

# 第1章 判別器の構成ブロック

## 1.1 はじめに

この章では、判別器を構成する各種ブロックの詳細を説明する。各有音程楽器ピアノロールの特徴抽出を行う Piano Block, Guitar Block, Bass Block, Strings Block を説明する。続いて、有音程楽器間に共通する特徴を抽出する Low-band Block, Mid-band Block, High-band Block, Tonal Block を説明する。さらに、打楽器ピアノロールの特徴抽出を行う Percussion Block, Drums Block, Other Percussion Block を説明する。また、有音程楽器の補助的な特徴抽出を行う Chroma Block, Piano Polyphonicity Block, Guitar Polyphonicity Block, Bass Polyphonicity Block, Strings Polyphonicity Block と各 Block の出力を統合し特徴抽出を行う Merged Block 1, Merged Block 2, Merged Block 3 について説明する。

判別器への入力はマルチトラックピアノロール  $\mathbf{X} = (\mathbf{X}_{Pi}, \mathbf{X}_{Gu}, \mathbf{X}_{Base}, \mathbf{X}_{St}, \mathbf{X}_{Pe})$  である。Conv 関数での計算のために各ピアノロールを 4 階テンソルへの変換を行う。まず、ピアノトラック  $\mathbf{X}_{Pi} \in \mathbb{R}^{4 \times 96 \times 84}$  を

$$\mathbf{X}_{Piano} = \text{Reshape}_{expand}(\mathbf{X}_{Pi}) \in \mathbb{R}^{1 \times 4 \times 96 \times 84} \quad (1.1)$$

へと変換する。ここで、 $\text{Reshape}_{expand}$  は入力の  $4 \times 96 \times 84$  の三階テンソルを  $1 \times 4 \times 96 \times 84$  の 4 階テンソルへと拡張する関数である。同様に、ギタートラック  $\mathbf{X}_{Gu} \in \mathbb{R}^{4 \times 96 \times 84}$  を

$$\mathbf{X}_{Guitar} = \text{Reshape}_{expand}(\mathbf{X}_{Gu}) \in \mathbb{R}^{1 \times 4 \times 96 \times 84} \quad (1.2)$$

へと変換する。ベーストラック  $\mathbf{X}_{Ba} \in \mathbb{R}^{4 \times 96 \times 84}$  を

$$\mathbf{X}_{Bass} = \text{Reshape}_{expand}(\mathbf{X}_{Ba}) \in \mathbb{R}^{1 \times 4 \times 96 \times 84} \quad (1.3)$$

へと変換する。ストリングストラック  $\mathbf{X}_{St} \in \mathbb{R}^{4 \times 96 \times 84}$  を

$$\mathbf{X}_{Strings} = \text{Reshape}_{expand}(\mathbf{X}_{St}) \in \mathbb{R}^{1 \times 4 \times 96 \times 84} \quad (1.4)$$

へと変換する。打楽器トラック  $\mathbf{X}_{Pe} \in \mathbb{R}^{4 \times 96 \times 84}$  を

$$\mathbf{X}_{Percussion} = \text{Reshape}_{expand}(\mathbf{X}_{Pe}) \in \mathbb{R}^{1 \times 4 \times 96 \times 84} \quad (1.5)$$

へと変換する。

## 1.2 Individual Blocks

### 1.2.1 Piano Block

従来手法である BinaryMuseGAN における判別器の Piano Block は、3 個のサブブロックで構成されていた。本研究では Piano Block を、6 個のサブブロックで構成する。それらを Piano Block 1, Piano Block 2, Piano Block 3, Piano Block 4, Piano Block 5, および Piano Block 6 と呼ぶことにする。順に説明する。

Piano Block 1 は、表 1.1 に示す 2 層の畳み込み層で構成される。第 1 畳み込み層は  $\mathbf{X}_{Piano} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Piano}, \mathcal{L}_{PiB1}^{(1)}, \mathbf{t}_{PiB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 7} \quad (1.6)$$

を抽出する。ここで、 $\mathcal{L}_{PiB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{PiB1}^{(1)} = (1, 1, 12)$  である。この段階で第 4 モードのサイズが 7 になり、オクターブごとの特徴抽出を開始する。次に、第 2 畳み込み層が  $\mathbf{Y}_{PiB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 7}$  から

$$\mathbf{Y}_{PiB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB1}^{(1)}, \mathcal{L}_{PiB1}^{(2)}, \mathbf{t}_{PiB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.7)$$

を抽出する。ここで、 $\mathcal{L}_{PiB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{PiB1}^{(2)} = (1, 6, 1)$  である。この段階で第 3 モードのサイズが 16 になり、微細なリズム構造の特徴抽出を開始する。以上のように、Piano Block 1 では、第 4 モードの特徴抽出を行ってから第 3 モードの特徴抽出を行う。すなわち、音高方向の特徴抽出の後に時間方向の特徴抽出を行う。

表 1.1: Piano Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 12}$	$\mathbf{t}_{PiB1}^{(1)} = (1, 1, 12)$
第 2 層	$\mathcal{L}_{PiB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{PiB1}^{(2)} = (1, 6, 1)$

一方、Piano Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する。すなわち、第 3 モードの時間方向特徴量を抽出してから第 4 モードである音高方向の特徴抽出を行う。まず、第 1 畳み込み層が  $\mathbf{X}_{Piano} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Piano}, \mathcal{L}_{PiB2}^{(1)}, \mathbf{t}_{PiB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.8)$$

を抽出する。ここで、 $\mathcal{L}_{PiB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{PiB2}^{(1)} = (1, 6, 1)$  である。この段階で第 3 モードのサイズが 16 になり、微細なリズム構造の特徴抽出を開始する。次に、第 2 畳み込み層が  $\mathbf{Y}_{PiB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{PiB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB2}^{(1)}, \mathcal{L}_{PiB2}^{(2)}, \mathbf{t}_{PiB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.9)$$

を抽出する．ここで， $\mathcal{L}_{PiB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{PiB2}^{(2)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Piano Block 2 の畳み込み層における各種パラメータを表 1.2 にまとめて示す．

表 1.2: Piano Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$	$\mathbf{t}_{PiB2}^{(1)} = (1, 6, 1)$
第 2 層	$\mathcal{L}_{PiB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{PiB2}^{(2)} = (1, 1, 12)$

続いて，Piano Block 3 は表 1.3 に示す 3 層の畳み込み層で構成される．まず，第 1 畳み込み層が  $\mathbf{X}_{Piano} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB3}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Piano}, \mathcal{L}_{PiB3}^{(1)}, \mathbf{t}_{PiB3}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.10)$$

を抽出する．ここで，畳み込みは Same 畳み込みであり， $\mathcal{L}_{PiB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$ ， $\mathbf{t}_{PiB3}^{(1)} = (1, 1, 1)$  である．この段階で第4モードのサイズが3のカーネルを用いることにより，隣接 3 半音の特徴抽出を開始する．この第 1 畳み込み層は，隣接 3 半音からの不協和音検出を目的としている．次に，第 2 畳み込み層は  $\mathbf{Y}_{PiB3}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB3}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB3}^{(1)}, \mathcal{L}_{PiB3}^{(2)}, \mathbf{t}_{PiB3}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.11)$$

を抽出する．ここで， $\mathcal{L}_{PiB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{PiB3}^{(2)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．次に，第 3 畳み込み層が  $\mathbf{Y}_{PiB3}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{PiB3}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB3}^{(2)}, \mathcal{L}_{PiB3}^{(3)}, \mathbf{t}_{PiB3}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.12)$$

を抽出する．ここで， $\mathcal{L}_{PiB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{PiB3}^{(3)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Piano Block 3 の畳み込み層における各種パラメータを表 1.3 にまとめて示す．

表 1.3: Piano Block 3 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{PiB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$	$\mathbf{t}_{PiB3}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{PiB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{PiB3}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{PiB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{PiB3}^{(3)} = (1, 1, 12)$

続いて、Piano Block 4 は表 1.3 に示す 3 層の畳み込み層で構成される。Piano Block 3 では第一層で隣接 3 半音から特徴を抽出したが、Piano Block 4 では隣接 2 半音から特徴抽出を行う。まず、第 1 畳み込み層が  $\mathbf{X}_{Piano} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB4}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Piano}, \mathcal{L}_{PiB4}^{(1)}, \mathbf{t}_{PiB4}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.13)$$

を抽出する。ここで、畳み込みは Same 畳み込みであり、 $\mathcal{L}_{PiB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 2}$ 、 $\mathbf{t}_{PiB4}^{(1)} = (1, 1, 1)$  である。この段階で第 4 モードのサイズが 2 のカーネルを用いることにより、隣接 2 半音の特徴抽出を開始する。この第 1 畳み込み層は、隣接 2 半音からの不協和音検出を目的としている。次に、第 2 畳み込み層は  $\mathbf{Y}_{PiB4}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{PiB4}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB4}^{(1)}, \mathcal{L}_{PiB4}^{(2)}, \mathbf{t}_{PiB4}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.14)$$

を抽出する。ここで、 $\mathcal{L}_{PiB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{PiB4}^{(2)} = (1, 6, 1)$  である。この段階で第 3 モードのサイズが 16 になり、微細なリズム構造の特徴抽出を開始する。次に、第 3 畳み込み層が  $\mathbf{Y}_{PiB4}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{PiB4}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB4}^{(2)}, \mathcal{L}_{PiB4}^{(3)}, \mathbf{t}_{PiB4}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.15)$$

を抽出する。ここで、 $\mathcal{L}_{PiB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{PiB4}^{(3)} = (1, 1, 12)$  である。この段階で第 4 モードのサイズが 7 になり、オクターブごとの特徴抽出を開始する。Piano Block 4 の畳み込み層における各種パラメータを表 1.4 にまとめて示す。

表 1.4: Piano Block 4 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{PiB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$	$\mathbf{t}_{PiB4}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{PiB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{PiB4}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{PiB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{PiB4}^{(3)} = (1, 1, 12)$

Piano Block 5 では、Piano Block 3 と Piano Block 4 の出力から特徴抽出を以下のように行う。まず、入力された特徴  $\mathbf{Y}_{PiB3}^{(3)}$ 、 $\mathbf{Y}_{PiB4}^{(3)}$  を

$$\mathbf{Y}_{PiB5}^{(0)} = \text{Concat}(1, \mathbf{Y}_{PiB3}^{(3)}, \mathbf{Y}_{PiB4}^{(3)}) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.16)$$

によって結合する。 $\mathbf{Y}_{PiB3}^{(3)}$  と  $\mathbf{Y}_{PiB4}^{(3)}$  の 2 テンソルを結合することにより、二種類の不協和音検出の結果を統合している。続いて、畳み込み層が  $\mathbf{Y}_{PiB5}^{(0)} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{PiB5}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB5}^{(0)}, \mathcal{L}_{PiB5}^{(1)}, \mathbf{t}_{PiB5}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.17)$$

を生成する．ここで， $\mathcal{L}_{PiB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$  であり， $\mathbf{t}_{PiB5}^{(1)} = (1, 1, 1)$  である．Piano Block 5 の畳み込み層における各種パラメータを表 1.5 にまとめて示す．

表 1.5: Piano Block 5 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$	$\mathbf{t}_{PiB5}^{(1)} = (1, 1, 1)$

Piano Block 6 では，Piano Block 1，Piano Block 2，Piano Block 5 の出力から特徴抽出を以下のように行う．まず，入力された特徴  $\mathbf{Y}_{PiB1}^{(2)}$ ， $\mathbf{Y}_{PiB2}^{(2)}$ ， $\mathbf{Y}_{PiB5}^{(1)}$  を

$$\mathbf{Y}_{PiB6}^{(0)} = \text{Concat}(1, \mathbf{Y}_{PiB1}^{(2)}, \mathbf{Y}_{PiB2}^{(2)}, \mathbf{Y}_{PiB5}^{(1)}) \in \mathbb{R}^{160 \times 4 \times 16 \times 7} \quad (1.18)$$

によって結合する．続いて，畳み込み層が  $\mathbf{Y}_{PiB6}^{(0)} \in \mathbb{R}^{160 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{PiB6}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiB6}^{(0)}, \mathcal{L}_{PiB6}^{(1)}, \mathbf{t}_{PiB6}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.19)$$

を生成する．ここで， $\mathcal{L}_{PiB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$  であり， $\mathbf{t}_{PiB6}^{(1)} = (1, 1, 1)$  である．Piano Block 6 の畳み込み層における各種パラメータを表 1.6 にまとめて示す．Piano Block は処理をこれで終了し， $\mathbf{Y}_{PiB6}^{(1)}$  を Piano Block の出力  $\mathbf{Y}_{PiB}$  として出力する： $\mathbf{Y}_{PiB} = \mathbf{Y}_{PiB6}^{(1)}$ ．これが，ピアノトラックの特徴量であり，Band Blocks に入力される．

表 1.6: Piano Block 6 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$	$\mathbf{t}_{PiB6}^{(1)} = (1, 1, 1)$

## 1.2.2 Guitar Block

従来手法である BinaryMuseGAN における判別器の Guitar Block は，3 個のサブブロックで構成されていた．本研究では Guitar Block を Piano Block と同様に 6 個のサブブロックで構成する．それらを Guitar Block 1，Guitar Block 2，Guitar Block 3，Guitar Block 4，Guitar Block 5，および Guitar Block 6 と呼ぶことにする．順に説明する．

Guitar Block 1 は，表 1.7 に示す 2 層の畳み込み層で構成される．第 1 畳み込み層は  $\mathbf{X}_{Guitar} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Guitar}, \mathcal{L}_{GB1}^{(1)}, \mathbf{t}_{GB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 7} \quad (1.20)$$

を抽出する．ここで， $\mathcal{L}_{GB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{GB1}^{(1)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．次に，第2畳み込み層が  $\mathbf{Y}_{GB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 7}$  から

$$\mathbf{Y}_{GB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB1}^{(1)}, \mathcal{L}_{GB1}^{(2)}, \mathbf{t}_{GB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.21)$$

を抽出する．ここで， $\mathcal{L}_{GB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{GB1}^{(2)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．以上のように，Guitar Block 1 では，第4モードの特徴抽出を行ってから第3モードの特徴抽出を行う．すなわち，音高方向の特徴抽出の後に時間方向の特徴抽出を行う．

表 1.7: Guitar Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{GB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{GB1}^{(1)} = (1, 1, 12)$
第2層	$\mathcal{L}_{GB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{GB1}^{(2)} = (1, 6, 1)$

一方，Guitar Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する．すなわち，第3モードの時間方向特徴量を抽出してから第4モードである音高方向の特徴抽出を行う．まず，第1畳み込み層が  $\mathbf{X}_{Guitar} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Guitar}, \mathcal{L}_{GB2}^{(1)}, \mathbf{t}_{GB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.22)$$

を抽出する．ここで， $\mathcal{L}_{GB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{GB2}^{(1)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．次に，第2畳み込み層が  $\mathbf{Y}_{GB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{GB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB2}^{(1)}, \mathcal{L}_{GB2}^{(2)}, \mathbf{t}_{GB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.23)$$

を抽出する．ここで， $\mathcal{L}_{GB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{GB2}^{(2)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Guitar Block 2 の畳み込み層における各種パラメータを表 1.8 にまとめて示す．

表 1.8: Guitar Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{GB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$	$\mathbf{t}_{GB2}^{(1)} = (1, 6, 1)$
第2層	$\mathcal{L}_{GB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{GB2}^{(2)} = (1, 1, 12)$

続いて、Guitar Block 3 は表 1.9 に示す 3 層の畳み込み層で構成される。まず、第 1 畳み込み層が  $\mathbf{X}_{Guitar} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB3}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Guitar}, \mathcal{L}_{GB3}^{(1)}, \mathbf{t}_{GB3}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.24)$$

を抽出する。ここで、畳み込みは Same 畳み込みであり、 $\mathcal{L}_{GB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$ 、 $\mathbf{t}_{GB3}^{(1)} = (1, 1, 1)$  である。この段階で第 4 モードのサイズが 3 のカーネルを用いることにより、隣接 3 半音の特徴抽出を開始する。この第 1 畳み込み層は、隣接 3 半音からの不協和音検出を目的としている。次に、第 2 畳み込み層は  $\mathbf{Y}_{GB3}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB3}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB3}^{(1)}, \mathcal{L}_{GB3}^{(2)}, \mathbf{t}_{GB3}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.25)$$

を抽出する。ここで、 $\mathcal{L}_{GB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{GB3}^{(2)} = (1, 6, 1)$  である。この段階で第 3 モードのサイズが 16 になり、微細なリズム構造の特徴抽出を開始する。次に、第 3 畳み込み層が  $\mathbf{Y}_{GB3}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{GB3}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB3}^{(2)}, \mathcal{L}_{GB3}^{(3)}, \mathbf{t}_{GB3}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.26)$$

を抽出する。ここで、 $\mathcal{L}_{GB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{GB3}^{(3)} = (1, 1, 12)$  である。この段階で第 4 モードのサイズが 7 になり、オクターブごとの特徴抽出を開始する。Guitar Block 3 の畳み込み層における各種パラメータを表 1.9 にまとめて示す。

表 1.9: Guitar Block 3 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{GB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$	$\mathbf{t}_{GB3}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{GB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{GB3}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{GB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{GB3}^{(3)} = (1, 1, 12)$

続いて、Guitar Block 4 は表 1.9 に示す 3 層の畳み込み層で構成される。Guitar Block 3 では第一層で隣接 3 半音から特徴を抽出したが、Guitar Block 4 では隣接 2 半音から特徴抽出を行う。まず、第 1 畳み込み層が  $\mathbf{X}_{Guitar} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB4}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Guitar}, \mathcal{L}_{GB4}^{(1)}, \mathbf{t}_{GB4}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.27)$$

を抽出する。ここで、畳み込みは Same 畳み込みであり、 $\mathcal{L}_{GB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$ 、 $\mathbf{t}_{GB4}^{(1)} = (1, 1, 1)$  である。この段階で第 4 モードのサイズが 2 のカーネルを用いることにより、隣接 2 半音の特徴抽出を開始する。この第 1 畳み込み層は、隣接 2 半音からの不協和音検出を目的としている。次に、第 2 畳み込み層は  $\mathbf{Y}_{GB4}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{GB4}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB4}^{(1)}, \mathcal{L}_{GB4}^{(2)}, \mathbf{t}_{GB4}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.28)$$

を抽出する．ここで， $\mathcal{L}_{GB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{GB4}^{(2)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．次に，第3畳み込み層が  $\mathbf{Y}_{GB4}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{GB4}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB4}^{(2)}, \mathcal{L}_{GB4}^{(3)}, \mathbf{t}_{GB4}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.29)$$

を抽出する．ここで， $\mathcal{L}_{GB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{GB4}^{(3)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Guitar Block 4 の畳み込み層における各種パラメータを表 1.10 にまとめて示す．

表 1.10: Guitar Block 4 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層 (Same)	$\mathcal{L}_{GB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$	$\mathbf{t}_{GB4}^{(1)} = (1, 1, 1)$
第2層	$\mathcal{L}_{GB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{GB4}^{(2)} = (1, 6, 1)$
第3層	$\mathcal{L}_{GB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{GB4}^{(3)} = (1, 1, 12)$

Guitar Block 5 では，Guitar Block 3 と Guitar Block 4 の出力から特徴抽出を以下のように行う．まず，入力された特徴  $\mathbf{Y}_{GB3}^{(3)}$ ， $\mathbf{Y}_{GB4}^{(3)}$  を

$$\mathbf{Y}_{GB5}^{(0)} = \text{Concat}(1, \mathbf{Y}_{GB3}^{(3)}, \mathbf{Y}_{GB4}^{(3)}) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.30)$$

によって結合する． $\mathbf{Y}_{GB3}^{(3)}$  と  $\mathbf{Y}_{GB4}^{(3)}$  の2テンソルを結合することにより，二種類の不協和音検出の結果を統合している．続いて，畳み込み層が  $\mathbf{Y}_{GB5}^{(0)} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{GB5}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB5}^{(0)}, \mathcal{L}_{GB5}^{(1)}, \mathbf{t}_{GB5}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.31)$$

を生成する．ここで， $\mathcal{L}_{GB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$  であり， $\mathbf{t}_{GB5}^{(1)} = (1, 1, 1)$  である．Guitar Block 5 の畳み込み層における各種パラメータを表 1.11 にまとめて示す．

表 1.11: Guitar Block 5 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{GB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$	$\mathbf{t}_{GB5}^{(1)} = (1, 1, 1)$

Guitar Block 6 では，Guitar Block 1，Guitar Block 2，Guitar Block 5 の出力から特徴抽出を以下のように行う．まず，入力された特徴  $\mathbf{Y}_{GB1}^{(2)}$ ， $\mathbf{Y}_{GB2}^{(2)}$ ， $\mathbf{Y}_{GB5}^{(1)}$  を

$$\mathbf{Y}_{GB6}^{(0)} = \text{Concat}(1, \mathbf{Y}_{GB1}^{(2)}, \mathbf{Y}_{GB2}^{(2)}, \mathbf{Y}_{GB5}^{(1)}) \in \mathbb{R}^{160 \times 4 \times 16 \times 7} \quad (1.32)$$



によって結合する．続いて，畳み込み層が  $\mathbf{Y}_{GB6}^{(0)} \in \mathbb{R}^{160 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{GB6}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GB6}^{(0)}, \mathcal{L}_{GB6}^{(1)}, \mathbf{t}_{GB6}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.33)$$

を生成する．ここで， $\mathcal{L}_{GB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$  であり， $\mathbf{t}_{GB6}^{(1)} = (1, 1, 1)$  である．Guitar Block 6 の畳み込み層における各種パラメータを表 1.12 にまとめて示す．Guitar Block は処理をこれで終了し， $\mathbf{Y}_{GB6}^{(1)}$  を Guitar Block の出力  $\mathbf{Y}_{GB}$  として出力する： $\mathbf{Y}_{GB} = \mathbf{Y}_{GB6}^{(1)}$ ．これが，ギタートラックの特徴量であり，Band Blocks に入力される．

表 1.12: Guitar Block 6 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	スライドベクトル
第 1 層	$\mathcal{L}_{GB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$	$\mathbf{t}_{GB6}^{(1)} = (1, 1, 1)$

### 1.2.3 Bass Block

従来手法である BinaryMuseGAN における判別器の Bass Block は，3 個のサブブロックで構成されていた．本研究では Bass Block を Piano Block と同様に 6 個のサブブロックで構成する．それらを Bass Block 1, Bass Block 2, Bass Block 3, Bass Block 4, Bass Block 5, および Bass Block 6 と呼ぶことにする．順に説明する．

Bass Block 1 は，表 1.13 に示す 2 層の畳み込み層で構成される．第 1 畳み込み層は  $\mathbf{X}_{Bass} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Bass}, \mathcal{L}_{BaB1}^{(1)}, \mathbf{t}_{BaB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 7} \quad (1.34)$$

を抽出する．ここで， $\mathcal{L}_{BaB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{BaB1}^{(1)} = (1, 1, 12)$  である．この段階で第 4 モードのサイズが 7 になり，オクターブごとの特徴抽出を開始する．次に，第 2 畳み込み層が  $\mathbf{Y}_{BaB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 7}$  から

$$\mathbf{Y}_{BaB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB1}^{(1)}, \mathcal{L}_{BaB1}^{(2)}, \mathbf{t}_{BaB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.35)$$

を抽出する．ここで， $\mathcal{L}_{BaB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{BaB1}^{(2)} = (1, 6, 1)$  である．この段階で第 3 モードのサイズが 16 になり，微細なリズム構造の特徴抽出を開始する．以上のように，Bass Block 1 では，第 4 モードの特徴抽出を行ってから第 3 モードの特徴抽出を行う．すなわち，音高方向の特徴抽出の後に時間方向の特徴抽出を行う．

表 1.13: Bass Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{BaB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{BaB1}^{(1)} = (1, 1, 12)$
第 2 層	$\mathcal{L}_{BaB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{BaB1}^{(2)} = (1, 6, 1)$

一方, Bass Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する. すなわち, 第 3 モードの時間方向特徴量を抽出してから第 4 モードである音高方向の特徴抽出を行う. まず, 第 1 畳み込み層が  $\mathbf{X}_{Bass} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Bass}, \mathcal{L}_{BaB2}^{(1)}, \mathbf{t}_{BaB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.36)$$

を抽出する. ここで,  $\mathcal{L}_{BaB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$  であり,  $\mathbf{t}_{BaB2}^{(1)} = (1, 6, 1)$  である. この段階で第 3 モードのサイズが 16 になり, 微細なリズム構造の特徴抽出を開始する. 次に, 第 2 畳み込み層が  $\mathbf{Y}_{BaB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{BaB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB2}^{(1)}, \mathcal{L}_{BaB2}^{(2)}, \mathbf{t}_{BaB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.37)$$

を抽出する. ここで,  $\mathcal{L}_{BaB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$  であり,  $\mathbf{t}_{BaB2}^{(2)} = (1, 1, 12)$  である. この段階で第 4 モードのサイズが 7 になり, オクターブごとの特徴抽出を開始する. Bass Block 2 の畳み込み層における各種パラメータを表 1.14 にまとめて示す.

表 1.14: Bass Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{BaB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$	$\mathbf{t}_{BaB2}^{(1)} = (1, 6, 1)$
第 2 層	$\mathcal{L}_{BaB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{BaB2}^{(2)} = (1, 1, 12)$

続いて, Bass Block 3 は表 1.15 に示す 3 層の畳み込み層で構成される. まず, 第 1 畳み込み層が  $\mathbf{X}_{Bass} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB3}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Bass}, \mathcal{L}_{BaB3}^{(1)}, \mathbf{t}_{BaB3}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.38)$$

を抽出する. ここで, 畳み込みは Same 畳み込みであり,  $\mathcal{L}_{BaB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$ ,  $\mathbf{t}_{BaB3}^{(1)} = (1, 1, 1)$  である. この段階で第 4 モードのサイズが 3 のカーネルを用いることにより, 隣接 3 半音の特徴抽出を開始する. この第 1 畳み込み層は, 隣接 3 半音からの不協和音検出を目的としている. 次に, 第 2 畳み込み層は  $\mathbf{Y}_{BaB3}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB3}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB3}^{(1)}, \mathcal{L}_{BaB3}^{(2)}, \mathbf{t}_{BaB3}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.39)$$

を抽出する．ここで， $\mathcal{L}_{BaB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{BaB3}^{(2)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．次に，第3畳み込み層が  $\mathbf{Y}_{BaB3}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{BaB3}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB3}^{(2)}, \mathcal{L}_{BaB3}^{(3)}, \mathbf{t}_{BaB3}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.40)$$

を抽出する．ここで， $\mathcal{L}_{BaB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{BaB3}^{(3)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Bass Block 3 の畳み込み層における各種パラメータを表 1.15 にまとめて示す．

表 1.15: Bass Block 3 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層 (Same)	$\mathcal{L}_{BaB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$	$\mathbf{t}_{BaB3}^{(1)} = (1, 1, 1)$
第2層	$\mathcal{L}_{BaB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{BaB3}^{(2)} = (1, 6, 1)$
第3層	$\mathcal{L}_{BaB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{BaB3}^{(3)} = (1, 1, 12)$

続いて，Bass Block 4 は表 1.15 に示す3層の畳み込み層で構成される．Bass Block 3 では第一層で隣接3半音から特徴を抽出したが，Bass Block 3 では隣接2半音から特徴抽出を行う．まず，第1畳み込み層が  $\mathbf{X}_{Bass} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB4}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Bass}, \mathcal{L}_{BaB4}^{(1)}, \mathbf{t}_{BaB4}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.41)$$

を抽出する．ここで，畳み込みはSame畳み込みであり， $\mathcal{L}_{BaB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$ ， $\mathbf{t}_{BaB4}^{(1)} = (1, 1, 1)$  である．この段階で第4モードのサイズが2のカーネルを用いることにより，隣接2半音の特徴抽出を開始する．この第1畳み込み層は，隣接2半音からの不協和音検出を目的としている．次に，第2畳み込み層は  $\mathbf{Y}_{BaB4}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{BaB4}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB4}^{(1)}, \mathcal{L}_{BaB4}^{(2)}, \mathbf{t}_{BaB4}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.42)$$

を抽出する．ここで， $\mathcal{L}_{BaB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり， $\mathbf{t}_{BaB4}^{(2)} = (1, 6, 1)$  である．この段階で第3モードのサイズが16になり，微細なリズム構造の特徴抽出を開始する．次に，第3畳み込み層が  $\mathbf{Y}_{BaB4}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{BaB4}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB4}^{(2)}, \mathcal{L}_{BaB4}^{(3)}, \mathbf{t}_{BaB4}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.43)$$

を抽出する．ここで， $\mathcal{L}_{BaB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{BaB4}^{(3)} = (1, 1, 12)$  である．この段階で第4モードのサイズが7になり，オクターブごとの特徴抽出を開始する．Bass Block 4 の畳み込み層における各種パラメータを表 1.16 にまとめて示す．

表 1.16: Bass Block 4 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{BaB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 2}$	$\mathbf{t}_{BaB4}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{BaB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{BaB4}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{BaB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{BaB4}^{(3)} = (1, 1, 12)$

Bass Block 5 では, Bass Block 3 と Bass Block 4 の出力から特徴抽出を以下のように行う. まず, 入力された特徴  $\mathbf{Y}_{BaB3}^{(3)}$ ,  $\mathbf{Y}_{BaB4}^{(3)}$  を

$$\mathbf{Y}_{BaB5}^{(0)} = \text{Concat}(1, \mathbf{Y}_{BaB3}^{(3)}, \mathbf{Y}_{BaB4}^{(3)}) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.44)$$

によって結合する.  $\mathbf{Y}_{BaB3}^{(3)}$  と  $\mathbf{Y}_{BaB4}^{(3)}$  の 2 テンソルを結合することにより, 二種類の不協和音検出の結果を統合している. 続いて, 畳み込み層が  $\mathbf{Y}_{BaB5}^{(0)} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{BaB5}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB5}^{(0)}, \mathcal{L}_{BaB5}^{(1)}, \mathbf{t}_{BaB5}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.45)$$

を生成する. ここで,  $\mathcal{L}_{BaB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$  であり,  $\mathbf{t}_{BaB5}^{(1)} = (1, 1, 1)$  である. Bass Block 5 の畳み込み層における各種パラメータを表 1.17 にまとめて示す.

表 1.17: Bass Block 5 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{BaB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$	$\mathbf{t}_{BaB5}^{(1)} = (1, 1, 1)$

Bass Block 6 では, Bass Block 1, Bass Block 2, Bass Block 5 の出力から特徴抽出を以下のように行う. まず, 入力された特徴  $\mathbf{Y}_{BaB1}^{(2)}$ ,  $\mathbf{Y}_{BaB2}^{(2)}$ ,  $\mathbf{Y}_{BaB5}^{(1)}$  を

$$\mathbf{Y}_{BaB6}^{(0)} = \text{Concat}(1, \mathbf{Y}_{BaB1}^{(2)}, \mathbf{Y}_{BaB2}^{(2)}, \mathbf{Y}_{BaB5}^{(1)}) \in \mathbb{R}^{160 \times 4 \times 16 \times 7} \quad (1.46)$$

によって結合する. 続いて, 畳み込み層が  $\mathbf{Y}_{BaB6}^{(0)} \in \mathbb{R}^{160 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{BaB6}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaB6}^{(0)}, \mathcal{L}_{BaB6}^{(1)}, \mathbf{t}_{BaB6}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.47)$$

を生成する. ここで,  $\mathcal{L}_{BaB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$  であり,  $\mathbf{t}_{BaB6}^{(1)} = (1, 1, 1)$  である. Bass Block 6 の畳み込み層における各種パラメータを表 1.18 にまとめて示す. Bass Block は処理をこれで終了し,  $\mathbf{Y}_{BaB6}^{(1)}$  を Bass Block の出力  $\mathbf{Y}_{BaB}$  として出力する:  $\mathbf{Y}_{BaB} = \mathbf{Y}_{BaB6}^{(1)}$ . これが, ベーストラックの特徴量であり, Band Blocks に入力される.

表 1.18: Bass Block 6 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{BaB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1}$	$\mathbf{t}_{BaB6}^{(1)} = (1, 1, 1)$

## 1.2.4 Strings Block

従来手法である BinaryMuseGAN における判別器の Strings Block は、3 個のサブブロックで構成されていた。本研究では Strings Block を Piano Block と同様に 6 個のサブブロックで構成する。それらを Strings Block 1, Strings Block 2, Strings Block 3, Strings Block 4, Strings Block 5, および Strings Block 6 と呼ぶことにする。順に説明する。

Strings Block 1 は、表 1.19 に示す 2 層の畳み込み層で構成される。第 1 畳み込み層は  $\mathbf{X}_{Strings} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Strings}, \mathcal{L}_{SB1}^{(1)}, \mathbf{t}_{SB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 7} \quad (1.48)$$

を抽出する。ここで、 $\mathcal{L}_{SB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{SB1}^{(1)} = (1, 1, 12)$  である。この段階で第 4 モードのサイズが 7 になり、オクターブごとの特徴抽出を開始する。次に、第 2 畳み込み層が  $\mathbf{Y}_{SB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 7}$  から

$$\mathbf{Y}_{SB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB1}^{(1)}, \mathcal{L}_{SB1}^{(2)}, \mathbf{t}_{SB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.49)$$

を抽出する。ここで、 $\mathcal{L}_{SB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{SB1}^{(2)} = (1, 6, 1)$  である。この段階で第 3 モードのサイズが 16 になり、微細なリズム構造の特徴抽出を開始する。以上のように、Strings Block 1 では、第 4 モードの特徴抽出を行ってから第 3 モードの特徴抽出を行う。すなわち、音高方向の特徴抽出の後に時間方向の特徴抽出を行う。

表 1.19: Strings Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{SB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{SB1}^{(1)} = (1, 1, 12)$
第 2 層	$\mathcal{L}_{SB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{SB1}^{(2)} = (1, 6, 1)$

一方、Strings Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する。すなわち、第 3 モードの時間方向特徴量を抽出してから第 4 モードである音高方向の

特徴抽出を行う。まず、第1畳み込み層が  $\mathbf{X}_{Strings} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Strings}, \mathcal{L}_{SB2}^{(1)}, \mathbf{t}_{SB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.50)$$

を抽出する。ここで、 $\mathcal{L}_{SB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{SB2}^{(1)} = (1, 6, 1)$  である。この段階で第3モードのサイズが16になり、微細なリズム構造の特徴抽出を開始する。次に、第2畳み込み層が  $\mathbf{Y}_{SB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{SB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB2}^{(1)}, \mathcal{L}_{SB2}^{(2)}, \mathbf{t}_{SB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.51)$$

を抽出する。ここで、 $\mathcal{L}_{SB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{SB2}^{(2)} = (1, 1, 12)$  である。この段階で第4モードのサイズが7になり、オクターブごとの特徴抽出を開始する。Strings Block 2 の畳み込み層における各種パラメータを表 1.20 にまとめて示す。

表 1.20: Strings Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{SB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 6 \times 1}$	$\mathbf{t}_{SB2}^{(1)} = (1, 6, 1)$
第2層	$\mathcal{L}_{SB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{SB2}^{(2)} = (1, 1, 12)$

続いて、Strings Block 3 は表 1.21 に示す3層の畳み込み層で構成される。まず、第1畳み込み層が  $\mathbf{X}_{Strings} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB3}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Strings}, \mathcal{L}_{SB3}^{(1)}, \mathbf{t}_{SB3}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.52)$$

を抽出する。ここで、畳み込みはSame畳み込みであり、 $\mathcal{L}_{SB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 3}$ 、 $\mathbf{t}_{SB3}^{(1)} = (1, 1, 1)$  である。この段階で第4モードのサイズが3のカーネルを用いることにより、隣接3半音の特徴抽出を開始する。この第1畳み込み層は、隣接3半音からの不協和音検出を目的としている。次に、第2畳み込み層は  $\mathbf{Y}_{SB3}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB3}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB3}^{(1)}, \mathcal{L}_{SB3}^{(2)}, \mathbf{t}_{SB3}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.53)$$

を抽出する。ここで、 $\mathcal{L}_{SB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり、 $\mathbf{t}_{SB3}^{(2)} = (1, 6, 1)$  である。この段階で第3モードのサイズが16になり、微細なリズム構造の特徴抽出を開始する。次に、第3畳み込み層が  $\mathbf{Y}_{SB3}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{SB3}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB3}^{(2)}, \mathcal{L}_{SB3}^{(3)}, \mathbf{t}_{SB3}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.54)$$

を抽出する。ここで、 $\mathcal{L}_{SB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり、 $\mathbf{t}_{SB3}^{(3)} = (1, 1, 12)$  である。この段階で第4モードのサイズが7になり、オクターブごとの特徴抽出を開始する。Strings Block 3 の畳み込み層における各種パラメータを表 1.21 にまとめて示す。

表 1.21: Strings Block 3 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{SB3}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3}$	$\mathbf{t}_{SB3}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{SB3}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{SB3}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{SB3}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{SB3}^{(3)} = (1, 1, 12)$

続いて, Strings Block 4 は表 1.21 に示す 3 層の畳み込み層で構成される. Strings Block 3 では第一層で隣接 3 半音から特徴を抽出したが, Strings Block 3 では隣接 2 半音から特徴抽出を行う. まず, 第 1 畳み込み層が  $\mathbf{X}_{Strings} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB4}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Strings}, \mathcal{L}_{SB4}^{(1)}, \mathbf{t}_{SB4}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 84} \quad (1.55)$$

を抽出する. ここで, 畳み込みは Same 畳み込みであり,  $\mathcal{L}_{SB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$ ,  $\mathbf{t}_{SB4}^{(1)} = (1, 1, 1)$  である. この段階で第 4 モードのサイズが 2 のカーネルを用いることにより, 隣接 2 半音の特徴抽出を開始する. この第 1 畳み込み層は, 隣接 2 半音からの不協和音検出を目的としている. 次に, 第 2 畳み込み層は  $\mathbf{Y}_{SB4}^{(1)} \in \mathbb{R}^{32 \times 4 \times 96 \times 84}$  から

$$\mathbf{Y}_{SB4}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB4}^{(1)}, \mathcal{L}_{SB4}^{(2)}, \mathbf{t}_{SB4}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 84} \quad (1.56)$$

を抽出する. ここで,  $\mathcal{L}_{SB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$  であり,  $\mathbf{t}_{SB4}^{(2)} = (1, 6, 1)$  である. この段階で第 3 モードのサイズが 16 になり, 微細なリズム構造の特徴抽出を開始する. 次に, 第 3 畳み込み層が  $\mathbf{Y}_{SB4}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 84}$  から

$$\mathbf{Y}_{SB4}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB4}^{(2)}, \mathcal{L}_{SB4}^{(3)}, \mathbf{t}_{SB4}^{(3)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.57)$$

を抽出する. ここで,  $\mathcal{L}_{SB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$  であり,  $\mathbf{t}_{SB4}^{(3)} = (1, 1, 12)$  である. この段階で第 4 モードのサイズが 7 になり, オクターブごとの特徴抽出を開始する. Strings Block 4 の畳み込み層における各種パラメータを表 1.22 にまとめて示す.

表 1.22: Strings Block 4 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層 (Same)	$\mathcal{L}_{SB4}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 2}$	$\mathbf{t}_{SB4}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{SB4}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 6 \times 1}$	$\mathbf{t}_{SB4}^{(2)} = (1, 6, 1)$
第 3 層	$\mathcal{L}_{SB4}^{(3)} \in \mathbb{R}^{32 \times 32 \times 1 \times 1 \times 12}$	$\mathbf{t}_{SB4}^{(3)} = (1, 1, 12)$

Strings Block 5 では, Strings Block 3 と Strings Block 4 の出力から特徴抽出を以下のように行う. まず, 入力された特徴  $\mathbf{Y}_{SB3}^{(3)}$ ,  $\mathbf{Y}_{SB4}^{(3)}$  を

$$\mathbf{Y}_{SB5}^{(0)} = \text{Concat}(1, \mathbf{Y}_{SB3}^{(3)}, \mathbf{Y}_{SB4}^{(3)}) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.58)$$

によって結合する.  $\mathbf{Y}_{SB3}^{(3)}$  と  $\mathbf{Y}_{SB4}^{(3)}$  の 2 テンソルを結合することにより, 二種類の不協和音検出の結果を統合している. 続いて, 畳み込み層が  $\mathbf{Y}_{SB5}^{(0)} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{SB5}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB5}^{(0)}, \mathcal{L}_{SB5}^{(1)}, \mathbf{t}_{SB5}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 7} \quad (1.59)$$

を生成する. ここで,  $\mathcal{L}_{SB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$  であり,  $\mathbf{t}_{SB5}^{(1)} = (1, 1, 1)$  である. Strings Block 5 の畳み込み層における各種パラメータを表 1.23 にまとめて示す.

表 1.23: Strings Block 5 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{SB5}^{(1)} \in \mathbb{R}^{32 \times 64 \times 1 \times 1 \times 1}$	$\mathbf{t}_{SB5}^{(1)} = (1, 1, 1)$

Strings Block 6 では, Strings Block 1, Strings Block 2, Strings Block 5 の出力から特徴抽出を以下のように行う. まず, 入力された特徴  $\mathbf{Y}_{SB1}^{(2)}$ ,  $\mathbf{Y}_{SB2}^{(2)}$ ,  $\mathbf{Y}_{SB5}^{(1)}$  を

$$\mathbf{Y}_{SB6}^{(0)} = \text{Concat}(1, \mathbf{Y}_{SB1}^{(2)}, \mathbf{Y}_{SB2}^{(2)}, \mathbf{Y}_{SB5}^{(1)}) \in \mathbb{R}^{160 \times 4 \times 16 \times 7} \quad (1.60)$$

によって結合する. 続いて, 畳み込み層が  $\mathbf{Y}_{SB6}^{(0)} \in \mathbb{R}^{160 \times 4 \times 16 \times 7}$  から

$$\mathbf{Y}_{SB6}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SB6}^{(0)}, \mathcal{L}_{SB6}^{(1)}, \mathbf{t}_{SB6}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 16 \times 7} \quad (1.61)$$

を生成する. ここで,  $\mathcal{L}_{SB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$  であり,  $\mathbf{t}_{SB6}^{(1)} = (1, 1, 1)$  である. Strings Block 6 の畳み込み層における各種パラメータを表 1.24 にまとめて示す. Strings Block は処理をこれで終了し,  $\mathbf{Y}_{SB6}^{(1)}$  を Strings Block の出力  $\mathbf{Y}_{SB}$  として出力する:  $\mathbf{Y}_{SB} = \mathbf{Y}_{SB6}^{(1)}$ . これが, スtring ストラックの特徴量であり, Band Blocks に入力される.

表 1.24: Strings Block 6 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{SB6}^{(1)} \in \mathbb{R}^{64 \times 160 \times 1 \times 1 \times 1}$	$\mathbf{t}_{SB6}^{(1)} = (1, 1, 1)$



## 1.3 Band Blocks

Band Blocks は Individual Blocks の出力を一つにまとめ、帯域ごとに特徴量を抽出する。帯域は低域 2 オクターブ、中域 2 オクターブ、高域オクターブの三種類である。まず、入力された各有音程楽器特徴量  $\mathbf{Y}_{PiB} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$ ,  $\mathbf{Y}_{GB} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$ ,  $\mathbf{Y}_{BaB} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$ ,  $\mathbf{Y}_{SB} \in \mathbb{R}^{64 \times 4 \times 16 \times 7}$  を

$$\mathbf{Y}_{BB} = \text{Concat}(1, \mathbf{Y}_{PiB}, \mathbf{Y}_{GB}, \mathbf{Y}_{BaB}, \mathbf{Y}_{SB}) \in \mathbb{R}^{256 \times 4 \times 16 \times 7} \quad (1.62)$$

によって結合する。 $\mathbf{Y}_{BB}$  は、Low-band Block, Mid-band Block, High-band Block に入力され、帯域ごとの特徴量抽出が行われる。

### 1.3.1 Low-band Block

Low-band Block では、有音程楽器特徴量  $\mathbf{Y}_{BB} \in \mathbb{R}^{256 \times 4 \times 16 \times 7}$  から低域 2 オクターブ分の特徴量を抽出する。まず、 $\mathbf{Y}_{BB}$  から

$$\mathbf{Y}_{LB}^{(0)} = \text{Slice}(\mathbf{Y}_{BB}, 4, 1, 2) \in \mathbb{R}^{256 \times 4 \times 16 \times 2} \quad (1.63)$$

を抽出する。続いて、畳み込み層が  $\mathbf{Y}_{LB}^{(0)} \in \mathbb{R}^{256 \times 4 \times 16 \times 2}$  から

$$\mathbf{Y}_{LB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{LB}^{(0)}, \mathcal{L}_{LB}^{(1)}, \mathbf{t}_{LB}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.64)$$

を生成する。ここで、 $\mathcal{L}_{LB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 2}$  であり、 $\mathbf{t}_{LB}^{(1)} = (1, 4, 2)$  である。この段階で第 3 モードが 4、第 4 モードが 1 となり、四分音符単位の低域特徴抽出を開始する。Low-band Block は処理をこれで終了し、 $\mathbf{Y}_{LB}^{(1)}$  を Low-band Block の出力  $\mathbf{Y}_{LB}$  として出力する： $\mathbf{Y}_{LB} = \mathbf{Y}_{LB}^{(1)}$ 。これが、低域の特徴量であり、Tonal Blocks に入力される。Low-band Block のパラメータをまとめて表 1.25 に示す。

表 1.25: Low-band Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{LB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 2}$	$\mathbf{t}_{LB}^{(1)} = (1, 4, 2)$

### 1.3.2 Mid-band Block

Mid-band Block では、有音程楽器特徴量  $\mathbf{Y}_{BB} \in \mathbb{R}^{256 \times 4 \times 16 \times 7}$  から中域 2 オクターブ分の特徴量を抽出する。まず、 $\mathbf{Y}_{BB}$  から

$$\mathbf{Y}_{MiB}^{(0)} = \text{Slice}(\mathbf{Y}_{BB}, 4, 3, 4) \in \mathbb{R}^{256 \times 4 \times 16 \times 2} \quad (1.65)$$

を抽出する．続いて，畳み込み層が  $\mathbf{Y}_{MiB}^{(0)} \in \mathbb{R}^{256 \times 4 \times 16 \times 2}$  から

$$\mathbf{Y}_{MiB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{MiB}^{(0)}, \mathcal{L}_{MiB}^{(1)}, \mathbf{t}_{MiB}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.66)$$

を生成する．ここで， $\mathcal{L}_{MiB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 2}$  であり， $\mathbf{t}_{MiB}^{(1)} = (1, 4, 2)$  である．この段階で第3モードのサイズが4，第4モードのサイズが1となり，四分音符単位の中域特徴抽出を開始する．Mid-band Block は処理をこれで終了し， $\mathbf{Y}_{MiB}^{(1)}$  を Mid-band Block の出力  $\mathbf{Y}_{MiB}$  として出力する： $\mathbf{Y}_{MiB} = \mathbf{Y}_{MiB}^{(1)}$ ．これが，中域の特徴量であり，Tonal Blocks に入力される．Mid-band Block のパラメータをまとめて表 1.26 に示す．

表 1.26: Mid-band Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{MiB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 2}$	$\mathbf{t}_{MiB}^{(1)} = (1, 4, 2)$

### 1.3.3 High-band Block

High-band Block では，有音程楽器特徴量  $\mathbf{Y}_{BB} \in \mathbb{R}^{256 \times 4 \times 16 \times 7}$  から高域3オクターブ分の特徴量を抽出する．まず， $\mathbf{Y}_{BB}$  から

$$\mathbf{Y}_{HB}^{(0)} = \text{Slice}(\mathbf{Y}_{BB}, 4, 5, 7) \in \mathbb{R}^{256 \times 4 \times 16 \times 3} \quad (1.67)$$

を抽出する．続いて，畳み込み層が  $\mathbf{Y}_{HB}^{(0)} \in \mathbb{R}^{256 \times 4 \times 16 \times 3}$  から

$$\mathbf{Y}_{HB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{HB}^{(0)}, \mathcal{L}_{HB}^{(1)}, \mathbf{t}_{HB}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.68)$$

を生成する．ここで， $\mathcal{L}_{HB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 3}$  であり， $\mathbf{t}_{HB}^{(1)} = (1, 4, 3)$  である．この段階で第3モードのサイズが4，第4モードのサイズが1となり，四分音符単位の高域特徴抽出を開始する．High-band Block は処理をこれで終了し， $\mathbf{Y}_{HB}^{(1)}$  を High-band Block の出力  $\mathbf{Y}_{HB}$  として出力する： $\mathbf{Y}_{HB} = \mathbf{Y}_{HB}^{(1)}$ ．これが，高域の特徴量であり，Tonal Blocks に入力される．High-band Block のパラメータをまとめて表 1.27 に示す．

表 1.27: High-band Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第1層	$\mathcal{L}_{HB}^{(1)} \in \mathbb{R}^{64 \times 256 \times 1 \times 4 \times 3}$	$\mathbf{t}_{HB}^{(1)} = (1, 4, 3)$

## 1.4 Tonal Block

Tonal Block では、帯域別特徴量  $\mathbf{Y}_{LB} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$ 、 $\mathbf{Y}_{MiB} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$ 、 $\mathbf{Y}_{HB} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から有音程楽器共通の特徴量を以下のように生成する。まず、入力された特徴  $\mathbf{Y}_{LB}$ 、 $\mathbf{Y}_{MiB}$ 、 $\mathbf{Y}_{HB}$  を

$$\mathbf{Y}_{TB}^{(0)} = \text{Concat}(1, \mathbf{Y}_{LB}, \mathbf{Y}_{MiB}, \mathbf{Y}_{HB}) \in \mathbb{R}^{192 \times 4 \times 4 \times 1} \quad (1.69)$$

によって結合する。続いて、2 段の畳み込み層によって処理を継続する。まず、第 1 畳み込み層が  $\mathbf{Y}_{TB}^{(0)} \in \mathbb{R}^{192 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{TB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{TB}^{(0)}, \mathcal{L}_{TB}^{(1)}, \mathbf{t}_{TB}^{(1)})) \in \mathbb{R}^{128 \times 4 \times 4 \times 1} \quad (1.70)$$

を生成する。ここで、 $\mathcal{L}_{TB}^{(1)} \in \mathbb{R}^{128 \times 192 \times 1 \times 1 \times 1}$  であり、 $\mathbf{t}_{TB}^{(1)} = (1, 1, 1)$  である。次に、第 2 畳み込み層が  $\mathbf{Y}_{TB}^{(1)} \in \mathbb{R}^{128 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{TB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{TB}^{(1)}, \mathcal{L}_{TB}^{(2)}, \mathbf{t}_{TB}^{(2)})) \in \mathbb{R}^{128 \times 4 \times 1 \times 1} \quad (1.71)$$

を生成する。ここで、 $\mathcal{L}_{TB}^{(2)} \in \mathbb{R}^{128 \times 128 \times 1 \times 4 \times 1}$  であり、 $\mathbf{t}_{TB}^{(2)} = (1, 4, 1)$  である。この段階で第 3 モードのサイズが 1 になり、全音符単位の有音程楽器共通特徴抽出を開始する。Tonal Block は処理をこれで終了し、 $\mathbf{Y}_{TB}^{(2)}$  を Tonal Block の出力  $\mathbf{Y}_{TB}$  として出力する： $\mathbf{Y}_{TB} = \mathbf{Y}_{TB}^{(2)}$ 。これが、有音程楽器共通特徴量であり、Merged Block 1 に入力される。Tonal Block のパラメータをまとめて表 1.28 に示す。

表 1.28: Tonal Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{TB}^{(1)} \in \mathbb{R}^{128 \times 192 \times 1 \times 1 \times 1}$	$\mathbf{t}_{TB}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{TB}^{(1)} \in \mathbb{R}^{128 \times 128 \times 1 \times 4 \times 1}$	$\mathbf{t}_{TB}^{(1)} = (1, 4, 1)$

## 1.5 Chroma Block

Chroma Block は、有音程楽器全体のクローマからハーモニー特徴量を抽出する。クローマは以下のように求める。まず、入力された特徴  $\mathbf{X}_{Piano} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$ 、 $\mathbf{X}_{Guitar} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$ 、 $\mathbf{X}_{Bass} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$ 、 $\mathbf{X}_{Strings} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  を

$$\mathbf{Y}_{Chroma}^{(0)} = \text{Concat}(1, \mathbf{X}_{Piano}, \mathbf{X}_{Guitar}, \mathbf{X}_{Bass}, \mathbf{X}_{Strings}) \in \mathbb{R}^{4 \times 4 \times 96 \times 84} \quad (1.72)$$

によって結合する。続いて、

$$\mathcal{Y}_{Chroma}^{(1)} = \text{Reshape}_{Chroma}(\mathbf{Y}_{Chroma}^{(0)}) \in \mathbb{R}^{4 \times 4 \times 8 \times 12 \times 12 \times 7} \quad (1.73)$$

を生成する．ここで， $\text{Reshape}_{\text{Chroma}}$  は入力テンソルを  $4 \times 4 \times 8 \times 12 \times 12 \times 7$  の6階テンソルに並び替える関数である．これは，入力テンソルの第3モードを8分音符単位で12分割し，第4モードをオクターブ単位で7分割する操作である．続いて， $\mathcal{Y}_{\text{Chroma}}^{(1)}$  の第1モードのインデックスを  $i$ ，第4モードのインデックスを  $j$ ，第6モードのインデックスを  $k$  とすると，以下のように

$$\mathbf{Y}_{\text{Chroma}} = \frac{1}{4 \cdot 12 \cdot 7} \sum_{i=1}^4 \sum_{j=1}^{12} \sum_{k=1}^7 \mathcal{Y}_{\text{Chroma}}^{(1)}[i, :, :, j, :, k] \in \mathbb{R}^{4 \times 8 \times 12} \quad (1.74)$$

を得る．これが有音程楽器全体のクローマである．Chroma Block は，入力のクローマ  $\mathbf{Y}_{\text{Chroma}} \in \mathbb{R}^{4 \times 8 \times 12}$  を

$$\mathbf{Y}_{\text{CB}}^{(0)} = \text{Reshape}_{\text{CB}}(\mathbf{Y}_{\text{Chroma}}) \in \mathbb{R}^{1 \times 4 \times 8 \times 12} \quad (1.75)$$

によって変形する．ここで  $\text{Reshape}_{\text{CB}}$  は入力テンソルを  $1 \times 4 \times 8 \times 12$  の4階テンソルに並び替える関数である．続いて，4段の畳み込み層によって処理を継続する．まず，第1畳み込み層が  $\mathbf{Y}_{\text{CB}}^{(0)} \in \mathbb{R}^{1 \times 4 \times 8 \times 12}$  から

$$\mathbf{Y}_{\text{CB}}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{\text{CB}}^{(0)}, \mathcal{L}_{\text{CB}}^{(1)}, \mathbf{t}_{\text{CB}}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.76)$$

を生成する．ここで， $\mathcal{L}_{\text{CB}}^{(1)} \in \mathbb{R}^{64 \times 1 \times 1 \times 1 \times 12}$  であり， $\mathbf{t}_{\text{CB}}^{(1)} = (1, 1, 12)$  である．この段階で第4モードのサイズが1になり，和音構造の特徴抽出を開始する．次に，第2畳み込み層が  $\mathbf{Y}_{\text{CB}}^{(1)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{\text{CB}}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{\text{CB}}^{(1)}, \mathcal{L}_{\text{CB}}^{(2)}, \mathbf{t}_{\text{CB}}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.77)$$

を生成する．ここで， $\mathcal{L}_{\text{CB}}^{(2)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{\text{CB}}^{(2)} = (1, 2, 1)$  である．この段階で第3モードのサイズが4になり，4分音符単位の和音進行の特徴抽出を開始する．次に，第3畳み込み層が  $\mathbf{Y}_{\text{CB}}^{(2)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{\text{CB}}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{\text{CB}}^{(2)}, \mathcal{L}_{\text{CB}}^{(3)}, \mathbf{t}_{\text{CB}}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.78)$$

を生成する．ここで， $\mathcal{L}_{\text{CB}}^{(3)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{\text{CB}}^{(3)} = (1, 2, 1)$  である．この段階で第3モードのサイズが2になり，2分音符単位の和音進行の特徴抽出を開始する．次に，第4畳み込み層が  $\mathbf{Y}_{\text{CB}}^{(3)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{\text{CB}}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{\text{CB}}^{(3)}, \mathcal{L}_{\text{CB}}^{(4)}, \mathbf{t}_{\text{CB}}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.79)$$

を生成する．ここで， $\mathcal{L}_{\text{CB}}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{\text{CB}}^{(4)} = (1, 2, 1)$  である．この段階で第3モードのサイズが1になり，全音符単位の和音進行の特徴抽出を開始する．

Chroma Block は処理をこれで終了し、 $\mathbf{Y}_{CB}^{(4)}$  を Chroma Block の出力  $\mathbf{Y}_{CB}$  として出力する： $\mathbf{Y}_{CB} = \mathbf{Y}_{CB}^{(4)}$ 。これが、有音程楽器全体のクローマから抽出したハーモニー特徴量であり、Merged Block 1 に入力される。Chroma Block のパラメータをまとめて表 1.29 に示す。

表 1.29: Chroma Block の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{CB}^{(1)} \in \mathbb{R}^{64 \times 1 \times 1 \times 12}$	$\mathbf{t}_{CB}^{(1)} = (1, 1, 12)$
第 2 層	$\mathcal{L}_{CB}^{(2)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{CB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{CB}^{(3)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{CB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{CB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{CB}^{(4)} = (1, 2, 1)$

## 1.6 Merged Block 1

Merged Block 1 は、ピッチ情報を持つ複数の特徴量を一つの特徴量に変換する。まず、入力された特徴  $\mathbf{Y}_{TB} \in \mathbb{R}^{128 \times 4 \times 1 \times 1}$ 、 $\mathbf{Y}_{CB} \in \mathbb{R}^{64 \times 4 \times 1 \times 1}$  を

$$\mathbf{Y}_{M1}^{(0)} = \text{Concat}(1, \mathbf{Y}_{TB}, \mathbf{Y}_{CB}) \in \mathbb{R}^{192 \times 4 \times 1 \times 1} \quad (1.80)$$

によって結合する。続いて、畳み込み層が  $\mathbf{Y}_{M1}^{(0)} \in \mathbb{R}^{192 \times 4 \times 1 \times 1}$  から

$$\mathbf{Y}_{MB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{M1}^{(0)}, \mathcal{L}_{MB1}^{(1)}, \mathbf{t}_{MB1}^{(1)})) \in \mathbb{R}^{128 \times 1 \times 1 \times 1} \quad (1.81)$$

を生成する。ここで、 $\mathcal{L}_{MB1}^{(1)} \in \mathbb{R}^{128 \times 192 \times 4 \times 1 \times 1}$  であり、 $\mathbf{t}_{MB1}^{(1)} = (4, 1, 1)$  である。この段階で第 2 モードのサイズが 1 になり、4 小節全体の特徴抽出を開始する。Merged Block 1 Block は処理をこれで終了し、 $\mathbf{Y}_{MB1}^{(1)}$  を Merged Block 1 の出力  $\mathbf{Y}_{MB1}$  として出力する： $\mathbf{Y}_{MB1} = \mathbf{Y}_{MB1}^{(1)}$ 。これが、ピッチ情報を持つ特徴量であり、Merged Block 3 に入力される。Merged Block 1 のパラメータをまとめて表 1.30 に示す。

表 1.30: Merged Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{MB1}^{(1)} \in \mathbb{R}^{128 \times 192 \times 4 \times 1 \times 1}$	$\mathbf{t}_{MB1}^{(1)} = (4, 1, 1)$

## 1.7 Polyphonicity Blocks

Polyphonicity Blocks は Individual Blocks と同様にトラック別の処理を行うブロックであり, Piano Polyphonicity Block, Guitar Polyphonicity Block, Bass Polyphonicity Block, Strings Polyphonicity Block の 4 ブロックからなる. これらの 4 ブロックは, 入力ピアノロールから同時発音数テンソルを作成し, 特徴抽出を行う.

### 1.7.1 Piano Polyphonicity Block

Piano Polyphonicity Block は, ピアノトラックの同時発音数テンソルから特徴抽出を行う. 生成器同様, 低域 3 オクターブを扱う Piano Low Polyphonicity Block と高域 4 オクターブを扱う Piano High Polyphonicity Block を用いて処理を行う. 順に説明する.

Piano Low Polyphonicity Block は, ピアノ低域部の同時発音数テンソル  $\mathbf{Y}_{PianoLPoly}$  から特徴抽出を行う.  $\mathbf{Y}_{PianoLPoly}$  をは  $\mathbf{X}_{Piano}$  の第 4 モードのインデックスを  $i$  とし, 以下のように求める.

$$\mathbf{Y}_{PianoLPoly} = \frac{1}{36} \sum_{i=1}^{36} \mathbf{X}_{Piano}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.82)$$

$$\mathbf{Y}_{PianoLPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{PianoLPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.83)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の 4 階テンソルに並び替える関数である. 続いて, 6 段の畳み込み層によって処理を継続する. まず, 第 1 畳み込み層が  $\mathbf{Y}_{PianoLPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PianoLPoly}, \mathcal{L}_{PiLPB}^{(1)}, \mathbf{t}_{PiLPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.84)$$

を抽出する. ここで,  $\mathcal{L}_{PiLPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり,  $\mathbf{t}_{PiLPB}^{(1)} = (1, 3, 1)$  である. この段階で 32 分音符の 3 連符のリズムの特徴抽出を開始し, 第 3 モードのサイズが 32 になる. 次に, 第 2 畳み込み層が  $\mathbf{Y}_{PiLPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiLPB}^{(1)}, \mathcal{L}_{PiLPB}^{(2)}, \mathbf{t}_{PiLPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.85)$$

を抽出する. ここで,  $\mathcal{L}_{PiLPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり,  $\mathbf{t}_{PiLPB}^{(2)} = (1, 2, 1)$  である. この段階で 32 分音符単位のリズムの特徴抽出を開始し, 第 3 モードのサイズが 16 になる. 次に, 第 3 畳み込み層が  $\mathbf{Y}_{PiLPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiLPB}^{(2)}, \mathcal{L}_{PiLPB}^{(3)}, \mathbf{t}_{PiLPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.86)$$

を抽出する．ここで、 $\mathcal{L}_{PiLPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{PiLPB}^{(3)} = (1, 2, 1)$  である．この段階で 16 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 8 になる．次に、第 4 畳み込み層が  $\mathbf{Y}_{PiLPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiLPB}^{(3)}, \mathcal{L}_{PiLPB}^{(4)}, \mathbf{t}_{PiLPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.87)$$

を抽出する．ここで、 $\mathcal{L}_{PiLPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{PiLPB}^{(4)} = (1, 2, 1)$  である．この段階で 8 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 4 になる．次に、第 5 畳み込み層が  $\mathbf{Y}_{PiLPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiLPB}^{(4)}, \mathcal{L}_{PiLPB}^{(5)}, \mathbf{t}_{PiLPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.88)$$

を抽出する．ここで、 $\mathcal{L}_{PiLPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{PiLPB}^{(5)} = (1, 2, 1)$  である．この段階で 4 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 2 になる．次に、第 6 畳み込み層が  $\mathbf{Y}_{PiLPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{PiLPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiLPB}^{(5)}, \mathcal{L}_{PiLPB}^{(6)}, \mathbf{t}_{PiLPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.89)$$

を抽出する．ここで、 $\mathcal{L}_{PiLPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{PiLPB}^{(6)} = (1, 2, 1)$  である．この段階で 2 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 1 になる．Piano Low Polyphonicity Block は処理をこれで終了する．Piano Low Polyphonicity Block のパラメータをまとめて表 1.31 に示す．

表 1.31: Piano Low Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiLPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{PiLPB}^{(1)} = (1, 3, 1)$
第 2 層	$\mathcal{L}_{PiLPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiLPB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{PiLPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiLPB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{PiLPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiLPB}^{(4)} = (1, 2, 1)$
第 5 層	$\mathcal{L}_{PiLPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiLPB}^{(5)} = (1, 2, 1)$
第 6 層	$\mathcal{L}_{PiLPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiLPB}^{(6)} = (1, 2, 1)$

Piano High Polyphonicity Block は、ピアノ高域部の同時発音数テンソル  $\mathbf{Y}_{PianoHPoly}$  から特徴抽出を行う． $\mathbf{Y}_{PianoHPoly}$  は  $\mathbf{X}_{Piano}$  の第 4 モードのインデックスを  $i$  とし、以下のように求める．

$$\mathbf{Y}_{PianoHPoly} = \frac{1}{48} \sum_{i=37}^{84} \mathbf{X}_{Piano}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.90)$$

$$\mathbf{Y}_{PianoHPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{PianoHPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.91)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の4階テンソルに並び替える関数である．続いて，6段の畳み込み層によって処理を継続する．まず，第1畳み込み層が  $\mathbf{Y}_{PianoHPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PianoHPoly}, \mathcal{L}_{PiHPB}^{(1)}, \mathbf{t}_{PiHPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.92)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(1)} = (1, 3, 1)$  である．この段階で32分音符の3連符のリズムの特徴抽出を開始し，第3モードのサイズが32になる．次に，第2畳み込み層が  $\mathbf{Y}_{PiHPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiHPB}^{(1)}, \mathcal{L}_{PiHPB}^{(2)}, \mathbf{t}_{PiHPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.93)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(2)} = (1, 2, 1)$  である．この段階で32分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが16になる．次に，第3畳み込み層が  $\mathbf{Y}_{PiHPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiHPB}^{(2)}, \mathcal{L}_{PiHPB}^{(3)}, \mathbf{t}_{PiHPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.94)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(3)} = (1, 2, 1)$  である．この段階で16分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが8になる．次に，第4畳み込み層が  $\mathbf{Y}_{PiHPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiHPB}^{(3)}, \mathcal{L}_{PiHPB}^{(4)}, \mathbf{t}_{PiHPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.95)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(4)} = (1, 2, 1)$  である．この段階で8分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが4になる．次に，第5畳み込み層が  $\mathbf{Y}_{PiHPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiHPB}^{(4)}, \mathcal{L}_{PiHPB}^{(5)}, \mathbf{t}_{PiHPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.96)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(5)} = (1, 2, 1)$  である．この段階で4分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが2になる．次に，第6畳み込み層が  $\mathbf{Y}_{PiHPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{PiHPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{PiHPB}^{(5)}, \mathcal{L}_{PiHPB}^{(6)}, \mathbf{t}_{PiHPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.97)$$

を抽出する．ここで， $\mathcal{L}_{PiHPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{PiHPB}^{(6)} = (1, 2, 1)$  である．この段階で2分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが1になる．



Piano High Polyphonicity Block は処理をこれで終了する． Piano High Polyphonicity Block のパラメータをまとめて表 1.32 に示す．

表 1.32: Piano High Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{PiHPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{PiHPB}^{(1)} = (1, 3, 1)$
第 2 層	$\mathcal{L}_{PiHPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiHPB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{PiHPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiHPB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{PiHPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiHPB}^{(4)} = (1, 2, 1)$
第 5 層	$\mathcal{L}_{PiHPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiHPB}^{(5)} = (1, 2, 1)$
第 6 層	$\mathcal{L}_{PiHPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{PiHPB}^{(6)} = (1, 2, 1)$

最後に低域  $\mathbf{Y}_{PiLPB}^{(6)}$  と高域  $\mathbf{Y}_{PiHPB}^{(6)}$  を

$$\mathbf{Y}_{PiPB} = \text{Concat}(4, \mathbf{Y}_{PiLPB}^{(6)}, \mathbf{Y}_{PiHPB}^{(6)}) \in \mathbb{R}^{128 \times 4 \times 1 \times 1} \quad (1.98)$$

によって結合する．これが，ピアノトラックの同時発音数特徴量であり，Merged Block 2 に入力される．

## 1.7.2 Guitar Polyphonicity Block

Guitar Polyphonicity Block は，ギタートrackの同時発音数テンソル  $\mathbf{Y}_{GuitarPoly}$  から特徴抽出を行う．  $\mathbf{Y}_{GuitarPoly}$  は  $\mathbf{X}_{Guitar}$  の第 4 モードのインデックスを  $i$  とし，以下のように求める．

$$\mathbf{Y}_{GuitarPoly} = \frac{1}{84} \sum_{i=1}^{84} \mathbf{X}_{Guitar}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.99)$$

$$\mathbf{Y}_{GuitarPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{GuitarPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.100)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の 4 階テンソルに並び替える関数である．続いて，6 段の畳み込み層によって処理を継続する．まず，第 1 畳み込み層が  $\mathbf{Y}_{GuitarPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GuitarPoly}, \mathcal{L}_{GPB}^{(1)}, \mathbf{t}_{GPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.101)$$

を抽出する．ここで，  $\mathcal{L}_{GPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり，  $\mathbf{t}_{GPB}^{(1)} = (1, 3, 1)$  である．この段階で 32 分音符の 3 連符のリズムの特徴抽出を開始し，第 3 モードのサイズが 32 になる．次に，第 2 畳み込み層が  $\mathbf{Y}_{GPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GPB}^{(1)}, \mathcal{L}_{GPB}^{(2)}, \mathbf{t}_{GPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.102)$$

を抽出する．ここで， $\mathcal{L}_{GPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{GPB}^{(2)} = (1, 2, 1)$  である．この段階で 32 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 16 になる．次に，第 3 畳み込み層が  $\mathbf{Y}_{GPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GPB}^{(2)}, \mathcal{L}_{GPB}^{(3)}, \mathbf{t}_{GPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.103)$$

を抽出する．ここで， $\mathcal{L}_{GPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{GPB}^{(3)} = (1, 2, 1)$  である．この段階で 16 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 8 になる．次に，第 4 畳み込み層が  $\mathbf{Y}_{GPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GPB}^{(3)}, \mathcal{L}_{GPB}^{(4)}, \mathbf{t}_{GPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.104)$$

を抽出する．ここで， $\mathcal{L}_{GPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{GPB}^{(4)} = (1, 2, 1)$  である．この段階で 8 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 4 になる．次に，第 5 畳み込み層が  $\mathbf{Y}_{GPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GPB}^{(4)}, \mathcal{L}_{GPB}^{(5)}, \mathbf{t}_{GPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.105)$$

を抽出する．ここで， $\mathcal{L}_{GPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{GPB}^{(5)} = (1, 2, 1)$  である．この段階で 4 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 2 になる．次に，第 6 畳み込み層が  $\mathbf{Y}_{GPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{GPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{GPB}^{(5)}, \mathcal{L}_{GPB}^{(6)}, \mathbf{t}_{GPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.106)$$

を抽出する．ここで， $\mathcal{L}_{GPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{GPB}^{(6)} = (1, 2, 1)$  である．この段階で 2 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 1 になる．Guitar Polyphonicity Block は処理をこれで終了し， $\mathbf{Y}_{GPB}^{(6)}$  を Guitar Polyphonicity Block の出力  $\mathbf{Y}_{GPB}$  として出力する： $\mathbf{Y}_{GPB} = \mathbf{Y}_{GPB}^{(6)}$ ．これが，ギタートラックの同時発音数特徴量であり，Merged Block 2 に入力される．Guitar Polyphonicity Block のパラメータをまとめて表 1.33 に示す．

表 1.33: Guitar Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{GPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{GPB}^{(1)} = (1, 3, 1)$
第 2 層	$\mathcal{L}_{GPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{GPB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{GPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{GPB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{GPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{GPB}^{(4)} = (1, 2, 1)$
第 5 層	$\mathcal{L}_{GPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{GPB}^{(5)} = (1, 2, 1)$
第 6 層	$\mathcal{L}_{GPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{GPB}^{(6)} = (1, 2, 1)$

### 1.7.3 Bass Polyphonicity Block

Bass Polyphonicity Block は、ベーストラックの同時発音数テンソル  $\mathbf{Y}_{BassPoly}$  から特徴抽出を行う。  $\mathbf{Y}_{BassPoly}$  は  $\mathbf{X}_{Bass}$  の第4モードのインデックスを  $i$  とし、以下のよう求める。

$$\mathbf{Y}_{BassPoly} = \frac{1}{84} \sum_{i=1}^{84} \mathbf{X}_{Bass}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.107)$$

$$\mathbf{Y}_{BassPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{BassPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.108)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の4階テンソルに並び替える関数である。続いて、6段の畳み込み層によって処理を継続する。まず、第1畳み込み層が  $\mathbf{Y}_{BassPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BassPoly}, \mathcal{L}_{BaPB}^{(1)}, \mathbf{t}_{BaPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.109)$$

を抽出する。ここで、 $\mathcal{L}_{BaPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり、 $\mathbf{t}_{BaPB}^{(1)} = (1, 3, 1)$  である。この段階で32分音符の3連符のリズムの特徴抽出を開始し、第3モードのサイズが32になる。次に、第2畳み込み層が  $\mathbf{Y}_{BaPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaPB}^{(1)}, \mathcal{L}_{BaPB}^{(2)}, \mathbf{t}_{BaPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.110)$$

を抽出する。ここで、 $\mathcal{L}_{BaPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{BaPB}^{(2)} = (1, 2, 1)$  である。この段階で32分音符単位のリズムの特徴抽出を開始し、第3モードのサイズが16になる。次に、第3畳み込み層が  $\mathbf{Y}_{BaPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaPB}^{(2)}, \mathcal{L}_{BaPB}^{(3)}, \mathbf{t}_{BaPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.111)$$

を抽出する。ここで、 $\mathcal{L}_{BaPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{BaPB}^{(3)} = (1, 2, 1)$  である。この段階で16分音符単位のリズムの特徴抽出を開始し、第3モードのサイズが8になる。次に、第4畳み込み層が  $\mathbf{Y}_{BaPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaPB}^{(3)}, \mathcal{L}_{BaPB}^{(4)}, \mathbf{t}_{BaPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.112)$$

を抽出する。ここで、 $\mathcal{L}_{BaPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{BaPB}^{(4)} = (1, 2, 1)$  である。この段階で8分音符単位のリズムの特徴抽出を開始し、第3モードのサイズが4になる。次に、第5畳み込み層が  $\mathbf{Y}_{BaPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaPB}^{(4)}, \mathcal{L}_{BaPB}^{(5)}, \mathbf{t}_{BaPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.113)$$

を抽出する。ここで、 $\mathcal{L}_{BaPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{BaPB}^{(5)} = (1, 2, 1)$  である。この段階で4分音符単位のリズムの特徴抽出を開始し、第3モードのサイズが2になる。次に、第6畳み込み層が  $\mathbf{Y}_{BaPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{BaPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{BaPB}^{(5)}, \mathcal{L}_{BaPB}^{(6)}, \mathbf{t}_{BaPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.114)$$

を抽出する．ここで， $\mathcal{L}_{BaPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{BaPB}^{(6)} = (1, 2, 1)$  である．この段階で2分音符単位のリズムの特徴抽出を開始し，第3モードのサイズが1になる．Bass Polyphonicity Block は処理をこれで終了し， $\mathbf{Y}_{BaPB}^{(6)}$  を Bass Polyphonicity Block の出力  $\mathbf{Y}_{BaPB}$  として出力する： $\mathbf{Y}_{BaPB} = \mathbf{Y}_{BaPB}^{(6)}$ ．これが，ベーストラックの同時発音数特徴量であり，Merged Block 2 に入力される．Bass Polyphonicity Block のパラメータをまとめて表 1.34 に示す．

表 1.34: Bass Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	スライドベクトル
第1層	$\mathcal{L}_{BaPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{BaPB}^{(1)} = (1, 3, 1)$
第2層	$\mathcal{L}_{BaPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{BaPB}^{(2)} = (1, 2, 1)$
第3層	$\mathcal{L}_{BaPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{BaPB}^{(3)} = (1, 2, 1)$
第4層	$\mathcal{L}_{BaPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{BaPB}^{(4)} = (1, 2, 1)$
第5層	$\mathcal{L}_{BaPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{BaPB}^{(5)} = (1, 2, 1)$
第6層	$\mathcal{L}_{BaPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{BaPB}^{(6)} = (1, 2, 1)$

## 1.7.4 Strings Polyphonicity Block

Strings Polyphonicity Block は，ストリングストラックの同時発音数テンソルから特徴抽出を行う．生成器同様，低域3オクターブと高域4オクターブを別々のブロックで生成する．それぞれ，Strings Low Block と Strings High Block と呼ぶことにする．順に説明する．

Strings Low Block は，ストリングス低域部の同時発音数テンソル  $\mathbf{Y}_{StringsLPoly}$  から特徴抽出を行う． $\mathbf{Y}_{StringsLPoly}$  は  $\mathbf{X}_{Strings}$  の第4モードのインデックスを  $i$  とし，以下のように求める．

$$\mathbf{Y}_{StringsLPoly} = \frac{1}{36} \sum_{i=1}^{36} \mathbf{X}_{Strings}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.115)$$

$$\mathbf{Y}_{StringsLPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{StringsLPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.116)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の4階テンソルに並び替える関数である．続いて，6段の畳み込み層によって処理を継続する．まず，第1畳み込み層が  $\mathbf{Y}_{StringsLPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{StringsLPoly}, \mathcal{L}_{SLPB}^{(1)}, \mathbf{t}_{SLPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.117)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(1)} = (1, 3, 1)$  である．この段階で 32 分音符の 3 連符のリズムの特徴抽出を開始し，第 3 モードのサイズが 32 になる．次に，第 2 畳み込み層が  $\mathbf{Y}_{SLPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SLPB}^{(1)}, \mathcal{L}_{SLPB}^{(2)}, \mathbf{t}_{SLPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.118)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(2)} = (1, 2, 1)$  である．この段階で 32 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 16 になる．次に，第 3 畳み込み層が  $\mathbf{Y}_{SLPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SLPB}^{(2)}, \mathcal{L}_{SLPB}^{(3)}, \mathbf{t}_{SLPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.119)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(3)} = (1, 2, 1)$  である．この段階で 16 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 8 になる．次に，第 4 畳み込み層が  $\mathbf{Y}_{SLPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SLPB}^{(3)}, \mathcal{L}_{SLPB}^{(4)}, \mathbf{t}_{SLPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.120)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(4)} = (1, 2, 1)$  である．この段階で 8 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 4 になる．次に，第 5 畳み込み層が  $\mathbf{Y}_{SLPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SLPB}^{(4)}, \mathcal{L}_{SLPB}^{(5)}, \mathbf{t}_{SLPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.121)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(5)} = (1, 2, 1)$  である．この段階で 4 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 2 になる．次に，第 6 畳み込み層が  $\mathbf{Y}_{SLPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{SLPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SLPB}^{(5)}, \mathcal{L}_{SLPB}^{(6)}, \mathbf{t}_{SLPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.122)$$

を抽出する．ここで， $\mathcal{L}_{SLPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり， $\mathbf{t}_{SLPB}^{(6)} = (1, 2, 1)$  である．この段階で 2 分音符単位のリズムの特徴抽出を開始し，第 3 モードのサイズが 1 になる．Strings Polyphonicity Block は処理をこれで終了する．Strings Low Polyphonicity Block のパラメータをまとめて表 1.35 に示す．

表 1.35: Strings Low Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{SLPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{SLPB}^{(1)} = (1, 3, 1)$
第 2 層	$\mathcal{L}_{SLPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SLPB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{SLPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SLPB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{SLPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SLPB}^{(4)} = (1, 2, 1)$
第 5 層	$\mathcal{L}_{SLPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SLPB}^{(5)} = (1, 2, 1)$
第 6 層	$\mathcal{L}_{SLPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SLPB}^{(6)} = (1, 2, 1)$

Strings High Block は、ストリングス高域部の同時発音数テンソル  $\mathbf{Y}_{StringsHPoly}$  から特徴抽出を行う。  $\mathbf{Y}_{StringsHPoly}$  は  $\mathbf{X}_{Strings}$  の第 4 モードのインデックスを  $i$  とし、以下のように求める。

$$\mathbf{Y}_{StringsHPoly} = \frac{1}{36} \sum_{i=37}^{72} \mathbf{X}_{Strings}[:, :, :, i] \in \mathbb{R}^{1 \times 4 \times 96} \quad (1.123)$$

$$\mathbf{Y}_{StringsHPoly} = \text{Reshape}_{Poly}(\mathbf{Y}_{StringsHPoly}) \in \mathbb{R}^{1 \times 4 \times 96 \times 1} \quad (1.124)$$

ここで  $\text{Reshape}_{Poly}$  は入力テンソルを  $1 \times 4 \times 96 \times 1$  の 4 階テンソルに並び替える関数である。続いて、6 段の畳み込み層によって処理を継続する。まず、第 1 畳み込み層が  $\mathbf{Y}_{StringsHPoly} \in \mathbb{R}^{1 \times 4 \times 96 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{StringsHPoly}, \mathcal{L}_{SHPB}^{(1)}, \mathbf{t}_{SHPB}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 32 \times 1} \quad (1.125)$$

を抽出する。ここで、 $\mathcal{L}_{SHPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(1)} = (1, 3, 1)$  である。この段階で 32 分音符の 3 連符のリズムの特徴抽出を開始し、第 3 モードのサイズが 32 になる。次に、第 2 畳み込み層が  $\mathbf{Y}_{SHPB}^{(1)} \in \mathbb{R}^{32 \times 4 \times 32 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SHPB}^{(1)}, \mathcal{L}_{SHPB}^{(2)}, \mathbf{t}_{SHPB}^{(2)})) \in \mathbb{R}^{32 \times 4 \times 16 \times 1} \quad (1.126)$$

を抽出する。ここで、 $\mathcal{L}_{SHPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(2)} = (1, 2, 1)$  である。この段階で 32 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 16 になる。次に、第 3 畳み込み層が  $\mathbf{Y}_{SHPB}^{(2)} \in \mathbb{R}^{32 \times 4 \times 16 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(3)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SHPB}^{(2)}, \mathcal{L}_{SHPB}^{(3)}, \mathbf{t}_{SHPB}^{(3)})) \in \mathbb{R}^{64 \times 4 \times 8 \times 1} \quad (1.127)$$

を抽出する。ここで、 $\mathcal{L}_{SHPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(3)} = (1, 2, 1)$  である。この段階で 16 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 8 になる。次に、第 4 畳み込み層が  $\mathbf{Y}_{SHPB}^{(3)} \in \mathbb{R}^{64 \times 4 \times 8 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(4)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SHPB}^{(3)}, \mathcal{L}_{SHPB}^{(4)}, \mathbf{t}_{SHPB}^{(4)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.128)$$

を抽出する．ここで、 $\mathcal{L}_{SHPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(4)} = (1, 2, 1)$  である．この段階で 8 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 4 になる．次に、第 5 畳み込み層が  $\mathbf{Y}_{SHPB}^{(4)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(5)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SHPB}^{(4)}, \mathcal{L}_{SHPB}^{(5)}, \mathbf{t}_{SHPB}^{(5)})) \in \mathbb{R}^{64 \times 4 \times 2 \times 1} \quad (1.129)$$

を抽出する．ここで、 $\mathcal{L}_{SHPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(5)} = (1, 2, 1)$  である．この段階で 4 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 2 になる．次に、第 6 畳み込み層が  $\mathbf{Y}_{SHPB}^{(5)} \in \mathbb{R}^{64 \times 4 \times 2 \times 1}$  から

$$\mathbf{Y}_{SHPB}^{(6)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{SHPB}^{(5)}, \mathcal{L}_{SHPB}^{(6)}, \mathbf{t}_{SHPB}^{(6)})) \in \mathbb{R}^{64 \times 4 \times 1 \times 1} \quad (1.130)$$

を抽出する．ここで、 $\mathcal{L}_{SHPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$  であり、 $\mathbf{t}_{SHPB}^{(6)} = (1, 2, 1)$  である．この段階で 2 分音符単位のリズムの特徴抽出を開始し、第 3 モードのサイズが 1 になる．Strings Polyphonicity Block は処理をこれで終了する．Strings High Polyphonicity Block のパラメータをまとめて表 1.36 に示す．

表 1.36: Strings High Polyphonicity Block の転置畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{SHPB}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 3 \times 1}$	$\mathbf{t}_{SHPB}^{(1)} = (1, 3, 1)$
第 2 層	$\mathcal{L}_{SHPB}^{(2)} \in \mathbb{R}^{32 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SHPB}^{(2)} = (1, 2, 1)$
第 3 層	$\mathcal{L}_{SHPB}^{(3)} \in \mathbb{R}^{64 \times 32 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SHPB}^{(3)} = (1, 2, 1)$
第 4 層	$\mathcal{L}_{SHPB}^{(4)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SHPB}^{(4)} = (1, 2, 1)$
第 5 層	$\mathcal{L}_{SHPB}^{(5)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SHPB}^{(5)} = (1, 2, 1)$
第 6 層	$\mathcal{L}_{SHPB}^{(6)} \in \mathbb{R}^{64 \times 64 \times 1 \times 2 \times 1}$	$\mathbf{t}_{SHPB}^{(6)} = (1, 2, 1)$

最後に低域  $\mathbf{Y}_{SLPB}^{(1)}$  と高域  $\mathbf{Y}_{SHPB}^{(1)}$  を

$$\mathbf{Y}_{SPB} = \text{Concat}(4, \mathbf{Y}_{SLPB}^{(1)}, \mathbf{Y}_{SHPB}^{(1)}) \in \mathbb{R}^{128 \times 4 \times 1 \times 1} \quad (1.131)$$

によって結合する．これが、ストリングストラックの同時発音数特徴量であり、Merged Block 2 に入力される．

## 1.8 Merged Block 2

Merged Block 2 では、Polyphonicity Blocks で抽出した各トラックの同時発音数特徴量を一つにまとめ、特徴抽出を行う．まず、入力された特徴  $\mathbf{Y}_{PiPB} \in \mathbb{R}^{128 \times 4 \times 1 \times 1}$ 、

$\mathbf{Y}_{GPB} \in \mathbb{R}^{64 \times 4 \times 1 \times 1}$ ,  $\mathbf{Y}_{BaPB} \in \mathbb{R}^{64 \times 4 \times 1 \times 1}$ ,  $\mathbf{Y}_{SPB} \in \mathbb{R}^{128 \times 4 \times 1 \times 1}$  を

$$\mathbf{Y}_{MB2}^{(0)} = \text{Concat}(1, \mathbf{Y}_{PiPB}, \mathbf{Y}_{GPB}, \mathbf{Y}_{BaPB}, \mathbf{Y}_{SPB}) \in \mathbb{R}^{384 \times 4 \times 1 \times 1} \quad (1.132)$$

によって結合する．続いて，2 段の畳み込み層によって処理を継続する．まず，第 1 畳み込み層が  $\mathbf{Y}_{MB2}^{(0)} \in \mathbb{R}^{384 \times 4 \times 1 \times 1}$  から

$$\mathbf{Y}_{MB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{MB2}^{(0)}, \mathcal{L}_{MB2}^{(1)}, \mathbf{t}_{MB2}^{(1)})) \in \mathbb{R}^{128 \times 4 \times 1 \times 1} \quad (1.133)$$

を生成する．ここで， $\mathcal{L}_{MB2}^{(1)} \in \mathbb{R}^{128 \times 384 \times 1 \times 1 \times 1}$  であり， $\mathbf{t}_{MB2}^{(1)} = (1, 1, 1)$  である．次に，第 2 畳み込み層が  $\mathbf{Y}_{MB2}^{(1)} \in \mathbb{R}^{128 \times 4 \times 1 \times 1}$  から

$$\mathbf{Y}_{MB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{MB2}^{(1)}, \mathcal{L}_{MB2}^{(2)}, \mathbf{t}_{MB2}^{(2)})) \in \mathbb{R}^{128 \times 1 \times 1 \times 1} \quad (1.134)$$

を生成する．ここで， $\mathcal{L}_{MB2}^{(2)} \in \mathbb{R}^{128 \times 128 \times 4 \times 1 \times 1}$  であり， $\mathbf{t}_{MB2}^{(2)} = (4, 1, 1)$  である．この段階で小節ごとのリズム特徴抽出を開始し，第 2 モードのサイズが 1 になる．Merged Block 2 は処理をこれで終了し， $\mathbf{Y}_{MB2}^{(2)}$  を Merged Block 2 の出力  $\mathbf{Y}_{MB2}$  として出力する： $\mathbf{Y}_{MB2} = \mathbf{Y}_{MB2}^{(2)}$ ．これが，有音程楽器全体の同時発音数から抽出した特徴量であり，Merged Block 3 に入力される．Merged Block 2 のパラメータをまとめて表 1.37 に示す．

表 1.37: Merged Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{MB2}^{(1)} \in \mathbb{R}^{128 \times 384 \times 1 \times 1 \times 1}$	$\mathbf{t}_{MB2}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{MB2}^{(2)} \in \mathbb{R}^{128 \times 128 \times 4 \times 1 \times 1}$	$\mathbf{t}_{MB2}^{(2)} = (4, 1, 1)$

## 1.9 Individual Percussion Block

提案手法では，打楽器トラック  $\mathbf{X}_{Percussion} \in \mathbb{R}^{1 \times 4 \times 96 \times 84}$  をドラム打楽器  $\mathbf{X}_{Drums} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  とその他の打楽器  $\mathbf{X}_{OtherPercussion} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  に分離して処理を行う．この処理を，ドラム打楽器  $\mathbf{X}_{Drums}$  を抽出する関数  $\text{Exploit}_D$  とその他の打楽器  $\mathbf{X}_{OtherPercussion}$  を抽出する関数  $\text{Exploit}_{OP}$  を用いて，

$$\mathbf{X}_{Drums} = \text{Exploit}_D(\mathbf{X}_{Percussion}) \in \mathbb{R}^{1 \times 4 \times 96 \times 15} \quad (1.135)$$

$$\mathbf{X}_{OtherPercussion} = \text{Exploit}_{OP}(\mathbf{X}_{Percussion}) \in \mathbb{R}^{1 \times 4 \times 96 \times 15} \quad (1.136)$$



のように表す．ここで,

$$\mathbf{X}_{Percussion} = (\mathbf{X}_{Pe}^1, \mathbf{X}_{Pe}^2, \dots, \mathbf{X}_{Pe}^{84}),$$

$$\mathbf{X}_{Drums} = (\mathbf{X}_{Drums}^1, \mathbf{X}_{Drums}^2, \dots, \mathbf{X}_{Drums}^{15}),$$

$$\mathbf{X}_{OtherPercussion} = (\mathbf{X}_{OtherPercussion}^1, \mathbf{X}_{OtherPercussion}^2, \dots, \mathbf{X}_{OtherPercussion}^{15})$$

であり,  $\mathbf{X}_{Pe}, \mathbf{X}_{Drums}, \mathbf{X}_{OtherPercussion}$  の上付きの添え字は音高方向のインデックスを表す．表 1.38 が示す通りに, 関数  $\text{Exploit}_D$  は  $\mathbf{X}_{Percussion}$  を  $\mathbf{X}_{Drums}$  に変換し, 関数  $\text{Exploit}_{OP}$  は  $\mathbf{X}_{Percussion}$  を  $\mathbf{X}_{OtherPercussion}$  に変換する．なお, 判別器の入力が学習データの際, 表 1.38 の変換前に複数の  $\mathbf{X}_{Pe}$  が存在する場合がある．この時は, 変換前の複数の  $\mathbf{X}_{Pe}$  の和を算出する．こうして得られた  $\mathbf{X}_{Drums}$  は Drums Block に入力され,  $\mathbf{X}_{OtherPercussion}$  は Other Percussion Block に入力される．

表 1.38: 判別器による打楽器変換

変換前	変換後	変換前	変換後
$\mathbf{X}_{Pe}^{35}, \mathbf{X}_{Pe}^{36}$	$\mathbf{X}_{Drum}^1$	$\mathbf{X}_{Pe}^{39}$	$\mathbf{X}_{OtherPercussion}^1$
$\mathbf{X}_{Pe}^{37}$	$\mathbf{X}_{Drum}^2$	$\mathbf{X}_{Pe}^{54}$	$\mathbf{X}_{OtherPercussion}^2$
$\mathbf{X}_{Pe}^{38}, \mathbf{X}_{Pe}^{40}$	$\mathbf{X}_{Drum}^3$	$\mathbf{X}_{Pe}^{58}$	$\mathbf{X}_{OtherPercussion}^3$
$\mathbf{X}_{Pe}^{41}, \mathbf{X}_{Pe}^{43}$	$\mathbf{X}_{Drum}^4$	$\mathbf{X}_{Pe}^{60}, \mathbf{X}_{Pe}^{61}$	$\mathbf{X}_{OtherPercussion}^4$
$\mathbf{X}_{Pe}^{45}, \mathbf{X}_{Pe}^{47}$	$\mathbf{X}_{Drum}^5$	$\mathbf{X}_{Pe}^{62}, \mathbf{X}_{Pe}^{63}, \mathbf{X}_{Pe}^{64}$	$\mathbf{X}_{OtherPercussion}^5$
$\mathbf{X}_{Pe}^{48}, \mathbf{X}_{Pe}^{50}$	$\mathbf{X}_{Drum}^6$	$\mathbf{X}_{Pe}^{65}, \mathbf{X}_{Pe}^{66}$	$\mathbf{X}_{OtherPercussion}^6$
$\mathbf{X}_{Pe}^{44}$	$\mathbf{X}_{Drum}^7$	$\mathbf{X}_{Pe}^{67}, \mathbf{X}_{Pe}^{68}$	$\mathbf{X}_{OtherPercussion}^7$
$\mathbf{X}_{Pe}^{42}$	$\mathbf{X}_{Drum}^8$	$\mathbf{X}_{Pe}^{69}$	$\mathbf{X}_{OtherPercussion}^8$
$\mathbf{X}_{Pe}^{46}$	$\mathbf{X}_{Drum}^9$	$\mathbf{X}_{Pe}^{70}$	$\mathbf{X}_{OtherPercussion}^9$
$\mathbf{X}_{Pe}^{49}, \mathbf{X}_{Pe}^{57}$	$\mathbf{X}_{Drum}^{10}$	$\mathbf{X}_{Pe}^{71}, \mathbf{X}_{Pe}^{72}$	$\mathbf{X}_{OtherPercussion}^{10}$
$\mathbf{X}_{Pe}^{51}, \mathbf{X}_{Pe}^{59}$	$\mathbf{X}_{Drum}^{11}$	$\mathbf{X}_{Pe}^{73}, \mathbf{X}_{Pe}^{74}$	$\mathbf{X}_{OtherPercussion}^{11}$
$\mathbf{X}_{Pe}^{52}$	$\mathbf{X}_{Drum}^{12}$	$\mathbf{X}_{Pe}^{75}$	$\mathbf{X}_{OtherPercussion}^{12}$
$\mathbf{X}_{Pe}^{53}$	$\mathbf{X}_{Drum}^{13}$	$\mathbf{X}_{Pe}^{76}, \mathbf{X}_{Pe}^{77}$	$\mathbf{X}_{OtherPercussion}^{13}$
$\mathbf{X}_{Pe}^{55}$	$\mathbf{X}_{Drum}^{14}$	$\mathbf{X}_{Pe}^{78}, \mathbf{X}_{Pe}^{79}$	$\mathbf{X}_{OtherPercussion}^{14}$
$\mathbf{X}_{Pe}^{56}$	$\mathbf{X}_{Drum}^{15}$	$\mathbf{X}_{Pe}^{80}, \mathbf{X}_{Pe}^{81}$	$\mathbf{X}_{OtherPercussion}^{15}$

### 1.9.1 Drums Block

Drums Block は3個のサブブロックで構成される．それらを Drums Block 1, Drums Block 2, および Drums Block 3 と呼ぶことにする．Drums Block 1 は, 表 1.39 に示す 2 層の畳み込み層で構成される．第 1 畳み込み層は  $\mathbf{X}_{Drums} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  から

$$\mathbf{Y}_{DB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Drums}, \mathcal{L}_{DB1}^{(1)}, \mathbf{t}_{DB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 1} \quad (1.137)$$

を抽出する．ここで,  $\mathcal{L}_{DB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 15}$  であり,  $\mathbf{t}_{DB1}^{(1)} = (1, 1, 15)$  である．この段階で第 4 モードのサイズが 1 になり, 各音色の特徴抽出を開始する．次に, 第 2 畳み込み層は  $\mathbf{Y}_{DB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 4 \times 15}$  から

$$\mathbf{Y}_{DB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{DB1}^{(1)}, \mathcal{L}_{DB1}^{(2)}, \mathbf{t}_{DB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.138)$$

を抽出する．ここで,  $\mathcal{L}_{DB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 24 \times 1}$  であり,  $\mathbf{t}_{DB1}^{(2)} = (1, 24, 1)$  である．この段階で第 3 モードのサイズが 4 になり, 四分音符単位の特徴抽出を開始する．以上のように, Drums Block 1 では, 第 4 モードの特徴抽出を行ってから第 3 モードの特徴抽出を行う．

表 1.39: Drums Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{DB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 15}$	$\mathbf{t}_{DB1}^{(1)} = (1, 1, 15)$
第 2 層	$\mathcal{L}_{DB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 24 \times 1}$	$\mathbf{t}_{DB1}^{(2)} = (1, 24, 1)$

一方, Drums Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する．すなわち, 第 3 モードの特徴量を抽出してから第 4 モードの特徴抽出を行う．第 1 畳み込み層は  $\mathbf{X}_{Drums} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  から

$$\mathbf{Y}_{DB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{Drums}, \mathcal{L}_{DB2}^{(1)}, \mathbf{t}_{DB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 4 \times 15} \quad (1.139)$$

を抽出する．ここで,  $\mathcal{L}_{DB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 24 \times 1}$  であり,  $\mathbf{t}_{DB2}^{(1)} = (1, 24, 1)$  である．この段階で第 3 モードのサイズが 4 になり, 四分音符単位の特徴抽出を開始する．次に, 第 2 畳み込み層は  $\mathbf{Y}_{DB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 4 \times 15}$  から

$$\mathbf{Y}_{DB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{DB2}^{(1)}, \mathcal{L}_{DB2}^{(2)}, \mathbf{t}_{DB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.140)$$

を抽出する．ここで,  $\mathcal{L}_{DB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 15}$  であり,  $\mathbf{t}_{DB2}^{(2)} = (1, 1, 15)$  である．この段階で第 4 モードのサイズが 1 になり, 各音色の特徴抽出を開始する．

表 1.40: Drums Block 2 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{DB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 24 \times 1}$	$\mathbf{t}_{DB2}^{(1)} = (1, 24, 1)$
第 2 層	$\mathcal{L}_{DB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 15}$	$\mathbf{t}_{DB2}^{(2)} = (1, 1, 15)$

Drums Block 3 では, Drums Block 1 と Drums Block 2 の出力から特徴抽出を以下のように行う. まず, 入力された特徴  $\mathbf{Y}_{DB1}^{(2)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$ ,  $\mathbf{Y}_{DB2}^{(2)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  を

$$\mathbf{Y}_{DB3}^{(0)} = \text{Concat}(1, \mathbf{Y}_{DB1}^{(2)}, \mathbf{Y}_{DB2}^{(2)}) \in \mathbb{R}^{128 \times 4 \times 4 \times 1} \quad (1.141)$$

によって結合する. 続いて, 2 段の畳み込み層によって処理を継続する. まず, 第 1 畳み込み層が  $\mathbf{Y}_{DB3}^{(0)} \in \mathbb{R}^{128 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{DB3}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{DB3}^{(0)}, \mathcal{L}_{DB3}^{(1)}, \mathbf{t}_{DB3}^{(1)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.142)$$

を生成する. ここで,  $\mathcal{L}_{DB3}^{(1)} \in \mathbb{R}^{64 \times 128 \times 1 \times 1 \times 1}$  であり,  $\mathbf{t}_{DB3}^{(1)} = (1, 1, 1)$  である. 次に, 第 2 畳み込み層が  $\mathbf{Y}_{DB3}^{(1)} \in \mathbb{R}^{64 \times 4 \times 4 \times 1}$  から

$$\mathbf{Y}_{DB3}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{DB3}^{(1)}, \mathcal{L}_{DB3}^{(2)}, \mathbf{t}_{DB3}^{(2)})) \in \mathbb{R}^{128 \times 4 \times 1 \times 1} \quad (1.143)$$

を生成する. ここで,  $\mathcal{L}_{DB3}^{(2)} \in \mathbb{R}^{128 \times 64 \times 1 \times 4 \times 1}$  であり,  $\mathbf{t}_{DB3}^{(2)} = (1, 4, 1)$  である. この段階で全音符単位の特徴抽出を開始し, 第 3 モードのサイズが 1 になる. Drums Block 3 は処理をこれで終了し,  $\mathbf{Y}_{DB3}^{(2)}$  を Drums Block の出力  $\mathbf{Y}_{DB}$  として出力する:  $\mathbf{Y}_{DB} = \mathbf{Y}_{DB3}^{(2)}$ . これが, ドラム打楽器から抽出した特徴量であり, Percussion Block に入力される. Drums Block 3 のパラメータをまとめて表 1.41 に示す.

表 1.41: Drums Block 3 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{DB3}^{(1)} \in \mathbb{R}^{64 \times 128 \times 1 \times 1 \times 1}$	$\mathbf{t}_{DB3}^{(1)} = (1, 1, 1)$
第 2 層	$\mathcal{L}_{DB3}^{(2)} \in \mathbb{R}^{128 \times 64 \times 1 \times 4 \times 4}$	$\mathbf{t}_{DB3}^{(2)} = (1, 4, 1)$

## 1.9.2 Other Percussion Block

Other Percussion Block も Drums Block と同様に, 3 個のサブブロックで構成する. それらを Other Percussion Block 1, Other Percussion Block 2, および Other Percussion

Block 3 と呼ぶことにする。Other Percussion Block 1 は、表 1.42 に示す 2 層の畳み込み層で構成される。第 1 畳み込み層は  $\mathbf{X}_{OtherPercussion} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  から

$$\mathbf{Y}_{OPB1}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{OtherPercussion}, \mathcal{L}_{OPB1}^{(1)}, \mathbf{t}_{OPB1}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 96 \times 1} \quad (1.144)$$

を抽出する。ここで、 $\mathcal{L}_{OPB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 15}$  であり、 $\mathbf{t}_{OPB1}^{(1)} = (1, 1, 15)$  である。この段階で第 4 モードのサイズが 1 になり、各音色の特徴抽出を開始する。次に、第 2 畳み込み層は  $\mathbf{Y}_{OPB1}^{(1)} \in \mathbb{R}^{32 \times 4 \times 4 \times 15}$  から

$$\mathbf{Y}_{OPB1}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{OPB1}^{(1)}, \mathcal{L}_{OPB1}^{(2)}, \mathbf{t}_{OPB1}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.145)$$

を抽出する。ここで、 $\mathcal{L}_{OPB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 24 \times 1}$  であり、 $\mathbf{t}_{OPB1}^{(2)} = (1, 24, 1)$  である。この段階で第 3 モードのサイズが 4 になり、四分音符単位の特徴抽出を開始する。以上のように、Other Percussion 1 では、第 4 モードの特徴抽出を行ってから第 3 モードの特徴抽出を行う。

表 1.42: Other Percussion Block 1 の畳み込み層における各種パラメータ

畳み込み層番号	フィルタカーネル	ストライドベクトル
第 1 層	$\mathcal{L}_{OPB1}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 1 \times 15}$	$\mathbf{t}_{OPB1}^{(1)} = (1, 1, 15)$
第 2 層	$\mathcal{L}_{OPB1}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 24 \times 1}$	$\mathbf{t}_{OPB1}^{(2)} = (1, 24, 1)$

一方、Other Percussion Block 2 は同じサイズの特徴量テンソルを逆の手順で抽出する。すなわち、第 3 モードの特徴量を抽出してから第 4 モードの特徴抽出を行う。第 1 畳み込み層は  $\mathbf{X}_{OtherPercussion} \in \mathbb{R}^{1 \times 4 \times 96 \times 15}$  から

$$\mathbf{Y}_{OPB2}^{(1)} = \text{LeakyReLU}(\text{Conv}(\mathbf{X}_{OtherPercussion}, \mathcal{L}_{OPB2}^{(1)}, \mathbf{t}_{OPB2}^{(1)})) \in \mathbb{R}^{32 \times 4 \times 4 \times 15} \quad (1.146)$$

を抽出する。ここで、 $\mathcal{L}_{OPB2}^{(1)} \in \mathbb{R}^{32 \times 1 \times 1 \times 24 \times 1}$  であり、 $\mathbf{t}_{OPB2}^{(1)} = (1, 24, 1)$  である。この段階で第 3 モードのサイズが 4 になり、四分音符単位の特徴抽出を開始する。次に、第 2 畳み込み層は  $\mathbf{Y}_{OPB2}^{(1)} \in \mathbb{R}^{32 \times 4 \times 4 \times 15}$  から

$$\mathbf{Y}_{OPB2}^{(2)} = \text{LeakyReLU}(\text{Conv}(\mathbf{Y}_{OPB2}^{(1)}, \mathcal{L}_{OPB2}^{(2)}, \mathbf{t}_{OPB2}^{(2)})) \in \mathbb{R}^{64 \times 4 \times 4 \times 1} \quad (1.147)$$

を抽出する。ここで、 $\mathcal{L}_{OPB2}^{(2)} \in \mathbb{R}^{64 \times 32 \times 1 \times 1 \times 15}$  であり、 $\mathbf{t}_{OPB2}^{(2)} = (1, 1, 15)$  である。この段階で第 4 モードのサイズが 1 になり、各音色の特徴抽出を開始する。