# Improving Natural-Language-Based Audio Retrieval with Transfer Learning and Audio & Text Augmentations

Paul Primus[1], and Gerhard Widmer[1,2]

[1]Institute of Computational Perception, [2]LIT Artificial Intelligence Lab

JMU
JOHANNES KEPLER
UNIVERSITY LINZ

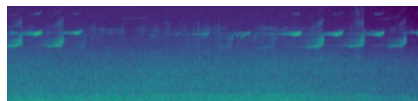Institute of
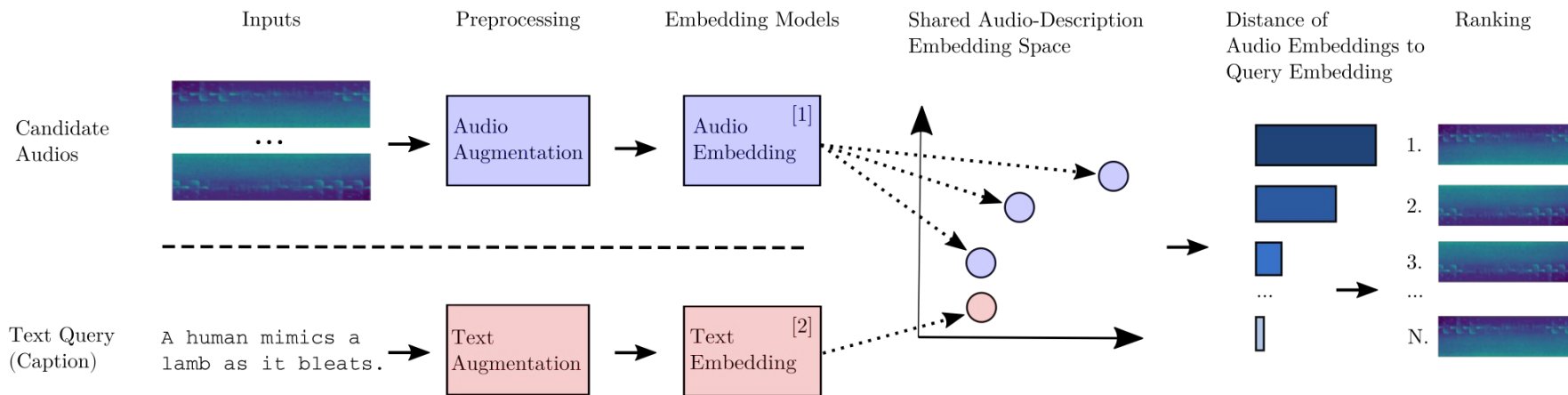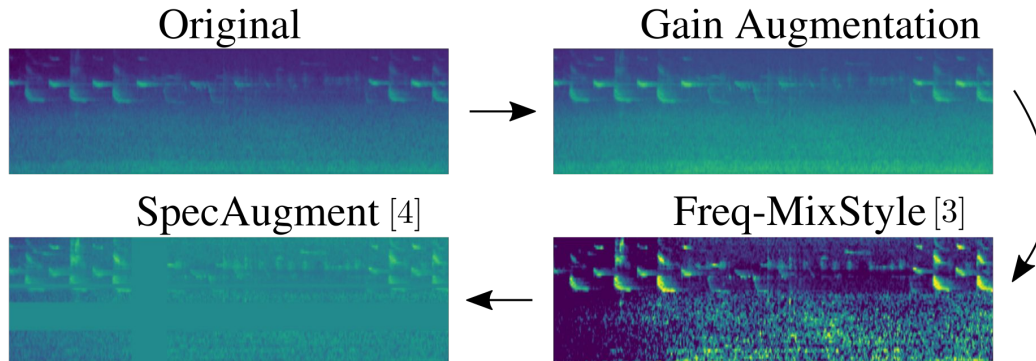Computational
Perception

# Natural-Language-Based Audio Retrieval

# Our Retrieval Framework

[1] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "PANNs: Large-scale pretrained audio neural networks for audio pattern recognition," IEEE ACM Trans. Audio Speech Lang. Process., 2020.

[2] Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pretraining of deep bidirectional transformers for language understanding," in Proc. of the North American Ch. of the Ass. for Computational Linguistics: Human Language Technologies, NAACL-HLT, 2019.

# Audio Augmentations



Original

Gain Augmentation

SpecAugment [4]

Freq-MixStyle [3]

[3]    B. Kim, S. Yang, J. Kim, H. Park, J. Lee, and S. Chang, "Domain generalization with relaxed instance frequency-wise normalization for multi-device acoustic scene classification," in 23rd Annual Conference of the International Speech Communication Association, Interspeech, 2022.

[4]    D. S. Park, W. Chan, Y. Zhang, C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," in 20th Annual Conf. of the Int. Speech Communication Association, Interspeech, 2019.

# Text Augmentations

| Augmentation | Caption |
|---|---|
| Original | `The rain pours down.` |
| Back Translation [5] | `It rains cats and dogs.` |
| Insert | `It tree rains cats and dogs.` |
| Delete | `It rains cats and dogs.` |
| Swap | `It and cats rains dogs.` |
| Synonym | `It drizzles cats and dogs.` |

[6]

[5] Sennrich, B. Haddow, and A. Birch, "Improving neural machine translation models with monolingual data," in Proc. of the 54th Annual Meeting of the Association for Computational Linguistics, ACL, 2016.

[6] W. Wei and K. Zou, "EDA: easy data augmentation techniques for boosting performance on text classification tasks," in Proc. of the 2019 Conf. on Empirical Methods in Natural Language Processing and the 9th Int. Joint Conf. on Natural Language Processing, EMNLP-IJCNLP, K. Inui, J. Jiang, V. Ng, and X. Wan, Eds., 2019

# Results

|  | R@1 | R@5 | R@10 | mAP@10 |
|---|---|---|---|---|
| DCASE baseline | 3.50 | 11.50 | 19.50 | $7.50 \pm 0.00$ |
| baseline | 6.63 | 20.06 | 31.52 | $12.53 \pm 0.08$ |
| + AudioSet pretraining | 13.18 | 35.30 | 48.61 | $22.80 \pm 0.29$ |
| + augmentations | 14.50 | 37.24 | 51.04 | $24.27 \pm 0.19$ |
| + AudioCaps pretraining | 14.34 | 38.12 | 52.04 | $24.57 \pm 0.15$ |

# Ablation Study Results

|                  | R@1   | R@5   | R@10  | mAP@10           |
| ---------------- | ----- | ----- | ----- | ---------------- |
| SMBO             | 14.50 | 37.24 | 51.04 | $24.27 \pm 0.19$ |
| no audio aug     | 13.88 | 36.94 | 51.06 | $23.74 \pm 0.16$ |
| no text aug      | 13.12 | 35.77 | 49.25 | $22.91 \pm 0.08$ |
| no SpeAugment    | 13.50 | 36.60 | 50.91 | $23.53 \pm 0.20$ |
| no FreqMixStyle  | 13.61 | 36.91 | 50.69 | $23.62 \pm 0.14$ |
| no Gain Augment  | 14.84 | 37.81 | 50.95 | $24.05 \pm 0.26$ |
| no BT            | 13.61 | 36.38 | 50.00 | $23.43 \pm 0.18$ |
| no EDA           | 13.33 | 37.02 | 49.94 | $23.27 \pm 0.03$ |

# Poster & Paper

# Contact

**email to paul.primus@jku.at**