

# 音楽のムード分類結果を利用した ホームビデオへの自動BGM付与・同期手法

Automatic Home Video BGM Selection and Synchronization  
Based on Music Mood Classification

小野 佑大†  
Yudai Ono

石先 広海‡  
Hiromi Ishizaki

帆足 啓一郎‡  
Keiichiro Hoashi

滝嶋 康弘‡  
Yasuhiro Takishima

甲藤 二郎†  
Jiro Katto

## 1. はじめに

近年、動画共有サイトが普及し、ホームビデオにBGMを付与して他者と共有するユーザが増えている。より魅力的なコンテンツを制作するために、映画やドラマの様に、動画のムードに合った楽曲を選定する事や、ミュージックビデオの様に、動画と楽曲とを連動させる事が重要である。しかしその作業は、多大な労力を要し、高い編集技術も必要なため、一般ユーザには難しい。そこで、その様なコンテンツの制作を支援する技術が必要である。

従来研究[1,2]では、動画編集ルールに基づき楽曲を選択したり、音楽構造と動画構造を合わせてミュージックビデオを自動生成したりする事で、ユーザのコンテンツ制作を支援している。しかし、これらには、ユーザの付加したいムードの楽曲を選定できないと言う課題がある。また、楽曲のサブ区間などをBGMとして使用しているが、それ以上に動画を効果的に演出する区間を付与する事で、より魅力的なコンテンツを生成できると考えられる。

そこで本稿では、この課題に対し、ホームビデオへのBGM付与の効率化を図るために、音楽のムード分類結果から得た印象語に基づき楽曲を選定し、動画と同期する区間を楽曲から自動抽出するBGM付与・同期手法を提案する。さらに本手法の有効性を被験者評価実験により検証する。

## 2. 提案手法

### 2.1 概要

本手法では、BGMとして使用する楽曲を選定するために、音楽のムード分類結果で得た印象語を利用する。そして、ユーザはその印象語を入力し、絞り込んだ楽曲群から選曲を行う。

次に、選択された楽曲から動画を効果的に演出する区間を抽出するために、動画の動きと楽曲の音量の変化するタイミングが一致する区間を楽曲から抽出して動画へ付与する。本手法では、心理学で定義される動画と音楽の調和に関する要因である時間的なアクセントの一致[3]に着目し、動画と調和する区間を楽曲から抽出する。

### 2.2 音楽のムード分類結果に基づく楽曲選定方法

本節では、文献[4]に基づいて楽曲をクラスタリングし、生成されたクラスタへ印象語を付与する。そして、ユーザが選択した印象語と一致するクラスタを選定し、楽曲を絞り込み、選曲を行う。音楽のムードは心理学で定義されるRussell空間[5]で表現する。これは人間の感情を活性-不活

性(Arousal)、快-不快(Valence)の二軸(以下、AV軸)から成る2次元平面上で定義し、その空間上の座標値(以下、AV値)で印象語を表現する。Arousalは音楽の音量やリズムの強さなど、Valenceは音楽の音色やリズム変動などに相関する。

まず初めに、音楽のムードに関連する29次元特徴量を抽出し、主成分分析を行う。そして、第一主成分に-1を乗じたものをArousal値、第二主成分をValence値としてRussell空間へ楽曲を写像する。次に、AV値を-1~1で正規化後、AV軸に対し階層的にk-means法を適用し、空間上の各象限に均一にクラスタを生成する。そして、各クラスタの重心に近い印象語でラベル付けする。最終的に、入力された印象語を含むクラスタを選定し、そのクラスタに属する楽曲群からBGMとなる楽曲を選択する。

### 2.3 動画を効果的に演出するBGM区間抽出法

本節では、動画を効果的に演出するBGM区間の抽出法について述べる。まず、動画と音楽から時系列の特徴量を抽出する。動画特徴はフレーム内の動きベクトルの総和で計算されるMotion Activityを、音楽特徴は信号の二乗平均平方根のフレーム間差分を用いる。

次に、各特徴量の変化を捉えやすくするために平滑化し、時系列解析で利用される特異スペクトル変換[6]によって変化を検出する。これは、時系列のある時点の変化に対して過去と未来の部分系列から特徴抽出をし、その非類似性で変化の度合い(変化度)を求める手法である。

最後に、検出した各特徴の時系列の変化度に対して式(1)で定義される相互相關関数を計算する。

$$(f * g)(\tau) = \sum_t f(t)g(\tau - t) \quad (1)$$

$f$ と $g$ は時刻 $t$ における動画と音楽の時系列の変化度を示し、ラグタイム $\tau$ 毎に $f$ を $g$ の先端から終端まで移動させる。そして、(1)の計算結果が高い区間を抽出する。

## 3. 評価実験

### 3.1 実験概要

提案手法の有効性を検証するために、2つの評価実験を行う。実験1では、クラスタリング結果に基づく楽曲選定方法を評価する。また、実験2では動画を効果的に演出するBGM区間抽出法を評価する。

### 3.2 実験1～楽曲選定方法の評価～

#### 3.2.1 実験方法

本実験では、Russell空間上の印象語を検索クエリとしてFlickr[7]から収集した500ファイルの動画と、文献[4]にて使用されているUSPOPの206曲を使用した。

†早稲田大学 理工学術院

‡株式会社 KDDI研究所

表1. 各クラスタが示すRussellの印象語[5](一部抜粋)

クラスタ1	クラスタ2	クラスタ3	クラスタ4
リラックスした	悲しい	楽しい	怒った
優しい	ほろ苦い	明るい	怖い
穏やかな	憂鬱な	嬉しい	イライラした
甘い	冷たい	興奮した	緊迫した
満足した	退屈な	驚いた	不安な

楽曲群に対して、クラスタ数を4に設定し、2.2節に記載される方法で、各クラスタにラベリングした。表1に各クラスタに付与された印象語を示す。そして、入力となる動画に、各クラスタから無作為に選定した楽曲のサビ部分をBGMとして付与する事で実験用の動画を作成した。これらとBGMのない動画の計5種類を1セットとし、計5セット用意した。

次に実験の流れについて述べる。まず、動画に付与された印象語を提示し、被験者は1セットずつ視聴する。そして、各動画から受けるムードが、提示した印象語にどれだけ近いか評価する。その評価基準は、10: とても近い~0: 全く近くない、の11段階を用いた。尚、被験者は大学生19名である。

### 3.2.2 実験結果

表2に被験者から得た評価値の平均を示す。これを見ると、動画の印象語が含まれるクラスタの楽曲をBGMとして用いた動画は、BGMなしの動画と比べ評価値が高い。この結果から、動画の印象語に合った楽曲を選定している事が分かり、本手法は有効であると言える。

しかし、「ほろ苦い」を見ると、動画の印象語を含むクラスタ1の評価値よりも、動画の印象語を含まないクラスタ2の評価値の方が僅かに高い。実際に、クラスタ1とクラスタ2それぞれの楽曲のAV値を見てみると、クラスタ1の楽曲のArousal値は-0.39, Valence値は0.13、クラスタ2の楽曲のArousal値は-0.5, Valence値は-0.11、であった。この2曲は異なるクラスタに属しているが、Russell空間上の距離が比較的近く、ムードが類似している可能性が高い。したがって、クラスタ1とクラスタ2のムードの違いが分かり辛く、被験者によって意見が分かれてしまったと言う事が原因と考えられる。

### 3.3 実験2~BGM区間抽出法の評価~

#### 3.3.1 実験方法

本実験では、実験1と同じ動画と音楽を用い、音楽は、動画の印象語が含まれるクラスタの楽曲のみを使用した。そして楽曲のサビ部分を付与した動画と、提案手法でBGMを付与した動画の2種類を1セットとして、計5セット用意した。

次に、実験の流れについて述べる。まず、被験者は1セットずつ視聴し、BGMが動画にどれだけ同期しているか評価する。その評価基準は、10: とても同期している~0: 全く同期していない、の11段階を用いた。そして全ての動画の評価が終了後、「動画の完成度を高めるためにBGMを同期させる事は効果的か」と言うアンケートを行った。その評価基準は、2: 効果的だ~0: 効果的でない、の3段階を用いた。尚、被験者は実験1と同様で、手法における順序効果は考慮してある。

表2. 実験1の評価値の平均

提示した印象語	BGMなし	クラスタ1	クラスタ2	クラスタ3	クラスタ4
怒り	5.28	3.89	3.28	4.56	6.00
冷たい	5.50	4.44	4.78	3.00	2.72
明るい	6.17	5.00	6.22	7.94	6.56
甘い	6.11	7.17	4.17	4.89	2.67
ほろ苦い	3.00	5.61	5.39	2.83	3.83

表3. 実験2の評価値の平均

動画の印象語	サビ部分	提案手法
怒り	6.00	7.72
冷たい	5.39	6.39
明るい	6.22	8.11
甘い	6.28	7.00
ほろ苦い	6.28	6.44

#### 3.3.2 実験結果

初めに全体アンケートの結果を見ると、19人中15人が効果的であると答えた。この事から、動画とBGMの同期を図る事は重要であると言える。次に、表3で示される被験者から得た評価値の平均を見ると、サビ部分をBGMとして加えた動画と比べ、提案手法の方が高い評価値を示している。この結果から、提案手法で付与したBGMは動画を効果的に演出できていると言える。

しかし、「ほろ苦い」の様に、評価値に差の無い結果も見られる。実際に「ほろ苦い」の楽曲のAV値を見てみると、Arousal値は-0.5, Valence値は-0.11で、Russell空間上の第3象限に存在している事が分かった。この様な楽曲は、リズムや音量変化等が小さい。そのため、動画と同期させた場合も被験者がこれらを認知し辛く、評価値に差が出なかつたと考えられる。

### 4.まとめ

本稿では、ホームビデオへのBGM付与の効率化を図るために、音楽のムード分類結果から得た印象語で楽曲を選定し、動画と同期する区間を楽曲から自動抽出するBGM付与・同期手法を提案した。そして、被験者による主観評価実験により提案手法の有効性を示した。今後は、BGM選曲の自動化とBGM区間抽出法の改善を試みる。

### 参考文献

- [1] Philippe Mulhem *et al.*, "Pivot Vector Space Approach for Audio-Video Mixing," *IEEE Multimedia*, Vol.10, No.2, pp.28–40, 2003.
- [2] Foote J *et al.*, "Creating music videos using automatic media analysis," *Proceedings of ACM multimedia*, New York, pp.553–560, 2002.
- [3] 岩宮眞一郎, “音楽と映像のマルチモーダル・コミュニケーション,”九州大学出版会, 2000.
- [4] 小野佑大ほか, “ホームビデオへの自動BGM付与のための心理学に基づく音楽分類手法,”第72回国情処全大, 1T-02, 2010.
- [5] J. A. Russell, *Journal of Personality and Social Psychology*, (6).
- [6] T. Id'e *et al.*, "Knowledge discovery from heterogeneous dynamic systems using changepoint correlations," In Proc. SIAM Intl. Conf. Data Mining, pp.571–575, 2005.
- [7] Flickr, available at <http://www.flickr.com/>