

ONOMA-TO-WAVE: ENVIRONMENTAL SOUND SYNTHESIS FROM ONOMATOPOEIC WORDS

Yuki Okamoto¹, Keisuke Imoto², Shinnosuke Takamichi³,
Ryosuke Yamanishi⁴, Takahiro Fukumori¹, Yoichi Yamashita¹

¹ Ritsumeikan University, Japan ² Doshisha University, Japan

³ The University of Tokyo, Japan ⁴ Kansai University, Japan

ABSTRACT

In this paper, we propose a framework for environmental sound synthesis from onomatopoeic words. As one way of expressing an environmental sound, we can use an onomatopoeic word, which is a character sequence for phonetically imitating a sound. An onomatopoeic word is effective for describing diverse sound features. Therefore, using onomatopoeic words for environmental sound synthesis will enable us to generate diverse environmental sounds. To generate diverse sounds, we propose a method based on a sequence-to-sequence framework for synthesizing environmental sounds from onomatopoeic words. We also propose a method of environmental sound synthesis using onomatopoeic words and sound event labels. The use of sound event labels in addition to onomatopoeic words enables us to capture each sound event's feature depending on the input sound event label. Our subjective experiments show that our proposed methods achieve higher diversity and naturalness than conventional methods using sound event labels.

Index Terms— Environmental sound synthesis, sound event, onomatopoeic word, sequence-to-sequence model

1. INTRODUCTION

Environmental sound synthesis is a research field of sound generation and is the task of generating natural environmental sounds. Many environmental sounds are used in the production of movies, games, and other contents [1]. However, there is a limit to the amount of environmental sound data that is openly available. In addition, there are cases where the environmental sound that exactly matches the required sound does not exist. Therefore, it is possible to solve these problems by using environmental sound synthesis. Moreover, environmental sound synthesis has great potential for many applications such as supporting movie and game production [1, 2, 3, 4], and data augmentation for sound event detection and scene classification [5, 6].

In recent years, some methods of environmental sound synthesis using deep learning approaches have been developed [2, 7, 8]. One of the methods of environmental sound

synthesis uses sound event labels as the input [7]. The method enables the generation of environmental sounds expressing sound events. In this method, since only sound event labels are input to the system, similar sounds are generated for the given sound event class; thus, the generated sounds are not sufficiently varied. Another possibility of environmental sound synthesis is to use onomatopoeic words, which are character sequences that phonetically imitate sounds. According to the studies of Lemaitre and Rocchesso [9] and Sundaram and Narayanan [10], onomatopoeic words are effective for expressing the features of audio samples. For example, when expressing *the sound of a whistle* using onomatopoeic words, we can distinguish the sounds with different durations and pitches using the length of the phoneme sequence, such as “py u” (short whistle) and “p i i i” (long whistle). Based on the idea of mapping onomatopoeic words to environmental sounds, Kawai developed KanaWave [11], software that generates environmental sounds from onomatopoeic words. KanaWave generates environmental sounds by simply connecting multiple sounds corresponding to the input onomatopoeic words, each of which is associated with a specific sound in a one-to-one correspondence. Therefore, the sounds generated by KanaWave do not have sufficient naturalness and diversity. To utilize environmental sounds in media content, such as in animation and movie production, an environmental sound synthesis method that can generate synthesized sounds with high naturalness and diversity is required.

In this paper, we propose environmental sound synthesis from onomatopoeic words using a statistical approach. Statistical methods make it possible to automatically learn the correspondence between environmental sounds and onomatopoeic words from large amounts of data with high diversity. Even if there is a large dataset, the diversity of generated sounds is limited because the conventional method generates sounds by combining sounds in a dataset. On the other hand, a statistical method enables us to generate more diverse synthesized sounds than conventional methods. In the proposed method, we utilize the sequence-to-sequence conversion framework (seq2seq framework) [12] to gener-

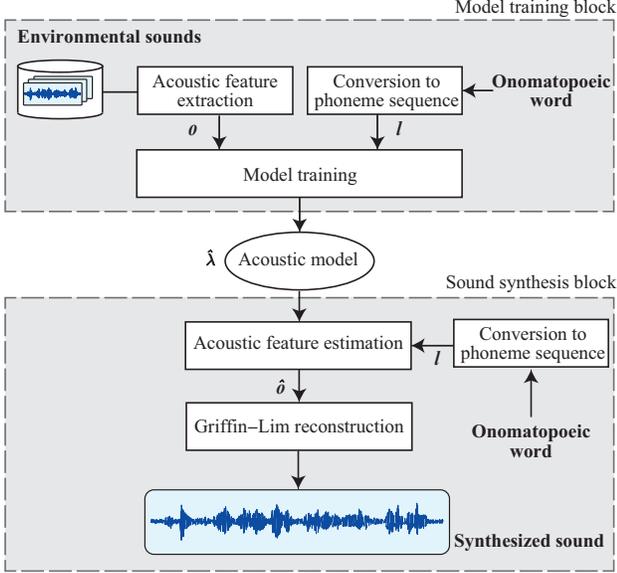


Fig. 1. Overview of environmental sound synthesis using onomatopoeia

ate environmental sounds from onomatopoeic words. The seq2seq framework is often used in sequence-to-sequence conversions, such as those in speech synthesis and neural machine translation, and has shown high performance in many studies [13, 14, 15]. The seq2seq framework uses several layers of recurrent neural network (RNN), which can model time-series information. Therefore, the seq2seq framework enables us to generate environmental sounds by considering the phoneme sequence for an onomatopoeic word. We also propose a method of environmental sound synthesis using sound event labels, which are used in the conventional method, and onomatopoeic words. The use of onomatopoeic words and sound event labels enables us to capture each sound event’s feature depending on the input sound event label.

The remainder of this paper is structured as follows. In Sec. 2, we describe the proposed methods of environmental sound synthesis from an onomatopoeic word. In Sec. 3, we report subjective experiments carried out to evaluate the performance of environmental sound synthesis from an onomatopoeic word. Finally, we summarize and conclude this paper in Sec. 4.

2. PROPOSED METHOD

2.1. Overview of Environmental Sound Synthesis from Onomatopoeic Words

Figure 1 shows the framework of environmental sound synthesis from onomatopoeic words. This approach consists of a model training block and a sound synthesis block. In

the model training block, acoustic feature sequence \mathbf{o} and phoneme sequence \mathbf{l} are extracted from environmental sounds and onomatopoeic words, respectively. Acoustic model parameter λ is estimated using extracted features \mathbf{o} and \mathbf{l} as follows:

$$\hat{\lambda} = \arg \max_{\lambda} P(\mathbf{o} | \mathbf{l}, \lambda). \quad (1)$$

In this paper, we propose two model training methods as follows.

- (I) Model training method using only onomatopoeic words as input to network (Sec. 2-B-1)
- (II) Model training method using onomatopoeic words and sound event labels as input to network (Sec. 2-B-2)

We will detail the model training methods in Sec. 2.2. In the sound synthesis block, phoneme sequence \mathbf{l} is converted from an input onomatopoeic word. Acoustic feature sequence \mathbf{o} is estimated from a phoneme sequence \mathbf{l} of the onomatopoeic word and acoustic model $\hat{\lambda}$ as follows:

$$\hat{\mathbf{o}} = \arg \max_{\mathbf{o}} P(\mathbf{o} | \mathbf{l}, \hat{\lambda}). \quad (2)$$

Finally, we reconstruct an environmental sound wave from estimated acoustic feature sequence $\hat{\mathbf{o}}$ using the Griffin-Lim algorithm [16].

2.2. Proposed Model Training Methods

2.2.1. Environmental Sound Synthesis Using Onomatopoeic Words

Figure 2 shows an overview of model training using onomatopoeic words. To synthesize environmental sounds from onomatopoeic words, we employ the seq2seq framework [12]. The seq2seq framework comprises an encoder and a decoder. Our method uses one-layered bidirectional long short-term memory (BiLSTM) as the encoder and two-layered long short-term memory (LSTM) as the decoder. As shown in Fig. 2, a phoneme sequence of the onomatopoeic word, $\mathbf{l} = \{l_1, \dots, l_T\}$, is input to the encoder. The encoder extracts feature vectors $\boldsymbol{\nu} = [\boldsymbol{\nu}^f, \boldsymbol{\nu}^b]$ from input sequence \mathbf{l} . Superscripts f and b indicate forward and backward networks, respectively. In unidirectional LSTM, the beginning features tend to be lost when the sequence is long. Therefore, using BiLSTM for the encoder, we can expect to extract a feature vector $\boldsymbol{\nu}$ that captures entire onomatopoeic words from past and future directions. The decoder estimates acoustic feature sequence $\mathbf{o} = \{\mathbf{o}_1, \dots, \mathbf{o}_{T'}\}$ from extracted feature vectors $\boldsymbol{\nu}$ in the encoder as follows:

$$p(\mathbf{o}_1, \dots, \mathbf{o}_{T'} | l_1, \dots, l_T) = \prod_{t=1}^{T'} p(\mathbf{o}_t | \boldsymbol{\nu}, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}). \quad (3)$$

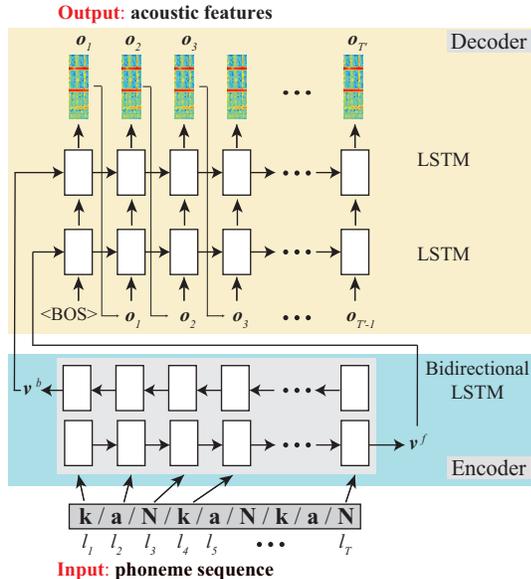


Fig. 2. Environmental sound synthesis from onomatopoeic words

Using two-layered LSTM for the decoder, we can expect to estimate acoustic features by considering features in the forward and backward directions of onomatopoeic words extracted by the encoder. The L1 norm between the estimated acoustic feature sequence \mathbf{o} and the target feature sequence at each time step is used as the loss function.

2.2.2. Environmental Sound Synthesis Using Onomatopoeic Words and Sound Event Labels

The method of environmental sound synthesis using only onomatopoeic words is expected to enable the control of the time-frequency structural features of synthesized sounds, such as sound duration. The method of environmental sound synthesis using only onomatopoeic words will generate diverse sounds. However, for example, the onomatopoeic word “p a N” could be considered to fit multiple sound events, such as *the sound of shooting guns* and *balloons breaking*. Therefore, we cannot control the frequency property associated with the sound categories using only onomatopoeic words. To control the frequency characteristics of sound events, we utilize sound event labels in addition to onomatopoeic words.

Figure 3 shows an overview of model training using onomatopoeic words and sound event labels. The method uses the seq2seq framework comprising one-layered BiLSTM as the encoder and two-layered LSTM as the decoder. The seq2seq-based intersequence conversion may involve conditioning on the decoder to control the decoder’s output features [17, 18, 19, 20]. In the proposed method, sound event labels \mathbf{c} represented as one-hot vectors and extracted feature vectors $\boldsymbol{\nu}$ are concatenated and given as the initial state of the

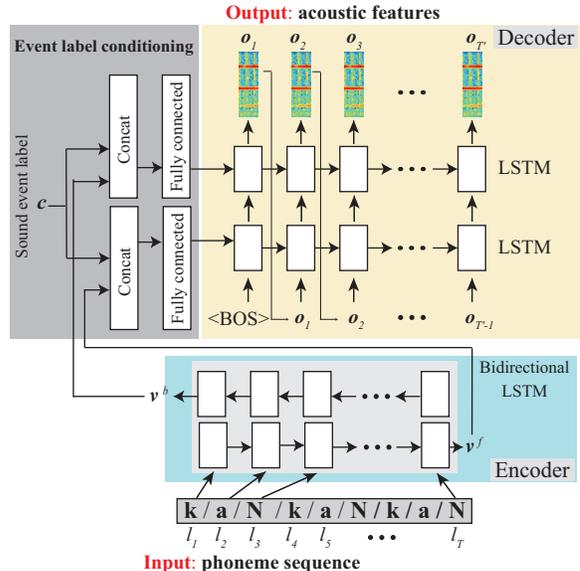


Fig. 3. Environmental sound synthesis from onomatopoeic words and sound event labels

decoder. The decoder estimates acoustic feature sequence $\mathbf{o} = \{\mathbf{o}_1, \dots, \mathbf{o}_{T'}\}$ from extracted feature vectors $\boldsymbol{\nu}$ in the encoder and sound event labels \mathbf{c} as follows:

$$p(\mathbf{o}_1, \dots, \mathbf{o}_{T'} \mid l_1, \dots, l_T) = \prod_{t=1}^{T'} p(\mathbf{o}_t \mid \boldsymbol{\nu}, \mathbf{o}_1, \dots, \mathbf{o}_{t-1}, \mathbf{c}). \quad (4)$$

The L1 norm between the estimated acoustic feature sequence \mathbf{o} and the target at each time step is used as the loss function.

3. EXPERIMENTS

The synthesized sounds must have high naturalness and diversity to use synthesized sounds as background sounds or sound effects in movies or games. From this viewpoint, we conducted two types of subjective test. For synthesized sounds, we evaluated their (I) naturalness and (II) sound diversity as environmental sounds. We aim to achieve the same level of quality as natural sound in terms of both the naturalness and diversity of the generated sound.

3.1. Experimental Conditions

For the evaluation, we used 10 types of sound event (*bell ringing, alarm clock, manual coffee grinder, cup clinking, drum, maracas, electric shaver, tearing paper, trash box banging, and whistle*) contained in the Real World Computing Partnership-Sound Scene Database (RWCP-SSD) [22]. We used a total of 1,000 samples (100 samples \times 10 sound events), in which 95 samples of each sound event were used

Table 1. Experimental conditions

Sound length	1–2 s
Sampling rate	16,000 Hz
Waveform encoding	16-bit linear PCM
Acoustic feature	log-amplitude spectrogram
Window length for FFT	0.128 s (2,048 samples)
Window shift for FFT	0.032 s (512 samples)
Encoder LSTM layers	1
# units in encoder LSTM	512
Decoder LSTM layers	2
# units in decoder LSTM	512, 512
Event label dimensions	10
Teacher forcing rate	0.6
Batch size	5
Optimizer	RAAdam [21]

for model training and the others were used for the subjective test. For the onomatopoeic words corresponding to each sound sample, we used the dataset in RWCP-SSD-Onomatopoeia [23]. Each sound sample has more than 15 onomatopoeic words, and we used 15 onomatopoeic words per audio sample for model training for a total of 14,250 onomatopoeic words (15 onomatopoeic words \times 950 audio samples). Table 1 shows the experimental conditions and parameters used for the proposed methods. In this study, we use the log-amplitude spectrogram as an acoustic feature. The generated audio samples are available on our web page¹.

3.2. Subjective Evaluations

Following the evaluation perspective described at the beginning of Sec. 3, we conducted the following two sets of experiments:

3.2.1. Experiment I: evaluation of naturalness for environmental sounds

The target sound of this paper is a sound that is comfortable as an environmental sound and that expresses the input onomatopoeic word. There are two perspectives of naturalness that should be satisfied. For this reason, we designed several experiments to evaluate each perspective. In Experiments I-1 and II-2, we presented environmental sounds and the onomatopoeic word used for the input, and evaluated how acceptable or expressive the presented sounds were in relation to the onomatopoeic word. In Experiment I-3, only the sound was presented to evaluate its naturalness as an environmental sound, and the sound itself was simply evaluated in terms of “quality as an environmental sound.”

Table 2. Number of synthesized sounds used for subjective test

Experiment	# labels	# samples in each label	# listeners	# total samples
Exp. I-1	10	10	30	3,000
Exp. I-2	10	10	30	3,000
Exp. I-3	10	5	30	1,500
Exp. II-1	5	5	30	750
Exp. II-2	10	2-3	50	1,300

Table 3. List of synthesis methods evaluated for each evaluation metric

Method	Exp. I-1	Exp. I-2	Exp. I-3	Exp. II-1	Exp. II-2
Natural sound	✓	✓	✓		
WaveNet			✓	✓	
KanaWave	✓	✓	✓		
Seq2seq (proposed)	✓	✓	✓		✓
Seq2seq + event label (proposed)	✓	✓	✓	✓	✓

- Experiment I-1: acceptance level of synthesized sounds for onomatopoeic words**
 We presented pairs of a sound (natural or synthesized) and an onomatopoeic word. The listener graded the acceptance level of the synthesized and natural sounds for onomatopoeic words on a scale of 1 (highly unacceptable) to 5 (highly acceptable).
- Experiment I-2: expressiveness of synthesized sounds for onomatopoeic words**
 We presented pairs of a sound (natural or synthesized) and an onomatopoeic word. The listener graded the expressive level of the synthesized and natural sounds for onomatopoeic words on a scale of 1 (very unexpressive) to 5 (very expressive).
- Experiment I-3: naturalness of environmental sounds**
 We presented a natural or synthesized sound. The listener graded the naturalness of the synthesized and natural sounds on a scale of 1 (very unnatural as an environmental sound) to 5 (very natural as an environmental sound).

3.2.2. Experiment II: evaluation of sound diversity for environmental sounds

To evaluate diversity for synthesized sounds, we conducted two types of subjective evaluation as follows:

¹https://y-okamoto1221.github.io/Onoma_to_wave_Demonstration/

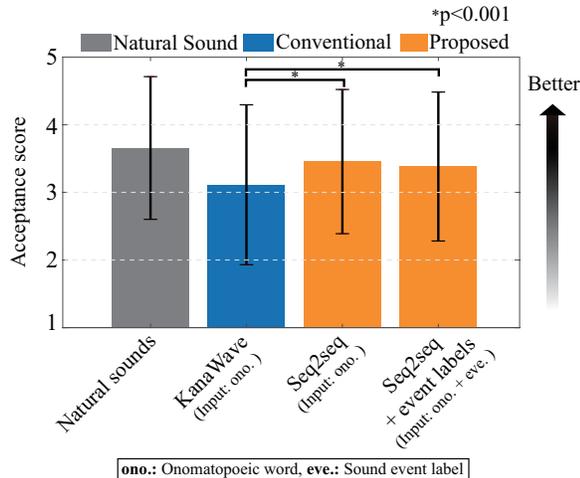


Fig. 4. Acceptance scores of natural and synthesized sounds

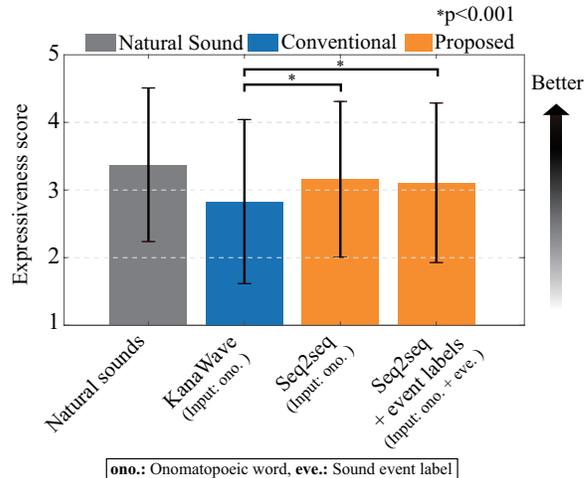


Fig. 5. Expressiveness scores of natural and synthesized sounds

- **Experiment II-1: diversity of synthesized sounds for each sound event**

We presented two sounds synthesized by the same method to listeners. In the proposed method, sounds are generated using randomly selected onomatopoeic words from the overall dataset as the input in each sound event. The listener graded the dissimilarity level between two presented sounds on a scale of 1 (very similar) to 5 (very dissimilar).

- **Experiment II-2: diversity of synthesized sounds for the same onomatopoeic words**

We presented listeners with a synthesized sound, and the listeners selected the best sound event label that represents the sound from ten choices. After listening to each sound synthesized by our methods presented randomly, the listener selected the sound event label that best represented the sound.

Each experiment was conducted using a crowdsourcing platform. Table 2 shows the numbers of audio samples and listeners in each experiment. To compare the synthesis methods, we evaluated the sounds synthesized by the conventional method using WaveNet [7] and KanaWave [11]. The conventional environmental sound synthesis method using WaveNet utilizes sound event labels as the input to the system to generate sounds. The conventional method using WaveNet [7] does not input onomatopoeic words. Therefore, we evaluated synthesized sounds by WaveNet in only experiments I-3 and II-1. KanaWave is the conventional non-statistical method of generating environmental sounds from only onomatopoeic words. The list of synthesis methods evaluated in each experiment is shown in Table 3.

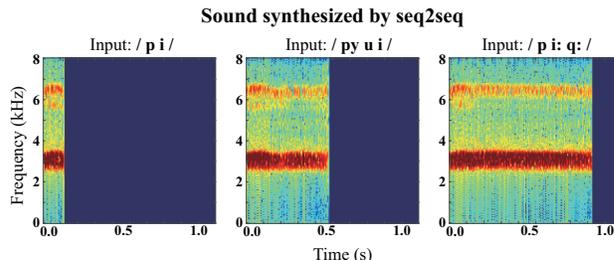


Fig. 6. Spectrograms of environmental sounds synthesized using only onomatopoeic words

3.3. Experimental Results and Discussion

3.3.1. Experiment I

Experiments I-1 and I-2: the average acceptance and expressiveness scores of synthesized and natural sounds for onomatopoeic words and their standard deviations are respectively shown in Figs. 4 and 5. The results show that our proposed methods can generate environmental sounds that are a better representation of onomatopoeic words than those generated by KanaWave.

Figure 6 shows spectrograms of sounds synthesized by our methods using only onomatopoeic words. As shown in Fig. 6, the proposed method can generate diverse environmental sounds. Also, the longest sound (right) is not the sound given by simply stretching the other sounds (left and center). Thus, onomatopoeic words are useful for generating diverse sounds with different characteristics, such as sound duration.

Figure 7 shows the spectrograms of sounds synthesized by KanaWave and the proposed method using both onomatopoeic words and sound event labels. In Fig. 7, each synthesized sound is generated from a phoneme sequence of the onomatopoeic word “b i i i i i” input to the system. In the

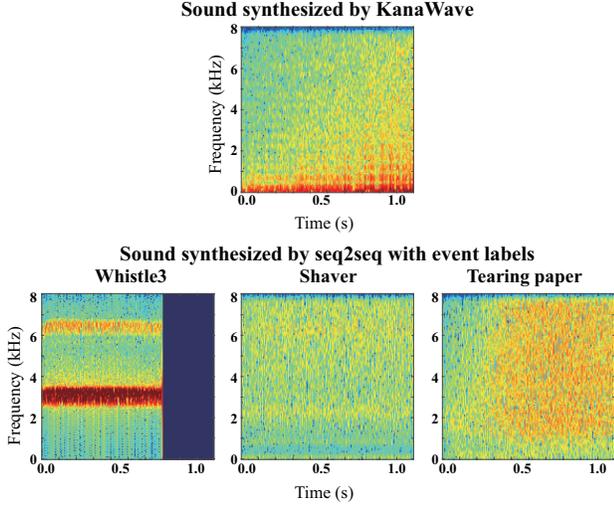


Fig. 7. Spectrograms of environmental sounds synthesized by KanaWave and the proposed method using onomatopoeic words and sound event labels

proposed method using both onomatopoeic words and sound event labels, we used sound event labels of *whistle*, *electric shaver*, and *tearing paper*. KanaWave can only generate one type of sound from the same onomatopoeic word. Therefore, the sound synthesized by KanaWave does not have diversity. On the other hand, the proposed method using onomatopoeic words and sound event labels can generate various sounds from the same onomatopoeic word by changing the input sound event labels.

Experiment I-3: the average MOS scores for the naturalness of synthesized and natural sounds, and their standard deviations are shown in Fig. 8. The results indicate that sounds synthesized by the proposed methods achieve higher naturalness than those synthesized by KanaWave. The experimental results also show that sounds synthesized by our methods had a similar sound quality to those synthesized by WaveNet. Thus, the proposed methods achieve environmental sound synthesis from onomatopoeic words without degrading the sound quality compared with conventional methods. In addition, natural sounds still have higher naturalness than sounds synthesized by the proposed methods. From these results, it is still necessary to develop a method of environmental sound synthesis that can provide quality equivalent to that of natural sounds.

3.3.2. Experiment II

Experiment II-1: the average dissimilarity score of the synthesized sound for each sound event is shown in Fig. 9. In this experiment, a high dissimilarity means that there is a rich diversity of synthesized sounds within the same type of event. The result indicates that the proposed method can generate

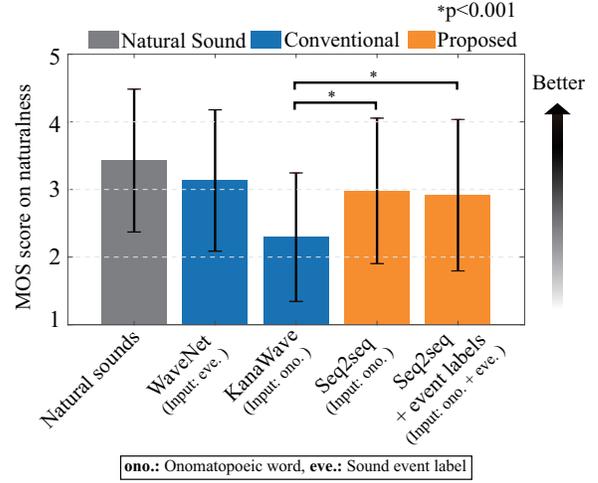


Fig. 8. MOS scores for naturalness of natural and synthesized sounds

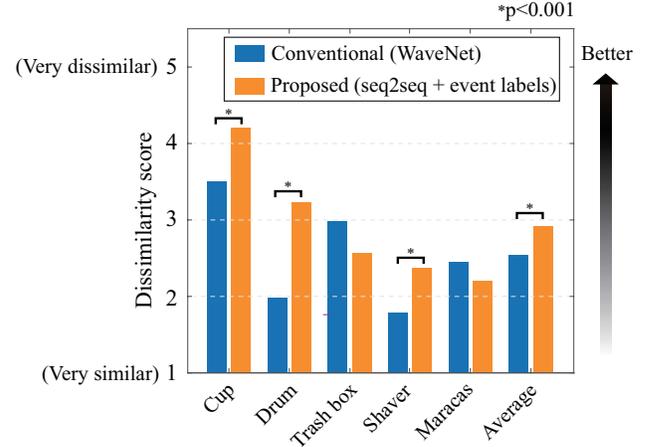


Fig. 9. Dissimilarity scores of synthesized sounds

synthesized sounds with richer diversity than the conventional method using WaveNet. In particular, the dissimilarity scores of the sound events *cup* and *drum* synthesized by the proposed methods are high, which indicates that the diversity of the synthesized sounds is richer than that of the sounds obtained by the conventional method. Thus, the proposed method enables us to generate diverse environmental sounds by using onomatopoeic words.

Experiment II-2: part of the distributions of sound event labels given to the synthesized sound from each onomatopoeic word are shown in Fig. 10. The sounds synthesized by our method using only onomatopoeic words tend to be given only one sound event label. On the other hand, the sounds synthesized by our method using both onomatopoeic words and sound event labels tend to be given various sound event labels. The entropies of the distribution of a given

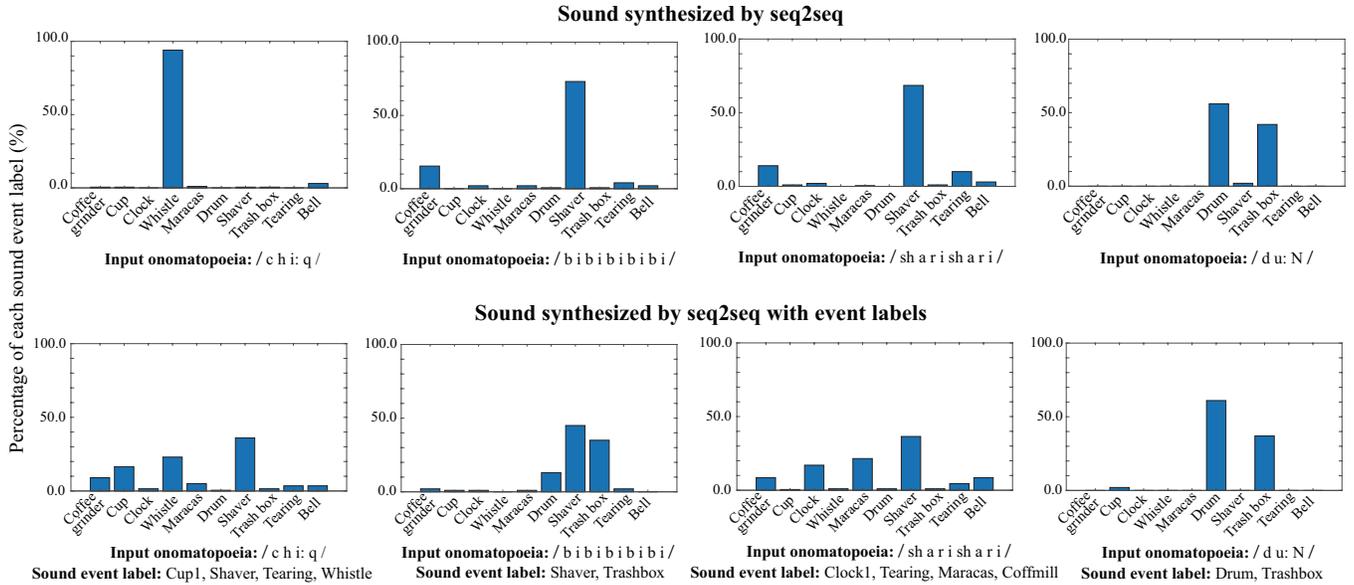


Fig. 10. Number of responses of sound event labels to each sound synthesized by our method using onomatopoeic words

acoustic event label are 1.70 bit for the method using only onomatopoeic words and 1.82 bit for the method using both onomatopoeic words and sound event labels. In this experiment, the maximum entropy is 3.02 bit when 10 types of sound event labels equally appear for each synthesized sound. This result shows that using both onomatopoeic words and sound event labels can represent multiple sound events for the same onomatopoeic word.

Figure 11 shows spectrograms of natural and synthesized sounds. In Fig. 11, each synthesized sound is generated from a phoneme sequence of the onomatopoeic word “b i: i q”. In the proposed method using both onomatopoeic words and sound event labels, we used the sound event labels of *whistle*, *electric shaver*, and *tearing paper*. As shown in Fig. 11, using only onomatopoeic words as an input generates sounds with similar features when the initial values of model parameters in model training are changed. On the other hand, using both onomatopoeic words and sound event labels, it is possible to generate sounds that capture each sound event’s feature depending on the input sound event label. These results also show that using sound event labels can control sound events of sound synthesized from onomatopoeic words.

4. CONCLUSION

In this paper, we proposed environmental sound synthesis from onomatopoeic words. We found that the proposed methods can generate sounds with high naturalness and diversity. Using sound event labels in addition to onomatopoeic words, we are also able to control not only the time-frequency

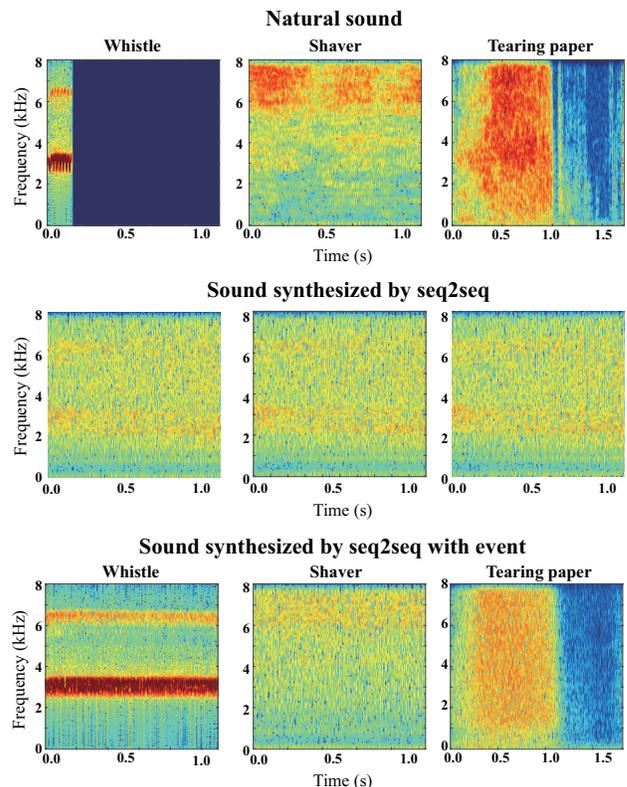


Fig. 11. Spectrograms of synthesized environmental sounds, which are generated from a phoneme sequence of the onomatopoeic word “b i: i q”, and natural sounds.

structure of the synthesized sounds but also the type of sound event. In the future, we will generate environmental sounds from onomatopoeic words using more types of sound event.

Acknowledgment

This work was supported by JSPS KAKENHI Grant Number JP19K20304 and ROIS NII Open Collaborative Research 2021 Grant Number 21S0502.

5. REFERENCES

- [1] Lloyd, D.B.; Raghuvanshi, N.; Govindaraju, K.: Sound synthesis for impact sounds in video games, in *Proc. Symposium on Interactive 3D Graphics and Games, ACM*, 2011, 55–61.
- [2] Kong, Q.; Xu, Y.; Iqbal, T.; Cao, Y.; Wang, W.; Plumbley, M.D.: Acoustic scene generation with conditional SampleRNN, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, 925–929.
- [3] Wang, K.; Cheng, H.; Liu, S.: Efficient sound synthesis for natural scenes, in *Proc. IEEE Virtual Reality (VR)*, 2017, 303–304.
- [4] Zhou, Y.; Wang, Z.; Fang, C.; Bui, T.; Berg, T.L.: Visual to sound: Generating natural sound for videos in the wild, in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, 3550–3558.
- [5] Salamon, J.; MacConnell, D.; Cartwright, M.; Li, P.; Bello, J.P.: Scaper: A library for soundscape synthesis and augmentation, in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2017, 344–348.
- [6] Gontier, F.; Lagrange, M.; Lavandier, C.; Petiot, J.F.: Privacy aware acoustic scene synthesis using deep spectral feature inversion, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, 886–890.
- [7] Okamoto, Y.; Imoto, K.; Komatsu, T.; Takamichi, S.; Yagyu, T.; Yamanishi, R.; Yamashita, Y.: Overview of tasks and investigation of subjective evaluation methods in environmental sound synthesis and conversion, *arXiv preprint arXiv:1908.10055*, 2019.
- [8] Liu, J.-Y.; Chen, Y.-H.; Yeh, Y.-C.; Yang, Y.-H.: Unconditional audio generation with generative adversarial networks and cycle regularization, *arXiv preprint arXiv:2005.08526*, 2020.
- [9] Lemaitre, G.; Rocchesso, D.: On the effectiveness of vocal imitations and verbal descriptions of sounds, *The Journal of the Acoustical Society of America*, 135 (2) (2014), 862–873.
- [10] Sundaram, S.; Narayanan, S.: Vector-based representation and clustering of audio using onomatopoeia words, in *Proc. American Association for Artificial Intelligence (AAAI) Symposium Series*, 2006, 55–58.
- [11] “KanaWave,” <https://www.vector.co.jp/soft/win95/art/se232653.html>.
- [12] Sutskever, I.; Vinyals, O.; Le, Q.V.: Sequence to sequence learning with neural networks, in *Proc. Advances in Neural Information Process. Systems (NIPS)*, 2014, 3104–3112.
- [13] Wang, Y.; Ryan, R.S.; Stanton, D.; Wu, Y.; Weiss, R.J.; Jaitly, N.; Yang, Z.; Xiao, Y.; Chen, Z.; Bengio, S.; Le, Q.; Agiomyrgiannakis, Y.; Clark, R.; Saurous, R.: Tacotron: Towards end-to-end speech synthesis, *arXiv preprint arXiv:1703.10135*, 2017.
- [14] Ikawa, S.; Kashino, K.: Generating sound words from audio signals of acoustic events with sequence-to-sequence models, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, 346–350.
- [15] Drossos, K.; Advanne, S.; Virtanen, T.: Automated audio captioning with recurrent neural networks, in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, 2017, 374–378.
- [16] Griffin, D.; Lim, J.: Signal estimation from modified short-time Fourier transform, *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 32 (2) (1984), 236–243.
- [17] Jia, Y.; Zhang, Y.; Weiss, R.J.; Wang, Q.; Shen, J.; Ren, F.; Chen, Z.; Nguyen, P.; Pang, R.; Moreno, I.L.; Wu, Y.: Transfer learning from speaker verification to multispeaker text-to-speech synthesis, in *Proc. Advances in Neural Information Process. Systems (NIPS)*, 2018, 4480–4490.
- [18] Ikawa, S.; Kashino, K.: Neural Audio captioning based on conditional sequence-to-sequence model, in *Proc. Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE)*, 2019, 99–103.
- [19] Cooper, E.; Lai, C.; Yasuda, Y.; Fang, F.; Wang, X.; Chen, N.; Yamagishi, J.: Zero-Shot multi-speaker text-to-speech with state-of-the-art neural speaker embeddings, in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, 6184–6188.

- [20] Park, J.; Zhao, K.; Peng, K.; Ping, W.: Multi-speaker end-to-end speech synthesis, *arXiv preprint arXiv:1907.04462*, 2019.
- [21] Liu, L.; Jiang, H.; He, P.; Chen, W.; Liu, X.; Gao, J.; Han, J.: On the variance of the adaptive learning rate and beyond, in *Proc. International Conference on Learning Representation (ICLR)*, 2020, 1–13.
- [22] Nakamura, S.; Hiyane, K.; Asano, F.; Endo, T.: Sound scene data collection in real acoustical environments, *The Journal of the Acoustic Society of Japan (E)*, 20 (3) (1999), 225–231.
- [23] Okamoto, Y.; Imoto, K.; Takamichi, S.; Yamashita, R.; Fukumori, T.; Yamashita, Y.: RWCP-SSD-Onomatopoeia: Onomatopoeic words dataset for environmental sound synthesis, in *Proc. Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE)*, 2020, 125–129.