

2種類のデータセットに対する ボトルネック特徴量を用いた環境音分類の検討

松原 拓未 (岡山大学 阿部研究室)

1. 研究背景・目的

◆ 環境音分類

- 収録された環境音から音源の種類を自動的に分類
 - 車の走行音, 鳥の鳴き声, 川の流れる音など
 - Neural Network (NN) を用いた研究
 - Convolutional NN (CNN) を利用した手法 [J.Guo+2017]
 - Recurrent NN (RNN) を利用した手法 [M. Zohrer+2017]
- 同一の機器, 条件で収録されたデータを用いる

◆ ボトルネック特徴量を用いた研究

- 汎化性能の向上を目的として用いられる
 - 性別に依存しない話者識別 [S. Ranjan+2017]
 - 雑音に対して頑健な話者識別 [H. Yu+2017]

目標

1つのモデルによる収録条件などによらない環境音分類

提案方式

CNN Autoencoder から得られる
ボトルネック特徴量を用いた環境音分類

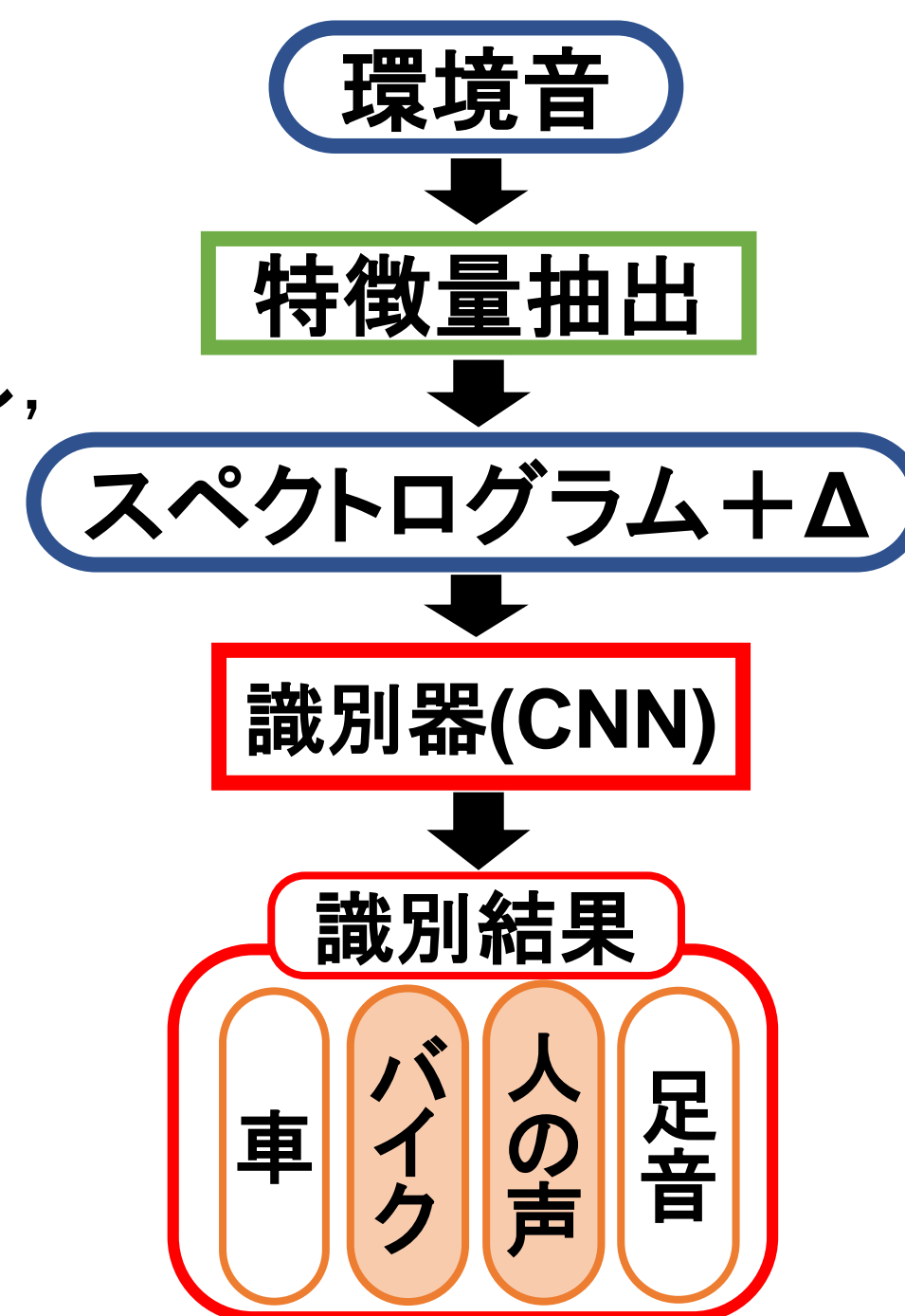
実験

収録条件が異なる2種類のデータセットの分類結果の比較

2. 環境音分類のための識別器

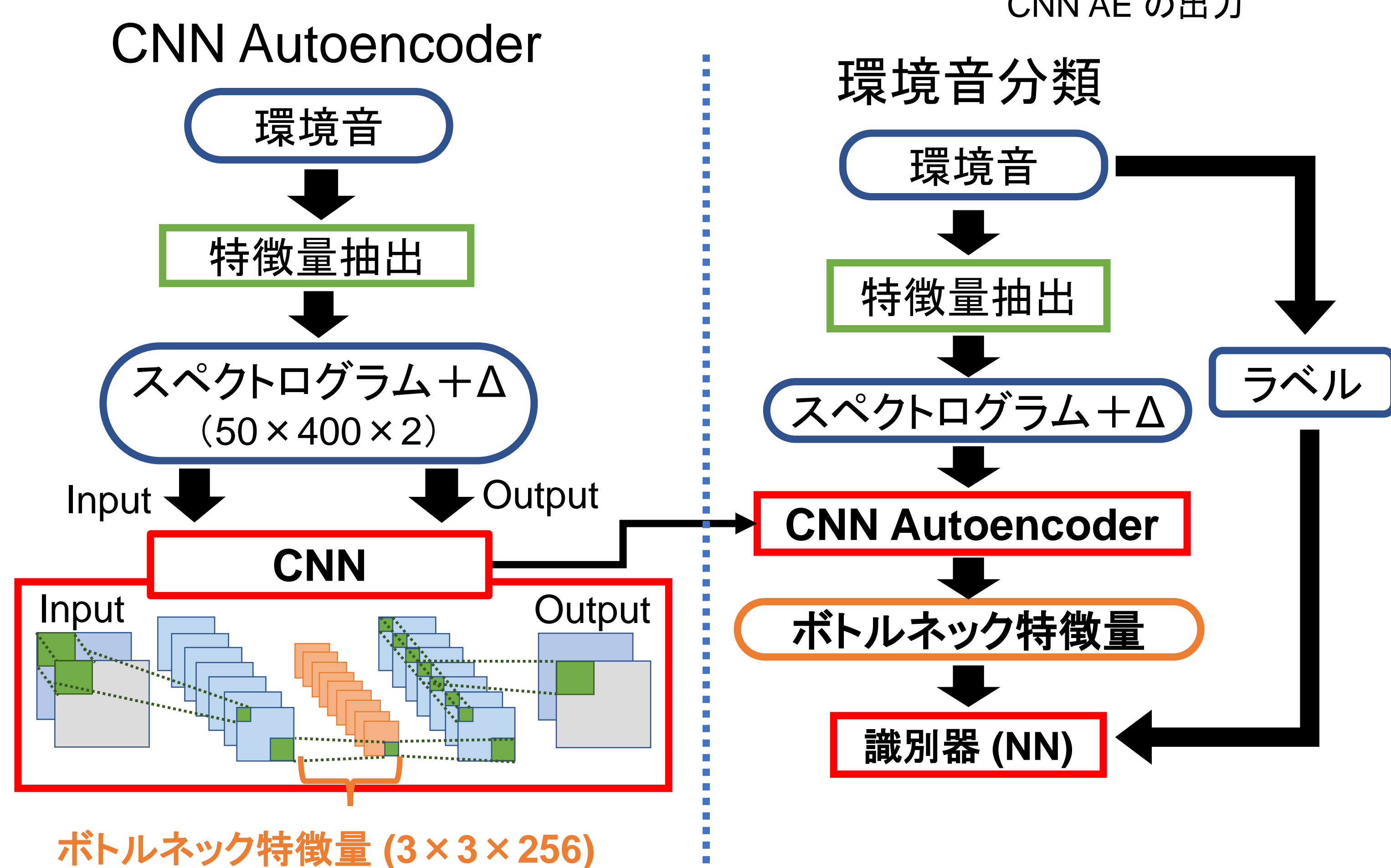
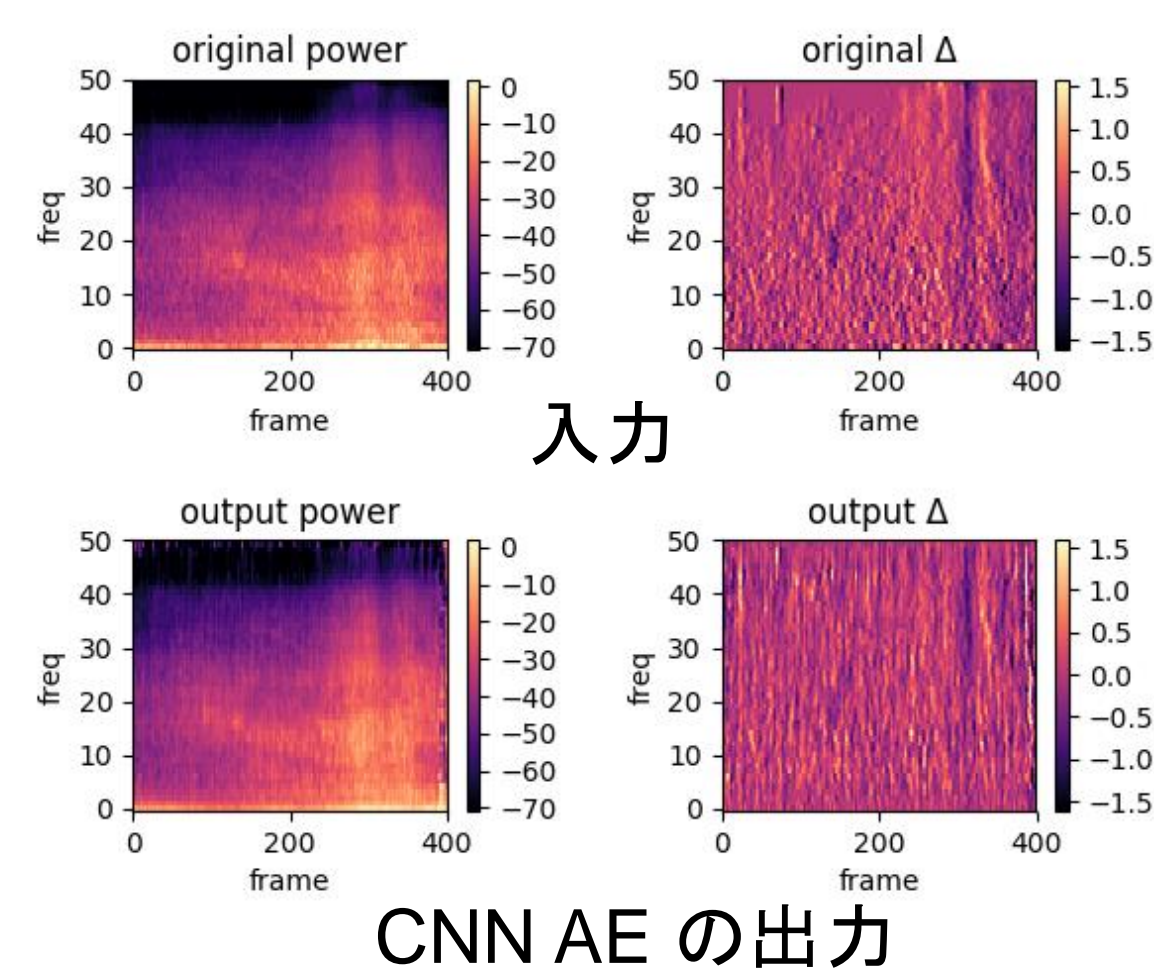
◆ 従来方式

- CNNを用いた環境音分類
 - スペクトログラムと動的特徴量(Δ)を入力し, 各ラベルの識別結果を出力
 - 識別機(CNN)
 1. 畳み込み層とプーリング層の組をいくつか通り入力を畳み込む
 2. 全結合層を通して各ラベルの識別結果を出力



◆ 提案方式

- CNN Autoencoder (AE) から抽出したボトルネック特徴量を用いた環境音分類
- ボトルネック特徴量には入力にとって重要な情報が含まれる
 - 分類をおこなうのに有用な情報
- 2種類のデータセットで学習することで両方のデータに適応



3. 使用データセット

◆ 2種類のデータセット

- TUT sound Events 2017 dataset (TUT) [A. Mesaros+ 2016]
 - 収録専用の機器で収録 (イヤーマイク, レコーダ)
 - 収録機器の操作などによるノイズが含まれていない
- 所属研究室で収録したデータセット (OU) [S. Hara+ 2017]
 - タブレット端末で収録
 - タップ音など収録機器の操作などによるノイズが含まれている

データ数 (TUT)	ラベル総数	break squeaking	car	children	large vehicle	people speaking	people walking
547	897	53	312	60	116	147	209

データ数 (OU)	ラベル総数	車	音楽	音響信号機	人	バイク	足音
1563	3131	880	438	362	566	535	350

4. 評価実験

◆ 分類性能の比較実験

- 1種類のデータセットの分類を学習 (青)
 - CNN AEの学習に用いるデータセットの違いによる性能比較
- 2種類のデータセットを用いて学習 (橙)

実験対象	分類ネットワーク		CNN AE の学習データセット	
	提案方式	従来方式	TUT	OU
AE_TUT+NN	○	×	○	×
AE_OU+NN	○	×	×	○
CNN_uni	×	○		
AE_multi+NN	○	×		
CNN_multi	×	○		

- 環境音の長さ: 10 秒
- 評価: 4-fold cross validation
- 評価指標: F-score

◆ 実験結果

- ラベル全体のF-score

F-score	AE_TUT+NN	AE_OU+NN	CNN_uni	AE_multi+NN	CNN_multi
TUT	0.737	0.737	0.715	0.731	0.693
OU	0.638	0.668	0.665	0.660	0.643

- 2種類のデータセットで学習した場合のF-scoreの変化量

	AE_multi+NN	CNN_multi
TUT	-0.006	-0.022
OU	-0.009	-0.023

5. まとめと今後の課題

◆ まとめ

- CNN Autoencoder から得られるボトルネック特徴量を用いた環境音分類手法について提案
- データセットごとに学習した場合に分類性能が向上した
- 2種類のデータセットの分類をまとめて学習した場合のF-scoreの低下を抑えられた

◆ 今後の課題

- ボトルネック特徴量を抽出するネットワークの検討
- より短い時間の分類に有効であるか調査