

SAMPLE-LEVEL CNN ARCHITECTURES FOR MUSIC AUTO-TAGGING USING RAW WAVEFORMS

Taejun Kim¹, Jongpil Lee², Juhan Nam²

¹School of Electrical and Computer Engineering, University of Seoul, Republic of Korea. ktj7147@uos.ac.kr

²Graduate School of Culture Technology, KAIST, Republic of Korea. {richter, juhanam}@kaist.ac.kr

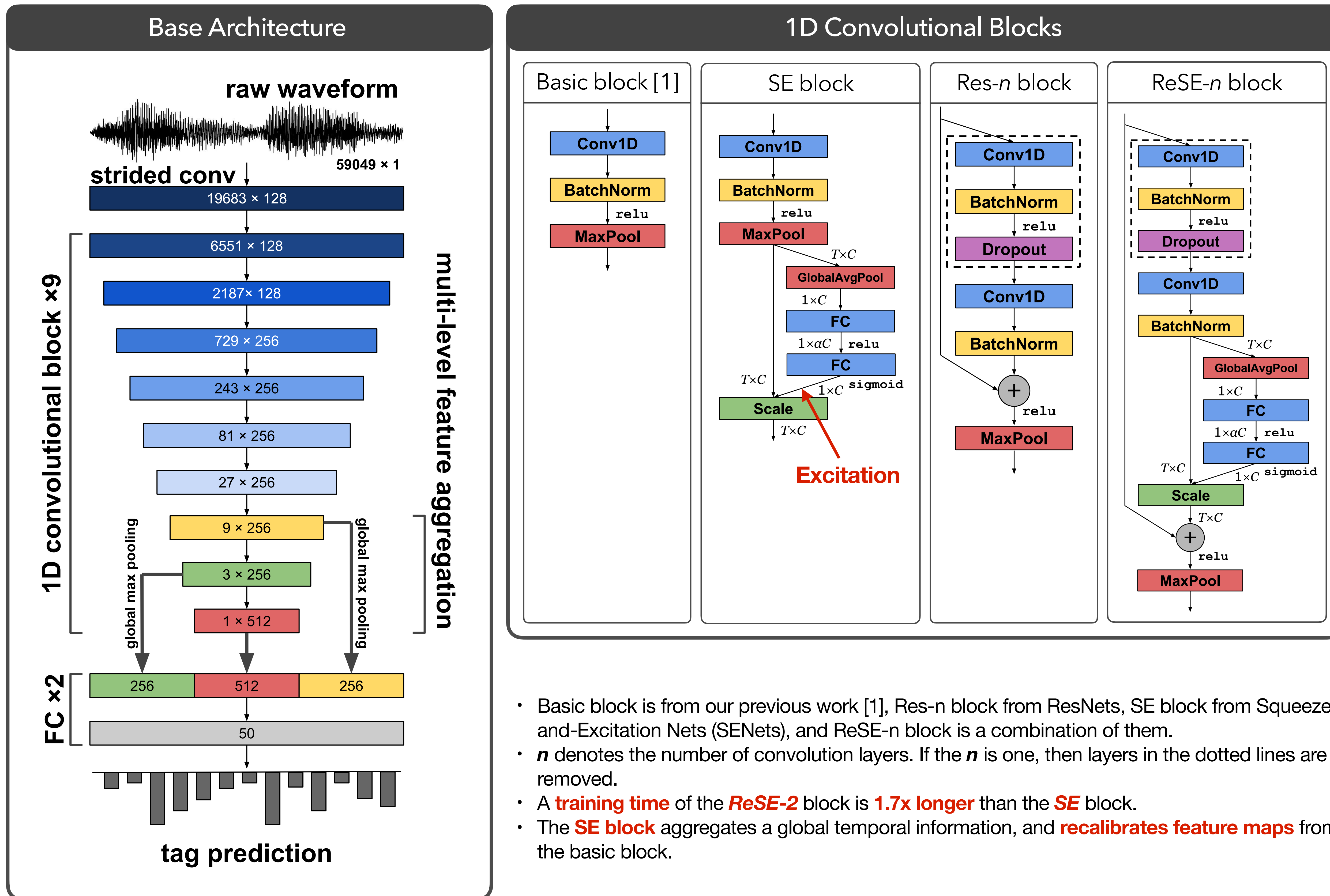
1

Summary

- We explore **end-to-end Convolutional Neural Network (CNN) architectures** for music auto-tagging tasks.
- We adopt architectures from state-of-the-art image classification networks (ResNets & SENets).
- Our models achieve **state-of-the-art results on MagnaTagATune dataset** and **comparable results on Million Song Dataset**.
- We analyze and visualize that the **SE blocks tend to normalize loudness of audios**.

2

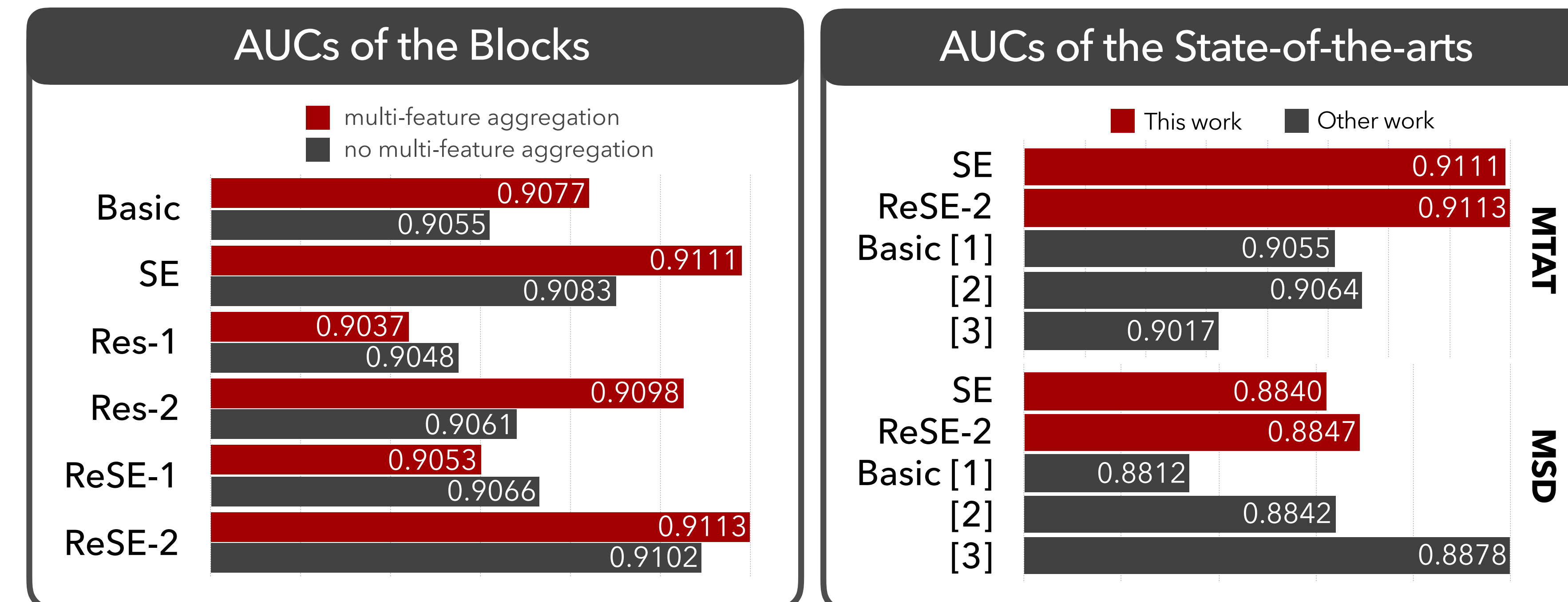
Architectures



- Basic block is from our previous work [1], Res-n block from ResNets, SE block from Squeeze-and-Excitation Nets (SEnets), and ReSE-n block is a combination of them.
- n denotes the number of convolution layers. If the n is one, then layers in the dotted lines are removed.
- A **training time** of the **ReSE-2** block is **1.7x longer** than the **SE** block.
- The **SE block** aggregates a global temporal information, and **recalibrates feature maps** from the basic block.

3

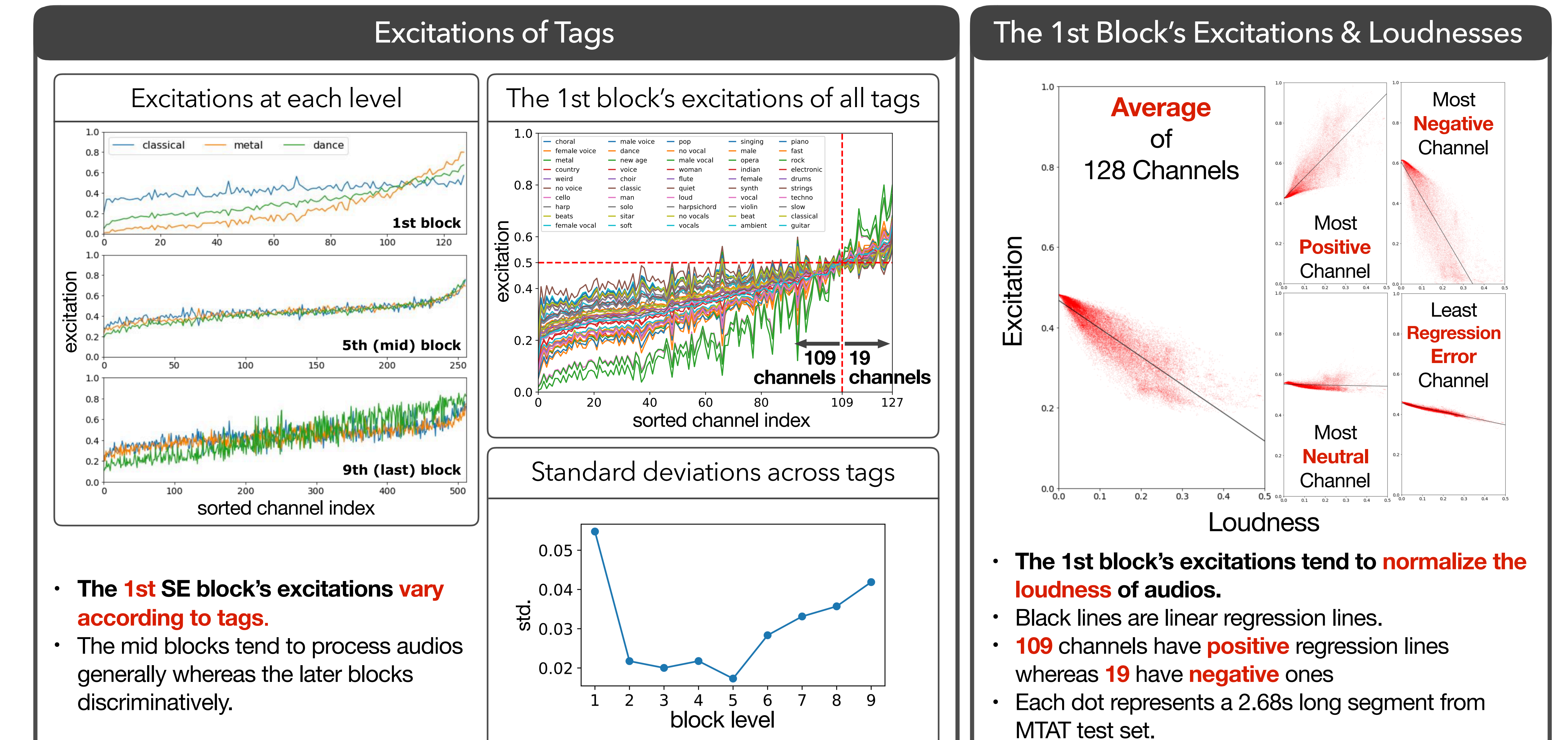
Experiments



- MagnaTagATune (MTAT)**
 - 21,105 songs
 - 170h long (29.1s for each song)
 - Split: 15,244 / 1529 / 4332
 - Top 50 most frequent tags
- Million Song Dataset (MSD)**
 - 241,889 songs
 - 1955h long (29.1s for each song)
 - Split: 201,680 / 11,774 / 28,435
 - 50 tags (Last.FM tag annotations)

4

Analysis of Excitations from the SE Blocks



References

- [1] Jongpil Lee, Jiyoung Park, Keunhyoung Luke Kim, and Juhan Nam, "Sample-level deep convolutional neural networks for music auto-tagging using raw waveforms," in Sound and Music Computing Conference (SMC), 2017.
- [2] Jongpil Lee and Juhan Nam, "Multi-level and multi-scale feature aggregation using sample-level deep convolutional neural networks for music classification," Machine Learning for Music Discovery Workshop, International Conference on Machine Learning (ICML), 2017.
- [3] Jongpil Lee and Juhan Nam, "Multi-level and multi-scale feature aggregation using pretrained convolutional neural networks for music auto-tagging," IEEE Signal Processing Letters, vol. 24, no. 8, pp. 1208–1212, 2017.

<https://github.com/tae-jun/resemul>

