音を与える意味では許容できると考えた。特に、人の体の重心のように手足の個別の動きではなく、それら複合体の動き全体に同期させる



図3. 実験で使用した背景画像

場合は、音と対応させる対象がまぎれてしまうため、この性能でも十分効果がある、と考えられる。また、背景画像がある場合のF尺度は0.690であり、背景画像の有無においてF尺度の差異は見られなかった。このことから、背景に動きがない場合、背景画像がない場合と同程度の性能で物体追跡が行われていることがわかった。

次に、動画画像毎の評価値をカメラの俯瞰角度毎に集計したものを表4に示す。同表のF尺度を見ると、俯瞰角度を45°に設定した場合のF尺度は0.723であり、俯瞰角度を0°にして撮影した場合より高い値を示している。これは、落下する球体を水平に撮影すると、動画内における物体の移動速度が早くなり、残像による球体の色調の変化により、正常な物体追跡が行えなかったからであると考えられる。

表3. 背景画像の有無毎の評価値

背景画像	適合率	再現率	F尺度
有	0.719	0.663	0.690
無	0.687	0.649	0.667

表4. カメラの俯瞰角度毎の評価値

カメラの俯瞰角度	適合率	再現率	F尺度
0°	0.706	0.563	0.627
45°	0.699	0.749	0.723

6. システムの応用と考察

構築した本システムを用いて、球体あるいは人間の運動する動画画像からのソニフィケーションを試み、現

段階のシステムの実用性を考察する。表5に入力として与えた動画の内容を示す。また、同表の各データから抜粋したフレーム画像と、それに対して追跡する物体の抽出を行った結果をそれぞれ図4、5に示す。データ#2では、図5aのフレーム画像内の、赤色のパーカーを追跡することにより、跳躍した人間が床に着地した時に対応した音響の合成を試みる。

表5. 入力する動画データの内容

データ	内容	追跡対象
#1	球体の跳ね返し合い	桃色の球体
#2	人間の跳躍（スキップ）	赤色のパーカー

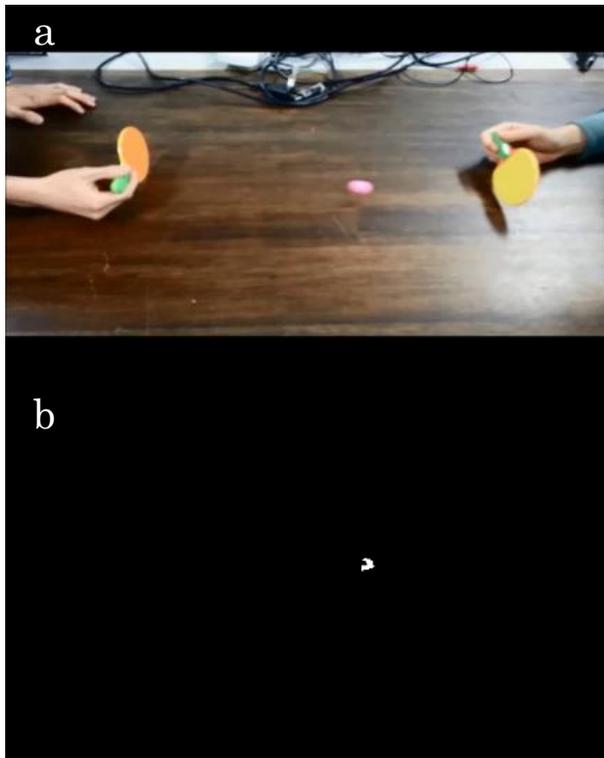


図4. データ#1のキャプチャ画像aと物体抽出の結果画像b

6.1. データ#1のソニフィケーション結果と考察

図4aのフレーム画像では、追跡する球体と似た色味を持つ人間の手が映っている。しかし、図4bを見ると、HSVにおける色抽出により、フレーム画像から追跡対象の球体のみ抽出できている。また、球体がラケットにより弾かれた時に、概ね衝突判定がされ、タイ

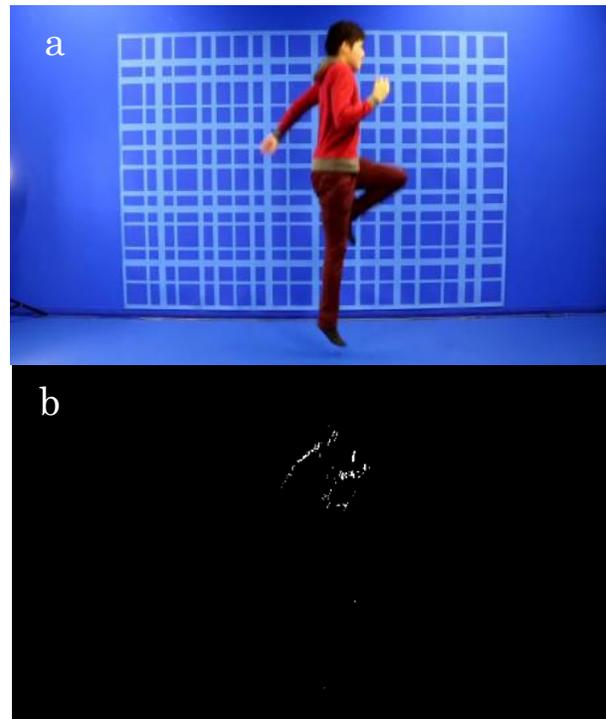


図5. データ#2のキャプチャ画像aと物体抽出の結果画像b

ミングに合わせて音響の付与が行われた。これは、5.2の実験に比べ、衝突がより抽出しやすい条件で行ったためと考えられる。音響に関しては、ラケットに衝突した球体の速度の変化に対応し、動画の演出を向上するような音響が合成されていると考えられる。

6.2. データ#2のソニフィケーション結果と考察

図5aのフレーム画像から追跡する物体の抽出を行った結果、赤色のパーカーだけ抽出することができた。しかし、図5bの抽出結果を見ると、元のフレーム画像の赤色のパーカーが映っている箇所のうち、一部しか抽出が行われていなかった。そのため、重心点の計算により求められる質点の座標の精度が低下し、動画内の人間の着地以外の場面において、衝突の判定が多発した。これは、赤色のパーカーの画面内の移動速度が、その色味を有している画素の面積の広さに対して遅いため、フレーム間差分の算出時に、除去されてしまったからであると考えられる。

7. まとめ

本稿では、動画内の物体の動作から質点の加速度を算出し、衝突判定時に音響をリアルタイム合成するシステムを構築した。衝突判定の物理評価実験の結果、背景画像が動いていない場合は、背景画像がない場合

と同程度の性能で物体の追跡が行えることがわかった。一方、移動速度の高い球体を撮影すると、各フレーム画像において球体の被写体ぶれが生じ、正常な物体追跡が行えなかったことも明らかになった。また、動画画像内の追跡する対象物を球体以外にし、動画画像からのソニフィケーションを行った結果、質点の座標の精度が低下し、衝突の誤判定が多発した。以上の課題点を改善するために、より高度な画像処理を導入することを検討している。

に開始した Media Project 他、コンピュータ音楽のコンサート企画多数。日本音響学会、電子情報通信学会、情報処理学会、ICMA、IEEE 日本電子音楽協会各会員。現在、東京電機大学 未来科学部教授、本会会長。

8. 参考文献

- [1] 森下沙耶, 岡部誠, & 尾内理紀夫. (2012). 動画への効果音付加支援システムの作成 (学生研究発表会). 映像情報メディア学会技術報告, 36(8), 119-122.
- [2] Puckette, M. (1996). Pure Data: another integrated computer music environment. Proceedings of the Second Intercollege Computer Music Concerts, 37-41.
- [3] OpenCV <http://opencv.jp/>
- [4] Wright, M., Freed, A., & Momeni, A. (2003, May). Opensound control: State of the art 2003. In Proceedings of the 2003 conference on NIME(New Interfaces for Musical Expression) (pp. 153-160). National University of Singapore.
- [5] 藤田紘久 Futta.NET <http://www.futta.net>

9. 著者プロフィール

原 一歩 (Kazuho HARA)

東京電機大学大学院未来科学研究科情報メディア学専攻音メディア表現研究室所属。高校2年からDTMによる作曲を行っており、大学院入学後にピアノを始める。現在は、動画と静止画からのソニフィケーションに関する研究を行っている。

小坂 直敏 (Naotoshi OSAKA)

昭51早大・理工・電気卒。昭53同大大学院修士課程了。同年日本電信電話公社(現NTT)入社。以来通話品質の研究、音声対話の研究、コンピュータ音楽あるいはマルチメディア創作のための音響研究などに従事。平6早大より博士(工学)。平8-14コミュニケーション科学基礎研究所音表現およびメディア表現研究グループリーダー、平成15東京電機大学工学部教授。メディアコンテンツのための音響情報処理の教育と研究に従事。また、音楽制作および発表活動も行う。2006