# 環境音埋め込みベクトル系列の類似度に基づく環境音合成の自動評価

Takamichi Laboratory

岸秀(きしみのる)†,阪井瞭介†,高道慎之介†‡,金森勇介‡,岡本悠希‡(†慶應義塾大学,‡東京大学)

#### 1. はじめに

#### 1.1 研究背景

- 環境音合成(text-to-audio; TTA)は, アニメや映画に必要な音の作成に使われる[1].
  - →環境音収録,検索コストを大幅に削減.
- TTAモデルの性能を上げるには何が必要? →合成音の適切な評価指標.

# "a dog is barking." TTA model

### 1.2 従来の評価指標

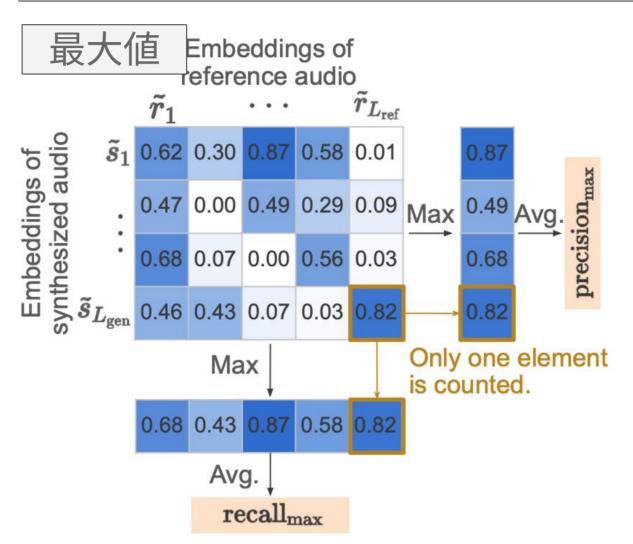
- <u>主観評価(MOS)</u>[2]:人間が1~5の5段階で評価した平均.
- OVL(overall quality):合成音の品質・自然さ
- REL (relavance) :入力テキストと合成音の関連度
- →課題:時間的・金銭的コストがかかる.
- 自動評価:合成音と{テキスト・正解音}を用いて合成音を自動で評価.
  - →課題:従来の自動評価 (MCD[3]等) は主観評価値と相関が低い[4].
  - →人間の評価に強く相関する自動評価指標を作りたい.

# 2. 提案法(AudioBERTScore):合成音,正解音それぞれの埋め込みベクトルの類似度計算によるスコア

#### 2.1 提案法概要 Text-to-audio Synthesized audio Reference audio Groundtruth pair Feature extractor Feature extractor Similarity-based High correlation Subjective score of synthesized audio

- 合成音と参照音(テキストと対の正解音) からそれぞれ埋め込みベクトル系列を獲得 →類似度計算によりスコアを算出.
- そのスコアと主観評価の相関が高くなる ことが望ましい.
- 主観評価値との相関を上げる際に重要な点
  - Feature extractorが適切に音を ベクトル系列で表現できるか.
  - 合成音と正解音ベクトルの類似度計算を どのように行うか.

# 2.2 類似度計算の詳細



p-norm

 $\tilde{s}_1$  0.62 0.30 0.87 0.58 0.01

 $\tilde{s}_{L_{\rm gen}}$  0.46 0.43 0.07 0.03

p-norm

0.07 0.00 0.56 0.03

0.57 0.26 0.50 0.43 0.41

- 類似度行列の計算
  - $\tilde{r}_*$  (行) は正解音ベクトル系列,  $\tilde{s}_*$  (列) は合成音ベクトル系列
  - 各要素は合成音と正解音ベクトルのコサイン類似度
- 2種類のノルム計算
  - 最大値:自然言語処理の既存手法 BERTScore [5] と同様. 局所フレームの対応を仮定.
- p-norm:時間遍在する環境音に対応するよう新たに定式化. p->∞で最大値ノルム.
- 3種類のスコア 0.47 0.00 0.49 0.29 0.09 p-norm 0.33 Avg.
  - o precision: 合成音の特徴を正解音がどのくらい含むか  $\operatorname{precision}_{\lambda,p} = \lambda \cdot \operatorname{precision}_{\max} + (1 - \lambda) \cdot \operatorname{precision}_{p}$

λ, pはハイパーパラメータ

- recall:正解音の特徴を合成音がどのくらい含むか  $\operatorname{recall}_{\lambda,p} = \lambda \cdot \operatorname{recall}_{\max} + (1 - \lambda) \cdot \operatorname{recall}_{p}$
- F1: precisionとrecallの調和平均

$$F1_{\lambda,p} = 2 \times \frac{\operatorname{precision}_{\lambda,p} \times \operatorname{recall}_{\lambda,p}}{\operatorname{precision}_{\lambda,p} + \operatorname{recall}_{\lambda,p}}$$

# 3. 実験的評価:提案法は人間の評価値とどのくらい相関する?

## 3.1 実験条件

- <u>評価セット</u> = PAM[5]テストセット332組
  - 参照音(正解音)
  - o テキスト 。 合成音

各合成音について

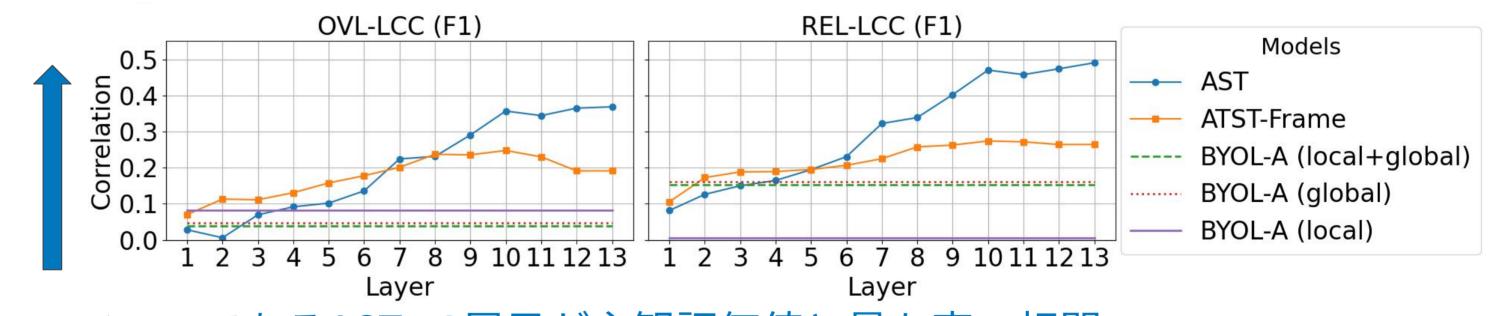
- 主観評価値:OVL, REL の5段階MOS値
- 特徴抽出器
  - BYOL-A[7], ATST-Frame[8], AST[9]
- 提案手法スコアと人間評価値の相関の計算
  - ピアソンの線形相関係数(LCC)
- スピアマンの順位相関係数(SRCC)

# 3.2 各特徴抽出器と各層の比較実験

All elements

are counted

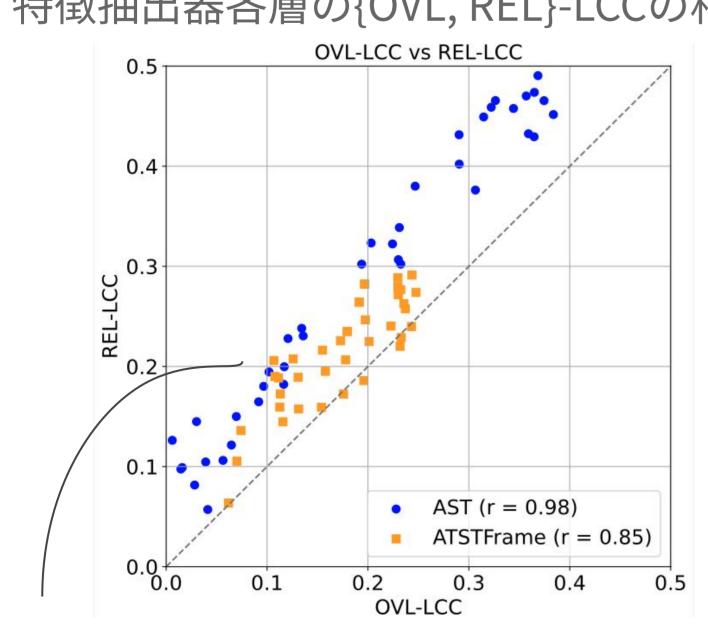
- 主観評価値と相関が高くなる特徴抽出器の探索.
- ●特徴抽出器はそれぞれ層を持つ、各層を用いて計算したスコアと主観評価値との相関。



- TransformerベースであるAST-13層目が主観評価値と最も高い相関.
- BYOL-A (CNNベース) → local層は音響特徴, global層は文脈情報をそれぞれ保持.
- AST, ATST-Frame(Transformerベース)→後半の層が文脈情報を含んでいる.

#### 3.3 OVLとRELの相関値の相関

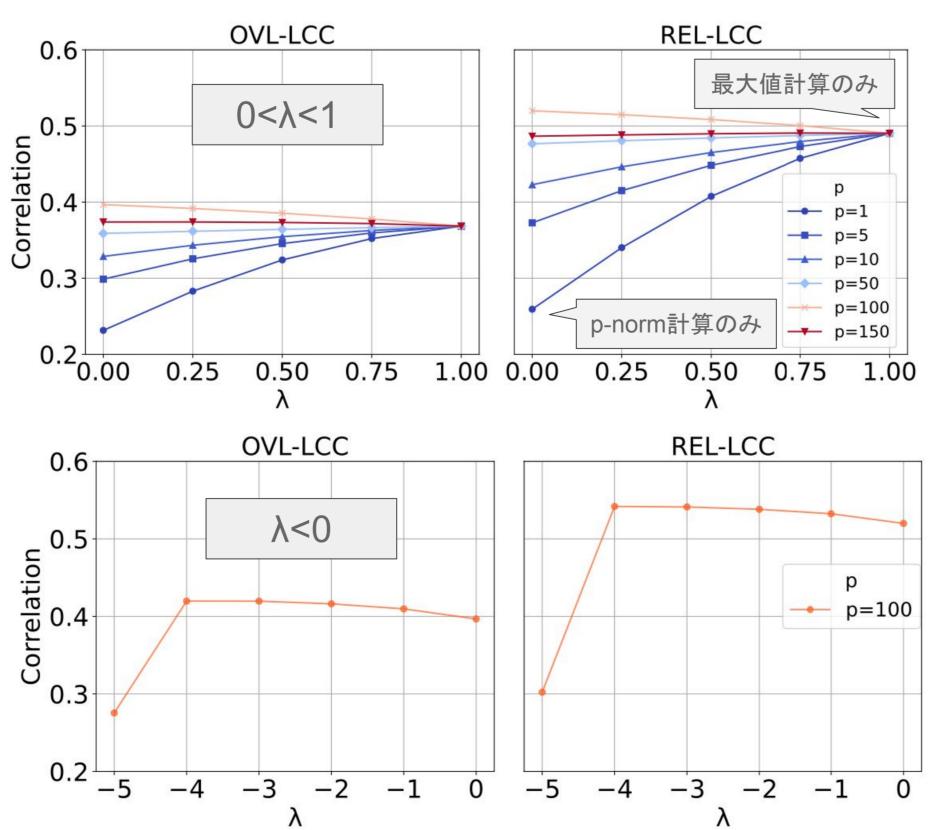
● 特徴抽出器各層の{OVL, REL}-LCCの相関.



- 全体的にOVL-LCC < REL-LCC.
  - →音質より関連度の定量化に優れる.
- 相関値の相関 > 0.8
  - →片方の性能向上でもう一方も向上

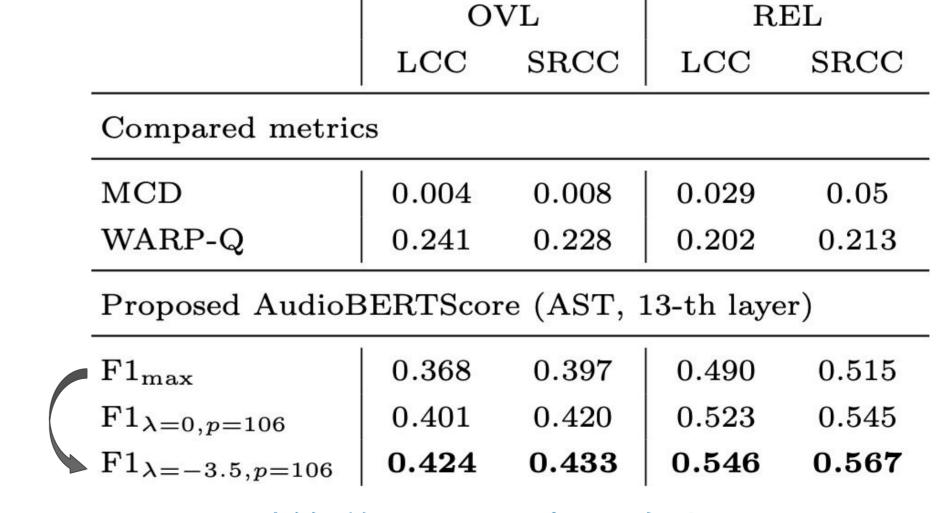
3.4 p, λの値の比較実験

• 主観評価値との相関が高くなるp, λの探索. ○ p=106, λ=-3.5で最も高い相関.



# 3.5 他の自動評価指標との比較

- 参照音を用いる従来指標MCD・WARP-Q[10] よりも高い相関.
  - →環境音の自動評価に強く貢献.



p-norm計算導入により相関向上.

Other metrics				
PAM	0.595	0.604	0.529	0.556
CLAPScore	0.337	0.323	0.487	0.475

#### Reference

- [1] T. Marrinan, arXiv, 2024
- [2] Y. Okamoto, "APSIPA ASC, 2022. [3] R. Kubichek, PACRIM, 1933
- [5] T. Zhang, ICLR, 2020 [6] T. Saeki, Interspeech, 2024

[4] 高野, 日本音響学会, 2025

- [7] D. Niizumi, TASLP, 2023
- [8] Y. Gong, Interspeech, 2021
- [9] X. Li, arXiv, 2023