# 楽譜情報を援用した多重奏音楽音響信号の音源分離と 調波・非調波統合モデルの制約付きパラメータ推定の同時実現

#### 糸 ш 克 **寿**†1 後 藤 直 **老**†2 駒 和 **節**†1 谷 博†1 也†1 形 哲 尾 ጉ

本論文では,多重奏の音楽音響信号とその楽曲に含まれるすべての単音の音高・音長・音量・発音 時刻・楽器の種類の組である楽譜情報を入力として,単音ごとの音響信号を出力する音源分離手法と, そのための制約付きモデルパラメータ推定手法について述べる.本分離手法では,Standard MIDI File (SMF)などから抽出された楽譜情報を用いることで混合音のパワースペクトルを単音ごとに分 離し,調波構造と非調波構造のそれぞれを表現する2つのモデルを統合した新たな重み付き混合モデ ルを用いることで,単音に複数の調波構造が含まれることを防ぎ,かつ音高を超えた楽器音の音色類 似性を考慮することを実現する.モデルパラメータは,楽譜情報に基づいて MIDI 音源から生成した テンプレート音によって初期化し,EM アルゴリズムを用いた最大事後確率推定により反復推定する. さらに,モデルの過学習を防ぎ,同一楽器の単音のモデルに類似した音色を持たせるための制約条件 も同時に用いる.ポピュラー音楽の SMF を用いた評価実験で,本手法により SNR が 0.4 – 0.9 dB 向上することを確認した.

# Simultanious Realization of Score-informed Sound Source Separation of Polyphonic Musical Signals and Constrained Parameter Estimation for Integrated Model of Harmonic and Inharmonic Structure

Katsutoshi Itoyama,<sup>†1</sup> Masataka Goto,<sup>†2</sup> Kazunori Komatani,<sup>†1</sup> Tetsuya Ogata<sup>†1</sup> and Hiroshi G. Okuno<sup>†1</sup>

This paper describes a sound sourse separation method for polyphonic sound mixtures of musical signals which include both harmonic instrument sounds and inharmonic instrument sounds, and a constrained parameter estimation method by using a score which includes pitch, duration, volume, onset time, and instrument of each note as prior information. We separate a power spectrum of sound mixtures into each musical note by using an integrated weighted-mixture model consisting of both harmonic-structure and inharmonic-structure tone models (generative models for the power spectrogram). The integrated model realize a parameter estimation method under a constraint of parameter similarity in the same musical instruments. We initialize model parameters using template sounds which are recorded from a MIDI tone generator. On the basis of the Maximum A Posteriori Probability estimation using the EM algorithm, we estimated all parameters of this integrated model under several original constraints for preventing over-training and maintaining intra-instrument consistency. Using standard MIDI files as prior information of the model parameters, we confirmed that the integrated model increased the SNR by  $0.4 - 0.9 \, dB$ .

1. はじめに

デジタルオーディオが普及し,価値観が多様化する 中で,より能動的に音楽を楽しみたいというユーザ の要求が現れてきた.これまでのオーディオ再生技術 は,受動的な音楽の楽しみ方をより豊かにする方向に 進歩し,ユーザの要求に応えてきた.たとえば,5.1 次元や7.1 次元などの大がかりなシステムで忠実な音 環境の再現を目指すというものや,アクティブノイズ キャンセルなどの簡便な装置で静かな音環境を作るこ とでどこでも手軽に音楽鑑賞を楽しむというものがあ る.一方,能動的な音楽の楽しみ方には作曲や編曲, 演奏などがある.一般的には能動的に音楽を楽しめる

<sup>†1</sup> 京都大学大学院情報学研究科

Graduate School of Informatics, Kyoto University †2 産業技術総合研究所

National Institute of Advanced Industrial Science and Technology (AIST)

のは技術や道具を持っている人に限られており,受動 的な楽しみと能動的な楽しみの間には大きなギャップ があった.

能動的な音楽鑑賞<sup>1)</sup> という要求に応える研究事例として, Yoshii らはドラムスを対象とした楽器音イコライザ INTER:D<sup>2)</sup> および Drumix<sup>3)</sup> を実現した.ユーザは Drumix を使って楽曲中のドラムスの音量を操作し, 音色を置き換え,また,ドラムパターンを編集でき,その結果能動的な音楽鑑賞がより簡便に可能となった.しかし,これらのシステムはドラムスだけを対象としており,一般の楽器音に対して適用するまでには至っていなかった.

これに対して我々の目的は, CD などによる音楽音響信号(混合音)中のあらゆる楽器パートに対して自由に音量を操作できる楽器音イコライザを実現することである.従来のグラフィックイコライザやパラメトリックイコライザでは,特定の周波数帯域ごとの音量を調整して周波数特性を変化させることはできたが,楽器ごとの音量を調整することはできなかった.

楽器音イコライザを実現するためには,楽曲中に含 まれるすべての楽器音を楽器パートごとに,もしくは 単音ごとに分離する必要があり,そのためにはどの楽 器が,どの時刻に,どの音高で演奏されているのかと いった「楽譜情報」が必要となる、そこで本論文では、 音楽音響信号とその楽曲の「楽譜情報」を入力とし, 音楽音響信号を楽器パートごとに分離して出力する手 法について論ずる.本論文での「楽譜情報」とは「楽 曲に含まれるすべての単音の音高・音長・音量・発音 時刻・楽器種類」である.インターネット上での標準 MIDI ファイル (Standard MIDI File; SMF)の販売 サービス<sup>4),5)</sup> などによって, 最新の楽曲であっても音 楽音響信号に対応した SMF を入手することが容易に なっており, これらの SMF から楽譜情報を抽出する ことが可能であるため, 音楽音響信号と楽譜情報の組 を得ることはさほど困難ではないと考える.ただし, 音響信号とSMFとは何らかの従来法<sup>6)-10)</sup>を用いて 同期がとられていると仮定し,本論文では同期の問題 については扱わない.

音楽音響信号の音源分離に関する従来研究は,以下 の2つに大別できる.

(1) 音高を明示的に扱うことで,調波構造を持ち, 音高に依存する楽器音を対象にするもの.人間の知 覚において,音高が変化しても変化を感じない成分 (高調波成分の強度比など)を適切に扱うことができ る.調波構造を表現する混合正弦波モデルを用いるも の<sup>11),12)</sup>, SMF を基に調波構造にフィルタをかけるも の<sup>13)</sup>,時間周波数平面上での調波構造を表現するモデ ルのフィッティングによるもの<sup>14)</sup>,高調波成分のパワー エンベロープの類似性を用いるもの<sup>15)</sup>,ステレオ信号 のパワーと位相の共通性を用いるもの<sup>16)</sup>などがある. これらは主に調波構造を含む楽器音のみを対象として おり,一般の楽器音を分離することは困難であった.

(2) 音高を明示的に扱わず,一般の楽器音を対象 にするもの.原理的には任意の楽器音を扱うことがで きるが,各々の従来研究ではピアノやバイオリンなどの 楽音とドラムスなどの噪音のいずれか一方を主として 扱ってきた.楽音を扱うものでは, Non-negative Matrix Factorization (NMF) やその拡張である Nonnegative Tensor Factorization (NTF)を用いるも の<sup>17),18)</sup> などが, 噪音を扱うものでは, Independent Component Analysis (ICA)を用いるもの<sup>19)</sup>, NMF を用いるもの<sup>20)</sup>,ドラム音検出の後にスペクトル変調 を行うもの<sup>21)</sup> などがある.また, ICA などの統計的手 法で楽音と噪音を同時に分離するもの22)-26)があるが, 発音区間や周波数成分のスパースネスが保証されない 複雑な音響信号を扱うには至っていない.このような 音響信号を分離するためには楽器音認識が不可欠で, 楽器音認識と音源分離を併用したもの<sup>27)</sup>などが研究 事例としてあげられるが,対象がドラム音のみであり, 調波構造を持つ音には適用されていなかったりした.

これらのアプローチは従来排他的で,双方の長所を あわせ持つ手法はこれまで存在しなかった.

本論文で分離の対象とする音は、「調波的な音」および「非調波的な音」、およびそれらを加算して得られる 音である.以下の性質を満たす音を調波的な音と呼ぶ.

- 調波構造を持つ.
- 各高調波成分の相対強度が時間の経過によって変化しない.

急激な F0 の変化を含まない.

具体的には,弦(ピアノ弦やギター弦など)や管内 の空気(フルートなど)の定常的な振動によって得ら れる音が相当する.歌声は各高調波成分の相対強度が 母音の遷移によって連続的に変化するため,本論文で は扱わない.また,パワースペクトルに調波構造を含 めた周波数方向への鋭いピークが存在しない音を非調 波的な音と呼ぶ.具体的には,ドラム音を想定してい る<sup>\*1</sup>.ピアノが発音時にハンマで弦を叩く音のような, 調波構造を含む音から調波構造を取り除くことによっ

<sup>\*1</sup> 理想的な膜振動から得られる信号には「整数倍でない倍音構造」 が含まれるため、ドラム音の中には周波数方向への鋭いピーク が存在するものがある.しかし本論文では「ドラム音には周波数 方向への鋭いピークは存在しない」と仮定し、ドラム音を扱う.

て得られる音も,周波数方向への鋭いピークがほとん ど存在しないと見なし非調波的な音に含める.

従来の音源分離に関する研究の多くは,前述のと おりこれらの2種類の音の一方のみに着目していた. Goto<sup>28)</sup>は,混合音中の最も優勢な調波構造を抽出す る手法について述べており,さらに調波構造モデルを 他の任意の関数の重み付き混合モデルに置き換えても 調波構造の場合と同様にモデルパラメータ推定を行う ことで,パワースペクトル上の任意の構造を扱うこと が理論的には可能であるとも述べているが,調波構造以 外の構造をどのようなモデルで扱えばよいか,そのよう なモデルを実現するうえでどのような問題点があるか, といった具体的な手法については述べていなかった.

我々は Goto<sup>28)</sup>の示唆を受け,調波構造モデルと非 調波構造モデルを統合した混合モデルを用いた音源分 離手法を設計し,実現した.調波構造モデルは,音高 を持つ楽器の単音の調波構造を表現するパラメトリッ クモデルに基づいており,発音時刻,音長,音量,音 高(F0)の時間変化,パワーエンベロープの時間変化, 各高調波成分の相対強度といったパラメータで表現さ れる.非調波構造モデルは,ノンパラメトリックモデ ルに基づいており,調波構造では表現が難しいドラム 音などのパワースペクトルをそのまま表現する.また 前述のように,ピアノやギターなどの調波構造を持つ 楽器音であっても,発音時には弦をハンマで叩くこと や弦を弾くことに由来する非調波成分を含んでいるの で,それらのパワースペクトルも非調波構造モデルで 表現する.

SMF などから抽出した各単音の音高,音長,音量, 発音時刻,楽器によってこれらのモデルのパラメータ を初期化し,モデルパラメータの最大事後確率推定を EM アルゴリズムを用いて実現する.このパラメータ 推定における問題点は,非調波構造モデルは大きな自 由度を持っており,あらゆるパワースペクトルを表現 できるため,パラメータ推定の結果,非調波構造モデ ルが調波成分も含めてすべての混合音を表現してし まうことである.この問題を解決するため,非調波構 造モデルの形状に関する制約や同一楽器のモデルパラ メータ類似性に関する制約を導入する.このようにし て得られた調波・非調波統合モデルを用いることで, 混合音のパワースペクトル上での分離が可能となる.

本論文の構成は以下のとおりである.まず,2章で 音源分離における2つの問題点を述べる.続いて,3章 では調波・非調波統合モデルの定式化,4章では統合モ デルに基づく音源分離処理,5章ではモデルパラメー タの推定処理について述べる.6章で評価実験につい て述べ,7章で本論文のまとめを行う.

2. 問題の所在と解決へのアプローチ

本研究の目標は,多重奏の音楽音響信号とその楽曲 の楽譜情報(各単音の音高・音長・音量・発音時刻・ 楽器種類)が与えられたとき,音響信号のパワースペ クトルを単音ごとに分離することである.いい換えれ ば,我々の目標は各楽器パートのすべての単音に対し て,単音に対応する調波構造モデルと非調波構造モデ ルの全パラメータを推定することである.

与えられた楽譜情報の各単音を個別に MIDI 音源で 演奏することで,音響信号中の各単音にある程度近い, 「音のサンプル」を作成できる.この音のサンプルを テンプレート音と呼ぶ.分離対象の音楽音響信号とそ の楽譜情報, MIDI 音源を用いて生成したテンプレー ト音が与えられたとき,我々が解くべき課題は以下の 2 点である.

(1) テンプレート音と入力音響信号とのずれの吸収.テンプレート音と入力信号との間には必ず音響的な違いがあるので,テンプレート音をそのまま用いたのでは完全に1つ1つの音を分離することはできない.そこで,テンプレート音と入力信号との音響的差異を吸収する手法が必要となる.

(2) 奏法に独立な楽器音の音色類似性の達成.あ る楽器を用いて,同一の音高,音長,音量を持つ単音を 複数回にわたって演奏したとしても,奏法(ビブラー トのかけ方など)が異なれば異なる音色を持つ音響信 号が生成される.単音ごとの音色の違いを表現するた めには,音色を表現するモデルを単音ごとに作成する 必要がある.しかし,単音ごとに独立したモデルを作 成すると,パラメータの自由度が大きくなりすぎてし まうため,混合音への適応によってモデルが過学習を 起こし,結果として分離性能が低下してしまう可能性 がある.これを防ぐためには,同一楽器の単音の間に 存在する音色の類似性を満たすような,音色の表現方 法を実現する必要がある.

これらの課題を,以下のアプローチで解決する.

(1) モデルパラメータ適応.テンプレート音で初期化した音モデルのパラメータを,モデルと入力音響信号とのパワースペクトル上での音響的差異を最小化するように更新する.すなわち,音モデルのパラメータを入力音響信号に適応させることによってテンプレート音と実演奏とのずれを吸収する.

(2) 同一楽器内パラメータ類似性に対する制約. 同一楽器の個々の単音モデル間のパラメータの類似性 を保ちつつも各単音の微小な違いを許容するような制 約の下でモデルパラメータの更新を行う.これは,同 一楽器に属する各単音のモデルパラメータの平均値 と現在着目している単音のモデルパラメータとの間の Kullback-Leibler ダイバージェンス(以下,KLDと 略す)を最小化することによる,モデルパラメータに 対する制約を与えることで達成できる.

# 2.1 問題の定義

本論文で扱う分離問題とは,入力混合音のパワース ペクトル X(c,t,f) を,k 番目の楽器,l 番目の単音 (以下, (k, l) 番目と記す)のパワースペクトルに分解 することである.ここで, $c \in \{1, \ldots, C\}$ は左右など のチャネルの番号,  $t \in [T_0, T_1]$ は時刻,  $f \in [F_0, F_1]$ は周波数を表す.入力された楽譜情報から,楽曲中で は K 種類の楽器が演奏されており, 各々の楽器は  $L_k$ 個の単音を持つものとする.すなわち, $k \in \{1, \ldots, K\}$ であり, 各々の k に対して  $l \in \{1, ..., L_k\}$  である.ま た, (k, l) 番目の単音を表すモデルを  $J_{kl}(c, t, f)$ , (k, l)番目のテンプレート音のパワースペクトルを  $Y_{kl}(t, f)$ とする.本論文で用いる楽譜情報にはチャネル間音圧 比は含まれていないため,テンプレート音にはチャネ ル間音圧比を設定せず,モノラル音響信号として生成  $farce{farce}{farce}$  ,  $Y_{kl}(t, f)$  には c がなく 1 チャネル となっている.また,

$$X_{0} = \frac{1}{C} \sum_{c} \iint X(c, t, f) dt df$$
$$= \sum_{k,l} \iint Y_{kl}(t, f) dt df$$
(1)

となるように ,  $Y_{kl}(t,f)$  のパワーを正規化してあるものとする .

# 3. 調波・非調波統合モデル

調波・非調波統合モデル  $J_{kl}(c,t,f)$  は, (k,l) 番目 の単音のパワースペクトルを表現するモデルである. 調波的な音のパワースペクトルを表現する調波構造モ デル  $H_{kl}(t,f)$  と非調波的な音のパワースペクトルを 表現する非調波構造モデル  $I_{kl}(t,f)$  との和に統合モデ ル全体の重み  $w_{kl}$  を乗じた  $J'_{kl}(t,f)$  に,さらにチャ ネルごとの重み  $r_{kl}(c)$  を乗じたもので,以下の式で 定義する.

$$J_{kl}(c,t,f) = r_{kl}(c)J'_{kl}(t,f)$$
(2)

$$J'_{kl}(t,f) = w_{kl} \left( H_{kl}(t,f) + I_{kl}(t,f) \right)$$
(3)

 $w_{kl}$  および  $r_{kl}(c)$  は以下の各条件を満たす.

$$\forall k, l, \ \sum_{k,l} w_{kl} = X_0 \tag{4}$$

$$\forall k, l, \sum_{c} r_{kl}(c) = \mathcal{C}$$
(5)

3.1 調波構造モデル

調波構造モデル  $H_{kl}(t,f)$  は,パラメトリックな基 底関数であるガウス分布関数の重み付き線形和として, パワーエンベロープを表現する関数  $E_{klm}(t)$  と各時 刻の調波構造を表現する関数  $F_{kln}(t,f)$  を用いて以下 の式で定義する.ただし,M,N は定数で,それぞれ パワーエンベロープを表現するガウス分布関数の数と 高調波成分を表現するガウス分布関数の数を表す.

$$H_{kl}(t,f) = \sum_{m,n} E_{klm}(t) F_{kln}(t,f)$$
(6)

$$E_{klm}(t) = \frac{u_{klm}}{\sqrt{2\pi}\phi_{kl}} e^{-\frac{(t-\tau_{kl}-m\phi_{kl})^2}{2\phi_{kl}^2}}$$
(7)

$$F_{kln}(t,f) = \frac{v_{kln}}{\sqrt{2\pi\sigma_{kl}}} e^{-\frac{(f-n\omega_{kl}(t))^2}{2\sigma_{kl}^2}}$$
(8)

このモデルは, Kameoka らの調波時間構造化クラス タリング(Harmonic Temporal Clustering; HTC)<sup>14)</sup> で用いられる音源モデルを参考に設計した.Kameoka らの HTC 音源モデルでは, $\omega_{kl}(t)$ は時間 t に関する 多項式として定義されていたが,多項式で表現可能な F0 時系列の集合は任意の F0 時系列の集合よりも小さ い.そこで本研究では,任意の F0 時系列を表現する ために,ノンパラメトリックな関数として $\omega_{kl}(t)$ を 定義した. $u_{klm}$ , $v_{kln}$ は以下の条件を満たす.

$$\forall k, l, n, \ \sum_{m} u_{klm} = 1 \tag{9}$$

$$\forall k, l, m, \ \sum_{n} u_{kln} = 1 \tag{10}$$

# 3.2 非調波構造モデル

非調波構造モデル  $I_{kl}(t, f)$  は, ノンパラメトリック な関数として,パワースペクトルの各時刻および周波 数における周波数成分の強度を直接表現するように以 下の式で定義する.

$$I_{kl}(t,f) = w_{kl}^{(I)} I'_{kl}(t,f)$$
(11)

ただし, $I_{kl}^{\prime}(t,f)$ , $w_{kl}^{(H)}$  および  $w_{kl}^{(I)}$  は以下の各条件を満たす.

$$\forall k, l, \iint I'_{kl}(t, f) \, dt \, df = 1 \tag{12}$$

$$\forall k, l, \ w_{kl}^{(H)} + w_{kl}^{(I)} = 1 \tag{13}$$

# 4. 音源分離

入力パワースペクトル X(c,t,f) を (k,l) 番目の 単音へと分離するためのパワースペクトル分配関数  $S_{kl}(c,t,f)$  を導入する.この関数は以下の条件を満 たす.

$$\forall c, t, f, \ \sum_{k,l} S_{kl}(c, t, f) = 1$$
(14)

分離された (k, l) 番目の単音のパワースペクトルは

$$X_{kl}^{(S)}(c,t,f) = S_{kl}(c,t,f)X(c,t,f)$$
(15)

#### で表される.

ここで,どのように  $S_{kl}(c,t,f)$  を定めると最も良 い分離が行えるかを考える.それには分離の良し悪し を計る尺度が必要なので, $X_{kl}^{(S)}(c,t,f)$  と  $J_{kl}(c,t,f)$ との Kullback-Leibler Divergence (KLD)  $Q'_{kl}$ :

$$Q'_{kl} = \sum_{c} \iint X^{(S)}_{kl}(c,t,f) \log \frac{X^{(S)}_{kl}(c,t,f)}{J_{kl}(c,t,f)} \, dt \, df$$
(16)

でこの尺度を定義する. KLD は距離の公理を満たさ ないが, あらゆる k, l, c, t, f に対して

$$X_{kl}^{(S)}(c,t,f) = J_{kl}(c,t,f)$$
(17)

となるときに限り最小値 0 をとるので , パワースペ クトル間の類似度として用いることができる .

このとき, $Q'_{kl}$ を最小化するような $S_{kl}(c,t,f)$ を 求めることができれば,それを用いた $X^{(S)}_{kl}(c,t,f)$ が $Q'_{kl}$ に基づく最適な分離結果となる.ただし,  $S_{kl}(c,t,f)$ は式(14)を満たさなければならないので, あらゆるk,lに関して同時に $Q'_{kl}$ を最小化する必要 がある.そこで, $Q'_{kl}$ をあらゆるk,lに関して足し 合わせ,式(14)の条件に対する未定乗数 $\lambda^{(S)}(c,t,f)$ による Lagrange の未定乗数項を加えたQ':

$$Q' = \sum_{k,l} Q'_{kl}$$
$$-\sum_{c} \iint \lambda^{(S)}(c,t,f)$$
$$\cdot \left(\sum_{k,l} S_{kl}(c,t,f) - 1\right) dt df$$
(18)

を最小化する . Q' は  $S_{kl}(c,t,f)$  に関する制約条件を 満たす空間において凸関数であるので,連立方程式

$$\frac{\partial Q'}{\partial S_{kl}(c,t,f)} = 0, \quad \frac{\partial Q'}{\partial \lambda^{(S)}(c,t,f)} = 0 \quad (19)$$

の解

$$S_{kl}(c,t,f) = \frac{J_{kl}(c,t,f)}{\sum_{k,l} J_{kl}(c,t,f)}$$
(20)

が, Q'を最小化する  $S_{kl}(c,t,f)$  であり, これにより 統合モデルに基づく分離が行われる.

図1に全体の処理の流れを示す.最初に楽譜情報と



図1 分離とモデル適応の処理の流れ Fig.1 Overview of separation and model adaptation.

MIDI 音源を用いて,テンプレート音の録音を行い, テンプレートを基にモデルパラメータを初期化する. その後,Q'をS<sub>kl</sub>(c,t,f) に関して最小化することに よるパワースペクトルの分離(Separation Process) と,分離パワースペクトルを用いたモデルパラメータ 推定(Model Adaptation)を交互に繰り返すことで 分離処理が進められる.モデルパラメータ推定につい ては次章で詳細を述べる.

# 5. パラメータ推定

分離の場合と同様に,推定されたパラメータの良し 悪しを計る尺度を $S_{kl}(c,t,f)X(c,t,f)$ と $J_{kl}(c,t,f)$ の間の KLD である $Q'_{kl}$ および制約条件として作用 するいくつかの追加コストの重み付き和で定め,この 尺度を最小化するようにパラメータ推定を行う.本章 では,パラメータ推定における各々の制約条件につい て述べ,その後パラメータ推定方法および推定手法と EM アルゴリズムとの関連について述べる.

5.1 テンプレート音との類似性

(k, l) 番目の単音のテンプレート音は, 混合音中の (k, l) 番目の単音とは楽器個体などが異なっているも のの, SMF から抽出した音高, 音長, 音量, 楽器に よって MIDI 音源から生成した音響信号であるので, 混合音中の(k, l) 番目の単音と「かなり」近いパワー スペクトルを持つと考えられる.つまり, モデルとテ ンプレートスペクトル Y<sub>kl</sub>(t, f) との差を同時に最小 化することで, 混合音中の(k, l) 番目の単音から推定 される理想的なパラメータに「かなり」近いパラメー タを推定することが可能になる.

分離スペクトルとモデルとの差を KLD で定義した ことと同様に,テンプレートスペクトルとモデルとの 差  $Q_{kl}^{(Y)}$  を KLD で定義する.ただし,テンプレート はモノラル音響信号であるので,モデルはチャネル間 音圧比を含まないものを用いる.

$$Q_{kl}^{(Y)} = \iint Y_{kl}(t,f) \log \frac{Y_{kl}(t,f)}{J'_{kl}(t,f)} dt df \quad (21)$$
  
5.2 調波構造モデルへの制約

調波構造モデルのパラメータ  $\omega_{kl}(t)$  は大きい自由 度を持つため,たとえば,分離音に他の楽器音が混入 しており,本来モデルが表現すべき楽器音の音高では なくそちらの楽器音の音高が推定されてしまった場合,  $\omega_{kl}(t)$  に不連続な区間が含まれてしまう.しかし,通 常はこのような F0 時系列は楽器音の単音として望ま しくなく,単音の F0 時系列は連続的に変化するべき と考える.そこで,以下の式で表現される新たなコス トを導入する.

$$Q_{kl}^{(\omega)} = \int \tilde{\omega}_{kl}(t) \log \frac{\tilde{\omega}_{kl}(t)}{\omega_{kl}(t)} dt$$
(22)

 $\tilde{\omega}_{kl}(t)$ は,  $\omega_{kl}(t)$ にガウシアンフィルタを畳み込む ことで時間方向に平滑化したもので,このコストは  $\tilde{\omega}_{kl}(t)$ と  $\omega_{kl}(t)$ の KLD で定義されている. $\tilde{\omega}_{kl}(t)$ は  $\omega_{kl}(t)$ よりも「滑らか」であるため, $Q'_{kl}$ を最小 化すると同時に  $Q^{(\omega)}_{kl}$ を最小化することで,推定され た  $\omega_{kl}(t)$ が「滑らか」であることを強制することが できる.

### 5.3 非調波構造モデルへの制約

1 章で述べたように,非調波構造モデル *I*<sub>kl</sub>(*t*, *f*) は 非常に大きい自由度を持つ.そのため,このモデルは 任意のパワースペクトルを表現することが可能で,入 カパワースペクトル,および分離パワースペクトルが 調波構造モデルを用いることなくこのモデルだけで表 現されてしまう可能性がある.しかし,調波構造モデ ルが表現すべき調波構造までも非調波構造モデルが表 現してしまうことは望ましくない.

我々は,パワースペクトル上の調波構造に関する最 大の特徴は周波数方向に複数のピークを持つこと,と 考える.逆に,周波数方向に強いピークを持たないパ ワースペクトルが調波構造を含むことはない.した がって,非調波構造モデルが周波数方向に強いピーク を持ちにくくさせる制約を与えることができれば,前 述の問題を解決できると考える.

この制約を実現するため,以下の式で表現される新 たなコストを導入する.

$$Q_{kl}^{(\tilde{I})} = \iint \tilde{I}_{kl}(t,f) \log \frac{\tilde{I}_{kl}(t,f)}{I'_{kl}(t,f)} dt df$$
(23)  
ここで,  $\bar{I}_{kl}(t,f)$ は  $I'_{kl}(t,f)$ にガウシアンフィルタを

畳み込むことで,周波数方向に平滑化したものであり,  $I'_{kl}(t,f)$ よりも周波数方向に滑らかであるため周波数 方向のピークを持たない,もしくは持っていたとして もピークにおける周波数成分のパワーは $I'_{kl}(t,f)$ の それより小さい.したがって,パラメータ推定時にこ のコスト $Q^{(\bar{I})}_{kl}$ も最小化することで,非調波構造モデ ルが調波構造を表現せず,それ以外の「非調波的な音」 のパワースペクトルだけを表現することを強制できる.

#### 5.4 同一楽器内のパラメータ類似性

2 章で述べたように,調波・非調波統合モデル *J<sub>kl</sub>(c,t,f)*のパラメータは,単音ごとの微小な誤差 を許容しつつ,かつ同一楽器ごとの類似性を満たす必 要がある.この性質を満たすようなパラメータを推定 するために,新たに2つの制約を導入する.

第1の制約は調波構造モデルの v<sub>kln</sub> に対する制約 で,以下の式で表されるコスト関数として与えられる.

$$Q_{kl}^{(v)} = \sum_{n} \bar{v}_{kn} \log \frac{\bar{v}_{kn}}{v_{kln}} \tag{24}$$

 $\bar{v}_{kn}$  は楽器ごとに  $v_{kln}$  を平均したものである.この コストは  $\bar{v}_{kn}$  と  $v_{kln}$  の KLD で定義されており,こ れを最小化することで  $v_{kln}$  を  $\bar{v}_{kn}$  へと近づけること が可能になる.すなわち, $v_{kln}$  の楽器ごとの類似性を 強制できる.

第2の制約は非調波構造モデル  $I'_{kl}(t,f)$  に対する もので,以下の式で表されるコスト関数として与えら れる.

$$Q_{kl}^{(\bar{I})} = \iint \bar{I}_k(t, f) \log \frac{\bar{I}_k(t, f)}{I'_{kl}(t, f)} \, dt \, df \qquad (25)$$

 $\bar{I}_k(t,f)$ は楽器ごとに  $I'_{kl}(t,f)$ を平均したものである. このコストは  $\bar{I}_k(t,f)$  と  $I'_{kl}(t,f)$ の KLD で定義され ており,これを最小化することで  $I'_{kl}(t,f)$ を  $\bar{I}_k(t,f)$ へと近づけることができる.すなわち, $I'_{kl}(t,f)$ の楽 器ごとの類似性を強制できる.

# 5.5 基底関数への分配関数

 $J_{kl}(c,t,f)$ は  $H_{kl}(t,f)$ と  $I_{kl}(t,f)$ の線形結合で定義されており, さらに  $H_{kl}(t,f)$ は基底関数であるガウス分布関数の線形結合で定義されている.そこで,パラメータ推定のために,  $X_{kl}^{(S)}(c,t,f)$ および  $Y_{kl}(t,f)$ を調波構造モデルの m 番目のパワーエンベロープのガウス分布, n 番目の調波構造のガウス分布(以下(m,n) 番目と記す)および非調波構造モデルへと分配する関数  $S_{klmn}^{(H)}(t,f)$ ,  $S_{kl}^{(I)}(t,f)$ を導入する.この場合,  $Q_{kl}'$ は以下の  $Q_{kl}''$ に置き換えられることになる.

$$Q_{kl}^{\prime\prime} = Q_{kl}^{\prime\prime(H)} + Q_{kl}^{\prime\prime(I)}$$
(26)

ただし,

$$= S_{kl}^{(I)}(c,t,f) S_{kl}(c,t,f) X(c,t,f)$$
(28)

$$J_{klmn}^{(m)}(c,t,f) = (c,t,f) =$$

$$= r_{kl}(c)w_{kl}w_{kl} + E_{klm}(t)F_{kln}(t,f)$$
(29)

$$J_{kl}^{(*)}(c,t,f) = r_{kl}(c)w_{kl}w_{kl}^{(*)}I_{kl}(t,f)$$
(30)

$$Q_{kl}^{\prime\prime(H)} = \sum_{c,m,n} \iint X_{klmn}^{(H)}(c,t,f)$$

$$X_{klmn}^{(H)}(c,t,f) \quad (a.b)$$

$$\cdot \log \frac{X_{klmn}^{(H)}(c,t,f)}{J_{klmn}^{(H)}(c,t,f)} dt df$$
(31)

$$Q_{kl}^{\prime\prime(I)} = \sum_{c} \iint X_{kl}^{(I)}(c,t,f) \\ \cdot \log \frac{X_{kl}^{(I)}(c,t,f)}{J_{kl}^{(I)}(c,t,f)} \, dt \, df$$
(32)

である.これらの分配関数の最適な値は,Lagrange の未定乗数項を加えることで *Skl*(*c*,*t*,*f*) と同様に導 出することが可能で,それぞれ

$$S_{klmn}^{(H)}(t,f) = \frac{w_{kl}w_{kl}^{(H)}E_{klm}(t)F_{kln}(t,f)}{J_{kl}(c,t,f)}$$
(33)  
$$S_{kl}^{(I)}(t,f) = \frac{w_{kl}w_{kl}^{(I)}I_{kl}(t,f)}{J_{kl}(c,t,f)}$$
(34)

となる.

5.6 コスト関数

推定されたパラメータの分離音に対する良し悪しを 計る尺度  $Q_{kl}^{''}$ , ここまでに述べた追加コスト  $Q_{kl}^{(Y)}$ ,  $Q_{kl}^{(\bar{D})}$ ,  $Q_{kl}^{(\bar{$ 

さらに,  $Q_{kl}$ をすべての (k,l)に対して合計したコ スト関数 Qを定義すると,分離とパラメータ推定の 両方をコスト関数 Q の最小化としてとらえることが できる.分離に関係がある項は  $Q_{kl}^{\prime\prime}$ だけであるので, 重み  $\alpha$ がかけられていたり,  $Q_{kl}^{(Y)}$ などのその他のコ ストが加えられていても Q を最小化して得られる分 配関数には影響しない.また,パラメータ推定に関し ては,(k,l) 番目の単音のモデルパラメータ推定に関 係がある項は  $Q_{kl}$  だけであるので, $Q_{kl}$  以外の Q の 項は  $J_{kl}(c,t,f)$  のパラメータ推定には影響しない.し たがって,分配関数およびモデルパラメータに関する Q の最小化を交互に繰り返すことで,分離とパラメー タ推定を行うことができる.

Qは,以下の式で定義される.

$$Q = \sum_{k,l} Q_{kl}$$

$$-\lambda^{(w)} \left( \sum_{k,l} w_{kl} - X_0 \right)$$

$$-\sum_c \iint \lambda^{(S)}(c,t,f)$$

$$\cdot \left( \sum_{k,l} S_{kl}(c,t,f) - 1 \right) dt df \qquad (35)$$

$$Q_{kl} = \alpha Q_{kl}'' + (1 - \alpha) Q_{kl}^{(Y)} + \beta_{\omega} Q_{kl}^{(\omega)}$$

$$+ \beta_{\tilde{I}} Q_{kl}^{(\tilde{I})} + \beta_v Q_{kl}^{(v)} + \beta_{\tilde{I}} Q_{kl}^{(\tilde{I})}$$

$$- \lambda_{kl}^{(r)} \left( \sum_c r_{kl}(c) - C \right)$$

$$- \lambda_{kl}^{(wHI)} \left( w_{kl}^{(H)} + w_{kl}^{(I)} - 1 \right)$$

$$- \lambda_{kl}^{(\omega)} \left( \sum_m u_{klm} - 1 \right)$$

$$- \lambda_{kl}^{(\omega)} \int (\omega_{kl}(t) - \tilde{\omega}_{kl}(t)) dt$$

$$- \lambda_{kl}^{(\omega)} \int (\omega_{kl}(t) - \tilde{\omega}_{kl}(t)) dt$$

$$- \lambda_{kl}^{(G)} \left( \iint I_{kl}'(t,f) dt df - 1 \right)$$

$$- \iint \lambda_{kl}^{(SHI)}(t,f)$$

$$\cdot \left( \sum_{m,n} S_{klmn}^{(H)}(t,f) + S_{kl}^{(I)}(t,f) - 1 \right) dt df$$

$$(36)$$

ただし、コスト $Q^{(\omega)}$ 、 $Q^{(\bar{I})}$ に関しては、これらを追 加すると分離とパラメータ推定の反復を行う過程でQのパラメータに関する極小点がパラメータの値によっ て変化してしまうため、パラメータの局所的な収束性 が保証されなくなる.しかしながら、本論文で実施し た実験においては非調波構造モデルが調波構造を表現 したり F0 時系列が不連続になったりすることはなく、 かつパラメータが発散している様子は確認できなかっ

₹Ⅰ	本調い	くで用い	16記号の-	一覧
Tab	ole 1	List o	of symbols	s.

記号	意味	備考
с	チャネル番号	$c \in \{1, \ldots, C\}$
t	時刻	$t \in [\mathrm{T}_0, \mathrm{T}_1]$
f	周波数	$f \in [F_0, F_1]$
k	楽器番号	$k \in \{1, \dots, \mathrm{K}\}$
l	単音番号	$\forall k, \exists \mathbf{L}_k, \ l \in \{1, \dots, \mathbf{L}_k\}$
m	パワーエンベロープを表現するガウス分布関数の番号	$m \in \{1, \dots, M\}$
n	調波構造を表現するガウス分布関数の番号(第 n 次倍音に相当)	$n \in \{1, \dots, \mathrm{N}\}$
$J_{kl}(c,t,f)$	(k,l) 番目の単音の調波・非調波統合モデル	
$J_{kl}^{\prime}(t,f)$	チャネル間音圧比 $r_{kl}(c)$ を取り除いた $J_{kl}(c,t,f)$	
$H_{kl}(t,f)$	調波構造モデル	
$E_{klm}(t)$	調波構造モデルのパワーエンベロープ関数	
$F_{kln}(t,f)$	調波構造モデルの調波構造関数	
$I_{kl}(t,f)$	非調波構造モデル	
X(c, t, f)	入力パワースペクトル	
$Y_{kl}(t,f)$	テンプレートスペクトル	
$S_{kl}(c,t,f)$	(k,l)番目の単音へのパワースペクトル分配関数	$\forall c, t, f, \sum_{k,l} S_{kl}(c, t, f) = 1$
$S_{klmn}^{(H)}(t,f)$	(k,l) 番目の単音, $(m,n)$ 番目の調波構造モデルの	$\forall k, l, t, f,$
	ガウス分布へのパワースペクトル分配関数	$\sum_{m=n} S_{klmn}^{(H)}(t,f) + S^{(I)}(t,f) = 1$
$S_{kl}^{\left( I\right) }(t,f)$	(k,l) 番目の単音の非調波構造モデルへのパワースペクトル分配関数	
$r_{kl}(c)$	チャネル間音圧比	$\forall k, l, \sum_{c} r_{kl}(c) = C$
$w_{kl}$	調波・非調波統合モデル全体の重み	$\sum_{k,l} \overline{w_{kl}} = 1$
$w_{kl}^{(H)}$	調波構造モデルの重み	$\overline{\forall k, l}, w_{kl}^{(H)} + w_{kl}^{(I)} = 1$
$u_{klm}$	パワーエンベロープを表現する $m$ 番目のガウス分布関数の重み係数	$\forall k, l, \sum_{m} u_{klm} = 1$
$v_{kln}$	n 次高調波成分の相対強度	$\forall k, l, \sum_{n=1}^{m} v_{kln} = 1$
$\tau_{kl}$	発音時刻	<i>it</i>
$\phi_{kl}$	パワーエンベロープを表現するガウス分布関数の広がりを表すパラメータ	
$\omega_{kl}(t)$	F0 時系列	
$\sigma_{kl}$	高調波成分を表現するガウス分布関数の広がりを表すパラメータ	
$w_{l,l}^{(I)}$	非調波構造モデルの重み	
$I_{kl}^{\prime \kappa \iota}(t,f)$	パワーを正規化した非調波構造モデル	$\forall k, l, \iint I'_{k,l}(t, f)  dt  df = 1$

たため,実験的には大きな問題にはならないと考える. 本論文で用いる記号の一覧を,表1に示す.

5.7 EM アルゴリズムとしての解釈

ここまでに述べた分離とパラメータ推定の繰返しは, Expectation-Maximization(EM)アルゴリズムを用 いた最大事後確率推定として解釈することもできる.

観測確率密度関数 p(c,t,f) が与えられたとき, p(c,t,f)を近似する確率密度分布  $p(k,l,c,t,f|\theta)$ の パラメータ  $\theta$ を推定する問題を考える.k, l に関す る隠れ変数の分布 p(k,l|c,t,f) および  $\theta$  の事前分布  $p(\theta)$ を導入すると,  $\theta \in \tilde{\theta}$  へと変化させたときの対 数事後確率の期待値の増加量  $Q(\theta, \tilde{\theta})$  は,

$$Q(\theta, \tilde{\theta}) = \sum_{k,l,c} \iint p(k, l|c, t, f, \theta) p(c, t, f)$$
$$\cdot \log p(k, l, c, t, f|\tilde{\theta}) dt df$$
$$+ \log p(\tilde{\theta})$$
(37)

と表される.

ここで,入力パワースペクトルX(c,t,f),パワース

ペクトル分配関数  $S_{kl}(c,t,f)$ , 調波・非調波統合モデ ル  $J_{kl}(c,t,f)$ をそれぞれ p(c,t,f),  $p(k,l|c,t,f,\theta)$ ,  $p(k,l,c,t,f|\theta)$ に対応させ<sup>\*1</sup>, さらにモデルに関する 制約を表現するコスト関数  $Q_{kl}^{(Y)}$ ,  $Q_{kl}^{(\omega)}$ ,  $Q_{kl}^{(I)}$ ,  $Q_{kl}^{(v)}$ ,  $Q_{kl}^{(I)}$ の総和を  $-\log p(\theta)$ に対応させると, EM アル ゴリズムにおける E ステップ, すなわち  $Q(\theta, \tilde{\theta})$ を最 大化する  $p(k,l|c,t,f,\theta)$ の推定はコスト関数 Qを最 小化する  $S_{kl}(c,t,f)$ の導出に対応する. 同様に, EM アルゴリズムにおける M ステップ, すなわち  $Q(\theta, \tilde{\theta})$ を最大化する  $\tilde{\theta}$ の推定はコスト関数を最小化するモ デルの各パラメータの導出に対応する.

# 6. 評価実験

本手法の性能を確認するため,評価実験を行った. 6.1 実験の目的 本実験の目的は,本論文で構築した調波・非調波統

<sup>\*1</sup> ただし、入力パワースペクトルとモデルはあらゆる変数に関し て積分した値が1になるように正規化したものを対応させる.

合モデルの有効性を確認することである.具体的には, 以下の3つの条件:

(1) 調波・非調波統合モデルを用いた場合(本手法)

(2) 調波構造モデルのみを用いた場合

(3) 非調波構造モデルのみを用いた場合

において楽曲の音響信号を楽器パートごとおよび単 音ごとに分離し,単音ごとにに分離したパワースペク トルとミックス前の各単音のパワースペクトルとの周 波数領域での SNR を用いて 3 つの条件を比較した. (*k*,*l*)番目の単音に関する周波数領域での SNR は,次 式で定義する.

$$\operatorname{SNR}_{kl} = \frac{1}{\operatorname{C}(\operatorname{T}_1 - \operatorname{T}_0)} \sum_{c} \int \operatorname{SNR}_{kl}(c, t) \, dt$$
(38)

$$SNR_{kl}(c,t) = \int \frac{X_{kl}^{(S)}(c,t,f)^2}{\left(X_{kl}^{(S)}(c,t,f) - X_{kl}^{(R)}(c,t,f)\right)^2} df \quad (39)$$

ただし, $X_{kl}^{(S)}(c,t,f)$ は(k,l)番目の単音の分離パワースペクトル, $X_{kl}^{(R)}(c,t,f)$ は(k,l)番目の単音の参照パワースペクトル,すなわちミックス前のパワースペクトルである.

本来,調波構造モデルでドラム音を表現することは 難しいが,本実験では用いたモデル以外の条件を同一 にするために,調波構造モデルのみを用いた場合でも ドラムパートを含む混合音の分離を行った.また,非 調波構造モデルのみを用いた場合は,式(23)の非調 波構造モデルを平滑化する制約を用いるとモデルが調 波構造を表現することができなくなるため,この制約 は用いていない.

6.2 実験データ

実験には, RWC 研究用音楽データベース:ポピュ ラー音楽(RWC-MDB-P-2001)<sup>29)</sup>から選んだ10曲 (No.1–10)を用いた.各楽曲は開始から30秒の区間 を利用した.楽曲ごとの K と L<sub>k</sub>の値は,各楽曲の SMF で用いられた楽器数,およびノートオンメッセー ジ数から求める.本実験では以下の理由により,MIDI 音源から録音した音響信号を分離対象とした.

本実験の目的に適した音楽データベースが入手
 困難である.目的に適したデータベースとは,楽曲の
 音響信号,それに同期がとられた SMF,さらに各楽器
 パートの音響信号をミックスする前のマスタートラックのすべてが利用できるものである.これは,ミックス前の音響信号に対する分離後の音響信号の SNRを
 測定し,定量的評価を行うために必要である.

一般的なポピュラー音楽には歌声が含まれてい

表 2 実験条件 Table 2 Experimental conditions.

Frequency analysis				
sampling rate	$44.1  \mathrm{kHz}$			
analysizing method	STFT			
STFT window	2,048 points Gaussian			
STFT shift	441 points			
Parameters				
# of channels: C	2			
# of kernels in $E_{klm}$ : M	10			
# of partials: N	20			
$\beta_v$	0.1			
$eta_\omega$	0.1			
$\beta_{ar{I}}$	3.5			
$\beta_{\tilde{I}}$	0.5			
MIDI sound generator				
test data	YAMAHA MU2000			
template sounds	Roland SD-90			

るが,1章で述べたように歌声は本手法で扱う対象に は含まれない.

テンプレート音と分離対象となる楽曲とは,異なる 楽器メーカの MIDI 音源で生成した.これは,分離対 象とは確実に音色(音源波形)が異なるテンプレート 音を用いるためである.また周波数解析,ハイパーパ ラメータ,用いた MIDI 音源に関する条件を表2に示 す.この表における β<sub>v</sub> などのパラメータは,実験的 に最適なものを求めたものである.

# 6.3 実験結果

表3に,統合モデルを用いて分離を行った場合(Integrated),調波構造モデルのみを用いて分離を行った 場合(Harmonic),および非調波構造モデルのみを用 いて分離を行った場合 (Inharmonic) に得た各単音ご との分離音を周波数領域での SNR で評価し,同種の 楽器ごとに平均した結果を示す. Piano, Guitar, ..., Drums は楽器の種類を, 1-8 などの番号は楽器に対応 する MIDI のプログラムナンバを表す. 楽器の種類は, MIDI のプログラムナンバを基準に,ピアノ(プログ ラムナンバ 1-8), ギター (プログラムナンバ 24-32) のように分類した.ドラムには,プログラムナンバー では分類しておらず, SMF でドラムトラックである ことを指定されたトラックに属する単音を分類した. ただし表 4 では,実験に用いた 10 曲中の 5 曲以上 に用いられている楽器パートについてのみ結果を示し ている.すべての楽器パートにおいて, Integrated の SNR が最も大きくなっている.

Integrated と Inharmonic に比べて Harmonic の SNR が全体的に小さな値になっており,特に Bass と Drums において SNR の低下が顕著である. Drums に 関しては調波構造モデルでドラム音を表現することは

表 3 楽器パートごとの平均 SNR Table 3 Average SNR of each instrument part.

Inst	Prog. #		SNR (dB)	
Dent	in MIDI	Interneted	Harmania	Inhomonio
Fart	III MIDI	Integrated	Harmonic	Innarmonic
Piano	1 - 8	48.87	25.67	48.17
Guitar	25 - 32	47.22	13.59	46.88
Bass	33 - 40	40.93	-22.19	40.53
Ensemble	49 - 56	48.11	31.45	47.45
Pipe	73 - 80	49.28	26.78	49.16
Drums		43.31	-18.02	42.56

表 4 楽器パートごとの平均 KLD Table 4 Average KLD of each instrument part.

Inst.		KLD	
Part	Initial	Estimated	Ideal
Piano	402.82	21.25	12.68
Guitar	342.92	28.12	13.72
Bass	1152.2	44.93	25.45
Ensemble	520.67	32.98	13.22
Pipe	504.57	23.73	13.24
Drums	842.32	56.41	25.39

困難なことが, Bass に関しては低音部を担当する楽 器が多く, その基本周波数に対して周波数解像度が粗 いために正しく調波構造を表現することができなかっ たことが,それぞれ原因であると考えられる.これは Integrated と Inharmonic の場合にも同様にあてはま る.この問題に対しては,STFT よりも周波数解像 度を詳細に制御可能な連続ウェーブレット変換などを 周波数解析手法に採用することで対処できると考え られる.また,BassとDrums以外の楽器パートでも HarmonicのSNR が低下しているが,上記2つの楽 器パートを正しく分離することができなかった結果, それ以外のパートにドラム音などが混入したためと考 えられる.

楽器パートごとに, Integrated と Harmonic, Integrated と Inharmonic のそれぞれの組において SNR の 平均値が等しいという帰無仮説を立てて t 検定を行っ たところ, Integrated と Inharmonic 間の Pipe 以外の SNR に関してはすべて有意水準 0.05 で帰無仮説を棄 却することができたため,この結果は SNR の向上と いう点において統合モデルを用いることの有用性を示 している.SNR に優位な差が現れなかった原因とし ては, Pipe のパートの多くの単音がフルートでかつ フルートがメロディ(ボーカル)に割り当てられてお り,単音の音量が大きくかつ他の楽器パートとの高調 波成分の重複が起こりにくいため,統合モデルを用い なくても単音のパワースペクトルを容易に推定できた ことがあげられる.

表4に,テンプレートでパラメータを初期化した直

後のモデル (Initial), それを混合音に適応させたモデ ル (Estimated), テストデータの各単音でパラメータ を初期化したモデル(Ideal)に関して,テストデータ の各単音とモデルとの KLD を同種の楽器ごとに平均 した結果を示す.楽器の分類は表3の場合と同様であ る. Estimated の KLD が Ideal のものに近ければ推定 されたパラメータが「良い」と考えられるが,それぞ れの単音に関しての KLD が非負となることは保証さ れていないため,必ずしも KLD は小さいほど良いと いうわけではなく, KLD の大小はあくまで良し悪し の目安であることに注意する必要がある.表4を見る と, Estimated の KLD は Initial のものよりもかなり Ideal に近い値を示しており,ある程度は妥当なパラ メータ推定が行われたと考えられる.詳細なデータは 省略するが,モデルとテストデータの単音間の KLD が150-200を超える、もしくは0を下回るとSNRが 大きく減少する傾向が見られた.

また,図2,図3に,統合モデル,調波構造モデル, 非調波構造モデルのそれぞれを用いて分離されたピア ノ単音,スネアドラム単音のパワースペクトル(それ ぞれ図 {2,3}-{b, c, d}),該当する区間の混合音およ びミックス前のピアノ単音のパワースペクトル(それ) ぞれ図 {2,3}-{a, e}) を示す.これらの図を見ると, いずれのモデルを用いて分離を行った場合でも調波構 造はある程度表現されていることが分かる.しかし, 調波構造モデルを用いた場合では発音時刻付近に現れ ている非調波成分が少ない点,非調波構造モデルを用 いた場合では調波構造がミックス前のものよりも早く 減衰している点,本来の調波構造以外の周波数成分が 含まれている点などが,ミックス前の単音とは異なっ ている.これに対して,統合モデルを用いた場合では, これら3つのモデルによる分離音の中では上記の点の ようなミックス前の単音とのずれは最も小さい.調波 構造の減衰が非調波構造モデルを用いた場合よりも遅 い原因は,調波構造モデルの高調波成分相対強度が時 間の経過によって変化しないように定義されており,高 次の高調波成分だけが先に減衰しないことがあげられ る.また,統合モデルや非調波構造モデルを用いた場合 には,発音時刻以外の時刻においてミックス前の単音 には見られない非調波成分が存在している.この非調 波成分は本来は他のピアノ音やドラム音の非調波成分 であるが,非調波構造モデルが混合音に過剰に適応し た結果混入したと考えられる.これは非調波構造モデ ルの自由度の高さに由来しているため,非調波構造モ デルに構造上の制約を持たせる,もしくは何らかの事 前分布を用いるなどによって解決されると考えられる.







7. おわりに

本論文では,調波的な音と非調波的な音の混合音を すべての楽器パートに分離するという問題に取り組ん だ.具体的な手法として,調波構造モデルと非調波構 造モデルを統合したモデルを用いた多重奏の音源分離 手法と,そのためのモデル適応手法について述べた. また,調波構造モデルと非調波構造モデルを統合する 際の問題点を非調波構造モデルへの制約という形でコ スト関数最小化によるパラメータ推定の枠組みを崩す ことなく解決した.本手法の性能を示すための評価実 験では,統合モデルを用いることの有効性を確認した.

また我々は,本手法で生成した楽器パートごとの分 離音を用いるアプリケーションとして,1章で述べた 楽器音イコライザを開発した.楽器音イコライザを用 いることで,ユーザは楽曲の楽器パートごとの音量バ ランスを自由に操作することができ,その結果として 能動的な音楽鑑賞<sup>1)</sup>が可能となる.

本論文で設計した調波・非調波統合モデルは,音源 分離への利用のみに限定されるものではない.たとえ ば,本モデルを用いて,多重音解析や楽器音認識へと 応用範囲を広げることが考えられる.また,任意のパ ワースペクトルを表現できる非調波構造モデルの性質 上,扱う対象も音楽音響信号に限定されるものではな く,音声や環境音へとその対象を広げることも可能で ある.今後は,歌声の分離など,扱える信号の対象を 増やし分離性能を改善することと,SMFから得られ る事前情報(各単音の音高,音長,音量,発音時刻, 楽器)が一部利用できない,もしくは利用はできるが 信頼性が低いといった状況における分離手法を構築す ることで,分離技術の汎用的を高めることに取り組ん でいく予定である.

謝辞 本研究の一部は,科学研究費補助金(基盤研 究(S),特定領域「情報爆発IT基盤」),21世紀 COE プログラム「知識社会基盤構築のための情報学拠点形 成」,科学技術振興機構 CrestMuse プロジェクトによ る支援を受けた.

# 参考文献

- Goto, M.: Active Music Listening Interfaces Based on Signal Processing, *Proc. ICASSP*, Vol.IV, pp.1441–1444 (2007).
- Yoshii, K., Goto, M. and Okuno, H.G.: IN-TER:D: A Drum Sound Equalizer for Controlling Volume and Timbre of Drums, *Proc. EWIMT*, pp.205–212 (2005).
- 3) Yoshii, K., Goto, M., Komatani, K., Ogata,

T. and Okuno, H.G.: Drumix: An Audio Player with Real-time Drum-part Rearrangement Functions for Active Music Listening, *IPSJ Journal*, Vol.48, No.3, pp.134–144 (2007).

- 4) ヤマハミュージックイークラブ.http://www. music-eclub.com/
- 5) internet MIDILINK. http://www.midilink. com/
- Cano, P., Loscos, A. and Bonada, J.: Score-Performance Matching Using HMMs, *Proc. ICMC*, pp.441–444 (1999).
- Dannenberg, R.B. and Hu, N.: Polyphonic Audio Matching for Score Following and Intelligent Audio Editors, *Proc. ICMC*, pp.27–33 (2003).
- Adams, N., Marquez, D. and Wakefield, G.: Iterative Deepening for Melody Alignment and Retrieval, *Proc. ISMIR*, pp.199–206 (2005).
- 9) Dixon, S. and Widmer, G.: MATCH: A Music Alignment Tool Chest, *Proc. ISMIR*, pp.492– 497 (2005).
- 10) Cont, A.: Realtime Audio to Score Alignment for Polyphonic Music Instruments Using Sparce Non-negative Constraints and Hierarchical HMMs, *Proc. ICASSP*, Vol.II, pp.641– 644 (2006).
- Virtanen, T. and Klapuri, A.: Separation of Harmonic Sound Sources Using Sinusoidal Modeling, *Proc. ICASSP*, Vol.II, pp.765–768 (2000).
- 12) Virtanen, T. and Klapuri, A.: Separation of Harmonic Sounds Using Linear Models for the Overtone Series, *Proc. ICASSP*, Vol.2, pp.1757–1760 (2002).
- Every, M. and Szymanski, J.: A Spectralfiltering Approach to Music Signal Separation, *Proc. DAFx*, pp.197–200 (2004).
- 14) Kameoka, H., Nishimoto, T. and Sagayama, S.: Harmonic-temporal Structured Clustering via Deterministic Annealing EM Algorithm for Audio Feature Extraction, *Proc. ISMIR*, pp.115–122 (2005).
- 15) Viste, H. and Evangelista, G.: A Method for Separation of Overlapping Partials Based on Similarity of Temporal Envelopes in Multichannel Mixtures, *IEEE Trans. Speech and Audio Processing*, Vol.14, No.3, pp.1051–1061 (2006).
- 16) Woodruff, J., Pardo, B. and Dannenberg, R.: Remixing Stereo Music with Score-informed Source Separation, *Proc. ISMIR*, pp.314–319 (2006).
- 17) Smaragdis, P. and Brown, J.C.: Non-negative Matrix Factorization for Polyphonic Music

Transcription, *Proc. WASPAA*, pp.177–180 (2003).

- 18) Fitzgerald, D., Cranitch, M. and Coyle, E.: Sound Source Separation using Shifted Nonnegative Tensor Factorization, *Proc. ICASSP*, Vol.V, pp.653–656 (2006).
- 19) Uhle, C., Dittmar, C. and Sporer, T.: Extraction of Drum Tracks from Polyphonic Music Using Independent Subspace Analysis, *Proc. ICA*, pp.834–848 (2003).
- 20) Helen, M. and Virtanen, T.: Separation of Drums from Polyphonic Music Using Nonnegative Matrix Factorization and Support Vector Machine, *Proc. EUSIPCO* (2005).
- 21) Barry, D., Fitzgerald, D., Coyle, E. and Lawlor, B.: Drum Source Separation Using Percussive Feature Detection and Spectral Modulation, *Proc. ISSC*, pp.13–17 (2005).
- 22) Casey, M. and Westner, A.: Separation of Mixed Audio Sources by Independent Subspace Analysis, *Proc. ICMC*, pp.154–161 (2000).
- 23) Dubnov, S.: Extracting Sound Objects by Independent Subspace Analysis, Proc. AES 22nd International Conference on Virtual, Synthetic and Entertainment Audio (2002).
- 24) FItzGerald, D., Coyle, E. and Lawlor, B.: Independent Subspace Analysis Using Locally Linear Embedding, *Proc. DAFx*, pp.13–17 (2003).
- 25) Brown, J.C. and Smaragdis, P.: Independent Component Analysis for Automatic Note Extraction from Musical Trills, J. Acoust. Soc. Am., Vol.115, No.5, pp.2295–2306 (2004).
- 26) Barry, D., FitzGerald, D. and Lawlor, B.: Single Channel Source Separation Using Shorttime Independent Component Analysis, *Proc. AES* (2005).
- 27) Yoshii, K., Goto, M. and Okuno, H.G.: Drum Sound Recognition for Polyphonic Audio Signals by Adaptation and Matching of Spectrogram Templates with Harmonic Structure Suppression, *IEEE Trans. Audio, Speech, and Language Processing*, Vol.15, No.1, pp.333–345 (2007).
- 28) Goto, M.: A Real-time Music-scene-description System: Predominant-F0 Estimation for Detecting Melody and Bass Lines in Real-world Audio Signals, Speech Communication (ISCA Journal), Vol.43, No.4, pp.311–329 (2004).
- 29) Goto, M., Hashiguchi, H., Nishimura, T. and Oka, R.: RWC Music Database: Popular, Classical and Jazz Music Databases, *Proc. ISMIR*, pp.287–288 (2002).

# 付 録

A.1 パラメータ更新式の導出

式 (35) のコスト関数を各パラメータで偏微分した ものの零点を求めることで, J が極小になるようにパ ラメータを更新する式を導出する.

A.1.1 r<sub>kl</sub>(c): 各チャネルの相対強度

$$\frac{\partial J}{\partial r_{kl}(c)} = \sum_{m,n} \iint \left( -\frac{X_{klmn}^{(H)}}{r_{kl}(c)} \right) dt \, df$$
$$+ \iint \left( -\frac{X_{kl}^{(I)}}{r_{kl}(c)} \right) dt \, df + w_{kl} - \lambda_{kl}^{(r)} = 0 \quad (40)$$
$$\frac{\partial J}{\partial \lambda_{kl}^{(r)}} = \sum_{c} r_{kl}(c) - C = 0 \quad (41)$$

この連立方程式を解き,以下を得る.

$$r_{kl}(c) = \frac{C \iint \left(\sum_{m,n} X_{klmn}^{(H)} + X_{kl}^{(I)}\right) dt \, df}{\sum_{c} \iint \left(\sum_{m,n} X_{klmn}^{(H)} + X_{kl}^{(I)}\right) dt \, df}$$
(42)

#### A.1.2 w<sub>kl</sub>: 統合モデル全体の重み

$$\frac{\partial J}{\partial w_{kl}} = \sum_{c,m,n} \iint \left( -\frac{X_{klmn}^{(H)}}{w_{kl}} \right) dt \, df + \sum_{c} \iint \left( -\frac{X_{kl}^{(I)}}{w_{kl}} \right) dt \, df + \mathcal{C} - \lambda^{(w)} = 0 \tag{43}$$
$$\frac{\partial J}{\partial \lambda^{(w)}} = \sum_{k,l} w_{kl} - X_0 = 0 \tag{44}$$

この連立方程式を解き,以下を得る.

$$w_{kl} = \iint \left( \sum_{m,n} X_{klmn}^{(H)} + X_{kl}^{(I)} \right) dt df \quad (45)$$

A.1.3  $w_{kl}^{(H)}$ ,  $w_{kl}^{(I)}$ : 調波・非調波構造モデルの 重み

$$\frac{\partial J}{\partial w_{kl}^{(H)}} = \sum_{c,m,n} \iint \left( -\frac{X_{klmn}^{(H)}}{w_{kl}^{(H)}} \right) + Cw_{kl} - \lambda_{kl}^{(wHI)} = 0$$
(46)

$$\frac{\partial J}{\partial w_{kl}^{(I)}} = \sum_{c} \iint \left( -\frac{X_{kl}^{(I)}}{w_{kl}^{(I)}} \right) dt \, df + \mathbf{C} w_{kl} - \lambda_{kl}^{(wHI)} = 0$$
(47)

$$\frac{\partial J}{\partial \lambda_{kl}^{(wHI)}} = w_{kl}^{(H)} + w_{kl}^{(I)} - 1 = 0$$
(48)

この連立方程式を解き,以下を得る.

$$w_{kl}^{(H)} = \frac{\sum_{c,m,n} \iint X_{klmn}^{(H)} \, dt \, df}{\sum_{c} \iint \left( \sum_{m,n} X_{klmn}^{(H)} + X_{kl}^{(I)} \right) \, dt \, df}$$
(49)

$$w_{kl}^{(I)} = \frac{\sum_{c} \iint X_{kl}^{(I)} dt df}{\sum_{c} \iint \left( \sum_{m,n} X_{klmn}^{(H)} + X_{kl}^{(I)} \right) dt df}$$
(50)

A.1.4  $\omega_{kl}(t)$ : F0 の軌跡

$$\frac{\partial J}{\partial \omega_{kl}(t)} = \sum_{c,m,n} \int -\frac{n\left(f - n\omega_{kl}(t)\right) X_{klmn}^{(H)}}{\sigma_{kl}^2} df + \beta_{\omega} \left(-\frac{\tilde{\omega}_{kl}(t)}{\omega_{kl}(t)} + 1\right) = 0$$
(51)

$$\Rightarrow \quad a_{\omega}\omega_{kl}(t)^2 + b_{\omega}\omega_{kl}(t) + c_{\omega} = 0 \tag{52}$$

$$\begin{cases} a_{\omega} = \sum_{c,m,n} \int n^2 X_{klmn}^{(H)} df \\ b_{\omega} = \sigma_{kl}^2 \beta_{\omega} - \sum_{c,m,n} \int nf X_{klmn}^{(H)} df \\ c_{\omega} = -\sigma_{kl}^2 \beta_{\omega} \tilde{\omega}_{kl}(t) \end{cases}$$
(53)

この方程式を解き,以下を得る.

$$\omega_{kl}(t) = \frac{-b_{\omega} + \sqrt{b_{\omega}^2 - 4a_{\omega}c_{\omega}}}{2a_{\omega}}$$
  
A.1.5  $u_{klm}$ :パワーエンベロープの概形

$$\frac{\partial J}{\partial u_{klm}} = \sum_{c,n} \iint -\frac{X_{klmn}^{(H)}}{u_{klm}} dt df - \lambda_{kl}^{(u)}$$
$$= 0 \tag{54}$$

$$\frac{\partial J}{\partial \lambda_{kl}^{(u)}} = \sum_{m} u_{klm} - 1 = 0 \tag{55}$$

# この連立方程式を解き,以下を得る.

$$u_{klm} = \frac{\sum_{c,n} \iint X_{klmn}^{(H)} dt \, df}{\sum_{c,m,n} \iint X_{klmn}^{(H)} dt \, df}$$
(56)

# A.1.6 *v<sub>kln</sub>*: *n* 次倍音成分の相対強度

$$\frac{\partial J}{\partial v_{kln}} = \sum_{c,m} \iint -\frac{X_{klmn}^{(H)}}{v_{kln}} dt df + \beta_v \left(-\frac{\bar{v}_{kn}}{v_{kln}} + 1\right) - \lambda_{kl}^{(v)} = 0 \quad (57)$$

$$\frac{\partial J}{\partial \lambda_{kl}^{(v)}} = \sum_{n} v_{kln} - 1 = 0 \tag{58}$$

# この連立方程式を解き,以下を得る.

$$v_{kln} = \frac{\beta_v \bar{v}_{kn} + \sum_{c,m} \iint X_{klmn}^{(H)} dt \, df}{\beta_v + \sum_{c,m,n} \iint X_{klmn}^{(H)} dt \, df} \qquad (59)$$

A.1.7 *τ<sub>kl</sub>*: 発音時刻

$$\frac{\partial J}{\partial \tau_{kl}} = \sum_{c,m,n} \iint -X^{(H)}_{klmn} \frac{t - \tau_{kl} - m\phi_{kl}}{\phi_{kl}^2} dt df$$
$$= 0 \tag{60}$$

# この方程式を解き,以下を得る.

$$\tau_{kl} = \frac{\sum_{c,m,n} \iint (t - m\phi_{kl}) X_{klmn}^{(H)} dt df}{\sum_{c,m,n} \iint X_{klmn}^{(H)} dt df}$$
(61)

A.1.8 
$$\mathbf{Y}\phi_{kl}$$
: 音長  

$$\frac{\partial J}{\partial \phi_{kl}} = \sum_{c,m,n} \iint X_{klmn}^{(H)}$$

$$\cdot \frac{(t - \tau_{kl}) (t - \tau_{kl} - m\phi_{kl}) - \phi_{kl}^2}{\phi_{kl}^3} dt df$$

$$= 0$$
(62)

$$\Rightarrow \quad a_{\phi}\phi_{kl}^2 + b_{\phi}\phi_{kl} + c_{\phi} = 0 \tag{63}$$

$$\begin{cases} a_{\phi} = \sum_{c,m,n} \iint X_{klmn}^{(H)} dt df \\ b_{\phi} = \sum_{c,m,n} \iint m (t - \tau_{kl}) X_{klmn}^{(H)} dt df \\ c_{\phi} = -\sum_{c,m,n} \iint (t - \tau_{kl})^2 X_{klmn}^{(H)} dt df \end{cases}$$
(64)

この方程式を解き,以下を得る.
$$\phi_{kl} = \frac{-b_{\phi} + \sqrt{b_{\phi}^2 - 4a_{\phi}c_{\phi}}}{2a_{\phi}}$$
(65)

A.1.9 *σ<sub>kl</sub>*:周波数方向の広がり

$$\frac{\partial J}{\partial \sigma_{kl}} = \sum_{c,m,n} \iint -X_{klmn}^{(H)}$$
$$-\frac{-n^2 \sigma_{kl}^2 + (f - n\omega_{kl}(t))^2}{n^2 \sigma_{kl}^3} dt df$$
$$= 0 \tag{66}$$

この方程式を解き,以下を得る.

$$\sigma_{kl} = \sqrt{\frac{\sum_{c,m,n} \iint (f - n\omega_{kl}(t))^2 X_{klmn}^{(H)} dt df}{\sum_{c,m,n} \iint X_{klmn}^{(H)} dt df}}$$
(67)
.1.10  $I'_{kl}(t, f)$ : 非調波構造モデル

A.1.10 
$$I'_{kl}(t, f)$$
: 非調波構造モデル

$$\frac{\partial J}{\partial I'_{kl}(t,f)} = \sum_{c} \left( -\frac{X^{(I)}_{kl}}{I'_{kl}(t,f)} + r_{kl}(c) \right)$$
$$+\beta_{\bar{I}} \left( -\frac{\bar{I}_{k}}{I'_{kl}(t,f)} + 1 \right)$$
$$+\beta_{\bar{I}} \left( -\frac{\tilde{I}_{kl}}{I'_{kl}(t,f)} + 1 \right) = 0 \quad (68)$$

$$\frac{\partial J}{\partial \lambda_{kl}^{(I)}} = \iint I_{kl}'(t,f) \, dt \, df - 1 = 0 \quad (69)$$

この方程式を解き,以下を得る.

$$I_{kl} = \frac{CX_{kl}^{(I)} + \beta_{\bar{I}}\bar{I}_k + \beta_{\bar{I}}\tilde{I}_{kl}}{C\iint X_{kl}^{(I)} dt \, df + \beta_{\bar{I}} + \beta_{\bar{I}}}$$
(70)

(平成 19 年 6 月 21 日受付)(平成 19 年 12 月 4 日採録)



糸山 克寿(学生会員)
 2006年京都大学工学部情報学科
 卒業,現在,京都大学大学院情報学
 研究科知能情報学専攻修士課程に在

籍中.2008年より日本学術振興会
 特別研究員(DC1).音楽情報処理,

音楽鑑賞インタフェース等の研究に従事.



後藤 真孝(正会員)

1998年早稲田大学大学院理工学 研究科博士後期課程修了.博士(工 学).同年電子技術総合研究所(2001 年に産業技術総合研究所に改組)に 入所し,現在,主任研究員.2000年

から 2003 年まで科学技術振興事業団さきがけ研究 21 「情報と知」領域研究員,2005 年から筑波大学大学院 准教授(連携大学院)を兼任.音楽情報処理,音声言 語情報処理等に興味を持つ.2000 年 WISS2000 論文 賞・発表賞,2001 年日本音響学会粟屋潔学術奨励賞・ ポスター賞,2003 年インタラクション 2003 ベスト ペーパー賞,2005 年情報処理学会論文賞,2007 年第 6 回ドコモ・モバイル・サイエンス賞基礎科学部門優 秀賞等 21 件受賞.電子情報通信学会,日本音響学会, 日本音楽知覚認知学会各会員.



駒谷 和範(正会員) 1998年京都大学工学部情報工学 科卒業.2000年同大学院情報学研 究科知能情報学専攻修士課程修了. 2002年同大学院博士後期課程修了. 京都大学博士(情報学).同年京都

大学情報学研究科助手.2007年より助教.音声対話 システムの研究に従事.情報処理学会平成16年度山 下記念研究賞,FIT2002 ヤングリサーチャー賞等受 賞.電子情報通信学会,言語処理学会,人工知能学会, ISCA,ACL 各会員.



尾形 哲也(正会員)
 1997年日本学術振興会特別研究
 員,1999年早稲田大学大学院博士
 後期課程単位修得退学,同年早稲田
 大学助手.2000年情緒交流ロボットの開発研究により博士(工学).2001

年理化学研究所脳科学総合研究センター研究員,2003 年京都大学大学院情報学研究科講師,2005年同大学 助教授を経て,2007年より同大学准教授.この間, 2001~2003年早稲田大学ヒューマノイド研究所客員 講師,2005~2007年同研究所客員助教授,2007年よ り同研究所准教授.2005年より理化学研究所脳科学総 合研究センター客員研究員.研究分野は神経回路モデ ルおよび人間とロボットのコミュニケーション発達を 考えるインタラクション創発システム情報学.2000年 日本機械学会論文賞,2005年国際会議 IEA-AIE 最優 秀論文賞等受賞.著書は,身体性とコンピュータ(共 立出版),メカノクリーチャ(コロナ社)等.



奥乃 博(正会員) 1972年東京大学教養学部基礎科学 科卒業.日本電信電話公社,NTT, JST,東京理科大学を経て,2001年 より京都大学大学院情報学研究科知 能情報学専攻教授.博士(工学).こ

の間,スタンフォード大学客員研究員,東京大学工学部 客員助教授.人工知能,音環境理解,ロボット聴覚,音 楽情報処理の研究に従事.1990年度人工知能学会論文 賞,IEA/AIE-2001,2005最優秀論文賞,IEEE/RSJ IROS-2001,IROS-2005 Best Paper Nomination Finalist,第2回船井情報科学振興賞等受賞.JSAI,RSJ, ACM,IEEE,AAAI等各会員.本学会英文図書出版 委員.