# EDAudio:
# Easy Data Augmentation Techniques for Audio Classification
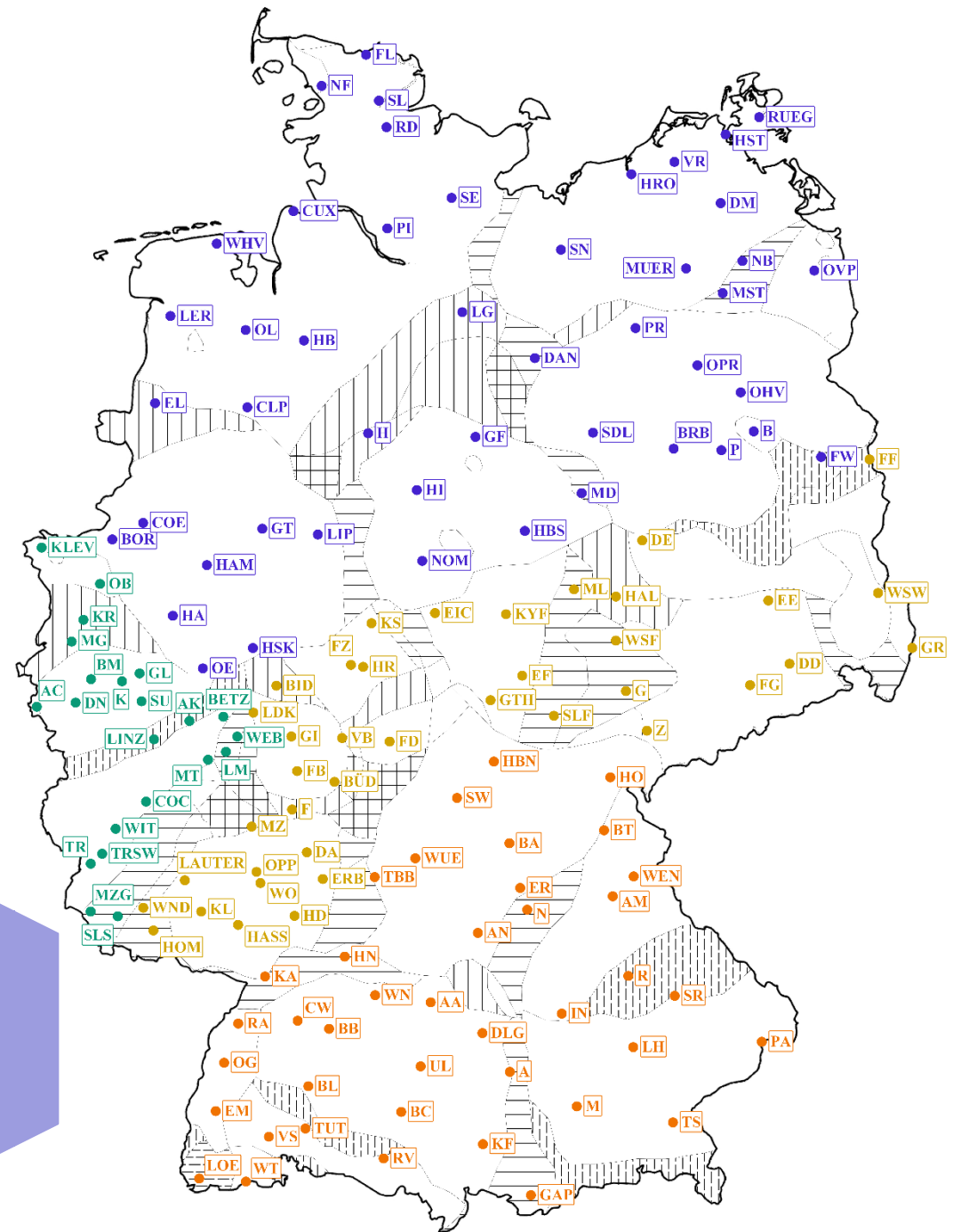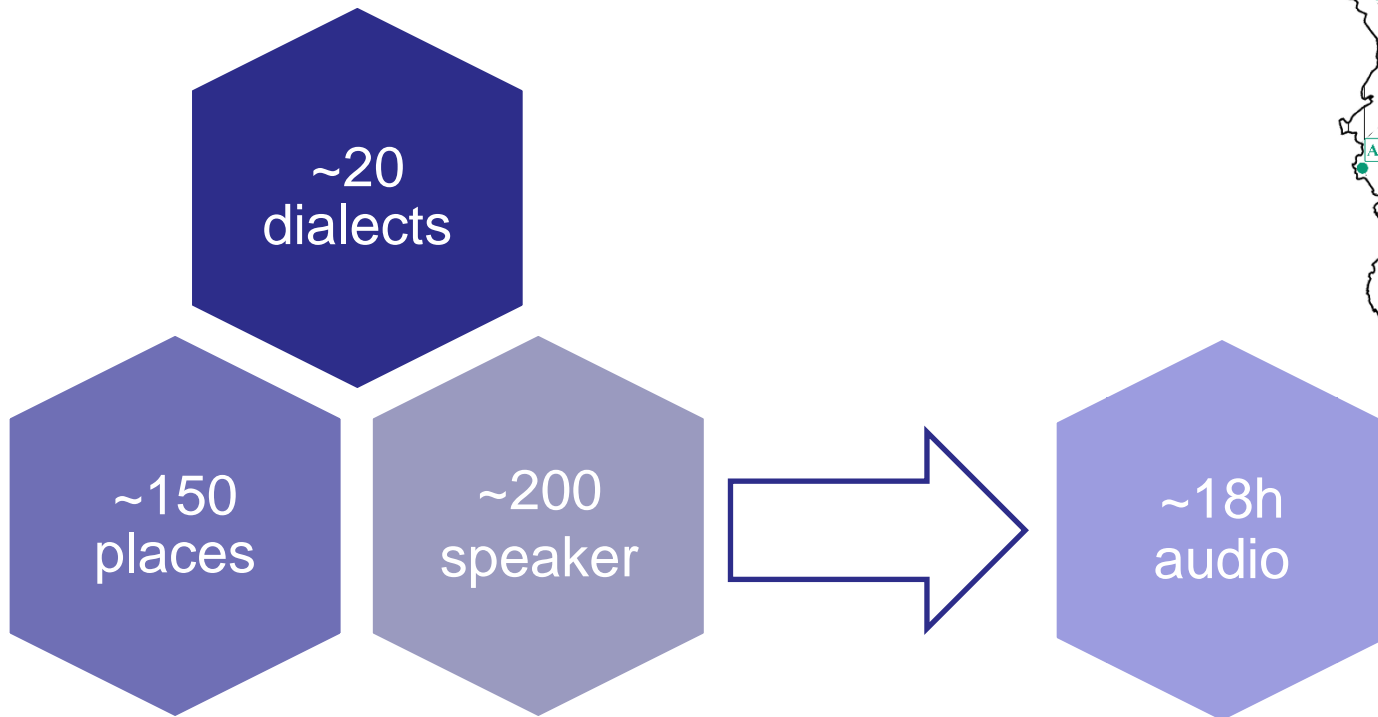
Alfred Lameli

Lea Fischbach

Caroline Kleen

Akbar Karimi

Lucie Flek

10th July 2024, ICLaVE - Vienna

GEFÖRDERT VOM

Bundesministerium für Bildung und Forschung

Finanziert von der Europäischen Union
NextGenerationEU

REDE
regionalsprache.de

Akademie der Wissenschaften und der Literatur | Mainz

Philipps Universität Marburg

Forschungszentrum Deutscher Sprachatlas
DSA

# Overview

- Recordings from Regionalsprache.de (2008-2012)
- Used Speaker are male and +65years
- Used Recordings contain dialectal „Wenker phrases"[1]



**~20 dialects**

**~150 places**

**~200 speaker**

**~18h audio**

[1]https://www.uni-marburg.de/en/fb09/dsa/research-documentation-center/wenkersaetze

# Why Augmentation?

- Increase in Dataset Size
- Regularization
- Balancing Classes

Improvement in Model Performance

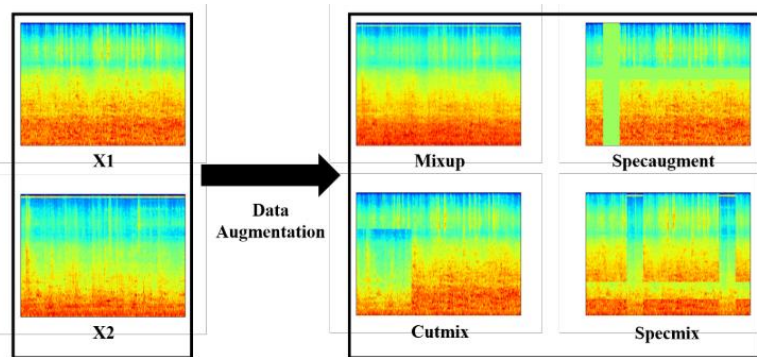# Experiment Motivation

## Spectrogram Data Augmentation (SDA)



Figure 1: *Overview of the data augmentation methods:Mixup, Cutmix, Specaugment and our Specmix.*

G. Kim, D. K. Han, and H. Ko, "Specmix: A mixed sample data augmentation method for training withtime-frequency domain features," arXiv preprint arXiv:2108.03020, 2021

D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "Specaugment: A simple data augmentation method for automatic speech recognition," arXiv preprint arXiv:1904.08779, 2019

## Environmental Sound Classification

J. Salamon and J. P. Bello, "Deep convolutional neural networks and data augmentation for environmental sound classification," IEEE Signal processing letters, vol. 24, no. 3, pp. 279–283, 2017.

- speed modification yields the most significant improvement
- noise addition contributes the least

Not true for dialect classification

## Automatic Speech Recognition (ASR)

T. Fukuda, R. Fernandez, A. Rosenberg, S. Thomas, B. Ram- abhadran, A. Sorin, and G. Kurata, "Data augmentation im- proves recognition of foreign accented speech." in Interspeech, no. September, 2018, pp. 2409–2413.

- Pitch Shift only method leading to improvement across all classes
- BN in charge of the least improvement

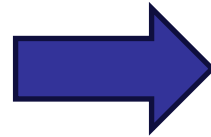Not true for dialect classification

# Experiment Setup

- Weighted f1-Score
- Cut Audio Files into 10-second Segments
- Fixed Speaker for training/validation/testing
- Run Model 50 times → get mean Score
- Starting with weighted f1_Score of 0.221

# Experiment Motivation

**Text Classification**

J. Wei and K. Zou, "Eda: Easy data augmentation techniques for boosting performance on text classification tasks," arXiv preprint arXiv:1901.11196, 2019

| Operation | Sentence |
|---|---|
| None | A sad, superior human comedy played out on the back roads of life. |
| SR<br>Synonym Replacement | A *lamentable*, superior human comedy played out on the *backward* road of life. |
| RI<br>Random Insertion | A sad, superior human comedy played out on *funniness* the back roads of life. |
| RS<br>Random Swap | A sad, superior human comedy played out on *roads* back *the* of life. |
| RD<br>Random Deletion | A sad, superior human out on the roads of life. |

Augmentation

Segment Removal

Time Masking

Segment Swap

Speaker Insertion

Shifting Pitch

Time Stretching

Volume Confusion

Speed Confusion

Time Reversing

Background Noise

Frequency Masking

Frequency Swapping

Frequency Insertion

14

# Experimental Setup



Original Segment → Starting Points for Intervals → Original Segment with marked Intervals → Function → Augmented Segment

# Experimental Setup

Praat

P. Boersma and D. Weenink, "Praat: doing phonetics by computer [Computer program]," Version 6.1.38, retrieved 2 January 2021 http://www.praat.org/, 2021.

$$n_{aug} = (\propto * l_{audio})/l_{aug}$$

$\propto$

| $l_{aug}$ | 0.1 | 0.3 | 0.5 | 1.0 |
|---|---|---|---|---|
| 0.3 | 3 | 10 | 16 | 33 |
| 1 | 1 | 3 | 5 | 10 |
| 4 | - | - | 1 | 2 |
| 5 | - | - | 1 | 2 |
| 10 | - | - | - | 1 |

$$n_{augFiles} = \{1,2,4,6\}$$

```python
def generate_intervals(length, times, total_len):
    result = []
    # Ensure there's enough space for intervals
    if times * length > total_len:
        raise ValueError("Not enough space for intervals in the given range.")

    # Generate 'times' random interval starting points
    end = 0
    for i in range(times):
        old_end = end
        start_tmp = random.randint(0, total_len - ((times-i) * length))
        start = start_tmp + old_end
        end = start + length
        result.append(start)
        # Adjust starting point for the next interval to avoid overlap
        total_len -= start_tmp + length
    return result
```
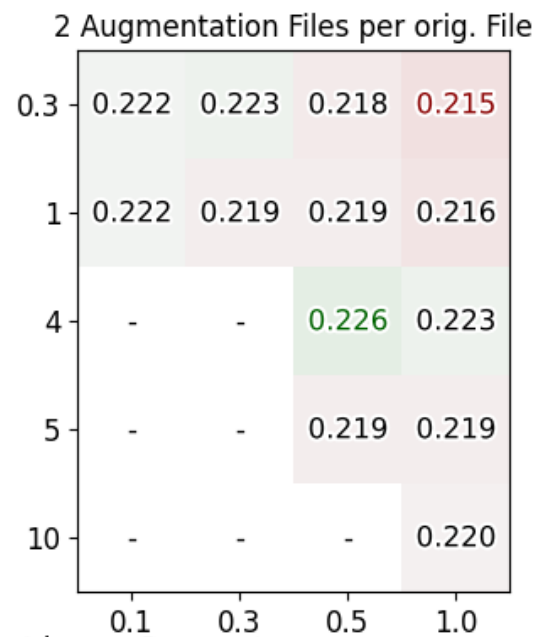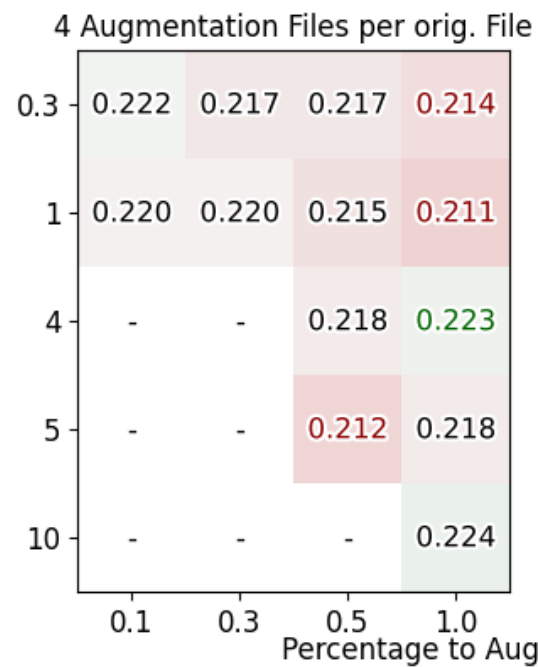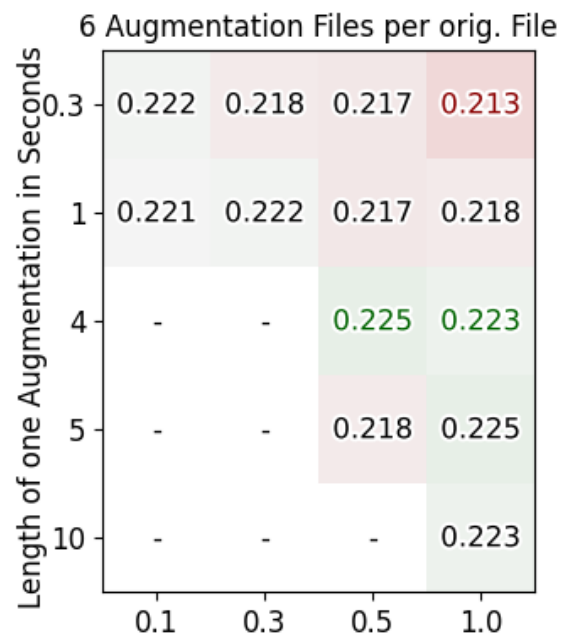
# EDA

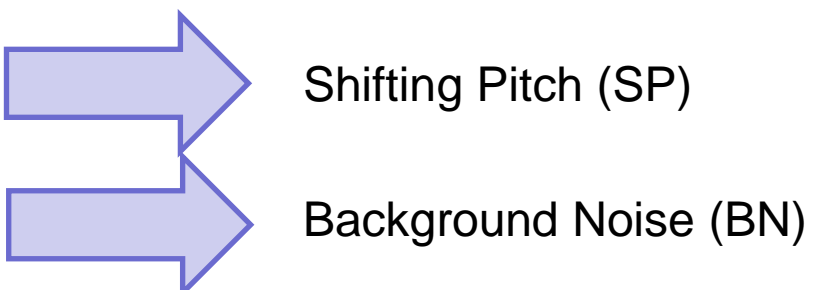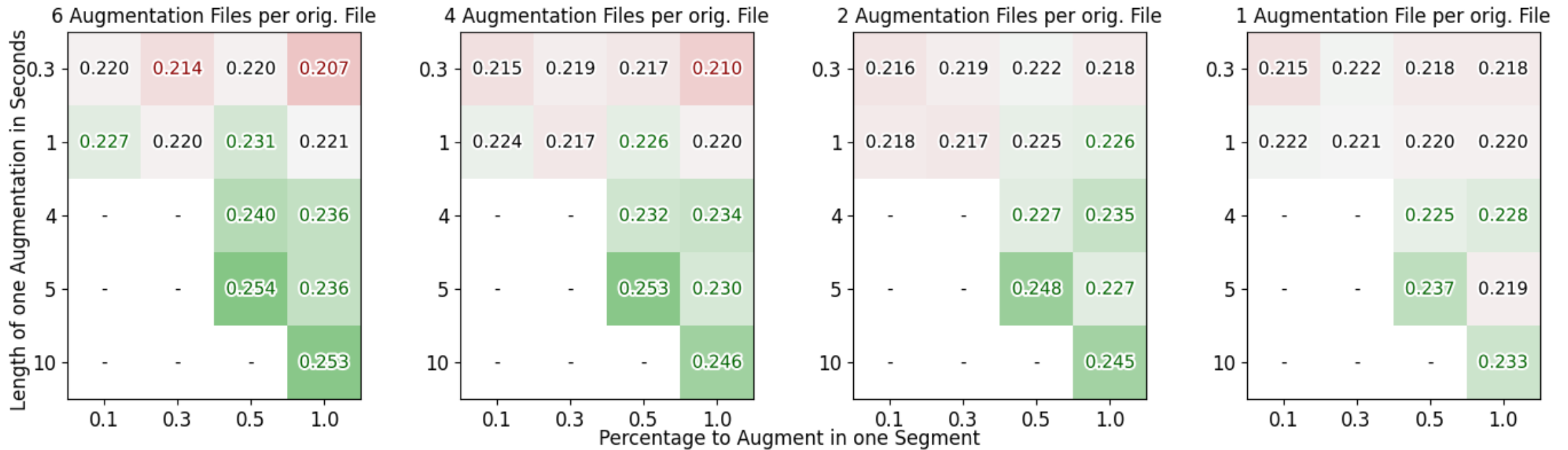| Operation | Sentence |
|---|---|
| None | A sad, superior human comedy played out on the back roads of life. |
| SR<br>Synonym Replacement | A *lamentable*, superior human comedy played out on the *backward* road of life. |
| RI<br>Random Insertion | A sad, superior human comedy played out on *funniness* the back roads of life. |
| RS<br>Random Swap | A sad, superior human comedy played out on *roads* back *the* of life. |
| RD<br>Random Deletion | A sad, superior human out on the roads of life. |

➡ Shifting Pitch (SP)

# Shifting Pitch

# Shifting Pitch



- Optimal: 2 Files per original File, 50% augmentation rate, length of 4 seconds each
- 0.5% enhancement compared to without augmentation

# EDA

| Operation | Sentence |
|---|---|
| None | A sad, superior human comedy played out on the back roads of life. |
| SR<br>Synonym Replacement | A *lamentable*, superior human comedy played out on the *backward* road of life. |
| RI<br>Random Insertion | A sad, superior human comedy played out on *funniness* the back roads of life. |
| RS<br>Random Swap | A sad, superior human comedy played out on *roads* back *the* of life. |
| RD<br>Random Deletion | A sad, superior human out on the roads of life. |

Shifting Pitch (SP)
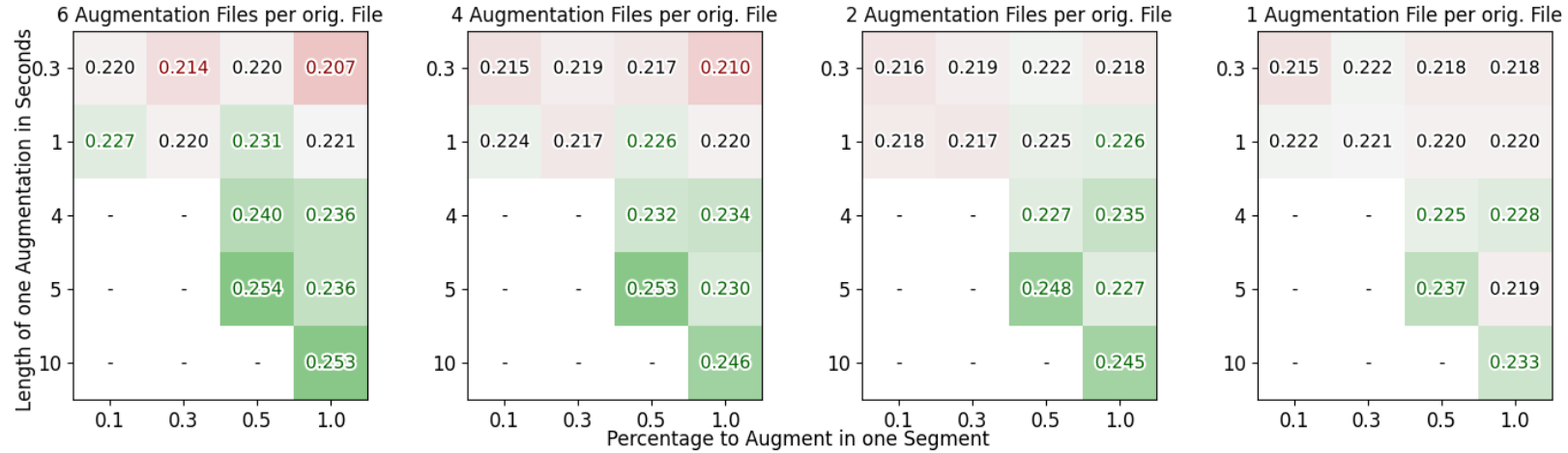
Background Noise (BN)

# Background Noise

## MUSAN

D. Snyder, G. Chen, and D. Povey, "Musan: A music, speech, and noise corpus," arXiv preprint arXiv:1510.08484, 2015.

- 929 noise files
- Total duration ~6h
- Technical noises such as
  – Dialtones
  – Fax machine
- Ambient sounds such as
  – Car idling
  – Thunder/wind/rain
  – Paper rustling
  – Animal noises

- Random part of random noise file
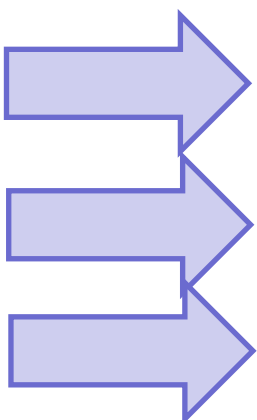- Scale noise that resulting SNRdB is in [0,30]

# Background Noise

# Background Noise



- inserting background noise works best when only one noise sound is inserted
- Optimal: 6 Files per original File, 50% augmentation rate, length of 5 seconds each
- 3.3% enhancement compared to without augmentation
- worse with shorter augmentation length
  - Because of general shorter length or the chosen noise sounds?
  - Test again for 6 Files per original File, 100% augmentation rate, length of 0.3 seconds each
  - Only use files from MULAN with >=5seconds
    - → significant better result, but still significant worse than without augmentation
    - → important to use the right noise file

# EDA

| Operation | Sentence |
|---|---|
| None | A sad, superior human comedy played out on the back roads of life. |
| SR<br>Synonym Replacement | A *lamentable*, superior human comedy played out on the *backward* road of life. |
| RI<br>Random Insertion | A sad, superior human comedy played out on *funniness* the back roads of life. |
| RS<br>Random Swap | A sad, superior human comedy played out on *roads* back *the* of life. |
| RD<br>Random Deletion | A sad, superior human out on the roads of life. |

Shifting Pitch (SP)

Background Noise (BN)

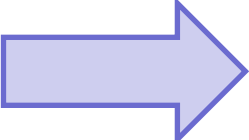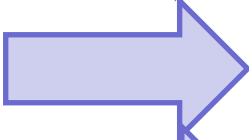Segment Swap (SeS)

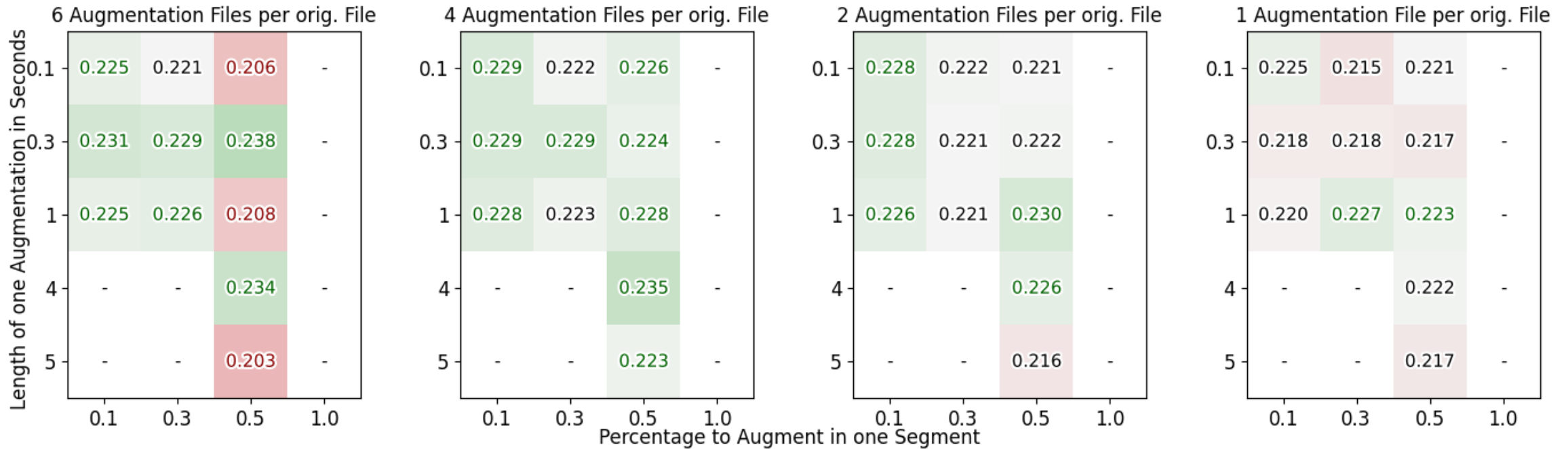# Segment Swap

# Segment Swap



- one parameter combination with a significantly poorer outcome:
  100% augmentation and 0.1 seconds
  - insufficient duration of 0.1 seconds
  - related to the length of the vowels and consonants (length <0.3seconds)
- Optimal: 1 File per original File, 30% augmentation rate, length of 0.3 seconds each
  - using 6 files (with same Score) may not justify the increased computational overhead
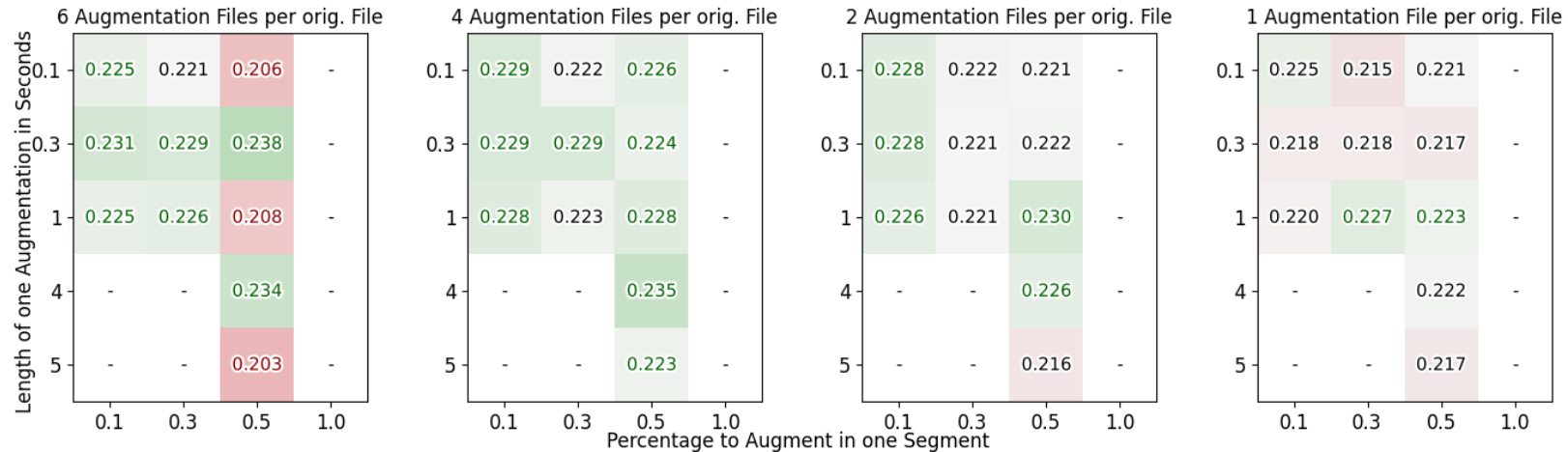- 1.1% enhancement compared to without augmentation

# EDA

| Operation | Sentence |
|---|---|
| None | A sad, superior human comedy played out on the back roads of life. |
| SR _Synonym Replacement_ | A **_lamentable_**, superior human comedy played out on the **_backward_** road of life. |
| RI _Random Insertion_ | A sad, superior human comedy played out on **_funniness_** the back roads of life. |
| RS _Random Swap_ | A sad, superior human comedy played out on **_roads_** back **_the_** of life. |
| RD _Random Deletion_ | A sad, superior human out on the roads of life. |

Shifting Pitch (SP)

Background Noise (BN)

Segment Swap (SeS)
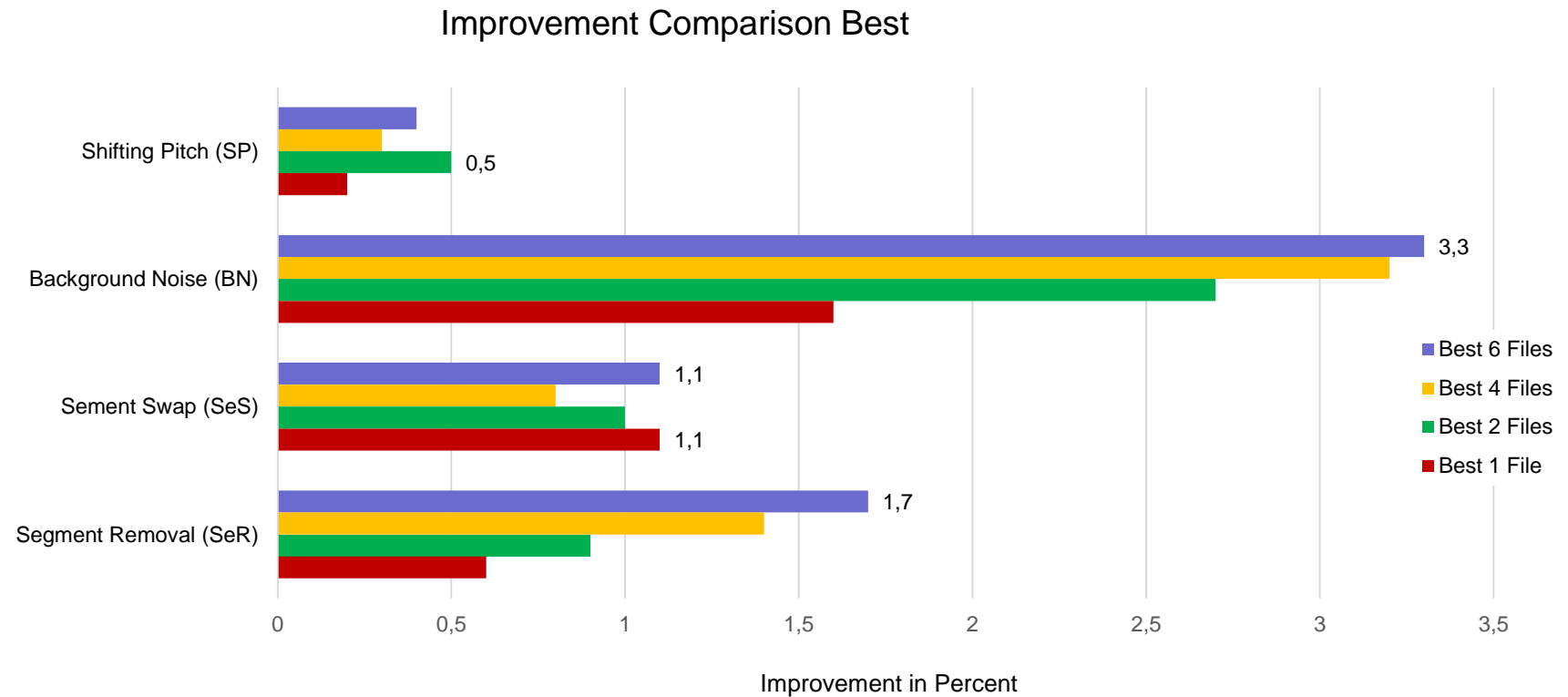
Segment Removal (SeR)

# Segment Removal

# Segment Removal



- Optimal: 6 Files per original File, 50% augmentation rate, length of 0.3 seconds each
- 1.7% enhancement compared to without augmentation
- For 6 Files per original File, 50% augmentation rate there are three significant worse results
  – Only results with significant worse performance
  – Conversely not for 4, 2 or 1 Files per original File

# Results Main Methods



Improvement Comparison Best

Segment Removal

Time Masking

Segment Swap

Speaker Insertion

Shifting Pitch

Time Stretching

Volume Confusion

Speed Confusion

Time Reversing

Background Noise

Frequency Masking

Frequency Swapping

Frequency Insertion

Augmentation

35

# Time Masking

- Similar to Segment Removal
- but the interval is not removed
- instead, it is replaced by zeros

- Used values for Hyperparameters:
  - $n_{augFiles} = 6$
  - $\propto = 0.5$
  - $l_{aug} = 0.3$

# Speaker Insertion

- The specific interval is replaced by another random interval from a different speaker
- Speaker is of the same class (hence, the same dialect).

- Used values for Hyperparameters:
  - $n_{augFiles} = 1$
  - $\propto = 0.3$
  - $l_{aug} = 0.3$

# Time Stretching & Speed Confusion

- Time Stretching
  - Intervals are time stretched within the range of [0.8, 1.2]
  - Pitch remains unchanged

- Speed Confusion
  - Similar to TS
  - But Pitch changes too
  - To archive that, the Interval gets resampled, but saved with the original sample rate
  - newSamplingRate=rate*oldSamplingRate,rate$\in$[0.8,1.2]

# Volume Confusion & Time Reversing

- Volume Confusion
  - Peak of the segment is set to a value within the range [0.2, 0.8]

- Time Reversing
  - Order of the samples in the Interval gets reversed

- Used values for Hyperparameters:
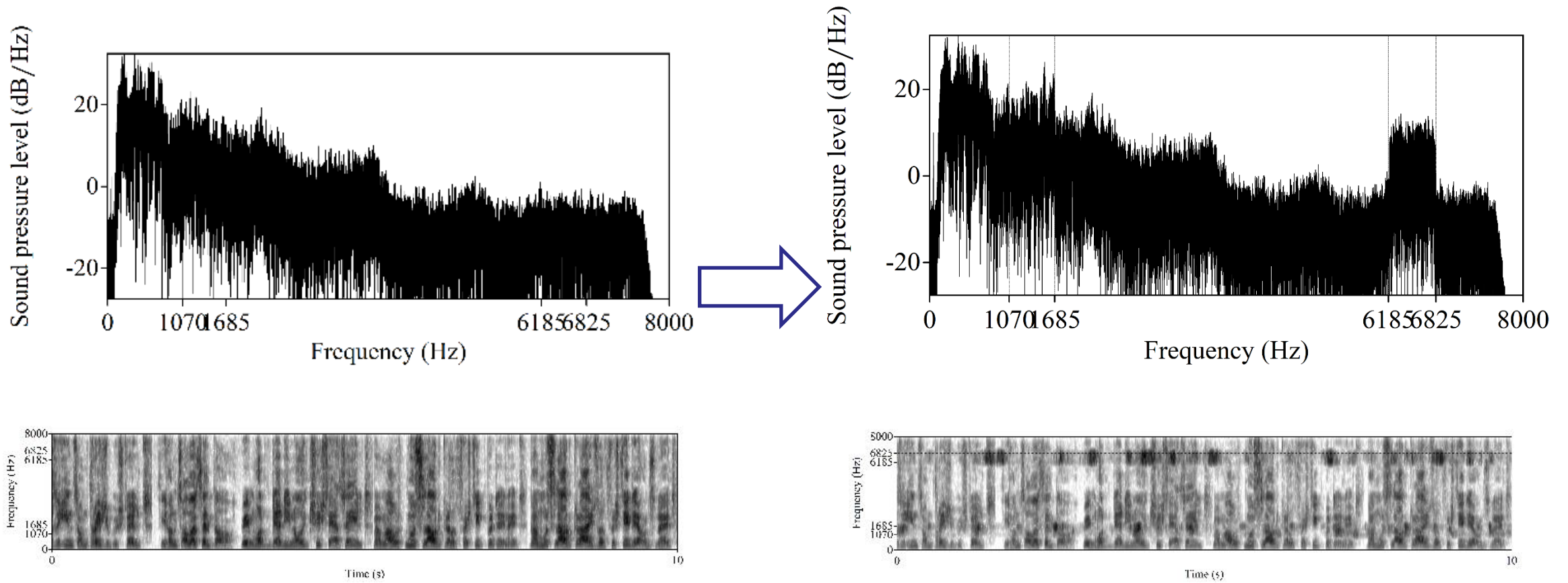  - $n_{augFiles} = 2$
  - $\propto = 0.5$
  - $l_{aug} = 4.0$

# Frequency manipulation

- Used values for Hyperparameters:
  - $n_{augFiles} = 6$
  - $\propto$ = not needed
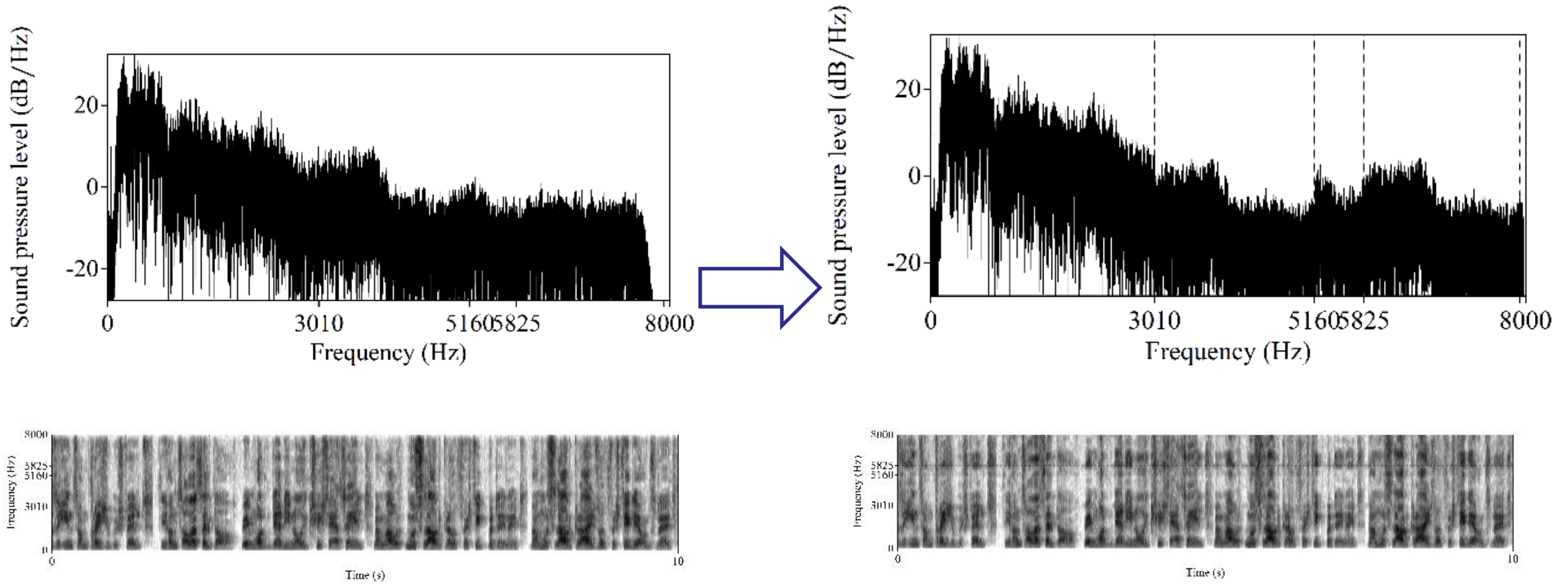  - $l_{aug}$ = not needed

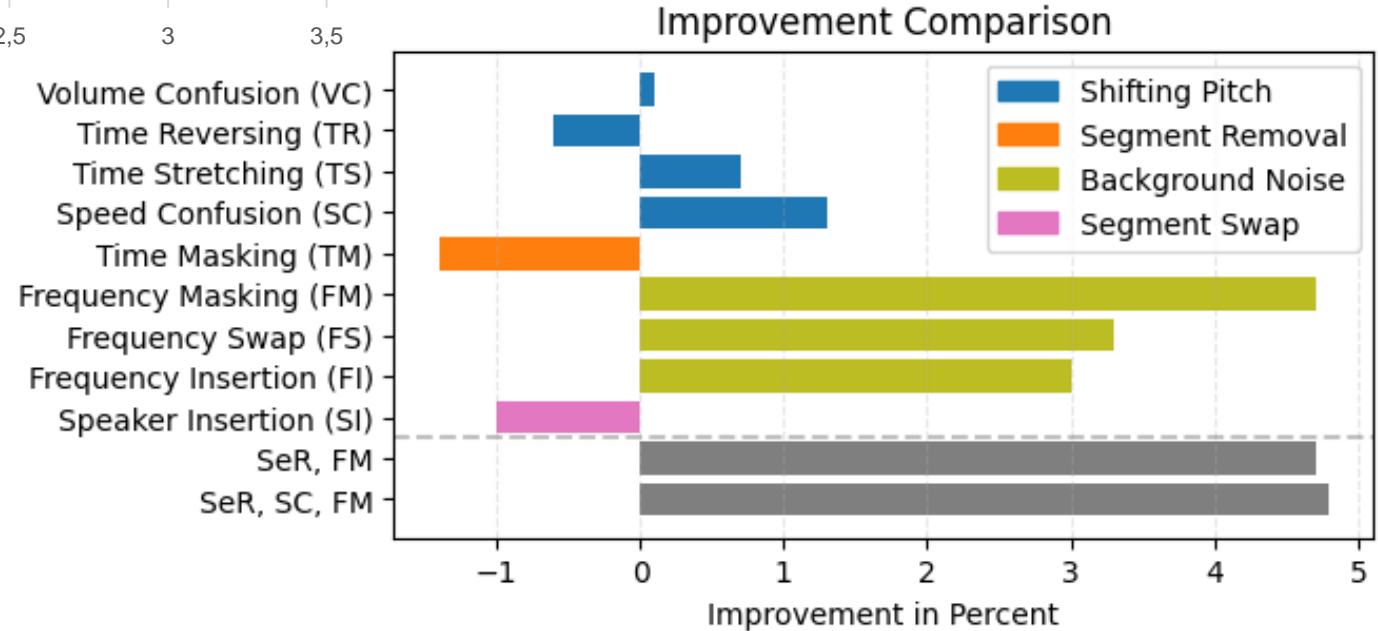# Frequency Masking
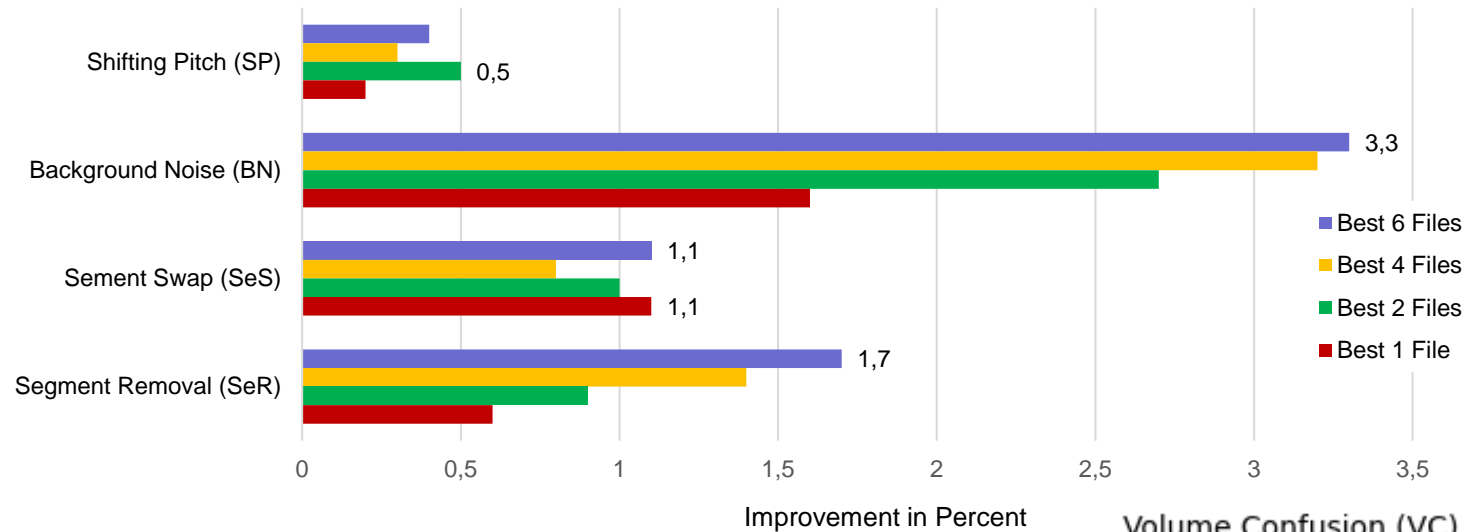
# Frequency Insertion

# Frequency Swapping

# Results of all Methods



Improvement Comparison Best

# Conlusion

- Best Method is Frequency Masking

- 4.7% better than without augmentation (from 0.221 to 0.268)

- Generally, all methods that are masking frequencies yield the best results

- Can add Segment Removal without performance loss to reduce computation effort

# Thanks for your attention!
## Any Questions?