

\*Keio University, Japan

## **1. Summary**



# 2. Background

Sound Event Detection (SED):



### Experimental setting

Pre-training Dataset:

• Audioset data w/ Strong labels (80k) [2]

**Formula-SED (ours) {50k, 100k, 1M}** 

Downstream Task: DCASE 2023, task 4 (DESED dataset) Evaluation metrics: PSDS1, 2, Event-F1, Inersection-F1

Table1. Accuracy comparison							lable2. The impact of pre-training dataset size						
Model	PSDS1	PSDS2	E-F1(%)	I-F1(%)	0.3	Man Martin Martin Martin	Model	Size	PSDS1	PSDS2	E-F1(%)	I-F1(%)	т Т Т
CRNN baseline [3] w/ Formula-SED (100k) w/ Audioset Strong (79k) [2]	0.352 0.405 0.387	0.579 <b>0.641</b> 0.618	45.7 <b>49.6</b> 47.7	65.8 <b>72.3</b> 70.7	SOS 0.2	Random	CRNN baseline CRNN baseline CRNN baseline	50k   100k   1M	0.380 0.405 <b>0.420</b>	0.620 0.641 <b>0.653</b>	49.4 49.6 <b>51.4</b>	70.8 72.3 <b>72.8</b>	Even
Paderborn CRNN [4] w/ Formula-SED (100k) w/ Audioset Strong (79k) [2]	0.262 0.278 0.355	0.506 0.539 <b>0.622</b>	34.9 35.3 <b>43.9</b>	57.3 59.4 <b>67.8</b>	0.1 0.0 0	Audioset Formula-SED 4000 8000 12000 step Training curve	Paderborn CRNN Paderborn CRNN Paderborn CRNN	50k   100k   1M	0.288 0.278 0.304	<b>0.553</b> 0.539 0.552	35.8 35.3 <b>37.1</b>	<b>62.4</b> 59.4 61.4	
5. Future work						6. Reference							
Evaluate the effectiveness of our Formula-SED for various tasks beyond sound event detection.						<ul> <li>[1] E. V. Bonilla, K. Chai, and C. Williams, "Multi-task gaussian process prediction," NIPS, vol. 20, 2007.</li> <li>[2] S. Hershey et al, "The benefit of temporally-strong labels in audio event classification," ICASSP, pp. 366–</li> <li>[3] M. Fuentes et al, DCASE 2023, Tampere, Finland: Tampere University, September 2023.</li> </ul>							

- Investigate the relationship between label thresholds and downstream accuracy.

# Formula-Supervised Sound Event Detection: Pre-Training Without Real Data Yuto Shibata\*<sup>†</sup>, Keitaro Tanaka<sup>†</sup><sup>‡</sup>, Yoshiaki Bando<sup>†</sup>, Keisuke Imoto<sup>†</sup>§, Hirokatsu Kataoka<sup>†</sup>¶, and Yoshimitsu Aoki

\*National Institute of Advanced Industrial Science and Technology (AIST), Japan

- Downstream accuracy is improved by pre-training using the proposed dataset (Tab 1).
- For CRNN baseline, our pre-training method outperformed real but noisy AudioSet.
- The training curve shows our pre-training method speeds up fine-tuning convergence.
- In the baseline CRNN, accuracy improved monotonically with the increase in the pre-training data scale (Tab 2).

-370, 2021. Yuto Shibata : yuto071508@keio.jp [4] J.EbbersandR.Haeb-Umbach, "Pre-trainingandself-trainingforsound event detection in domestic environments," Tech. Rep. of DCASE 2022 Challenge Task 4, 2022. [5] Mesaros et al. "Sound event detection: A tutorial." *IEEE Signal Processing Magazine* 38.5 (2021): 67-83.

<sup>‡</sup>Waseda University, Japan §Doshisha University, Japan

¶University of Oxford, United Kingdom

Project page

**Generated Data Samples** File 2, # Events = 2 File 4, # Events = 3 File 5, # Events = 3 File 7, # Events = 1 File 8, # Events = 1 File 11, # Events = 1 File 10, # Events = 1 \_\_\_\_\_\_ time [s] time [s]

### Pre-training Label Analysis

0.4

0.3

0.2

0.1

- Formula-SED consists of finite set of synthesis params. - We investigate which audio parameters are crucial. • F0-related params Reverberation

