

深層多重音検出を用いた音響信号から楽譜へのピアノ採譜

柴田 健太郎^{1,a)} 中村 栄太^{1,2,b)} 錦見 亮^{1,c)} 吉井 和佳^{1,d)}

概要: 本稿では多重音響信号から楽譜を推定するピアノ自動採譜システムについて述べる。最近の研究で、多重音検出とリズム量子化手法の統合による自動採譜システムの可能性が示されたが、多重音検出の誤りによって採譜精度が頭打ちになってしまう問題があった。一方、深層学習技術が飛躍的に進歩しているため、本研究では現在最高精度の深層多重音検出手法を統合した新たな採譜システムを構築し、その効果を検証する。具体的には、多重音を含む音響信号から音符の音高、発音時刻、消音時刻、ベロシティーを推定する新たな畳み込み型のニューラルネットワークを提案する。また、和音を含む声部を扱うことができる声部分離手法を新たに提案する。評価実験により提案システムの採譜精度が従来手法を大幅に上回ることを示すと同時に、提案する声部推定によって楽譜の可読性が向上することを示す。

1. はじめに

自動採譜は音楽情報処理の分野において長年の未解決問題である [1–3]。自動採譜の最終目標は音楽音響信号を完全な楽譜へ変換することであり、得られた楽譜は音楽演奏や記号領域での音楽内容解析に役立つ。多声音楽の自動採譜はその問題の難しさから、多重音検出とリズム量子化の2つの問題に分けて主に研究されてきた。多重音検出は、音響信号を半音単位の音高と発音・消音時刻で表される音符列、およびベロシティー（音強）の系列である演奏 MIDI 系列（ピアノロール）へと変換する [4–9]。リズム量子化は、演奏 MIDI 系列を発音・消音時刻が拍単位で記述される量子化 MIDI 系列へと変換する [10–13]。

自動採譜に関するの研究が多く行われるなか、音響信号から完全な楽譜まで推定する研究の数は限られる [14–17]。文献 [14] は多重音検出とリズム量子化を多段階処理の枠組みで統合したピアノ採譜手法を提案している。この手法は、当時最も性能が良い手法の1つであった確率的潜在成分分析（probabilistic latent component analysis; PLCA）を多重音検出に用いている。PLCA の出力は非常に多くの誤りを含むため、実用には程遠い採譜結果がしばしば得られた。一方、深層ニューラルネットワーク（Deep Neural Network; DNN） [7–9] を用いた多重音検出法が目覚ましい成果を挙げており、音響から楽譜への採譜における大幅な

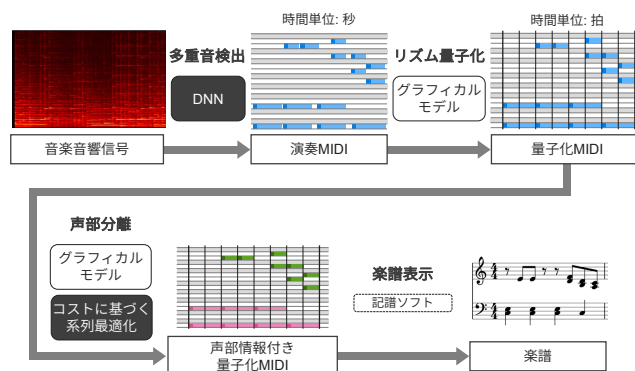


図1 深層多重音検出を用いた音響信号から楽譜へのピアノ採譜手法

性能向上が期待できる。

また、ピアノ採譜において声部分離は、読みやすい楽譜の生成に必要な、重要な問題の一つである [18–24]。ピアノ楽譜の声部構造は階層的である。まず、通常右手と左手に対応する二つのパートがある。さらにそれぞれのパートの中に複数の声部があり、それら声部は楽譜上では異なるレイヤーで表記される。既存手法の大半は各声部が単旋律だと仮定しているため、各声部が和音を含むことが多いピアノ楽譜の採譜には適していない。文献 [14] では記譜ソフトに実装されている声部分離手法を用いており、不要なタイや誤った場所に配置された休符によって出力楽譜の読みやすさが損なわれている。

本研究の目的は、最高精度のピアノ音響信号に対する自動採譜法のベンチマークを確立するとともに、現状での限界と残された主要な問題を示すことである。最新の多重音検出・リズム量子化・声部分離手法を統合することでピアノ採譜システムを構築する。現在フレーム単位の多重音推

¹ 京都大学 大学院情報学研究所

² 京都大学 白眉センター

a) shibata@sap.ist.i.kyoto-u.ac.jp

b) enakamura@sap.ist.i.kyoto-u.ac.jp

c) nishikimi@sap.ist.i.kyoto-u.ac.jp

d) yoshii@kuis.kyoto-u.ac.jp

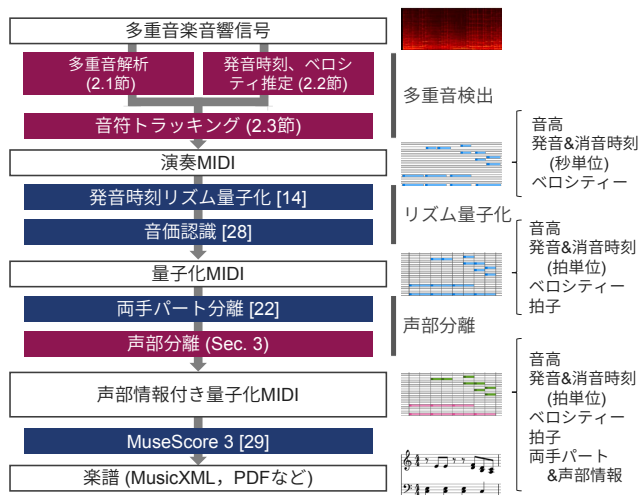


図2 提案システムの構成

定で最高性能を達成している手法 [9,25] に基づき、新たな多重音検出法を提案する。また、和音を含む声部を扱えるコストに基づく声部分離法を提案する。MAPS データセット [26] を用いて、文献 [14,27] で提案された評価尺度に基づく系統的な評価実験を行う。また、提案する声部推定手法の有効性を確かめるための主観評価実験を行う。

提案する採譜システムの構成は文献 [14] で提案されたものの拡張であり図2に示す通りである。我々は新たな多重音検出法を提案するとともに (2章)、代表的な手法 [8] との比較も行う。リズム量子化ステップでは、拍節隠れマルコフモデル (hidden Markov model; HMM) を用いて拍単位の発音時刻を推定する (拍子が常に4/4拍子に固定されることを除き、文献 [14] の4.1章と同様)。その後、拍単位の消音時刻をマルコフ確率場モデル [28] を用いて推定する。声部分離ステップでは、まず初めに両手パート分離手法 [22] を適用する。この手法は、一般的なMIDI系列を入力として受け取り、右手と左手に対応する2つのMIDI系列を出力する。その後、3章で述べる声部分離の新技术を適用する。最終ステップでは、公開ソフト MuseScore 3 [29] を使用して、楽譜のグラフィカルな表現を得る。

2. 多重音検出

提案する多重音検出法 (POVNet) は2つのDNNから構成される。1つは音高解析用のネットワーク (PitchNet)、もう1つは発音時刻とベロシティーの推定用のネットワーク (OnVelNet) である (図3)。これらのネットワークは、文献 [9] で初めて多重音解析に用いられた、DeepLabv3+ [30] と呼ばれる画像領域分類用の畳み込みニューラルネットワーク (convolutional neural network; CNN) に基づく。各ネットワークを別々に学習し、それらの出力を音符トラッキングステップで統合することで演奏MIDI系列が得られる。両ネットワークはHCFP (harmonic combined frequency and periodicity) 特徴量 $\mathbf{Z} \in \mathbb{R}_+^{2H \times F \times T}$ [9] を入力とする。

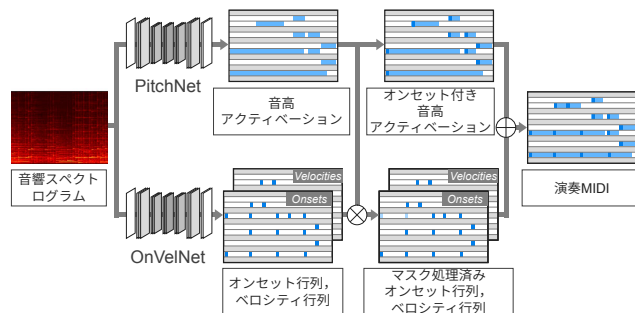


図3 多重音検出ネットワーク (POVNet) の構成図。

ここで、 H, F, T はそれぞれ調和音の数、周波数ビンの数、時間フレームの数である。本稿では文献 [9] と同様に $H = 6, F = 352 = 88 \times 4$ とした。

2.1 多重音解析ネットワーク (PitchNet)

HCFP 特徴量 \mathbf{Z} を与えると、PitchNet は確率行列 $\mathbf{P}_p \in [0, 1]^{F \times T}$ を出力する。ここで、 $\mathbf{P}_p(f, t)$ はフレーム t における周波数 f のサリエンス (顕著さ) を表す。入力が $2H = 12$ チャンネル、出力が1チャンネルであることを除いて、PitchNet は文献 [9] で提案されたものと同じネットワークである。最終層でシグモイド関数を用いて \mathbf{P}_p が得られる。PitchNet は以下の損失関数を最小化するように学習される。

$$\mathcal{L}_p = -\frac{1}{F \times T} \sum_{f=1}^F \sum_{t=1}^T \hat{\mathbf{P}}_p(f, t) \log \mathbf{P}_p(f, t) \quad (1)$$

ここで、 $\hat{\mathbf{P}}_p \in \{0, 1\}^{F \times T}$ は正解MIDIデータ (サスティンペダルイベントを含む) をアップサンプルして得られたバイナリ行列である。最後に、 \mathbf{P}_p を2値化し、ダウンサンプルすることで音高アクティベーション行列 $\mathbf{D}_p \in \{0, 1\}^{M \times T}$ が得られる。ここで、 $\mathbf{D}_p(m, t)$ はフレーム t における半音単位の音高 m の音の有無を表し、 $M = 88$ はピアノの鍵盤の個数を表す。

2.2 発音時刻・ベロシティー推定ネットワーク (OnVelNet)

HCFP 特徴量 \mathbf{Z} を与えると、OnVelNet はオンセット確率行列 $\mathbf{P}_o \in [0, 1]^{F \times T}$ と音強行列 $\mathbf{P}_v \in [0, 1]^{F \times T}$ を出力する。ここで、 $\mathbf{P}_o(f, t)$ と $\mathbf{P}_v(f, t)$ はそれぞれ、フレーム t 、周波数 f におけるオンセットサリエンスと音強を表す。また、音強はMIDI規格でベロシティーとして定められた0から127まで整数値をリスケールした0から1までの値をとる。出力がオンセット行列用と音強行列用の2チャンネルであることを除いて、OnVelNet は PitchNet と同じ構造である。最終層では、シグモイド関数を用いて \mathbf{P}_o が、 $[0, 1]$ 区間のクリップ関数を用いて \mathbf{P}_v がそれぞれ得られる。このネットワークは以下で定義されるバイナリクロスエントロピー損失と \mathcal{L}_o と平均二乗誤差損失 \mathcal{L}_v の重み付き和を最小化するように学習される。

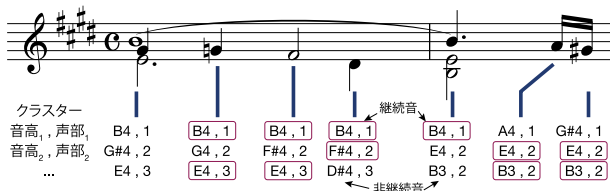


図4 声部配置の例

$$\mathcal{L}_0 = -\frac{1}{F \times T} \sum_{f=1}^F \sum_{t=1}^T \hat{\mathbf{P}}_0(f, t) \log \mathbf{P}_0(f, t), \quad (2)$$

$$\mathcal{L}_v = \frac{1}{F \times T} \sum_{f=1}^F \sum_{t=1}^T \hat{\mathbf{P}}_0(f, t) \left(\hat{\mathbf{P}}_v(f, t) - \mathbf{P}_v(f, t) \right)^2, \quad (3)$$

$$\mathcal{L}_{ov} = w_0 \mathcal{L}_0 + w_v \mathcal{L}_v, \quad (4)$$

ここで、 $\hat{\mathbf{P}}_0 \in \{0, 1\}^{F \times T}$ は正解音符のオンセットの有無を表すバイナリ行列、 $\hat{\mathbf{P}}_v \in [0, 1]^{F \times T}$ は正解音符のオンセットにおける音強を表す実数行列である。 $\hat{\mathbf{P}}_0$ と $\hat{\mathbf{P}}_v$ は 2.1 節における $\hat{\mathbf{P}}_p$ と同様の方法で正解 MIDI から得られる。発音時刻のわずかな揺らぎを許容するため、 $\hat{\mathbf{P}}_0(f, t) = 1$ ならば $\hat{\mathbf{P}}_0(f, t \pm 1) = 1$ とした。文献 [8] と同様に \mathcal{L}_v は正解のオンセットでマスクされた後で計算される。

最後に、 \mathbf{P}_0 を二値化しダウンサンプルすることでオンセット行列 $\mathbf{D}_0 \in \{0, 1\}^{M \times T}$ が得られる。 \mathbf{D}_0 において時間方向に連続した 1 の要素は中央の要素を除いて 0 に設定する。ベロシティー行列 $\mathbf{D}_v \in \{0, \dots, 127\}^{M \times T}$ は、 \mathbf{P}_v に対してスケールリング、整数化した後ダウンサンプルすることで得られる。

2.3 音符トラッキング

演奏 MIDI 系列 \mathbf{D}_n は音符の音高、発音時刻、消音時刻とベロシティーで構成され、上述の音高アクティベーション行列 \mathbf{D}_p 、オンセット行列 \mathbf{D}_o 、ベロシティー行列 \mathbf{D}_v から求められる。 \mathbf{D}_n のオンセットは $\mathbf{D}_p(f, t) = 1$ かつ $\mathbf{D}_p(f, t-1) = 0$ の条件を満たす TF ビンを列挙することで得られる。 \mathbf{D}_o は同音連打オンセットを求めるためにだけに用いられるが、音高行列 \mathbf{D}_n との不整合が生じないように、 $\mathbf{D}_p(f, t) = 0$ となる TF ビンでは $\mathbf{D}_o(f, t)$ は 0 となるようにマスクをかける。最終的に \mathbf{D}_n は以下の規則によって生成される。

- \mathbf{D}_p によって得られるオンセットとマスクされた \mathbf{D}_o から得られるオンセットが 50 ms 以内に入っている場合、先行するオンセットのみを採用する。
- 最終的な音符継続長が 30 ms 以内の音符は除去する。

3. ピアノ採譜のための声部分離

我々の提案する声部分離手法は、声部の適切な構造と、入力 MIDI 系列に対する声部の適合度を表現するコスト関

数に基づいた系列最適化手法である。入力 MIDI 系列は右手左手パートに分離済みとし、本手法では右手左手パートに対して独立に声部分離を行う。系列最適化では発音時刻が同じである音符の組を単位として扱う。発音時刻が同じである音符の組に対してクラスターを構成し、それよりも発音時間が前で音がオーバーラップする音符（継続音と呼ぶ）もそのクラスターに加える。クラスター内の音符に対する声部を整数のラベル $1, 2, \dots, V_{\max}$ で記述する。ここで、声部数の上限値 V_{\max} は任意に設定できる変数である。各クラスター C_k に対して、クラスター内の音符 $n \in C_k$ の声部ラベル $S_k = (s_n)$ の組を声部配置と呼ぶ。声部分離の探索空間は全てのクラスターに対する全ての可能な声部配置の集合である。

コスト関数は、以下で定義される垂直コストと水平コストの和で構成される。垂直コスト $V(S_k)$ は、あるクラスターに対する声部配置の適切さを記述し、以下の 4 つの因子の和で定義される。

- 各音符 $n \in C_k$ に s_n を与える（不要な声部にペナルティを与える）。
- 声部の順番と音高の順番が逆になっている音符のペアそれぞれに λ_2 を与える（声部の交差にペナルティを与える）。
- 声部ラベルが同じで消音時刻が異なる音符のペアそれぞれに λ_3 を与える。
- 同じ声部ラベルを持つ継続音と継続音でない音符のペアそれぞれに λ_4 を与える。

水平コスト $H(S_{k-1}, S_k)$ は隣接するクラスターの声部配置の適切なつながりを記述し、以下の 3 つの因子の和で定義される。

- 声部ラベルが一貫していない継続音に λ_5 を与える。
- 同じ声部ラベルを持ちながら時間の開きがある隣接する音符のペアそれぞれに λ_6 を与える。（休符にペナルティを与える）。
- 同じ声部ラベルを持ちながら時間的オーバーラップがある隣接する音符のペアそれぞれに λ_7 を与える。

系列最適化は Viterbi アルゴリズムによって解くことができる。声部内の和音の消音時刻は一致せねばならず、かつその声部の次の音符の発音時刻と同じかそれより前でないといけないという声部の制約を満たすために、声部推定後に、消音時刻は得られた声部にしがたって補正される。上記のコスト関数の重み付けはいくつかの試行により $(\lambda_2, \dots, \lambda_7) = (3, 1, 1, 5, 0.2, 1)$ としたが、系統的なパラメータ最適化の余地が残っている。

4. 評価

4.1 実験設定

多重音検出の精度を評価するために MAPS データベース [26] の “ENSTDkCl” と “ENSTDkAm” のラベルが付い

手法	フレーム単位			ノート単位		
	\mathcal{P}_f	\mathcal{R}_f	\mathcal{F}_f	\mathcal{P}_n	\mathcal{R}_n	\mathcal{F}_n
PLCA [14]	—	—	—	77.9	68.9	72.8
Onsets and Frames [33]	92.9	78.5	84.9	87.5	85.6	86.4
DeepLabv3+ [9]	87.5	86.3	86.7	—	—	—
PitchNet のみ	89.3	84.4	86.6	91.1	68.4	77.5
POVNet (提案手法)	89.3	85.7	87.3	89.7	84.1	86.7

表1 多重音検出精度 (%)

た60曲を用いた。また、提案システムの最終的な出力である楽譜の採譜精度を評価するためには“ENSTDkCl”のラベルが付いた30曲を用いた。従来手法[14]と同様に、楽譜の評価にはそれぞれの録音の最初の30sを評価の対象とした。PitchNetとOnVelNetの学習にはMAPSデータセットのうち上記の60曲を除く全ての曲を用いた。ネットワークのパラメータは最適手法RAdam[31]を用いて更新し、学習率の初期値は0.001とした。OnVelNetの損失関数の重みはそれぞれ $w_o = 0.9$, $w_v = 0.1$ とした。

2章の多重音検出の性能評価には文献[32]で導入されているフレーム単位及び音符単位の評価尺度を用いた。フレーム単位の尺度では10msの時間分解能で適合率 \mathcal{P}_f 、再現率 \mathcal{R}_f 、F値 \mathcal{F}_f を定義する。音符単位の尺度では音符の発音時刻に基づいて評価され、推定音符の発音時刻が正解音符の発音時刻の ± 50 ms以内に入っていれば正解とみなされる。この基準に基づき音符単位の適合率 \mathcal{P}_n 、再現率 \mathcal{R}_n 、F値 \mathcal{F}_n を定義する。採譜システム全体の性能評価には文献[14]にて導入された編集距離に基づく評価尺度と文献[27]にて導入された評価尺度MV2Hを用いた。前者では以下の誤り率を定義する：音高誤り率 E_p 、不足音符率 E_m 、余分音符率 E_e 、発音時刻誤り率 E_{on} 、消音時刻誤り率 E_{off} とこれら5つの尺度の平均の総合誤り率 E_{all} である。評価尺度MV2Hでは、以下の精度を定義する：多重音検出 \mathcal{F}_p 、声部分離 \mathcal{F}_{voi} 、拍節配位 \mathcal{F}_{met} 、音価推定 \mathcal{F}_{val} 、和声解析 \mathcal{F}_{harm} とこれらの平均 \mathcal{F}_{MV2H} 。ここで、文献[27]では \mathcal{F}_{harm} がコードラベル推定精度と調推定精度の重み付き和として定義されているが、提案システムではコードラベルは出力しないので、調推定精度を意味する。

4.2 結果

表1に多重音検出手法の精度を示す。提案手法(POVNet)との比較として、文献[14]のPLCAに基づく手法、文献[33]のDNNに基づく代表的な手法(MAESTROデータベースで学習したOnsets and Frames)、および文献[9]の手法(DeepLabv3+)による結果を示す。さらに、OnVelNetを用いないPitchNetのみによる結果も示すが、これはDeepLabv3+とほぼ等価である。表より、我々の提案するPOVNetがフレーム単位と音符単位両方のF値において比較手法を上回ることが見て取れる。POVNetとOnsets and Framesは \mathcal{F}_n において同等な精度となっているが、PLCAに基づく手法

に対しては大きな差が見られる。PitchNetのみとPOVNetの \mathcal{R}_n の間に大きな差があることから、OnVelNetによって同音連打を扱えるようになった効果が見て取れる。

表2に採譜システム全体の評価の結果を示す。深層音高検出によって採譜精度が大幅に改善していることが見て取れる。編集距離に基づく評価尺度上では、多重音検出機構にPONVnetを用いたシステムがOnsets and Frames[33]を用いたシステムより総じて精度がよい。POVNet+RQとPOVNet+RQ+VSの結果を見比べることで、提案した声部推定法がわずかながら編集距離基準の誤り率を低下させることが分かる。声部推定手法を用いないPOVNet+RQが \mathcal{F}_{voi} においてPLCA+RQ+VS上回っていることは注目に値する。これに関しては、評価尺度 \mathcal{F}_{voi} では単旋律の声部のみを仮定しているため、ピアノ楽曲のように和音を含む声部分離の評価には適していない可能性も考えられる。

PLCA+RQ+VSの平均誤り率がNakamura *et al.* [14]よりも悪い理由として、後者は余分な音符を省く機構を持つNoisy拍節HMMを用いていることが考えられる。この余分な音を取り除く機構を現在のシステムに導入することも考えられるが、深層音高検出を用いたモデルでは適合率 \mathcal{P}_n が高いため、効果は限定的だと思われる。表2の一番下の行はリズム量子化と声部分離を正解演奏MIDIに対して適用した結果を示している。この結果から、多重音検出の精度は、音高推定に関連する尺度のみならず、リズム採譜に関する尺度にも依然として大きな影響を与えることが見て取れる。

推定楽譜の一例を図5に示す。深層音高検出を用いた手法によって得られた楽譜は従来手法[14]による楽譜より大幅に改善していることは明らかである。POVNet+RQとPOVNet+RQ+VSの結果を比較することで、声部分離手法が余分なタイや休符を効果的に減らしていることが見て取れる。また、全ての推定結果において図5の青枠で示す通り、後打音は正しく認識されておらず、これは我々のリズム量子化手法の現在の限界を示している。

さらに、声部推定が楽譜の可読性に与える影響を検証するために、主観評価実験を行った。POVNet+RQ+VSとPOVNet+RQによる推定楽譜をどちらによるものか隠して2枚同時に被験者に見せ、「演奏するにあたっての楽譜の読みやすさ」が高い方、もしくは「どちらとも言えない」の選択肢を選んでもらった。実験に用いたプラットフォームはリンク先[34]から確認できるようになっている。実験には、MAPSデータセットのうち“ENSTDkCl”のラベルが付いた30曲を用いた。実験にはピアノ楽譜を読む21人が参加し、そのうち約60%の人が10年以上のピアノ演奏経験者であった。得られた630の回答のうち、POVNet+RQ+VSが496回(78.7%)、POVNet+RQが54回(8.57%)、「どちらとも言えない」が80回(12.7%)選ばれた。この結果は、提案した声部推定法により楽譜の可読性が向上することを

手法	E_p	E_m	E_e	E_{on}	E_{off}	E_{all}	\mathcal{F}_p	\mathcal{F}_{voi}	\mathcal{F}_{met}	\mathcal{F}_{val}	\mathcal{F}_{harm}	\mathcal{F}_{MV2H}
Nakamura ら [14]	2.92	28.1	11.4	19.4	41.8	20.7	68.1	45.8	31.1	80.5	57.2	56.5
PLCA [14]+RQ+VS	3.29	23.5	13.4	25.1	41.8	21.4	67.2	46.7	32.5	81.3	57.1	57.0
Onsets and Frames [33]+RQ+VS	0.96	8.23	7.36	12.8	28.1	11.5	82.8	46.7	41.6	86.6	74.6	66.4
POVNet+RQ	0.74	8.19	5.83	12.9	28.6	11.3	83.9	48.2	43.6	86.5	71.1	66.7
POVNet+RQ+VS (提案手法)	0.74	8.11	5.75	12.8	27.9	11.0	84.3	46.1	43.4	87.3	71.1	66.4
Ground truth MIDI+RQ+VS	0.45	2.60	2.46	5.36	21.3	6.44	90.4	49.2	47.4	86.7	77.8	70.3

表2 MAPS-ENSTDkCl データセットに対する採譜の誤り率 (%) と精度 (%). '+RQ' はリズム量子化法を適用したことを表し, '+VS' は我々が提案した声部分離法を適用したことを表す. '+VS' が無い場合, 声部分離には MuseScore 3 を用いたことを表す.

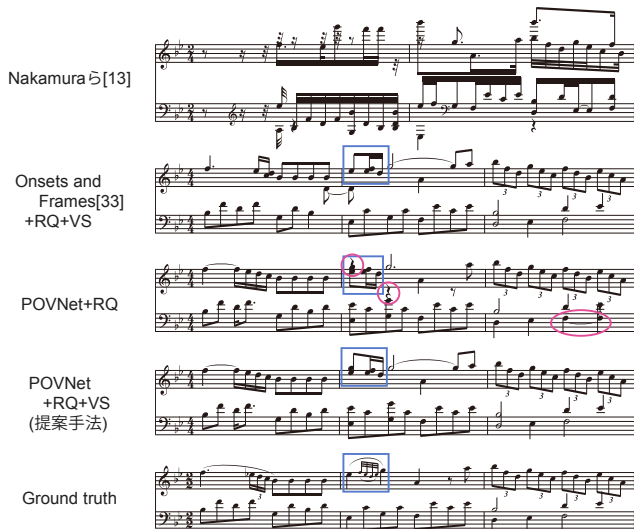


図5 採譜結果の例 (MAPS-ENSTDkCl データセット, Mozart: Piano Sonata K. 333).

示している.

4.3 考察

実験により, 我々の提案するピアノ採譜システムの楽譜推定精度は, 従来の最高性能のシステムを大幅に上回ることを示した. デモページ^{*1}に掲載された例から見て取れる様に, 提案システムによって推定される楽譜はピアノ演奏のための楽譜として実用的なレベルに達していることも多い. しかし, まだ我々の提案するシステムには限界も多く, 改善の余地が残されている. 表2から見て取れる様に, 多重音検出の改善によって最終的な楽譜の推定精度は更に向上する. 同様にリズム量子化, 特に消音時刻推定精度 E_{off} は正解 MIDI からの採譜結果でさえ 20% を越えており, 改善により推定楽譜の品質は大きく向上すると思われる. この問題への解決策として, 深層学習の技術を消音時刻推定に適用することと, 声部配置と消音時刻を同時推定することが見込みがあると考えられる. さらに, 消音時刻をより正しく推定するにはサステインペダルイベントの検出も重要である.

トリル, アルペジオ, 前打音, 後打音, グリッサンドと

^{*1} <http://sap.ist.i.kyoto-u.ac.jp/members/shibata/sigmus125/>

いった装飾音はこれまで自動採譜の研究において扱われることは稀であったが, リズム量子化の精度を向上させるためには明示的なモデル化が不可欠である. 完璧な楽譜の推定を目指す上で, ペダルイベント, 強弱記号, スラー, アーティキュレーションや運指番号の推定も重要である. 本研究ではクラシック音楽のみで構成される MAPS データセットのみを用いたが, ジャズやポピュラー音楽を用いてモデルの学習及び検証を行うことも今後の課題である.

5. 結論

本稿では, 最新の多重音検出, リズム量子化手法, および声部分離手法を統合することで, 現在最高精度のピアノ音響信号に対する自動採譜システムを提案した. 評価実験により, 提案システムが既存システムの採譜精度を大きく上回ることを確認すると同時に, 現在の自動採譜の限界と今後考えられるの研究の方向性を論じた. 提案した手法により, 実用レベルのピアノ自動採譜の実現が期待される. 今後の研究でリズム量子化と声部推定に深層学習の技術を適用することで, さらに精度向上を目指す予定である. また, 今回提案したシステムをウェブサービスとして公開し, YouTube 等のピアノ演奏動画や任意のピアノ音響信号から楽譜を推定して提供する枠組みを構築する予定である.

6. 謝辞

本研究の一部は, JSPS 科研費 No. 16H01744, No. 19H04137, 19K20340 および JST ACCEL No. JPMJAC1602 の支援を受けた.

参考文献

- [1] Klapuri, A. and Davy, M.: *Signal Processing Methods for Music Transcription*, Springer Science & Business Media (2006).
- [2] Benetos, E., Dixon, S., Giannoulis, D., Kirchhoff, H. and Klapuri, A.: Automatic Music Transcription: Challenges and Future Directions, *Journal of Intelligent Information Systems*, Vol. 41, No. 3, pp. 407–434 (2013).
- [3] Benetos, E., Dixon, S., Duan, Z. and Ewert, S.: Automatic Music Transcription: An Overview, *IEEE Signal Processing Magazine*, Vol. 36, No. 1, pp. 20–30 (2018).
- [4] Benetos, E. and Weyde, T.: An efficient temporally-

- constrained probabilistic model for multiple-instrument music transcription, *Proc. ISMIR*, pp. 701–707 (2015).
- [5] Cheng, T., Mauch, M., Benetos, E., Dixon, S. et al.: An attack/decay model for piano transcription, *Proc. ISMIR*, pp. 584–590 (2016).
- [6] Sigtia, S., Benetos, E. and Dixon, S.: An end-to-end neural network for polyphonic piano music transcription, *IEEE/ACM Transactions on Audio, Speech and Language Processing*, Vol. 24, No. 5, pp. 927–939 (2016).
- [7] Bittner, R. M., McFee, B., Salamon, J., Li, P. and Bello, J. P.: Deep salience representations for F0 estimation in polyphonic music, *Proc. ISMIR*, pp. 63–70 (2017).
- [8] Hawthorne, C., Elsen, E., Song, J., Roberts, A., Simon, I., Raffel, C., Engel, J., Oore, S. and Eck, D.: Onsets and Frames: Dual-objective piano transcription, *Proc. ISMIR*, pp. 50–57 (2018).
- [9] Wu, Y.-T., Chen, B. and Su, L.: Polyphonic Music Transcription with Semantic Segmentation, *Proc. ICASSP*, IEEE, pp. 166–170 (2019).
- [10] Cemgil, A. T., Desain, P. and Kappen, B.: Rhythm Quantization for Transcription, *Comp. Mus. J.*, Vol. 24, No. 2, pp. 60–76 (2000).
- [11] Raphael, C.: A Hybrid Graphical Model for Rhythmic Parsing, *Artificial Intelligence*, Vol. 137, pp. 217–238 (2002).
- [12] Hamanaka, M., Goto, M., Asoh, H. and Otsu, N.: A Learning-Based Quantization: Unsupervised Estimation of the Model Parameters, *Proc. ICMC*, pp. 369–372 (2003).
- [13] Nakamura, E., Yoshii, K. and Sagayama, S.: Rhythm transcription of polyphonic piano music based on merged-output HMM for multiple voices, *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, Vol. 25, No. 4, pp. 794–806 (2017).
- [14] Nakamura, E., Benetos, E., Yoshii, K. and Dixon, S.: Towards complete polyphonic music transcription: Integrating multi-pitch detection and rhythm quantization, *Proc. ICASSP*, IEEE, pp. 101–105 (2018).
- [15] Carvalho, R. G. C. and Smaragdis, P.: Towards End-to-End Polyphonic Music Transcription: Transforming Music Audio Directly to a Score, *Proc. WASPAA*, pp. 151–155 (2017).
- [16] Román, M. A., Pertusa, A. and Calvo-Zaragoza, J.: An End-to-end Framework for Audio-to-Score Music Transcription on Monophonic Excerpts., *Proc. ISMIR*, pp. 34–41 (2018).
- [17] Nishikimi, R., Nakamura, E., Fukayama, S., Goto, M. and Yoshii, K.: Automatic Singing Transcription Based on Encoder-Decoder Recurrent Neural Networks with a Weakly-Supervised Attention Mechanism, *Proc. ICASSP*, pp. 161–165 (2019).
- [18] Chew, E. and Wu, X.: Separating voices in polyphonic music: A contig mapping approach, *Proc. CMMR*, pp. 1–20 (2004).
- [19] Cambouropoulos, E.: Voice and stream: Perceptual and computational modeling of voice separation, *Music Perception: An Interdisciplinary Journal*, Vol. 26, No. 1, pp. 75–94 (2008).
- [20] Temperley, D.: A unified probabilistic model for polyphonic music analysis, *Journal of New Music Research*, Vol. 38, No. 1, pp. 3–18 (2009).
- [21] Duane, B. and Pardo, B.: Streaming from MIDI Using Constraint Satisfaction Optimization and Sequence Alignment, *Proc. ICMC*, pp. 1–8 (2009).
- [22] Nakamura, E., Ono, N. and Sagayama, S.: Merged-Output HMM for Piano Fingering of Both Hands, *Proc. ISMIR*, pp. 531–536 (2014).
- [23] McLeod, A. and Steedman, M.: HMM-based voice separation of MIDI performance, *Journal of New Music Research*, Vol. 45, No. 1, pp. 17–26 (2016).
- [24] de Valk, R. and Weyde, T.: Deep neural networks with voice entry estimation heuristics for voice separation in symbolic music representations, *Proc. ISMIR*, pp. 281–288 (2018).
- [25] 中村栄太, Benetos, E., 吉井和佳, Dixon, S.: 多重音検出とリズム量子化の統合による多声音楽の自動採譜, *Proc. SIGMUS*, Vol. 2018-MUS-120, No. 19, pp. 1–6 (2018).
- [26] Emiya, V., Badeau, R. and David, B.: Multipitch estimation of piano sounds using a new probabilistic spectral smoothness principle, *IEEE Transactions on Audio, Speech, and Language Processing*, Vol. 18, No. 6, pp. 1643–1654 (2009).
- [27] McLeod, A. and Steedman, M.: Evaluating Automatic Polyphonic Music Transcription., *Proc. ISMIR*, pp. 42–49 (2018).
- [28] Nakamura, E., Yoshii, K. and Dixon, S.: Note Value Recognition for Piano Transcription Using Markov Random Fields, *IEEE/ACM Transactions on Audio, Speech and Language Processing*, Vol. 25, No. 9, pp. 1542–1554 (2017).
- [29] : MuseScore3, <https://musescore.org/en>.
- [30] Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F. and Adam, H.: Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation, *Proc. ECCV*, pp. 801–818 (2018).
- [31] Liu, L., Jiang, H., He, P., Chen, W., Liu, X., Gao, J. and Han, J.: On the Variance of the Adaptive Learning Rate and Beyond, *arXiv preprint arXiv:1908.03265* (2019).
- [32] Bay, M., Ehmann, A. F. and Downie, J. S.: Evaluation of Multiple-F0 Estimation and Tracking Systems, *Proc. ISMIR*, pp. 315–320 (2009).
- [33] Hawthorne, C., Stasyuk, A., Roberts, A., Simon, I., Huang, C.-Z. A., Dieleman, S., Elsen, E., Engel, J. and Eck, D.: Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset, *Proc. ICLR*, (online), available from (<https://openreview.net/forum?id=r11YRjC9F7>) (2019).
- [34] : 主観実験プラットフォーム, <https://forms.gle/iTnKncKGBN3CFkv9>.