

双方向RNNによるMFCC及びラウドネスからの 振幅スペクトログラム予測

川口 翔也¹ 北村 大地¹

概要: 変分自己符号化器 (variational autoencoder: VAE) は入力データの潜在的な特徴量を教師無しで学習できる DNN であり, 潜在特徴量に複数クラスの相対関係を表す構造を導入することで, 一定の解釈性を持たせることができる. 例えば, 複数楽器音の音色特徴量を表すメル周波数ケプストラム係数 (mel-frequency cepstrum coefficient: MFCC) の潜在特徴量を VAE で学習することで, 各楽器音の特徴量を併せ持つような新しい音響信号の MFCC を新たに生成できる. 我々は現在, VAE を用いて楽器音の音色を変換するシステムの構築を目指しており, その一例として, 入力信号の MFCC を VAE で生成した MFCC に置き換える音色変換を検討している. このようなシステムでは, 基本周波数, MFCC, 及び音量変化の 3 つの音響特徴量を入力とすることを想定しているが, MFCC を置き換えた後に音響信号に戻すためには, 前述の 3 つの音響特徴量からスペクトログラムを生成する必要がある, これは解析的な処理ではない. そこで本稿では, 基本周波数, MFCC, 及び音量から振幅スペクトログラムを予測する手法について検討する. 特に, 前述の音響特徴量を入力とする双方向再帰型ニューラルネットワークを用いた振幅スペクトログラムの予測について実験的に調査する. ピアノ及びギターを用いた実験では, 両楽器において比較的高精度に振幅スペクトログラムが予測可能であることを示す.

Amplitude Spectrogram Prediction from MFCC and Loudness Using Bidirectional RNN

KAWAGUCHI SHOYA¹ KITAMURA DAICHI¹

1. はじめに

生成モデル系の深層ニューラルネットワーク (deep neural network: DNN) に基づく音色変換は様々なものが研究されている. 例えば, 潜在的な特徴量を教師無しで学習できる生成モデル系 DNN である変分自己符号化器 (variational auto-encoder: VAE) を用いた楽器音解析や生成が提案されている [1, 2]. Fig. 1 に, 数字の手書き画像に VAE を適用したものを示す. Fig. 1 のように, 潜在空間と呼ばれる空間にそれぞれの数字の集合を見ることができる. さらに空間上では, 「7」と「9」の集合の間に相当する潜在変数値を入力することで「7」と「9」の中間の手書き数字が出力される. 本稿では, このような各集合の相対関係の学習

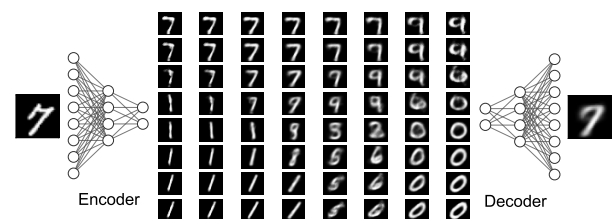


Fig. 1 The latent space of VAE trained with images of handwritten numbers.

を, 楽器音の音色集合に対して適用することで, VAE を用いた新しい音色変換アルゴリズムの構築を目指す. 以後, このシステムを「提案音生成システム」と呼ぶ. 提案音生成システムを用いることで, ギターとピアノの中間の音色等の新しい音色を持つ音響信号を生成できる. さらに新しい芸術及び音楽の発展に寄与できると考える.

提案音生成システムを実現には, 音高, 音色, 及び音量の

¹ 香川高等専門学校
National Institute of Technology, Kagawa College

3つの特徴量から振幅スペクトログラムの予測が必要となるが、実現可能かは不明である。そこで本稿では、提案音生成システムにおいて必要な部分システムとして、前述の3つの特徴量からDNNデコーダを用いて振幅スペクトログラムを予測する手法を提案する。DNNデコーダとして、多層パーセプトロン (multi-layer perceptron: MLP) 及び再帰型ニューラルネットワーク (recurrent neural network: RNN) の1つであるゲート付き再帰型ユニット (gated recurrent unit: GRU) [3] を用いた双方向再帰型ニューラルネットワーク (bidirectional RNN using GRU: BiGRU) の2種類のネットワーク構造を比較し、どのようなネットワーク構造が高精度な振幅スペクトログラムの予測に効果的かについて、実験的に調査する。

2. 要素技術と類似研究

2.1 メル周波数ケプストラム係数

メル周波数ケプストラム係数 (mel-frequency cepstrum coefficient: MFCC) [4] とは、時間周波数領域で表現された音声及び楽器音等の音色の特徴量である。音色のみの特徴量を可能な限り抽出しているため、音高や音量はあまり反映されない特徴量である。

MFCCを説明する前に、まずメル周波数について説明する。メル周波数とは、1000 Hzの純音に対応する音高尺度を1000 mel (メル周波数と呼ぶ) と定義し、これを基準として人間が知覚する音高を線形な一次元軸に対応付けた尺度である。周波数 f Hz とメル周波数 m mel の対応関係は次式で定義される [5]。

$$m = 1127 \ln \left(1 + \frac{f}{700} \right) \quad (1)$$

$$= 2595 \log_{10} \left(1 + \frac{f}{700} \right) \quad (2)$$

このメル周波数軸上で等間隔に三角状のバンドパスフィルタを複数個 (K 個) 構成したものをメルフィルタバンクと呼ぶ。

まず、音響信号に短時間 Fourier 変換 (short-time Fourier transform: STFT) を適用し複素スペクトログラム $\mathbf{Z} \in \mathbb{C}^{I \times J}$ を得る。次に、各短時間区間のパワースペクトル (パワースペクトログラム $|\mathbf{Z}|^2$ の各列ベクトル) に、 K 個のメルフィルタをそれぞれ畳み込むことでメルスペクトルと呼ばれる K 次元の音色の特徴量が得られる。ここで、 I 及び J はそれぞれ周波数ビン数及び時間フレーム数を表し、行列に対する演算子 $|\cdot|^q$ は要素毎の絶対値と q 乗を表す。パワースペクトログラム $|\mathbf{Z}|^2$ の全列ベクトルをメルスペクトルに変換したものをメルスペクトログラム $\mathbf{P} \in \mathbb{R}^{K \times J}$ と呼び、この変換を次式で表す。

$$\mathbf{P} = \text{MelFiltering}(|\mathbf{Z}|^2) \quad (3)$$

ここで、周波数は I から K に次元圧縮されている。

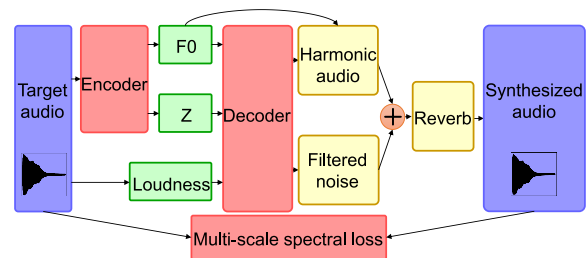


Fig. 2 DDSP autoencoder architecture. Red components are part of the neural network architecture, green components are the latent representation, and yellow components are deterministic synthesizers and effects [6].

MFCCは、メルスペクトログラム \mathbf{P} の各列ベクトル (各メルスペクトル) に離散コサイン変換 (discrete cosine transform: DCT) を適用して得られる実数係数である。MFCCを \mathbf{C} と表記し、この処理を次式で表す。

$$\mathbf{C} = \text{DCT}(\mathbf{P}) \in \mathbb{R}^{K \times J} \quad (4)$$

以上の変換により得られるMFCCは、音高や音量の影響を可能な限り排除した音色に関する K 次元の特徴量ベクトルとなる。これはパワースペクトルの包絡を良く表現した特徴量であり、音色に関する情報を強く保持している。

2.2 類似研究1: Differentiable digital signal processing

DNNを用いた楽器音信号の生成手法として、Differentiable digital signal processing (DDSP) [6]がある。Fig. 2にDDSPの概要図を示す。DDSPでは、入力された音響信号から、音高、音色、及び音量 (Fig. 2のF0, Z, 及び Loudness) の3つの特徴量をエンコーダで抽出する。この3つの特徴量はデコーダに入力され、事前に抽出済みのF0及びその整数倍の周波数からなる複数正弦波 (Fig. 2の Harmonic audio) を駆動する変数と、白色雑音に音色を与えるFIRフィルタの伝達関数 (Fig. 2の Filtered noise) を出力する。このようにして推論された Harmonic audio と Filtered noise を合成し、最後に必要に応じて残響を付与して合成音を出力する。すべての学習可能なパラメータは、出力の合成音と入力の音響信号間で計算される損失関数が小さくなるように最適化される。DDSPの損失関数には、multi-scale spectral (MSS) ロス [6] が用いられている。

DDSPはデコーダを通して得られるパラメータから Harmonic audio 及び Filtered noise を駆動しており、これは入力の音響信号を正弦波とノイズで人工的に合成していることに対応する。このような音響信号の合成は比較的頑健な音響信号の生成が可能である反面、どれだけパラメータを正しく推定しても、非常にリアルな楽器音の再現は難しく、人工的な音響信号が生成されてしまうデメリットがある。さらに、DDSPは入力音響信号の合成器 (シンセサイザ) を構成することが目的であり、結果的に音色変換等

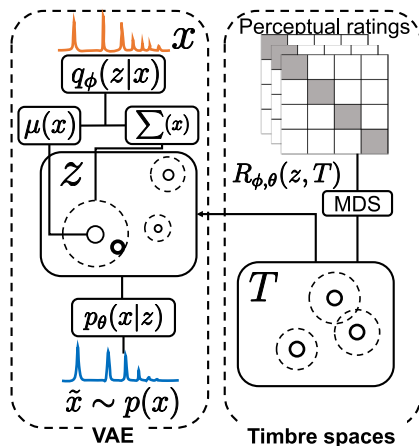


Fig. 3 Concept of regularizing VAE with perceptual metrics. Latent spaces and variables in VAE are forced to be matched to the perceptual timbre spaces structured by perceptual ratings [7].

へ応用することは可能であるが、音色変換が主目的ではない。MFCCもGRU等で非線形に変換されているため、本稿が目指すような音色の潜在特徴量の制御は難しい。

2.3 類似研究2：知覚的メトリクス正則化付きVAE

VAEは学習データに潜む構造を低次元の多次元正規分布に従う特徴量として潜在空間に埋め込むことができる。この潜在空間は多次元正規分布に従うことから、Fig. 1に示すように一定の解釈性や乱数からのデータ生成が可能であるが、依然として潜在空間は物理的又は知覚的な構造と結びついていないという問題がある。この問題を解決するために、楽器音の生成を目的として、知覚的メトリクスに基づく正則化付きVAEが提案されている[7]。Fig. 3に、この手法の概要図を示す。Fig. 3の右上に示すPerceptual ratingsは、過去の研究[8, 9]で蓄積された楽器音の音色に関する知覚メトリクスである。即ち、複数の楽器音間の音色の類似・相違を数値的にレーティングした表形式のデータである。この知覚的メトリクスを、多次元尺度構成法(multi-dimensional scaling: MDS)で低次元空間に変換することで、知覚的メトリクスに基づく音色空間 T を構成している。さらに、音響信号のスペクトルを入出力とするVAEにおいて、潜在空間の分布構造が先の音色空間 T の構造と類似するよう、VAEの学習に正則化項(Fig. 3中の $R_{\phi, \theta}(z, T)$)を付与している。このような正則化を導入することで、VAEの潜在空間が知覚的メトリクスと対応付けられた形で学習され、音色という知覚的情報と深く結びついた解釈性の高い生成モデルを得ることができる。そのため、例えばVAEの潜在空間上で各楽器音の相関を解析することや、ある楽器から別の楽器へ連続的に音響信号を変化させることなどが実現でき、1章で述べた本稿の目的が達成できる可能性がある。

しかしながら、文献[7]の手法は主観的な音色の尺度を

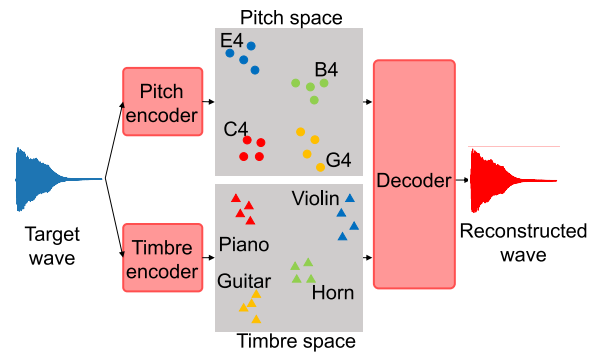


Fig. 4 VAE for learning disentangled pitch and timbre representations of music instrument sounds [11].

持つ知覚メトリクスに基づいており、十分な主観評価データを収集するコストは非常に大きくなるという困難性がある。さらに、主観的な音色の尺度は個人差を多く含んでいるため、正確な音色空間が構成できない可能性も残る。

2.4 類似研究3：VAEによる音色・音高の分離表現学習

本稿の目的に合致した既存手法として、VAEを用いて音色と音高が分離された表現(timbre-pitch-disentangled representations)の学習法が提案されている[1, 10, 11]。これらの手法ではいずれも、Fig. 4に示すように、楽器音の音色と音高が分離された潜在空間及び生成モデルを学習することを目的としている。このような生成モデルの学習においてVAEを用いることで、音色の潜在空間及び音高の潜在空間を独立に獲得でき、1章で述べた本稿の目的に合致したモデルを構築できる可能性がある。分離された音色と音高の潜在変数はデコーダの入力で結合され、楽器音信号が再構成される仕組みである。

しかしながら、前述のいずれの手法でも、DDSPでは用いられたラウドネス(時間的な音量変化)特徴量は分離されておらず、ラウドネスに関する特徴量は音色か音高のいずれか又は両方の潜在空間に押し込められていると考えられる。音を決定づける3要素は音色・音高・音量であるため、時間的な音量変化が人間の「その楽器音らしさ」という知覚に与える影響は無視できない可能性がある。

3. 提案手法

3.1 提案音生成システム全体の説明

提案音生成システムの全体図をFig. 5に示す。まず、入力となる音響信号をエンコーダに通し、音高、音色、及び音量の3つの特徴量を抽出する。この3つの特徴量には、基本周波数 f_0 、MFCC、及びラウドネスをそれぞれ用いる。さらに抽出されたMFCCのみをVAEに入力し、MFCCの潜在特徴量を学習する。出力部では、基本周波数、VAEで生成されたMFCC、及びラウドネスの3つの特徴量をデコーダに入力し、新しい音響信号を生成する。

以上が提案音生成システムの概要である。学習時はFig. 5に示す入力と出力(Original waveとGenerated wave)間

の損失が小さくなるように、VAE 及びデコーダをそれぞれ学習する。学習後は、任意の潜在変数 z から新しい楽器音の音響信号を生成できることが期待される。

3.2 本稿で扱う問題

前節で述べたように、Fig. 5 の提案音生成システムでは、音高、VAE から生成された MFCC、及びラウドネスの 3 つの特徴量から振幅スペクトログラムを生成する必要がある。しかしながら、MFCC は音色のみを低次元空間で表現した特徴量であることから、音高、MFCC、及びラウドネスの 3 つの特徴量から解析的に振幅スペクトログラムを求めることはできない。従って、3 つの特徴量から振幅スペクトログラムを予測する何らかの非線形な変換をデコーダに用いる必要が生じる。

本稿では、前述の問題を取り扱うことを主目的とし、Fig. 5 に示す提案音生成システムを実装する上で必要不可欠なデコーダを DNN で実現する方法について実験的に検討する。すなわち、音高、MFCC、及びラウドネスの 3 つの特徴量から振幅スペクトログラムを高精度に予測する DNN の構築を目指す。この本稿で取り扱う主目的を Fig. 6 に示す。なお、本稿で取り扱うデコーダは MFCC 及びラウドネスを入力とする DNN を想定している。音高の特徴量は離散的であることから、DNN の入力に与えるのではなく、各音高専用に学習した DNN を選択するために用いる。すなわち、予め学習された音高依存の DNN を複数用意し、いずれかの DNN が入力音高により選択される。このような方式を取ることで、DNN に基づくデコーダは音高に対する汎化性能を獲得する必要がなくなり、より高精度な振幅スペクトログラムの予測が可能になると考えられる。デコーダとして用いる DNN には、MLP 型 DNN 及び BiGRU 型 DNN の 2 種類を取り扱い、予測精度を実験的に比較する。

3.3 DNN に基づくデコーダ

提案 DNN デコーダの入力には、2.1 節の計算方法を用いた MFCC とラウドネスを用いる。この時のラウドネスは DDSP の文献中の方法 [6] よりも簡易的な抽出方法として、次式で計算する。

$$v_j = \sum_{i=1}^I |z_{ij}| \quad (5)$$

ここで、 z_{ij} は \mathbf{Z} の要素であり、 $i = 1, 2, \dots, I$ 及び $j = 1, 2, \dots, J$ はそれぞれ周波数ピンのインデクス及び時間フレームを示す。一方、MFCC はラウドネスの影響を排除するために、次式の時間フレーム毎の正規化を施したスペクトログラムから求める。

$$\hat{z}_{ij} = \frac{z_{ij}}{v_j} \quad (6)$$

即ち、式 (6) で得られる正規化済みパワースペクトログラム $|\hat{\mathbf{Z}}|^2$ を用いて式 (3) 及び (4) から MFCC を求める。この MFCC を $\hat{\mathbf{C}}$ と定義する。提案 DNN デコーダの入力は、MFCC とラウドネスを次式で結合した行列とする。

$$\mathbf{X} = \begin{bmatrix} \hat{\mathbf{C}} \\ \mathbf{v}^T \end{bmatrix} \in \mathbb{R}^{(K+1) \times J} \quad (7)$$

ここで、 $\mathbf{v} = [v_1, v_2, \dots, v_J]^T$ である。

前述の 2 種類の提案 DNN デコーダのうち BiGRU について詳細を説明する。BiGRU は双方向 RNN (bidirectional RNN: BiRNN) の一種であり、MFCC の時間方向のように連続的な系列の次元を持つ入力に対して、その系列方向の再帰性を考慮した学習ができる。過去から未来の方向 (順方向) の RNN の出力を過去側から順に $\mathbf{h}_1^{(\text{forward})}, \mathbf{h}_2^{(\text{forward})}, \dots, \mathbf{h}_J^{(\text{forward})}$ とし、未来から過去方向 (逆方向) の RNN の出力を未来側から順に $\mathbf{h}_1^{(\text{backward})}, \mathbf{h}_2^{(\text{backward})}, \dots, \mathbf{h}_J^{(\text{backward})}$ とする。提案 DNN デコーダでは、式 (7) の \mathbf{X} を Fig. 7 のように入力する。具体的には、Fig. 8 に示すように、 \mathbf{X} の時間 j における入力ベクトル \mathbf{x}_j から 4 つの GRU を通して、時刻 j における出力ベクトル \mathbf{h}_j を出力する。この時、BiGRU では順方向の出力ベクトル $\mathbf{h}_j^{(\text{forward})}$ 及び逆方向の出力ベクトル $\mathbf{h}_{j-j}^{(\text{backward})}$ が出力され、その要素毎の積を取ったベクトルを時間 j の出力ベクトル \mathbf{h}_j として扱う。なお、Fig. 8 の GRU の出力ベクトルの要素数は全て I と設定している。

提案 DNN デコーダの学習では、次式の MSS ロスを損失関数に用いる。

$$L_{\text{MSS}}(\mathbf{y}, \hat{\mathbf{y}}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_1 + \|\log \mathbf{y} - \log \hat{\mathbf{y}}\|_1 \quad (8)$$

ここで、 $\mathbf{y} \in \mathbb{R}_{\geq 0}^I$ 及び $\hat{\mathbf{y}} \in \mathbb{R}_{\geq 0}^I$ はそれぞれ入力と予測の振幅スペクトログラムをベクトル化したものであり、 $\|\cdot\|_1$ は L_1 ノルムを表す。

4. 振幅スペクトログラム予測実験

4.1 実験条件

本章では、提案 DNN デコーダの振幅スペクトログラムの予測精度を確認するための実験について述べる。本実験で用いる楽器音信号には、musical instrument digital interface (MIDI) 音源の Roland SVC に含まれる楽器音のうち、Table 1 に示すピアノ 4 種類及びギター 4 種類の計 8 種類の楽器音を用いた。この MIDI 音源 8 種類のそれぞれに対して 4 種類の方法で音色を変化させたデータを加算し、合計でピアノ 20 種類及びギター 20 種類の合計 40 種類の楽器音信号を生成した。音色変化の種類は、低音域強調、高周波強調、コーラス付加、及び残響付加の 4 種類とした。用意した 40 種類の音源の内、ピアノ 18 種類及びギター 18 種類の計 36 種類を学習データとして各 DNN の最適化に用い、残りのピアノ 2 種類及びギター 2 種類の計 4

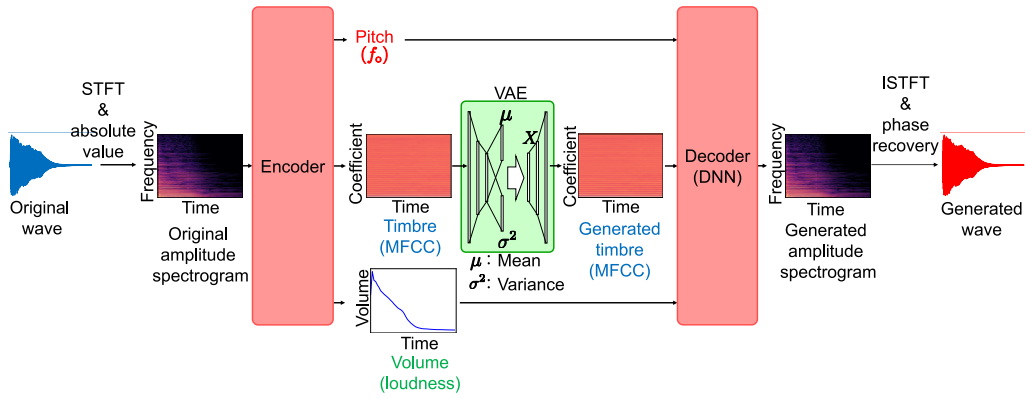


Fig. 5 Detailed process flow of the proposed sound generation system.

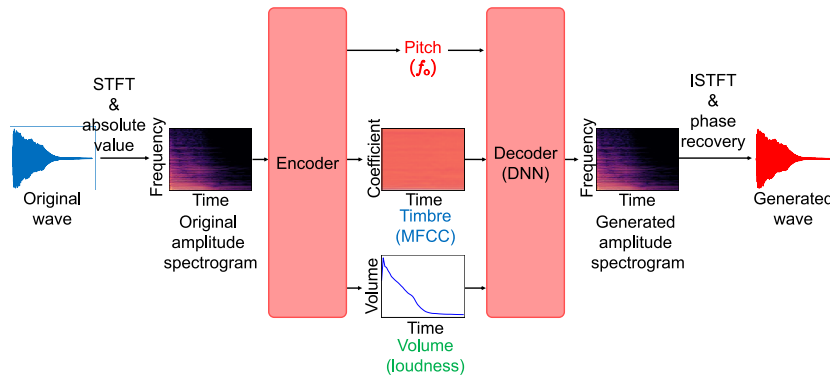


Fig. 6 Training process flow of the proposed DNN-based timbre decoder.

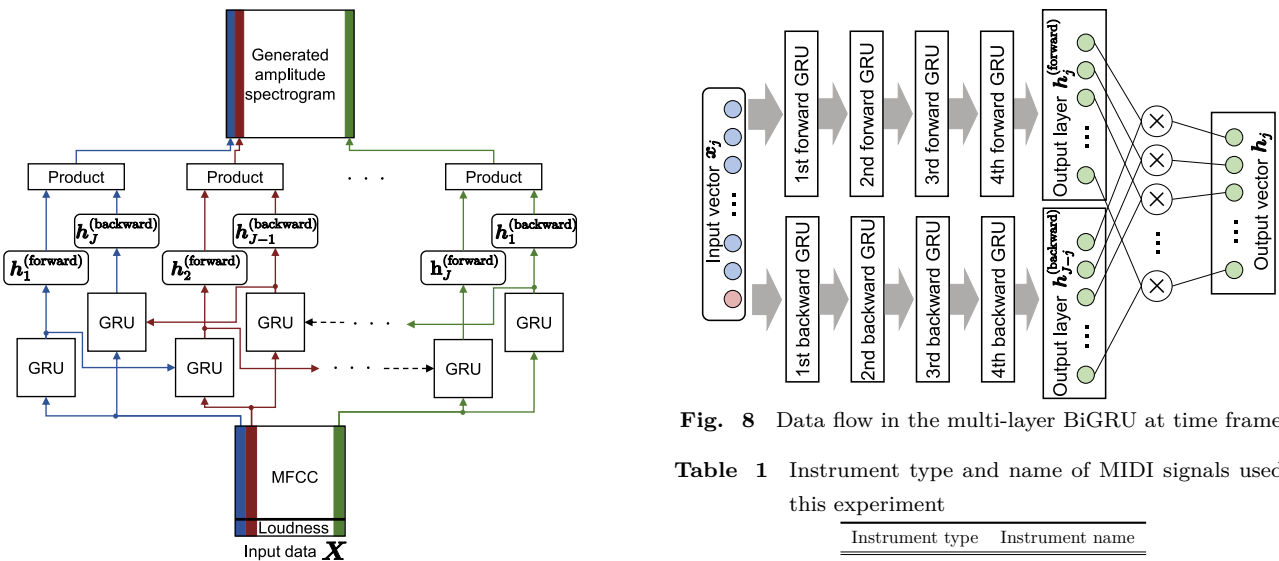


Fig. 7 Architecture of BiRNN used as the DNN decoder.

種類をテストデータに用いた。音源は C3 から B5 までの 36 音を MIDI 音源から作成したが、本稿では G4 音の結果のみを報告する。MIDI 音源から作成した音はサンプリング周波数が 44.1 kHz であり、テンポを 120 bpm とした際の 4 分音符 1 つで構成された 1.18 s の単音の音響信号である。なお、各音響信号に適用する STFT 及び MFCC への変換に用いた条件を Table 2 に示す。

MLP は隠れ層 3 層と出力層で構成され、隠れ層の次元数は入力側から 1024, 512, 及び 512 に設定した。また、MLP と BiGRU の全層において非線形関数に ReLU を用

Fig. 8 Data flow in the multi-layer BiGRU at time frame j .

Table 1 Instrument type and name of MIDI signals used in this experiment

Instrument type	Instrument name
Piano	Piano1
Piano	Piano2
Piano	Piano3
Piano	Honky-tonk
Guitar	Nylon-str. Gt
Guitar	Steel-str. Gt
Guitar	Jazz Gt
Guitar	Muted Gt

いた。最適化には Adam を使い、50000 エポックで各モデルを学習した。

4.2 実験結果

Figs. 9 及び 10 はそれぞれ、学習済みの MLP 及び BiGRU

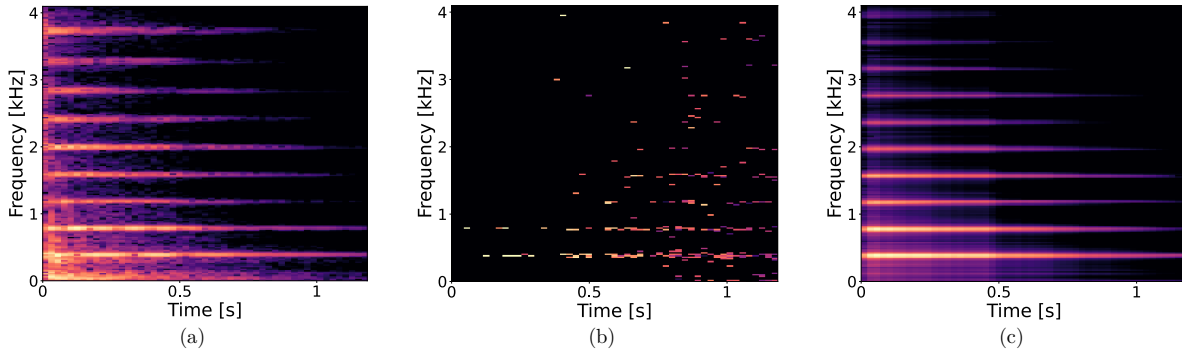


Fig. 9 Example of spectrograms of piano test data: (a) input, (b) predicted by MLP, and (c) predicted by BiGRU.

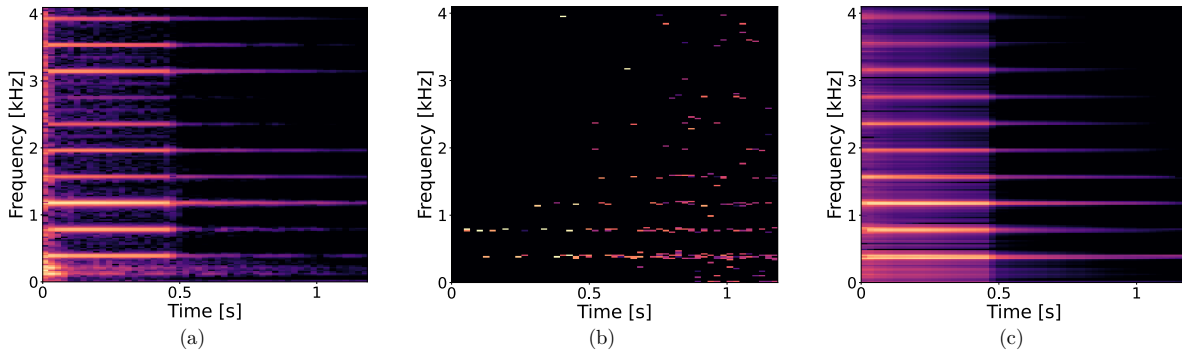


Fig. 10 Example of spectrograms of guitar test data: (a) input, (b) predicted by MLP, and (c) predicted by BiGRU.

Table 2 Conditions used in STFT and MFCC calculations

Window length in STFT	23.2 ms
Shift length in STFT	11.6 ms
Window function in STFT	Hann window
Maximum frequency of mel-filter bank	22.05 kHz
Minimum frequency of mel-filter bank	0.00 kHz
Number of mel-filters	$K = 64$

に対してテストデータ中のピアノ及びギターの G4 音の振幅スペクトログラムを入力した際の予測結果の例を示している。いずれの結果をみても、MLP を用いた振幅スペクトログラムの予測は失敗している。このような予測失敗の傾向は学習データに対しても同様であり、MLP を用いて MFCC とラウドネスのみから振幅スペクトログラムを予測することが非常に困難であることを示唆している。一方、BiGRU を用いた予測では、ピアノとギターの両楽器のテストデータで調波構造や時間的な推移を正確に予測できている。これは、BiGRU が MFCC 及びラウドネスの時間方向の連続性を考慮しつつ学習できたことに起因している。

5. おわりに

本稿では、VAE に基づく楽器音の音色変換を目的とした提案音生成システムについて説明し、その実現に必要な部分システムとして、音高、音色、及び音量の3つの特徴量から振幅スペクトログラムを予測する DNN デコーダを提案した。今後の課題として、生成された信号の歪みに関する客観評価及び提案音生成システムの構築が挙げられる。

謝辞 本研究の一部は公益信託小野音響学研究助成基金及び JSPS 科研費 22H03652 及びの助成を受けた。

参考文献

- [1] Y. J. Luo, K. Agres, and D. Herremans, "Learning disentangled representations of timbre and pitch for musical instrument sounds using Gaussian mixture variational autoencoders," in *Proc. ISMIR*, 2019.
- [2] R. Yu, "A tutorial on VAEs: From Bayes' rule to lossless compression," arXiv: 2006.10273, 2020.
- [3] K. Cho, B. V. Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," arXiv: 1406.1078, 2014.
- [4] F. Zheng, G. Zhang, and Z. Song, "Comparison of different implementations of MFCC," *J. Computer Science and Technology*, vol. 16, pp. 582–589, 2001.
- [5] D. O'Shaughnessy, *Speech Communication: Human and Machine*. Addison-Wesley, 1987.
- [6] J. Engel, L. Hantrakul, C. Gu, and A. Roberts, "Differentiable Digital Signal Processing," in *Proc. ICLR*, 2020.
- [7] P. Esling, A. Chemla-RomenuSantos, and A. Bitton, "Generative timbre spaces: regularizing variational autoencoders with perceptual metrics," in *Proc. DAFX*, 2018.
- [8] S. McAdams, S. Winsberg, S. Donnadieu, G. D. Seoete, and J. Krimphoff, "Perceptual scaling of synthesized musical timbres: Common dimensions, specificities, and latent subject classes," *Psychological research*, vol. 58, no. 3, pp. 177–192, 1995.
- [9] S. Lakatos, "A common perceptual space for harmonic and percussive timbres," *Perception & psychophysic*, vol. 62, no. 7, pp. 1426–1439, 2000.
- [10] Y. J. Luo, K. W. Cheuk, T. Nakano, M. Goto, and D. Herremans, "Unsupervised disentanglement of pitch and timbre for isolated musical instrument sounds," in *Proc. ISMIR*, pp. 700–707, 2020.
- [11] K. Tanaka, R. Nishikimi, Y. Bando, K. Yoshii, and S. Morishima, "Pitch-timbre disentanglement of musical instrument sounds based on VAE-based metric learning," in *Proc. ICASSP*, pp. 111–115, 2021.