

Abstract: As birds can be excellent indicators to assess the environmental quality and change, many projects are now addressing identifying the bird species by their songs. Converting the bird calls into spectrograms and inputting the images into the CNN-based classification model is the common method to identify the bird species. However, the transformer-based models have been gaining attention to strike the balance of the high performance and the low computational cost in the image classification task. In this study, I implemented the two types of 48 bird species identification systems, transformer-based and CNN-based, and compared the validation performance. The transformer-based and the CNN-based model achieved 75.8%, 74.3% in the F1-score metrics, respectively. The future task is to update the bird song sounds to increase the spectrograms.

Keywords: bird classification, acoustic monitoring, deep learning, transformer, convolutional neural network

1. 背景と目的

鳥類は、生息地の質や環境状態を理解するための重要な指標となる¹⁾。そのため、鳴き声から種を判別して鳥類の生息状況を把握する試みが世界的に行われている。鳥類の種判別は、鳥の鳴き声を含む音声をスペクトログラムへ変換し、CNN (Convolutional Neural Network) モデルで画像分類タスクとして種判別を行う手法が一般的である²⁾。しかし、画像分類の分野では高い精度と低い計算コストを両立するモデルとして Transformer をベースとした画像分類モデルが注目を集めている³⁾。そこで本研究では、従前の CNN の画像分類モデルに加えて、Transformer の画像分類モデルを用いて、2種類の鳥の種判別システムを構築し、画像分類モデル間の性能を比較し、種判別の性能を評価することを目的とした。

2. 研究手法

2.1 スペクトログラムのデータセットの構築

鳥の種判別システムの構築では、2003~2005年に千葉県立中央博物館が千葉県内の各地で録音した48種類の鳥のアノテーション付きの音声ファイル7,504を使用した (wave形式, sampling rate: 22.05 kHz)。アノテーションには、鳴き声の開始時間、鳴き声の継続時間、鳥の種類の情報が含まれている。アノテーションをもとに鳴き声の音声を抽出した後、音声の長さを全て5秒で統一するために、5秒以上継続する音声はその中からランダムに5秒を抽出、5秒未満の音声は不足分を zero padding した。5秒に調整された音声を短時間フーリエ変換でスペクトログラムに変換し、周波数をメル尺度 (N mels: 128) に変換した。鳥類48種と虫の鳴き声などのそれ以外の音 (Noise) の計8,841の音響イベントのスペクトログラムのデータセットを構築した。図1はウグイスのスペクトログラムのサンプルである。時間幅は5s、周波数帯は300 Hzから10,010 Hzを示している。色の濃淡は各周波数帯の音圧レベルを示す。

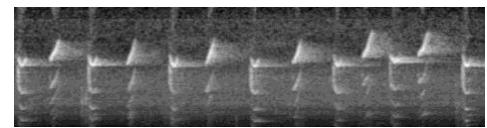


図1 スペクトログラム (ウグイス)
横軸: 時間情報, 縦軸: 周波数情報

2.2 交差検証によるモデルの訓練と種判別精度の評価

モデルには、CNNではResNet50d、TransformerではAST (Audio Spectrogram Transformer)⁴⁾を使用した。ASTは、スペクトログラムの分類問題でCNNの性能を上回るモデルである。ResNet50dはImageNet、ASTはImageNetとAudioSetでpre-trainedされたモデルを使用した。クラスごとにデータ数の80%を訓練用データ、20%を評価用データとした層化5分割交差検証を行い、評価用データに対して Overall accuracy, Precision, Recall, F1-score で交差検証精度を評価した (表1)。

表1 交差検証に用いたパラメータ

Parameters	Value	
	ResNet50d	AST
Input image size	224 x 224	128 x 512
Optimizer	Adam	
Scheduler	CosineAnnealingLR	
Loss function	BCEWithLogitsLoss	
Learning rate	8.0×10^{-5}	2.0×10^{-5}
Minimum learning rate	1.0×10^{-5}	1.0×10^{-6}
Epochs	32	50
Batch size	64	32

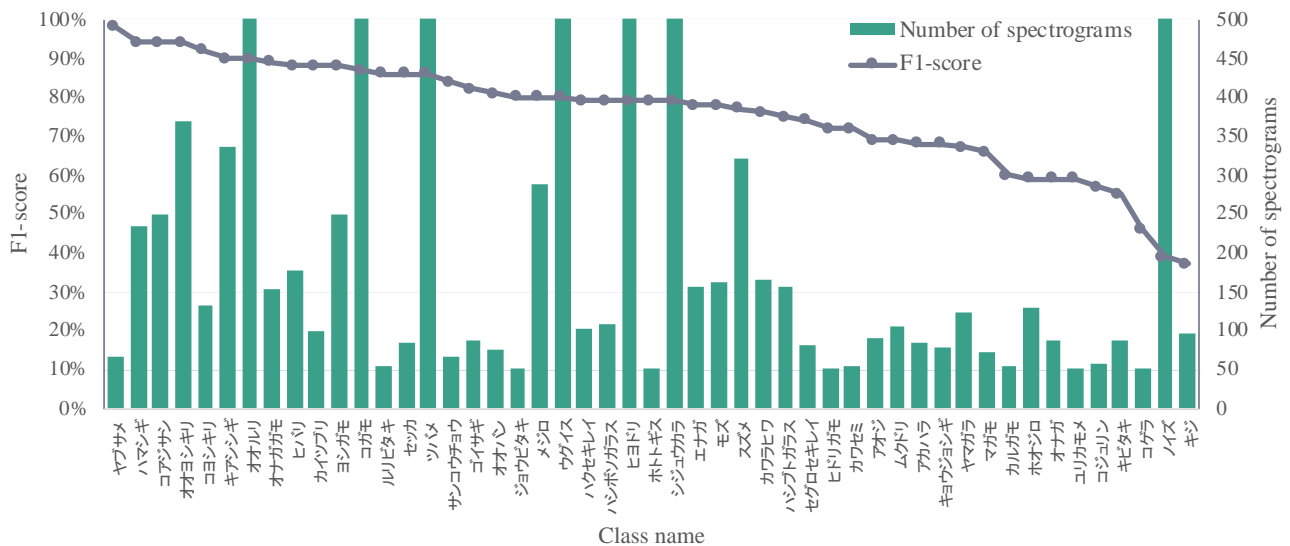


図2 ASTのクラスごとのF1-scoreとデータ数

3. 結果と考察

3.1 種判別の精度評価

図2にASTのクラスごとのF1-scoreとデータ数、表2に5分割交差検証の精度の結果を示す。まず Overall accuracy, Recall, F1-scoreではASTがResNet50dを上回る精度を示し、Transformerのモデルが鳥類の鳴き声の分類タスクでも高い判別性能を持つことがわかる。また鳥の種類別の精度では、ASTは14種の鳥で80%以上の精度を達成した。ここで8種の60%以下の精度の鳥はデータ数が少ない傾向があり、モデルが鳴き声の特徴を十分に捉えることができていない可能性がある。そのため、さらなる判別精度の向上のためにはデータ数の拡張が必要である。

3.2 種判別の精度向上に向けた考察

ASTの予測では、全スペクトログラムのうち誤判別は(1,872枚, 21.2%)となった。この誤判別は(1)鳥類の鳴き声とNoiseの誤判別(593枚, 6.7%), (2)鳥類の鳴き声間での誤判別(1279個, 14.5%)に分類することができる。(1)の誤判別の主な要因は、鳥類の鳴き声とアノテーションされた音声のSN比が低く、鳴き声の特徴量がスペクトログラムに明確に保持されず誤判別をしたと考えられる。(2)の誤判別の主な要因は、複数種が同時に鳴いていることが考えられる。本研究は1枚のスペクトログラムに1種類の音響イベントがあるタスクを想定したため、複数種の鳴き声の特徴量が1枚のスペクトログラムにある場合、モデルがどの種を予測すべきかを判断できなかったために誤判別が生じたと予想される。

4 今後の課題

今後の課題として、Noiseが含まれていても頑健な鳥の種判別を可能にするため、(1) Gaussian NoiseやPink Noise等を追加することで音声データへのData Augmentationを行うこと、(2) Noiseかどうかの判別をした後に鳥の種判別を行うという2段階の予測モデルへ変更すること、(3) 同時に複数種が鳴いた時でも判別を可能にするため、マルチラベル分類タスクに更新することが有効である可能性がある。

参考文献

- O'Connell, T.J., Jackson, L.E. and Brooks, R.P. (2000), BIRD GUILDS AS INDICATORS OF ECOLOGICAL CONDITION IN THE CENTRAL APPALACHIANS. *Ecological Applications*, 10: 1706-1721. [https://doi.org/10.1890/1051-0761\(2000\)010\[1706:BGAI0E\]2.0.CO;2](https://doi.org/10.1890/1051-0761(2000)010[1706:BGAI0E]2.0.CO;2)
- Stowell, D. (2021). Computational bioacoustics with deep learning: a review and roadmap. ArXiv, abs/2112.06725.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houlsby, N. (2021). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. ArXiv, abs/2010.11929.
- Gong, Y., Chung, Y., & Glass, J.R. (2021). AST: Audio Spectrogram Transformer. ArXiv, abs/2104.01778.

表2 5分割交差検証の精度の結果

	AST	ResNet50d
Overall accuracy	78.8 ± 0.3	75.7 ± 1.8
Precision	76.8 ± 0.6	77.2 ± 0.7
Recall	75.7 ± 0.5	72.8 ± 0.4
F1-score	75.8 ± 0.3	74.3 ± 0.5