

大規模な日本語笑い声コーパスを用いたテキストレス笑い声合成*

☆辛 徳泰, 高道 慎之介, 森松 亜依, 猿渡 洋 (東大院・情報理工)

1 Introduction

Human speech contains not only verbal but also nonverbal expressions like laughter, sobbing, and scream, etc. [1, 2], which can effectively convey internal affects [3, 4] of speakers in various languages and cultures [5]. Although recent advances in speech synthesis are able to synthesize natural verbal speech that is indistinguishable from human speech [6, 7, 8, 9], the progress in synthesizing nonverbal expressions is limited due to the lack of both data and technologies. In this work, we focus on a typical but important task in nonverbal-expression synthesis: laughter synthesis. Being able to synthesize laughter can intuitively improve the expressiveness and authenticity of a speech synthesis system. Such systems can be applied, for example, in virtual agents to smooth communication with users [10].

We present a method for laughter synthesis using pseudo phonetic tokens on a large-scale in-the-wild laughter corpus. In the proposed method, firstly a clustering model based on k-means [11] is trained on features extracted from the laughter utterances by a self-supervised learning (SSL) model called HuBERT [12]. The clustering model is then used to transcribe each utterance into a sequence of discrete tokens containing the phonetic information of the original laughter, which we call pseudo phonetic tokens (PPTs). A Text-to-speech (TTS) model is then trained by regarding PPTs as text inputs to synthesize laughter. Furthermore, we show it is possible to train a token language model (tLM) on the PPTs to enable unconditional laughter synthesis. Experimental results demonstrate that: (1) the proposed method significantly outperforms a baseline method that uses phonemes to represent laughter; (2) the proposed method can generate natural laughter unconditionally with the assistance of tLM. The contributions of this work are summarized as follows:

- We propose a large-scale in-the-wild Japanese laughter corpus. This corpus is, to our best knowledge, currently the largest laughter corpus that is suitable for laughter synthesis.
- We propose a method for laughter synthesis using pseudo phonetic tokens as the representation of laughter.
- We propose to train a token language model to generate PPTs and synthesize laughter unconditionally.
- We conduct comprehensive objective and subjective experiments to demonstrate the pro-

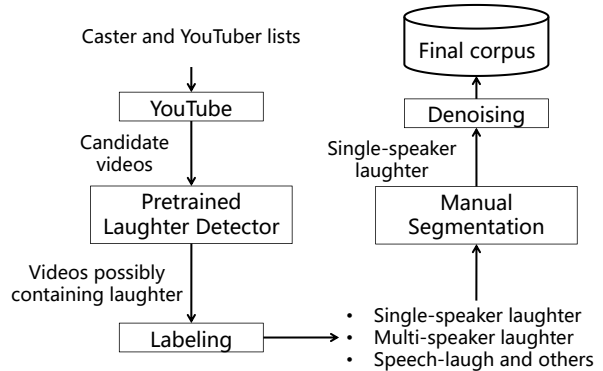


Fig. 1: Data-collection process for the proposed corpus.

posed method can synthesize natural laughter that is significantly better than a baseline method.

We publish the proposed corpus¹ and the code implementation² of the proposed method.

2 Laughter data collection

We aim to collect large-scale in-the-wild laughter utterances that are suitable for laughter synthesis. The general data-collection process is illustrated in Figure 1. We first use several lists of casters and YouTubers obtained from Wikipedia, e.g. en.wikipedia.org/wiki/List_of_YouTubers, to crawl candidate videos during June 2022 by searching the names in the list on YouTube, which results in about 10k videos. Using the lists ensures we only collect human speech instead of others like animal sounds. Second, we use an open-sourced pretrained laughter detection model (github.com/jrgillick/laughter-detection) to discover videos that possibly contain laughter, which results in about 1500 videos. However, we find that many of the detected videos include multi-speaker laughter or speech-laugh which are not suitable for synthesis. Therefore, we further conduct a listening test with crowdsourcing to label the detected videos. Specifically, we request about 1500 workers to label the videos with three categories: (1) single-speaker laughter; (2) multi-speaker laughter; (3) others including speech laugh.

After labeling, we manually segment laughter utterances from those videos that have at least one “single-speaker laughter” label. Note that, many videos contain background noises, and we only select those with non-speech noise to simplify the denois-

¹sites.google.com/site/shinnosuketakamichi/research-topics/laughter_corpus

²github.com/Aria-K-Alethia/laughter-synthesis

*Textless Laughter Synthesis with a Large-scale Japanese Laughter Corpus, Detai Xin, Shinnosuke Takamichi, Ai Morimatsu, Hiroshi Saruwatari (The University of Tokyo).

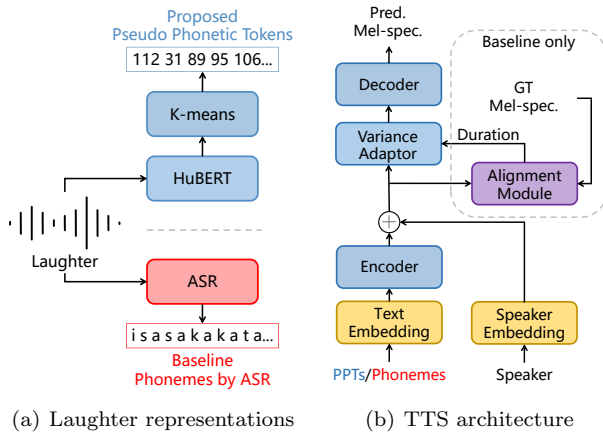


Fig. 2: Left: laughter representations used in this work; Right: architecture of the TTS model. For the baseline method an additional alignment module is used.

ing step. Besides, we discard non-Japanese videos in this step. Finally, to reduce noises in the utterances, we use a source separation model called Demucs³, which is a powerful source separation model based on DNNs, to extract the vocals from the videos. Specifically, we use the pretrained Demucs v3 (“hdemucs_mmi”) model [13] since we find it is more stable than the latest v4 model. The final corpus contains 7489 utterances of single-speaker laughter from 470 speakers. The total duration of the corpus is about 3.5 hours.

3 Laughter synthesis using pseudo phonetic tokens

3.1 The baseline method

As illustrated in the bottom half of Figure 2 (a), we use a pretrained multilingual ASR model based on wav2vec 2.0 [14] to transcribe laughter into phoneme sequences. We then adopt FastSpeech2 [8] to synthesize mel-spectrograms from the phoneme sequences. The original FastSpeech2 relies on an external alignment tool to get the duration information for each phoneme, but it is difficult to find an off-the-shelf alignment tool for the standard International Phonetic Alphabet (IPA) used by the multilingual ASR model in the baseline method. Therefore, we use an unsupervised alignment module inspired by Glow-TTS [7] that can be jointly trained with the TTS model. The module can be efficiently trained in an unsupervised manner using the connectionist temporal classification (CTC) loss [15].

3.2 Pseudo phonetic tokens

The proposed PPT is inspired by generative spoken language modeling [16, 17], which originally uses SSL models to discretize speech to do TTS in a textless manner. In this work, we further adapt this idea into nonverbal laughter. As shown in the top half of Figure 2 (a), the waveform is first fed into HuBERT [12] to convert it into continuous sequen-

tial features. Then, a k-means model [11] is trained upon the features, which can be used to convert the continuous features into discrete tokens (cluster indices). The obtained PPTs are then fed into a TTS model to synthesize laughter. The TTS model has all components used in the baseline method except for the alignment module. This is because the running length of each PPT can be regarded as its duration. For example, for a PPT sequence [21, 21, 34, 21], its duration sequence is [2, 1, 1]. Following original GSLM [16], we remove sequential repetitions (the sequence in the above example becomes [21, 34, 21] after removing) in all PPT sequences before inputting them to the phoneme encoder.

3.3 Token language model

PPT can be regarded as a symbolic representation of laughter. Thus, it is possible to train a token language model (tLM) on the PPTs of the proposed corpus. After training, one can generate laughter unconditionally by sampling from tLM.

4 Experiments

4.1 Setup

We downsample all waveforms into 16 kHz. Since the fps of HuBERT is 50, we set hop length to 320 to extract all acoustic features including pitch and mel-spectrograms. The pitch information of each utterance is extracted with WORLD vocoder [18].

We exclude utterances that are too long (over 20 s) or cannot get pitch values by the WORLD vocoder, which results in 7290 utterances. We split these utterances into train/validation/test sets with 7110/90/90 utterances, respectively. The test set consists of 30 speakers with 3 utterances per speaker, which are randomly selected from the speakers who have at least 10 utterances in the proposed corpus.

We use a pretrained multilingual wav2vec 2.0 model (XLSR) [14] fine-tuned on CommonVoice⁴ [19] as the multilingual ASR model used in the baseline method. The resulting transcriptions have 87 unique symbols in IPA.

We use the pretrained “hubert-base-ls960” model⁵ to extract the continuous sequential features used in the proposed method. We set the cluster number of k-means to 200, which means that there are 200 different PPTs used in the TTS model. We train several k-means models; most of them converge in about 250 iterations. After training, we convert all utterances into their PPT representations.

We use the same architecture of the original FastSpeech2 [8]. The dimension of the speaker embedding is set to 256. For the alignment module in the baseline method, we use exactly the same training

³github.com/facebookresearch/demucs

⁴huggingface.co/facebook/wav2vec2-xlsr-53-espeak-cv-ft

⁵huggingface.co/facebook/hubert-base-ls960

strategy used in RAD-TTS [20]. For all TTS models, the batch size is set to 16. Adam [21] is used as the optimizer with a scheduled learning rate proposed in [22]. All models converge in about 200k steps.

We use HiFi-GAN [9] as the vocoder to convert mel-spectrograms into time-domain waveforms. As the hop length of the officially released pre-trained models is not 320, we train a new HiFi-GAN vocoder from scratch on a multi-speaker Japanese corpus [23]. We use the official script⁶ to train the model.

We use fairseq [24] to train tLMs. We use the “transformer_lm” architecture, which is based on a 6-layer transformer [22]. Adam [21] is used as the optimizer with an initial learning rate of $5e-4$. The batch size is set to 16. All tLMs converge with about 30 epochs. After training, we generate 90 sequences of PPTs unconditionally for each tLM. The temperature is set to 0.7. These sequences are then inputted into the TTS model to synthesize laughter with the same speaker setting of the test set.

4.2 Objective metrics

We use several objective metrics computed on the test set or generated sequences of PPTs of laughter to evaluate the TTS models and tLMs:

- **Mel-cepstral distortion (MCD)** computed with dynamic time warping (DTW).
- **F0 root mean square error (F0-RMSE)** computed with DTW.
- **Perplexity (PPL)** defined as the normalized inverse probability on the test set of the tLM.
- **Self-BLEU [25]** defined as the average value of the n -gram (4-gram in this work) BLEU scores [26] between one generated sentence and the rest generated sentences for all generated sentences.

Here MCD and F0-RMSE reflect the quality of the synthesized laughter; PPL and Self-BLEU reflect the performance of the tLM and the diversity of the generated sentences, respectively. In particular, since each tLM has a unique set of PPTs, in this work we propose to use a normalized version of Self-BLEU that is defined as the ratio of the Self-BLEU of the generated sentences to the Self-BLEU of the test set: $\overline{\text{Self-BLEU}} = \text{Self-BLEU} / \text{Self-BLEU}_{\text{gt}}$. This metric has a value between $[0, 1]$, and can reflect how diverse the generated sentences are compared to the GT sentences.

4.3 Laughter synthesis

4.3.1 Layer selection of HuBERT

As the output of each layer of HuBERT is possible to be used as the features for PPTs, we train 12 proposed models and compute the objective metrics to select the best layer. The result is illustrated in

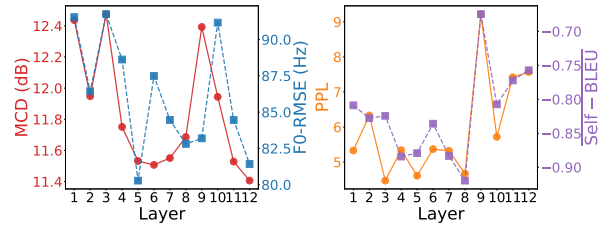


Fig. 3: Objective performance of the proposed method using the output of different layers of HuBERT as the feature for PPTs. Negative Self-BLEU scores are shown for the ease of comparison.

Table 1: Performance of all methods in the evaluations of laughter synthesis. **Bold** indicates the best score with $p < 1e-5$ comparing to Baseline.

Model	MCD(L)	F0-RMSE(L)	MOS(†)	SMOS(†)
GT	-	-	3.73	-
HiFi-GAN	6.68	53.70	3.31	4.74
Baseline	16.59	117.65	1.25	1.20
Baseline GT	10.74	85.69	-	-
Proposed-L5	11.53	80.28	3.00	3.07
Proposed-L8	11.69	82.81	2.98	3.17
Proposed-L12	11.41	81.43	2.96	3.22

Figure 3. Hereafter we use $L\{1, 2, \dots, 12\}$ to denote the proposed method using the corresponding layer of HuBERT for simplicity. It can be seen that the performances of the TTS models (left) and the tLMs (right) are not consistent. Specifically, L12 has the best performance among the TTS models, but L8 has the best performance among the tLMs. Besides, L5 has good performance in all metrics. Therefore, we use L5, L8, and L12 in the following evaluations. We also tried the proposed method with a fewer or larger cluster number but found no improvements in the preliminary experiments.

4.3.2 Comparison to the baseline

Next, we compare the proposed method to the baseline method. In addition to the objective metrics, we also use subjective mean opinion score (MOS) and similarity MOS (SMOS) to evaluate the naturalness and similarity of the synthesized laughter, respectively.

All results are shown in Table 1. First, the baseline method has poor performance in both the objective and subjective evaluations. To verify if this is because the model fails to learn from the inputted phonemes, we further use GT acoustic features (pitch and energy) to synthesize the test utterances. The corresponding model is denoted as “Baseline GT” in Table 1. It can be seen that the performance becomes comparable to the proposed method, which implies that the laughter representation makes the performance of the baseline method bad. Second, it can be seen that the 3 proposed models have significantly better performance than the baseline method in all metrics, which demonstrates the effectiveness of the proposed method using PPTs as the representation for laughter. Finally, we observe that L5 has the best naturalness and L12 has the best speaker similarity, which is consistent

⁶github.com/jik876/hifi-gan

Table 2: Subjective performance of the proposed method in the evaluation of unconditional laughter generation. **Bold** indicates the best score with $p < 1e-5$.

Model	MOS(↑)	SMOS(↑)
Proposed-L5	3.11	2.65
Proposed-L8	2.80	2.59
Proposed-L12	3.06	2.59

with the results in objective metrics as both of the two models have better objective performance than L8.

4.4 Unconditional laughter generation

Finally we evaluate the performance of unconditional laughter generation with a MOS test and a SMOS test. Given the poor performance of the baseline method shown in the previous section, we only use the three proposed models in this evaluation. 27 listeners join in the MOS test; each evaluates 33 utterances of which the first 3 are dummy samples. As a result, each utterance has 3 answers. The SMOS test is conducted using exactly the same setting of the MOS test.

The result is shown in Table 2. It can be seen that generally $L5 > L12 > L8$. This is quite different from the performance of tLMs shown in the right side of Figure 3, in which $L8 > L5 > L12$. We suppose this is because the quality of the synthesized laughter is mainly determined by the performance of the TTS models. However, it should be pointed out that MOS and SMOS cannot evaluate the diversity of the synthesized laughter subjectively. We leave this as future work. Combining this result with the result shown in Table 1, we conclude that layer 5 is the best layer of HuBERT for PPTs used in laughter synthesis.

謝辞：本研究は、JST 次世代研究者挑戦的研究プログラム JPMJSP2108、JSPS 科研費 JP23KJ0828、JST 創発的研究支援事業 JPMJFR22 の支援を受けたものです。

参考文献

- [1] K. R. Scherer, "Affect bursts," *Emotions: Essays on emotion theory*, vol. 161, p. 196, 1994.
- [2] J. Trouvain, K. P. Truong, "Comparing non-verbal vocalisations in conversational speech corpora," in *Proc. LREC*. Citeseer, 2012, pp. 36–39.
- [3] J. A. Hall et al., "Psychosocial correlates of interpersonal sensitivity: A meta-analysis," *Journal of nonverbal behavior*, vol. 33, no. 3, pp. 149–180, 2009.
- [4] K. R. Scherer, U. Scherer, "Assessing the ability to recognize facial and vocal expressions of emotion: Construction and validation of the emotion recognition index," *Journal of Nonverbal Behavior*, vol. 35, no. 4, pp. 305–326, 2011.
- [5] D. A. Sauter et al., "Cross-cultural recognition of basic emotions through nonverbal emotional vocalizations," *Proceedings of the National Academy of Sciences*, vol. 107, no. 6, pp. 2408–2412, 2010.
- [6] Y. Ren et al., "Fastspeech: Fast, robust and controllable text to speech," *Proc. NeurIPS*, vol. 32, 2019.

- [7] J. Kim et al., "Glow-tts: A generative flow for text-to-speech via monotonic alignment search," *Proc. NeurIPS*, vol. 33, pp. 8067–8077, 2020.
- [8] Y. Ren et al., "Fastspeech 2: Fast and high-quality end-to-end text to speech," in *Proc. ICLR*, 2021.
- [9] J. Kong et al., "Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis," *Proc. NeurIPS*, vol. 33, pp. 17 022–17 033, 2020.
- [10] K. El Haddad et al., "Speech-laugh: an hmm-based approach for amused speech synthesis," in *Proc. ICASSP*. IEEE, 2015, pp. 4939–4943.
- [11] J. MacQueen, "Classification and analysis of multivariate observations," in *5th Berkeley Symp. Math. Statist. Probability*. University of California Los Angeles LA USA, 1967, pp. 281–297.
- [12] W.-N. Hsu et al., "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [13] A. Défossez, "Hybrid spectrogram and waveform source separation," in *Proceedings of the ISMIR 2021 Workshop on Music Source Separation*, 2021.
- [14] Q. Xu et al., "Simple and effective zero-shot cross-lingual phoneme recognition," in *Proc. Interspeech*, 2022.
- [15] A. Graves et al., "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006, pp. 369–376.
- [16] K. Lakhotia et al., "On generative spoken language modeling from raw audio," *Transactions of the Association for Computational Linguistics*, vol. 9, pp. 1336–1354, 2021.
- [17] F. Kreuk et al., "Textless speech emotion conversion using decomposed and discrete representations," in *Proc. EMNLP*, 2022, pp. 11 200–11 214.
- [18] M. Morise et al., "World: a vocoder-based high-quality speech synthesis system for real-time applications," *IEICE TRANSACTIONS on Information and Systems*, vol. 99, no. 7, pp. 1877–1884, 2016.
- [19] R. Ardila et al., "Common voice: A massively-multilingual speech corpus," in *Proc. LREC*, 2020, pp. 4218–4222.
- [20] K. J. Shih et al., "Rad-tts: Parallel flow-based tts with robust alignment learning and diverse synthesis," in *ICML Workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models*, 2021.
- [21] D. P. Kingma, J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [22] A. Vaswani et al., "Attention is all you need," *Proc. NeurIPS*, vol. 30, 2017.
- [23] S. Takamichi et al., "JVS corpus: free Japanese multi-speaker voice corpus," *arXiv preprint arXiv:1908.06248*, 2019.
- [24] M. Ott et al., "fairseq: A fast, extensible toolkit for sequence modeling," in *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- [25] Y. Zhu et al., "Texygen: A benchmarking platform for text generation models," in *The 41st international ACM SIGIR conference on research & development in information retrieval*, 2018, pp. 1097–1100.
- [26] K. Papineni et al., "Bleu: a method for automatic evaluation of machine translation," in *Proc. ACL*, 2002, pp. 311–318.