



Inter-University Research Institute Corporation /
Research Organization of Information and Systems
National Institute of Informatics



科学を支え、未来へつなぐ
科学技術振興機構
Japan Science and Technology Agency



JSPS

A Preliminary Case Study on Long-Form In-the-Wild Audio Spoofing Detection

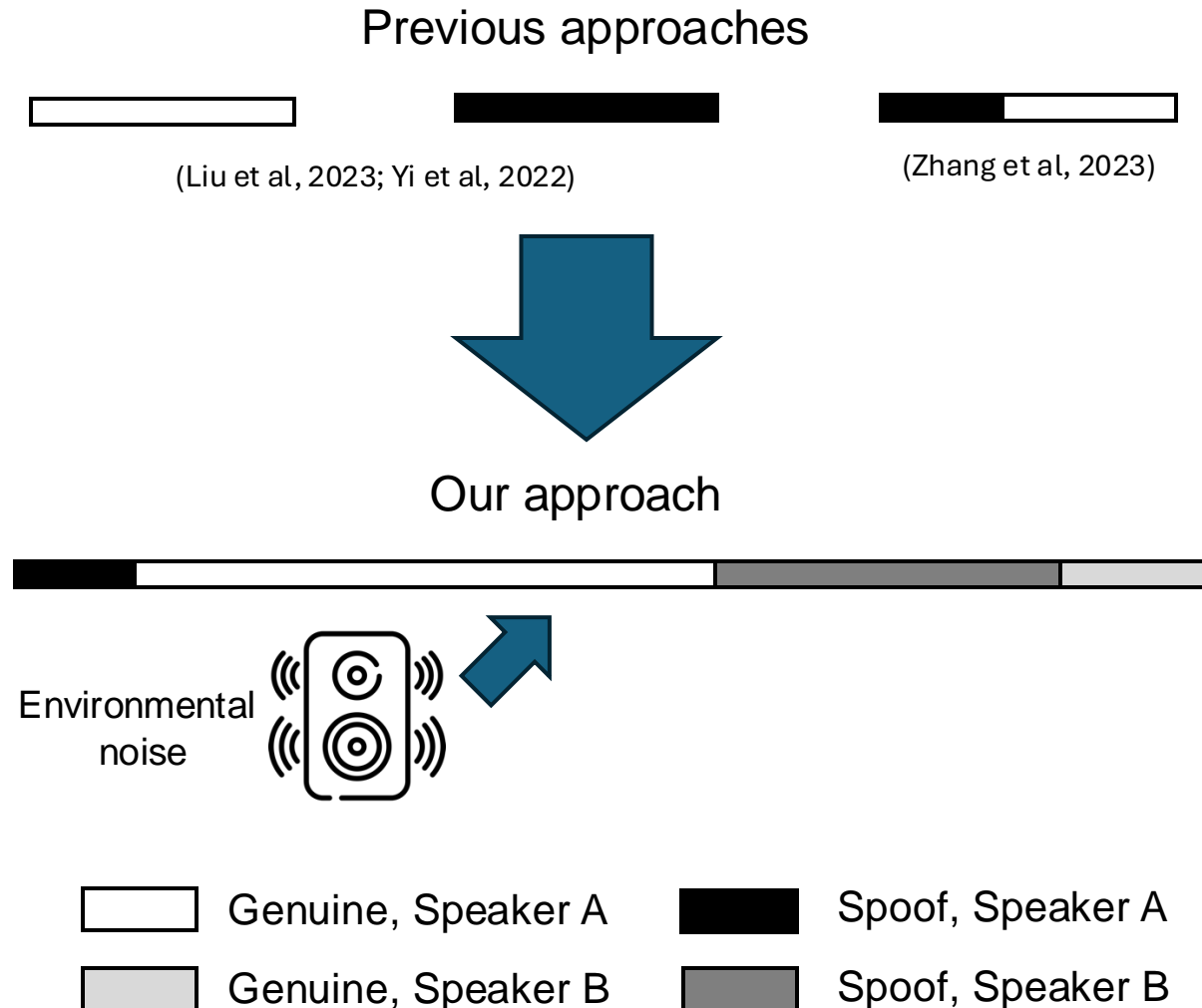
Xuechen Liu, Xin Wang, *Junichi Yamagishi*

BIOSIG 2024

2024-09-26

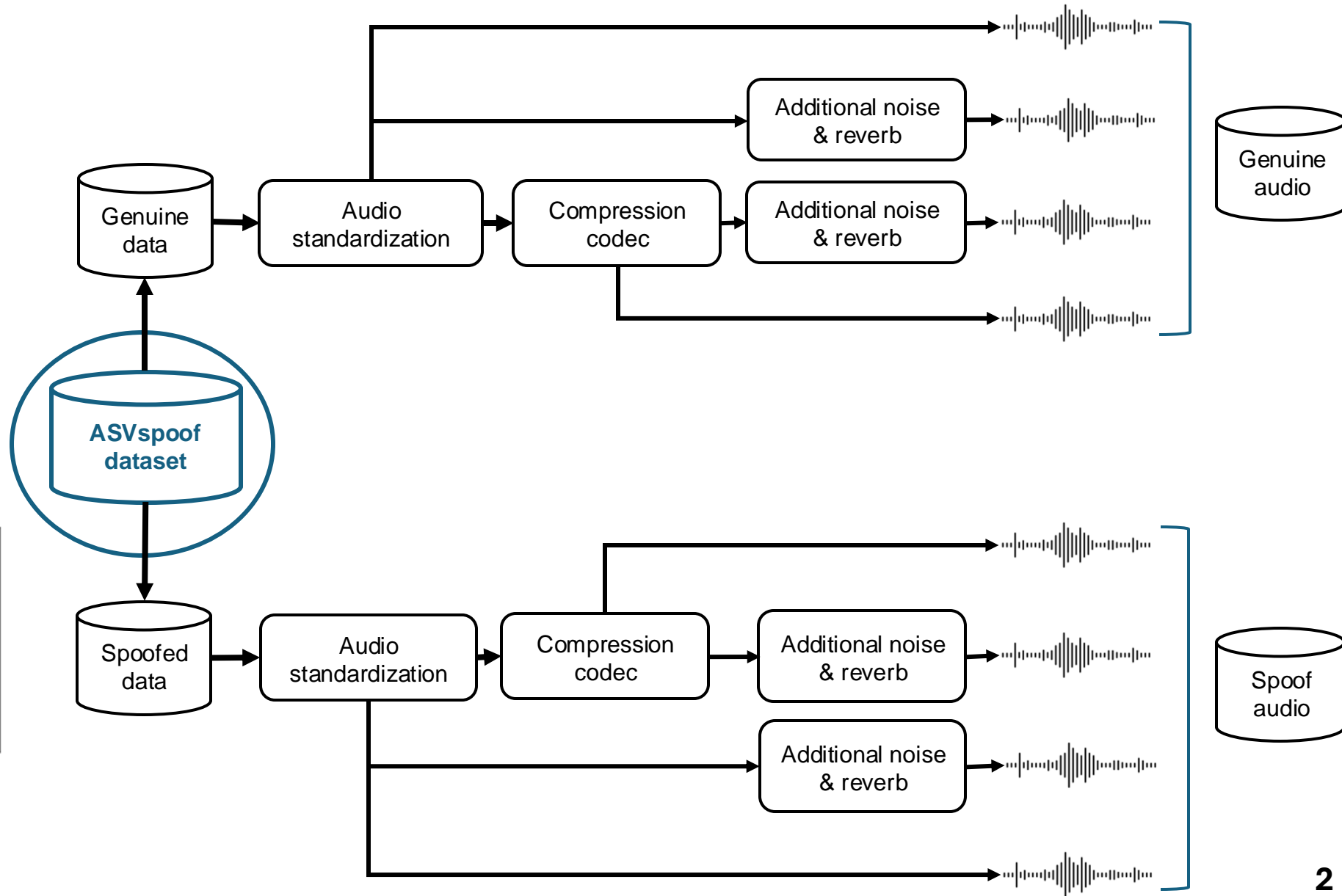
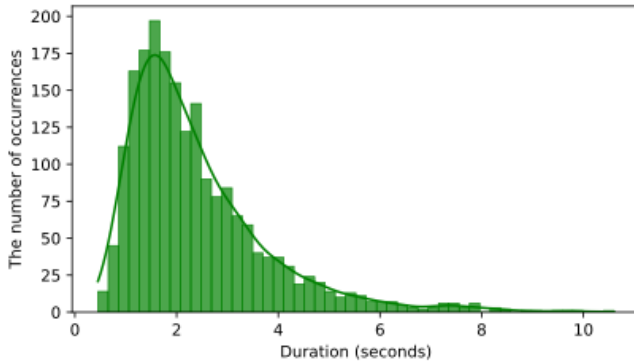
Our Objective

- Audio deepfakes pose a significant threat, with real-world cases of fraud and misinformation on the rise
- So far, we have focused on **short-duration (short-form), acoustically-clean and single-speaker** audio waveforms
- Longer durations and partially-spoofed audio have been addressed
- Our goal is to take a step towards realistic conditions, with **longer duration (long-form), acoustically-complex, and multi-speaker** audios



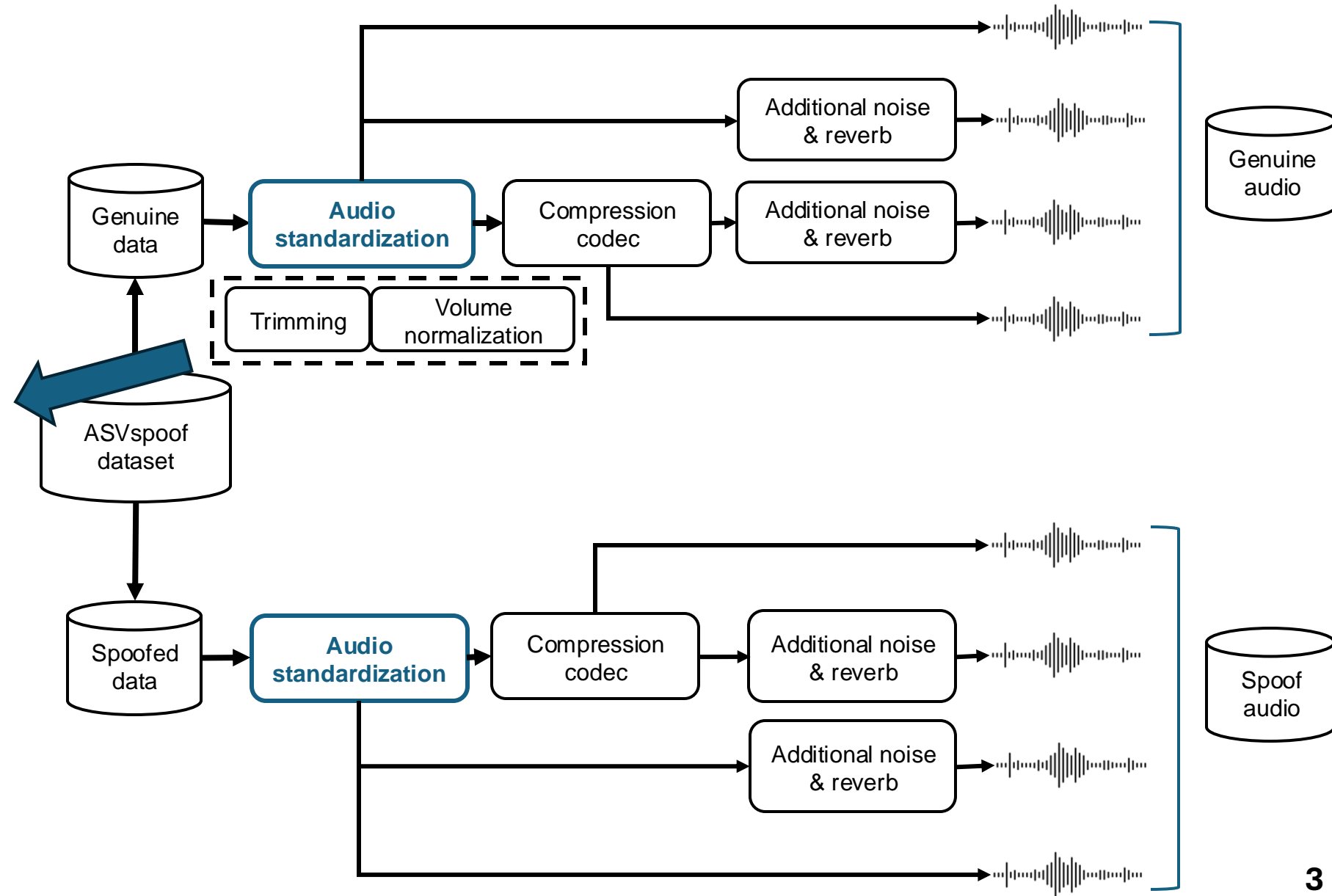
Long-form Spoofed Audio Generation

- 2-10 seconds
- Single speaker
- Clean condition



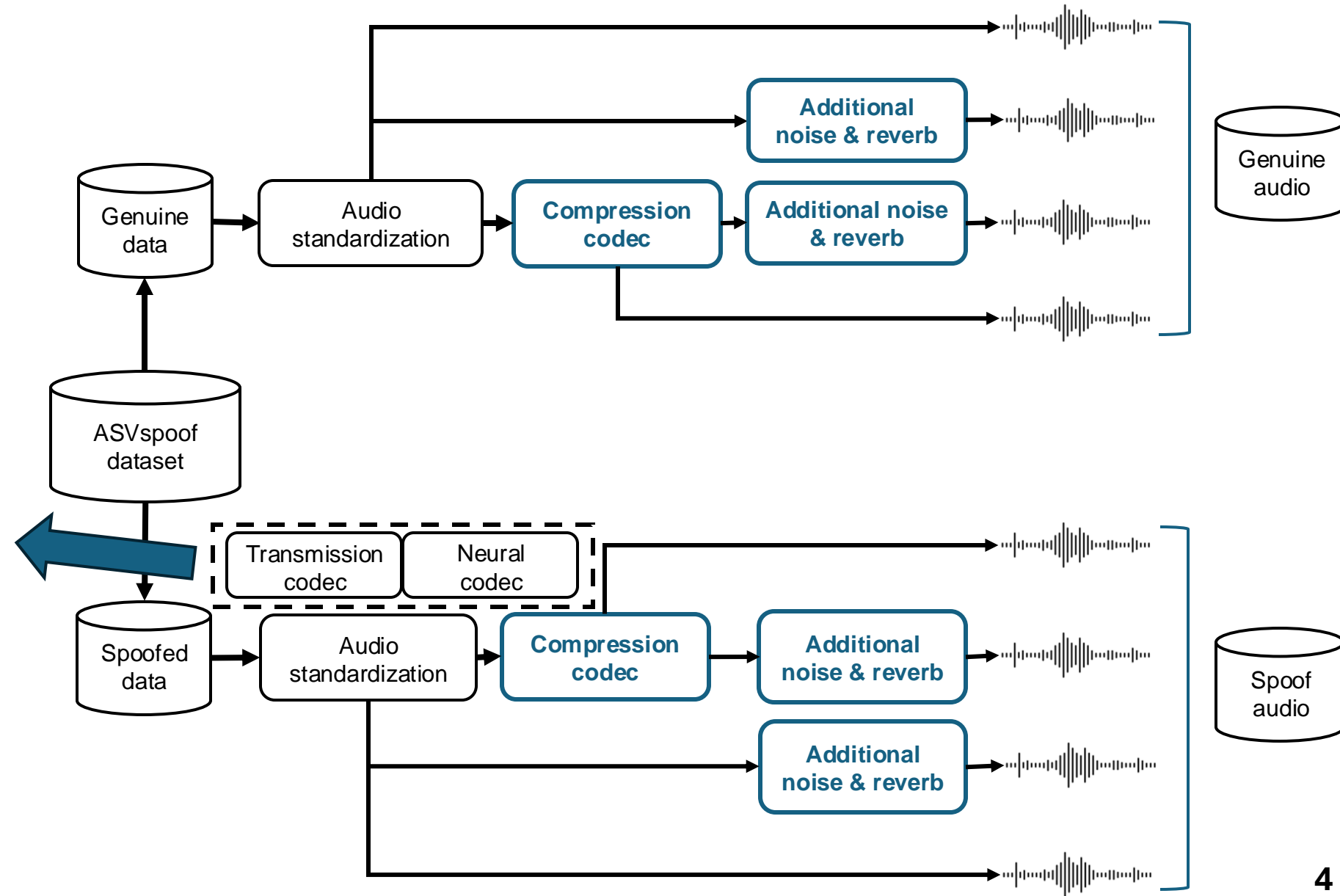
Long-form Spoofed Audio Generation

- Removing silence at beginning & end of the audio
- Normalize the speaker volume using sv56 toolkit



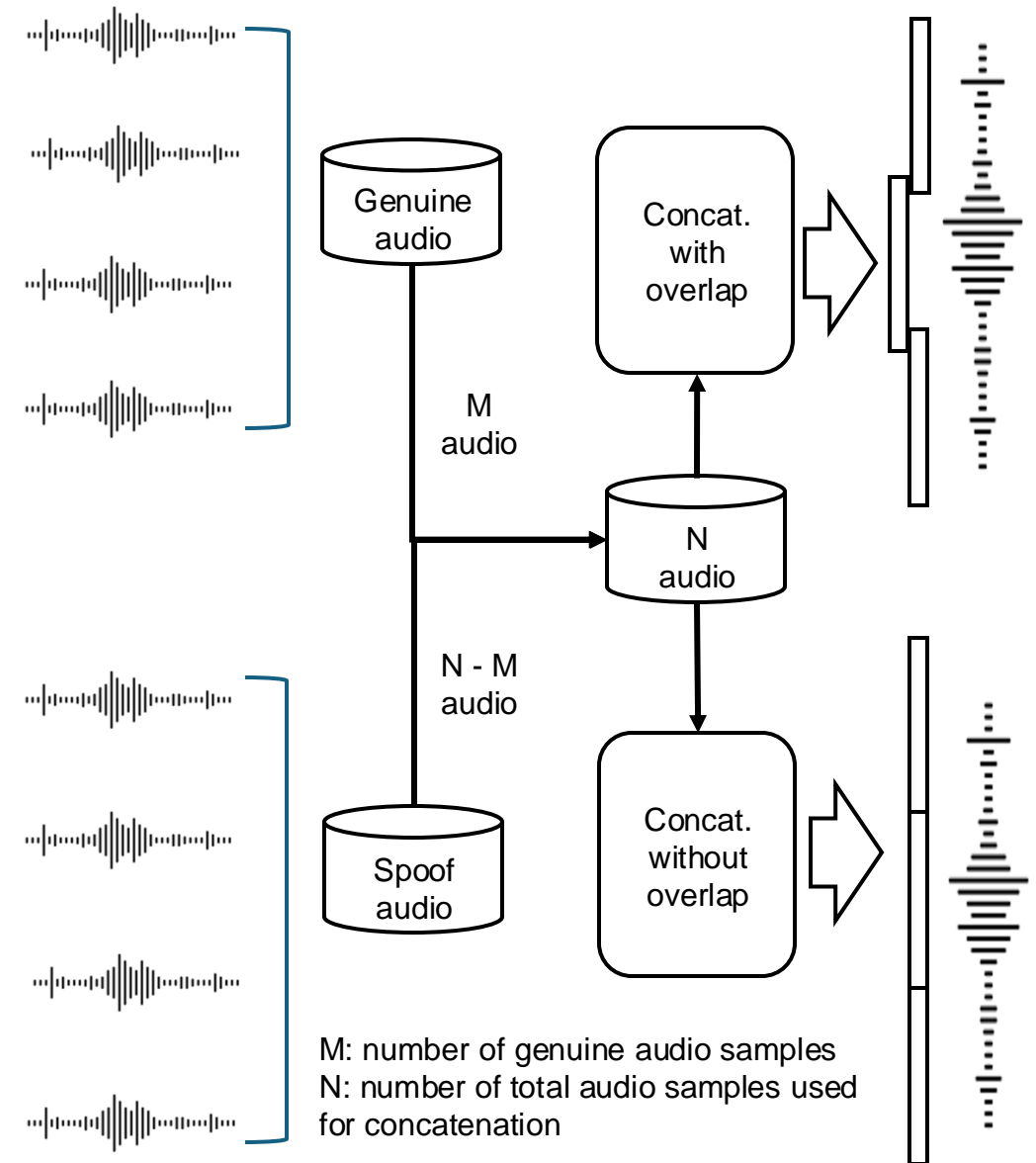
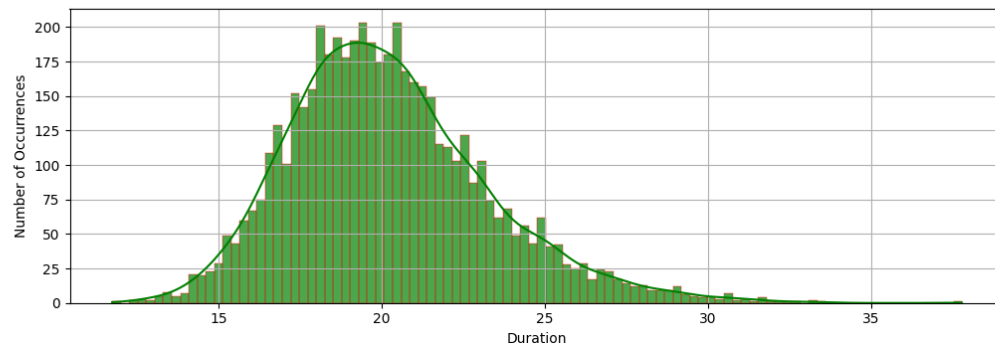
Long-form Spoofed Audio Generation

- We use various codec models as processing methods
- Transmission codec: mp3, opus, a-law.....
- Neural codec: EnCodec, Soundstream, FACodec.....
- Additions: RIR & MUSAN



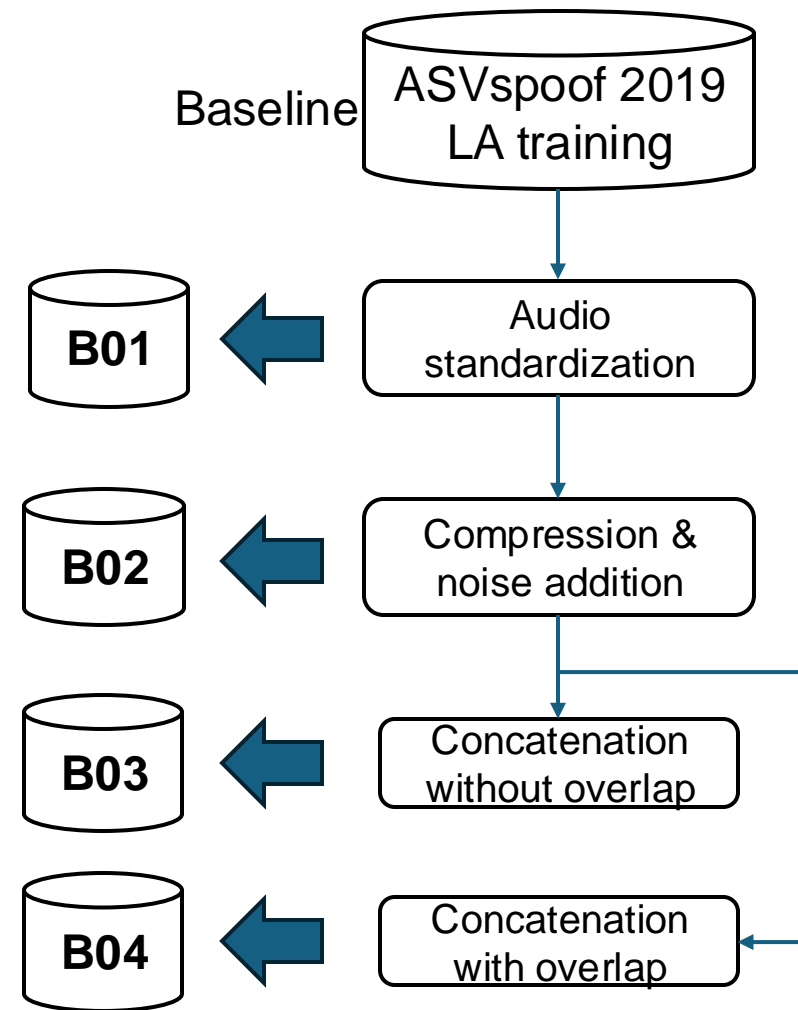
Long-form Spoofed Audio Generation

- Generated audio via concatenation, with optionally overlap
- Usually, 15-30 seconds long
- Multiple speakers in the same audio
- Mixed amount of genuine and spoofed content
- **Varying amount of genuine vs. spoof ratio**, so purely-genuine and purely-spoofed audios are also there



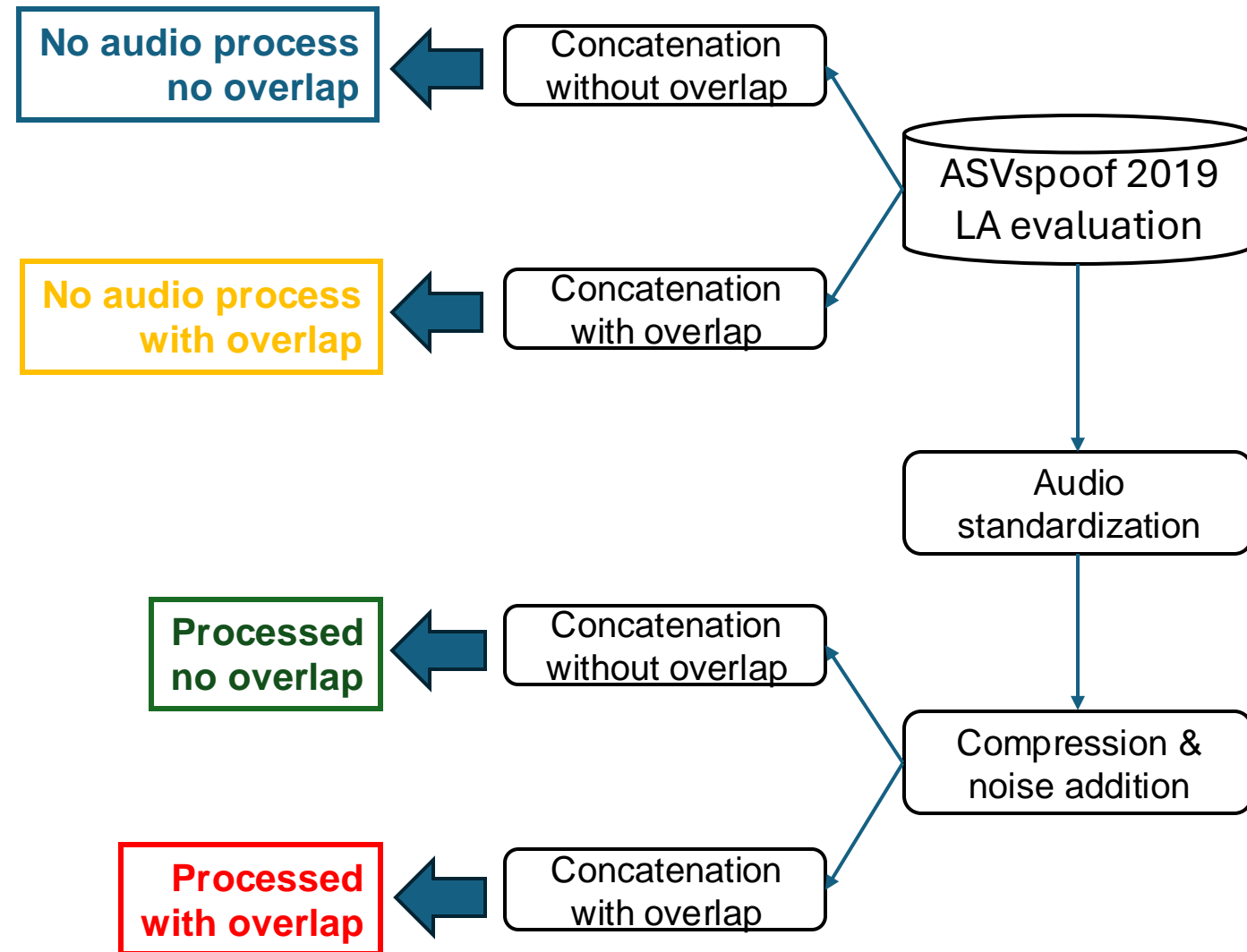
Experimental Setup

- We use state-of-the-art end-to-end audio deepfake detector called AASIST
- We developed multiple training sets
 - Baseline: The original ASVspoof 2019 LA training data (short-form, clean, single-speaker)
 - B01: Baseline + audio standardization
 - B02: B01 + audio processing
 - B03: B02 + long-form audio concatenation
 - B04: B02 + long-form audio concatenation, w/ overlap
- Evaluation audio are generated via same functions, with different placements
 - No audio process, no overlap
 - No audio process, with overlap
 - Processed, no overlap
 - Processed, with overlap



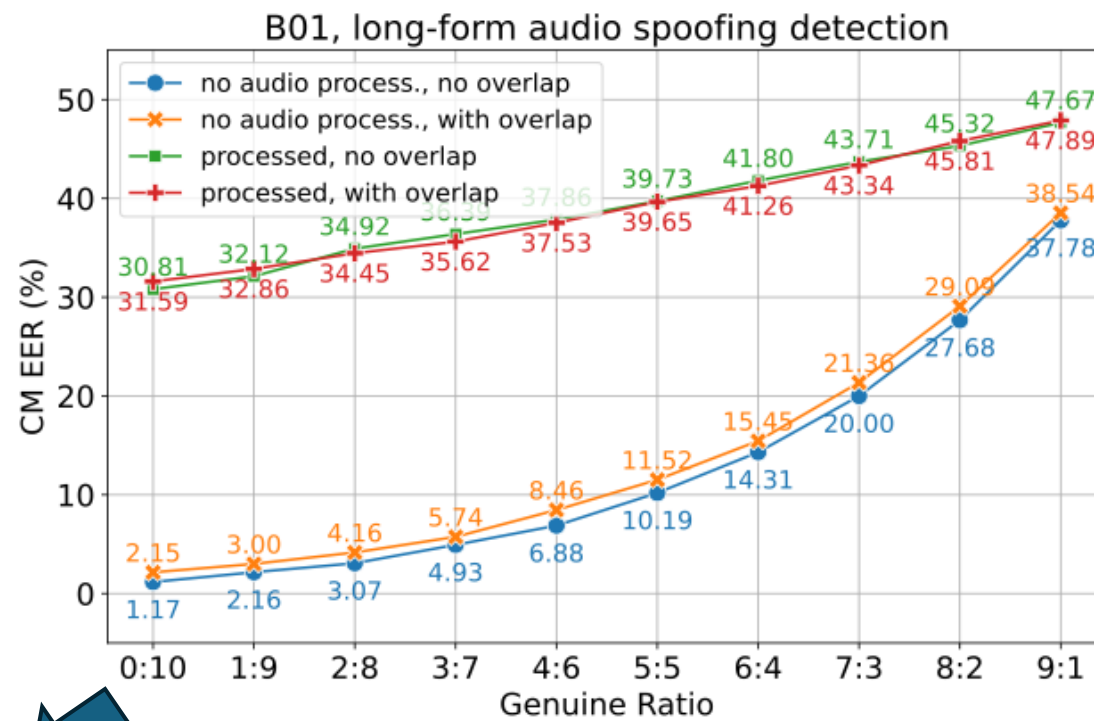
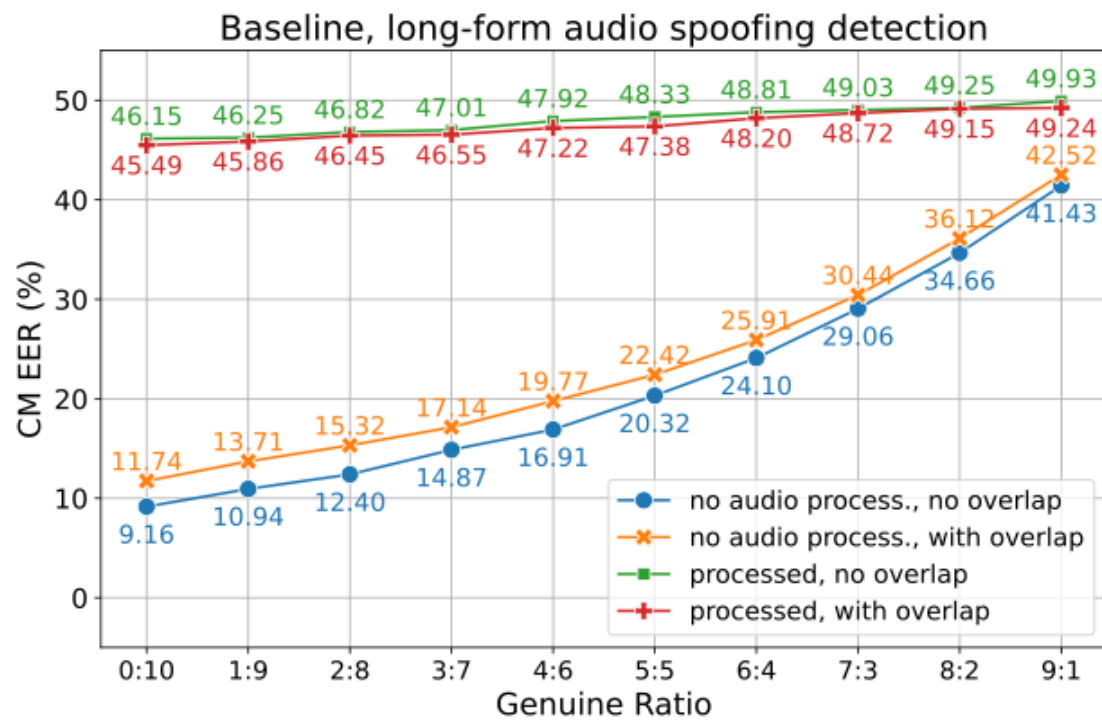
Experimental Setup

- We use state-of-the-art end-to-end audio deepfake detector called AASIST
- We developed multiple training sets
 - Baseline: The original ASVspoof 2019 LA training data (short-form, clean, single-speaker)
 - B01: Baseline + audio standardization
 - B02: B01 + audio processing (codec,
 - B03: B02 + long-form audio concatenation
 - B04: B02 + long-form audio concatenation, w/ overlap
- Evaluation audio are generated via same functions, with different placements
 - **No audio process, no overlap**
 - **No audio process, with overlap**
 - **Processed, no overlap**
 - **Processed, with overlap**



Results

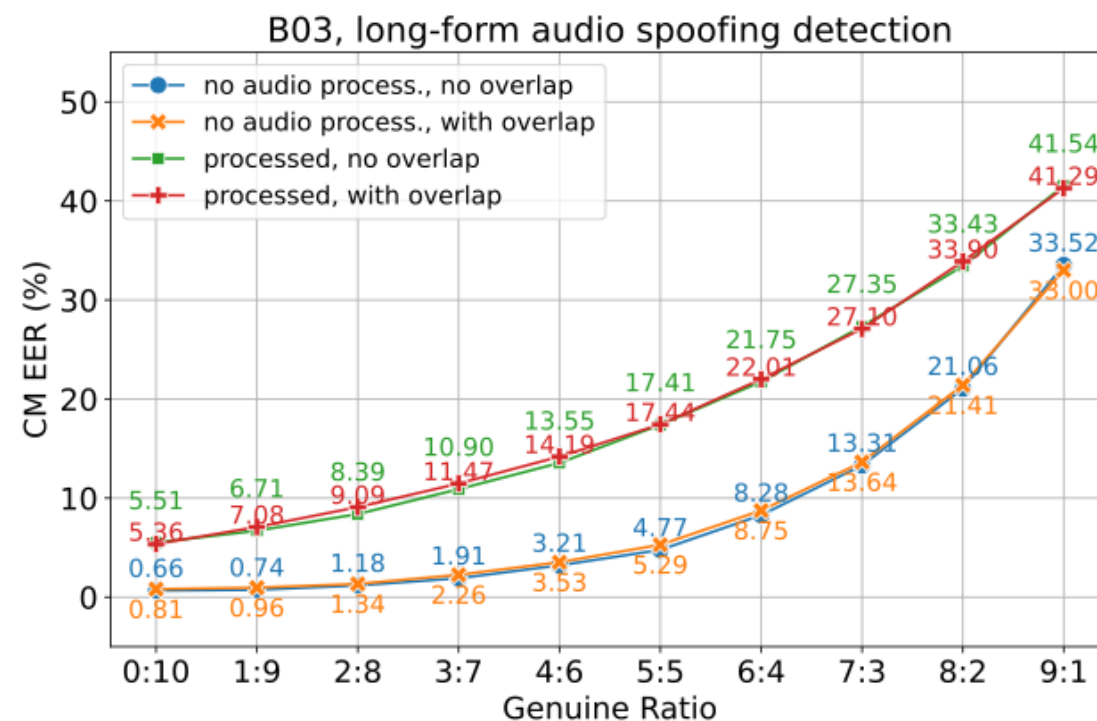
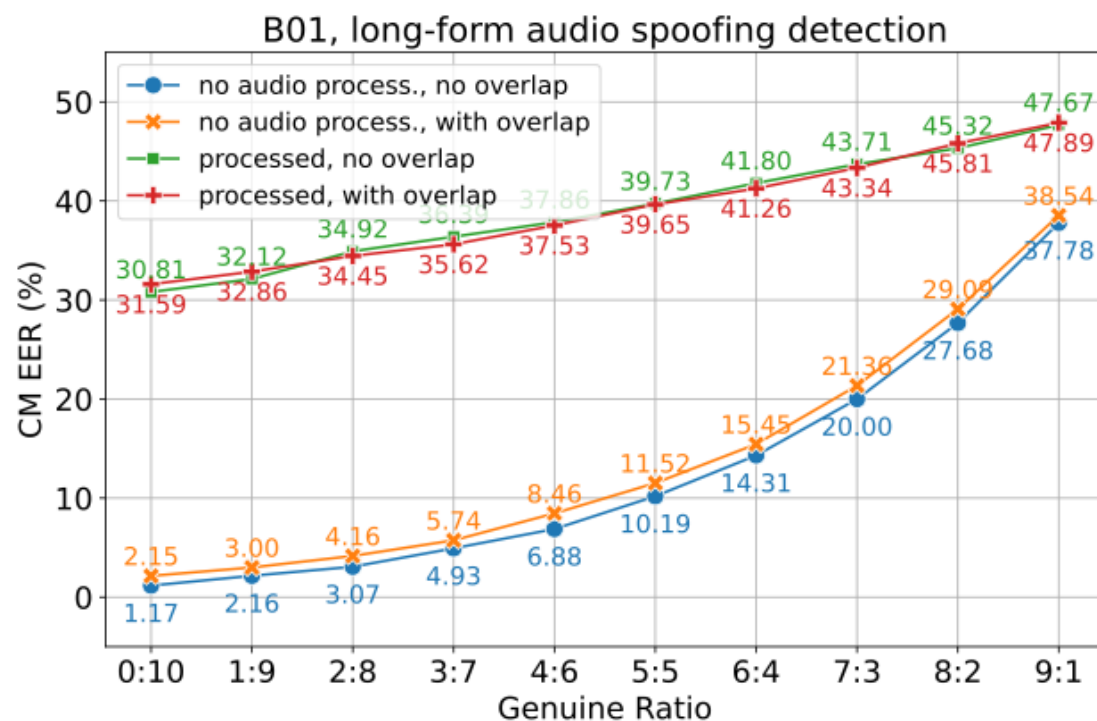
- Training on short, clean speech is not enough for detecting complex spoofed audios
- Performance drops as the amount of genuine content in the long-form audio increases



M : (N - M)

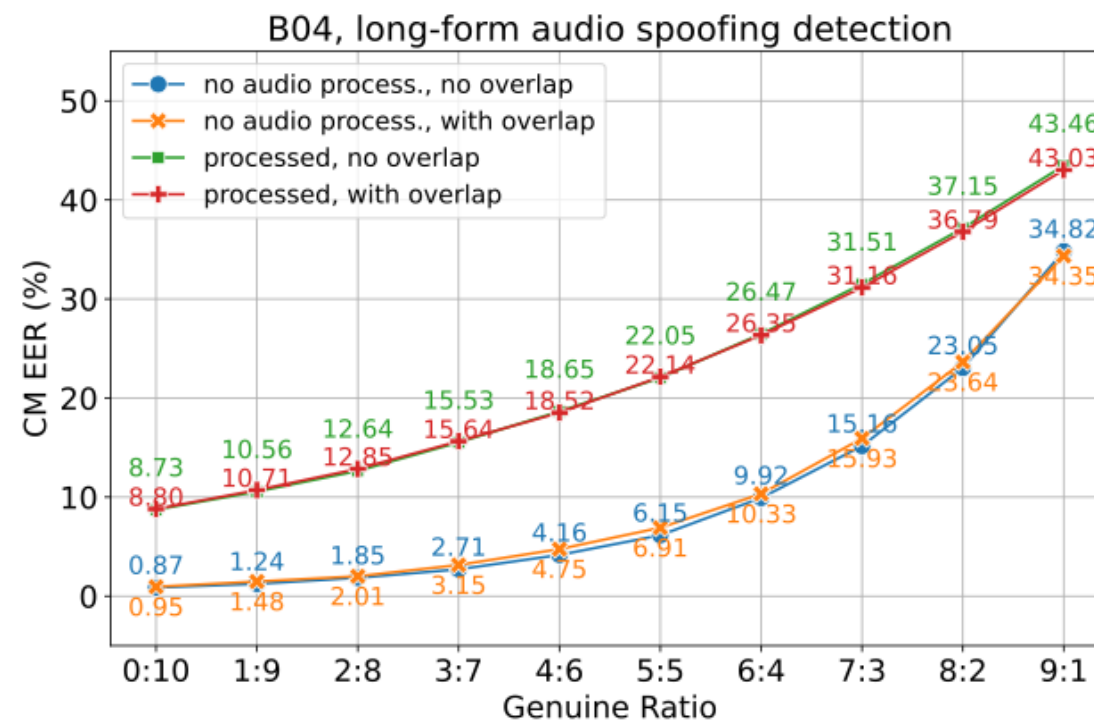
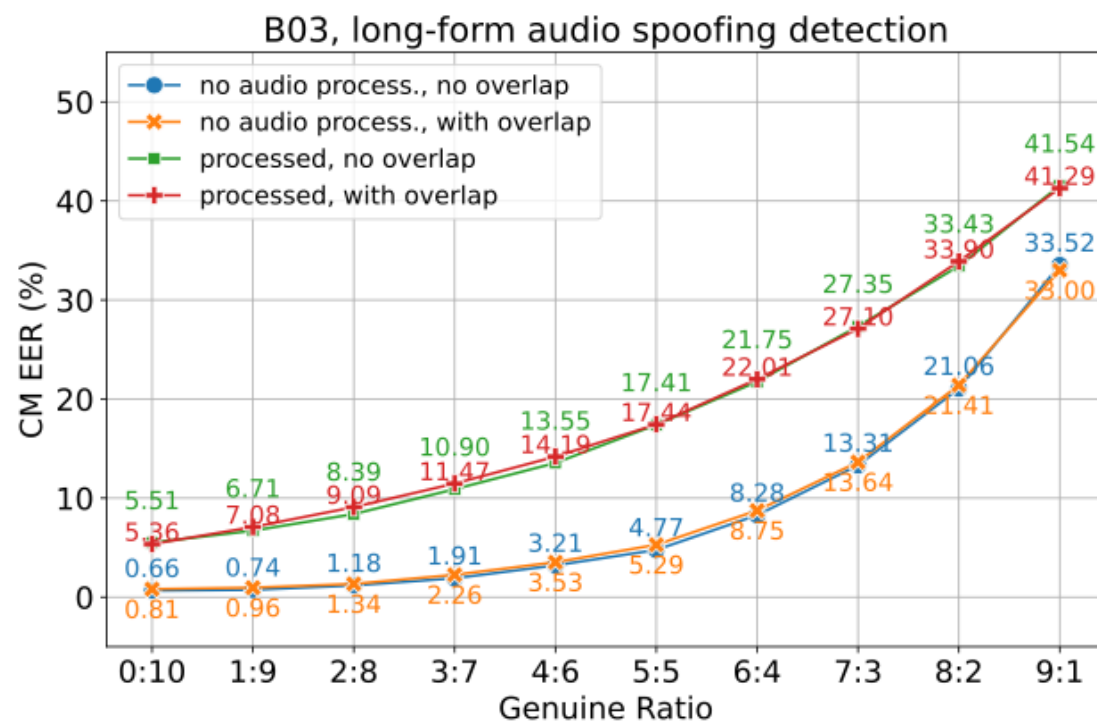
Results

- Further improvement on training set improves performance, especially for highly-spoofed audio



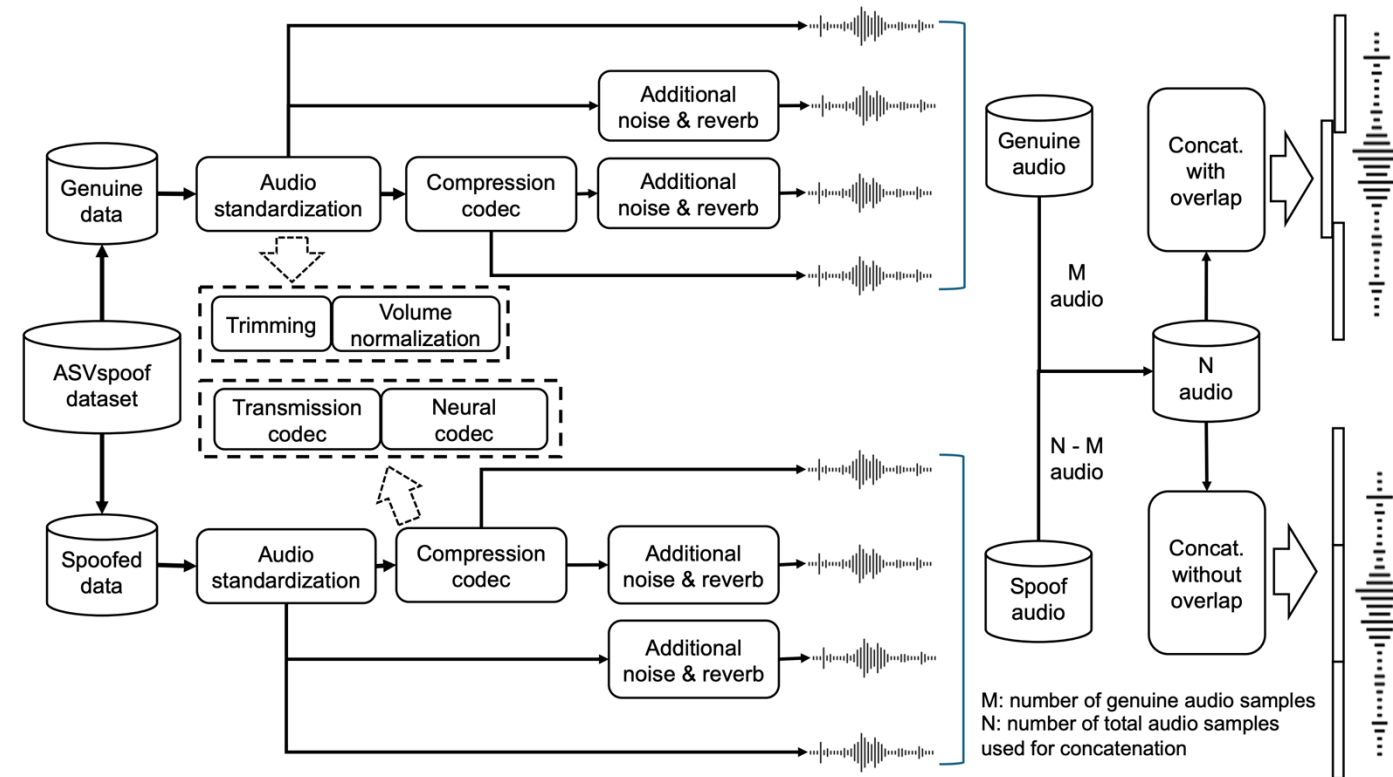
Results

- Overlapping the adjacent in long-form audio in the training set slightly worsens the performance



Main Take-Aways

- Short-form detection models are insufficient for more challenging scenarios posed by audio DeepFake technologies
- We have proposed a new pipeline on generating longer duration
 - 15-30 seconds
 - Varying amount of spoofed content
 - Multi-speaker presence
- Long-form training data improves detection accuracy for complex audio spoofing task
- Future work will explore larger datasets and more complex models





Inter-University Research Institute Corporation /
Research Organization of Information and Systems
National Institute of Informatics



科学を支え、未来へつなぐ
科学技術振興機構
Japan Science and Technology Agency



JSPS

Thanks for Listening!