

PAPER

A Method to Detect Chorus Sections in Lyrics Text*

Kento WATANABE^{†a)}, Nonmember and Masataka GOTO^{†b)}, Fellow

SUMMARY This paper addresses the novel task of detecting chorus sections in English and Japanese lyrics text. Although chorus-section detection using audio signals has been studied, whether chorus sections can be detected from text-only lyrics is an open issue. Another open issue is whether patterns of repeating lyric lines such as those appearing in chorus sections depend on language. To investigate these issues, we propose a neural-network-based model for sequence labeling. It can learn phrase repetition and linguistic features to detect chorus sections in lyrics text. It is, however, difficult to train this model since there was no dataset of lyrics with chorus-section annotations as there was no prior work on this task. We therefore generate a large amount of training data with such annotations by leveraging pairs of musical audio signals and their corresponding manually time-aligned lyrics; we first automatically detect chorus sections from the audio signals and then use their temporal positions to transfer them to the line-level chorus-section annotations for the lyrics. Experimental results show that the proposed model with the generated data contributes to detecting the chorus sections, that the model trained on Japanese lyrics can detect chorus sections surprisingly well in English lyrics, and that patterns of repeating lyric lines are language-independent.

key words: lyrics information processing, music information retrieval, natural language processing, lyrics structure analysis

1. Introduction

The digitization of lyrics collections has opened various areas of lyrics-based research, such as research on lyrics browsing [2]–[4], lyrics genre classification [5]–[7] and lyrics-to-audio synchronization [8]–[18]. Lyrics are usually plain text without any annotations, and some researchers have analyzed their structure, such as paragraph structure and topic transitions between paragraphs [19]–[23]. For example, Fell et al. [19] and Watanabe et al. [20] estimated section boundaries in lyrics text without empty lines but were not able to assign a section label such as verse or chorus to each estimated section. Chorus sections were not detected in lyrics text.

The goal of this paper is to achieve automatic chorus-section detection for lyrics text. This task has not been studied, though chorus-section detection, as well as music structure analysis, for audio signals has been a popular topic of research [24]–[42]. Since whether chorus sections can be

detected from text-only lyrics is an open issue, it is worth investigating this issue from an academic viewpoint. Moreover, a chorus-section detection method for lyrics text has potential applications. For example, when listeners want to find lyrics with a chorus section having a particular phrase such as “I love you” for the purpose of singing that section or reusing it in a short video clip, it is necessary for a lyrics search system to automatically detect which lines of the lyrics are included in chorus sections. The detected lyric lines of chorus sections could be used in a lyrics viewing function of music services displaying lyrics with those lines highlighted by a different color or typeface. Automatic lyric video generation technologies could give those lines more vivid animations.

Chorus sections are the most repeated and memorable portions of a song [40]. Since it is not easy to find such sections by exploiting heuristic rules, most existing chorus-section detection methods for audio signals have leveraged repetitive patterns of those sections within a song. In this paper, we propose a supervised model that can detect chorus sections in English and Japanese lyrics. Our model uses both structural features that represent patterns of repeating lyric lines, and linguistic features that are calculated either from word2vec [43] with context2vec [44] or from BERT [45]. To detect chorus sections using only plain text without any labels or even empty lines (i.e., section boundaries), we investigate a model and features effective for chorus-section detection. Experimental results show that our proposed model outperforms alternative baseline models and that combining structural and linguistic features contributes to better performance.

Although such a supervised model needs a large dataset of lyrics with line-level chorus-section annotations for its training, there was no such dataset as there was no prior work on chorus sections in lyrics text. To address this issue of lacking training data, we generated a dataset consisting of 9,313 English and 91,459 Japanese lyrics with chorus-section annotations by utilizing pairs of musical audio signals and their corresponding manually time-aligned lyrics. We first automatically detected chorus sections in audio signals of a song [40]. Then, since each lyric line had the corresponding start time within the song, we could find lyric lines that temporally correspond to the duration of each detected chorus section. We thus obtained the annotated dataset by assigning a chorus label to those lyric lines and a not-chorus label to the other lines. Experimental results show that the model trained with this large auto-

Manuscript received July 29, 2022.

Manuscript revised March 30, 2023.

Manuscript publicized June 2, 2023.

[†]The authors are with National Institute of Advanced Industrial Science and Technology (AIST), Tsukuba-shi, 305–8568 Japan.

*An earlier version of this paper was published at a conference [1].

a) E-mail: kento.watanabe@aist.go.jp

b) E-mail: m.goto@aist.go.jp

DOI: 10.1587/transinf.2022EDP7139

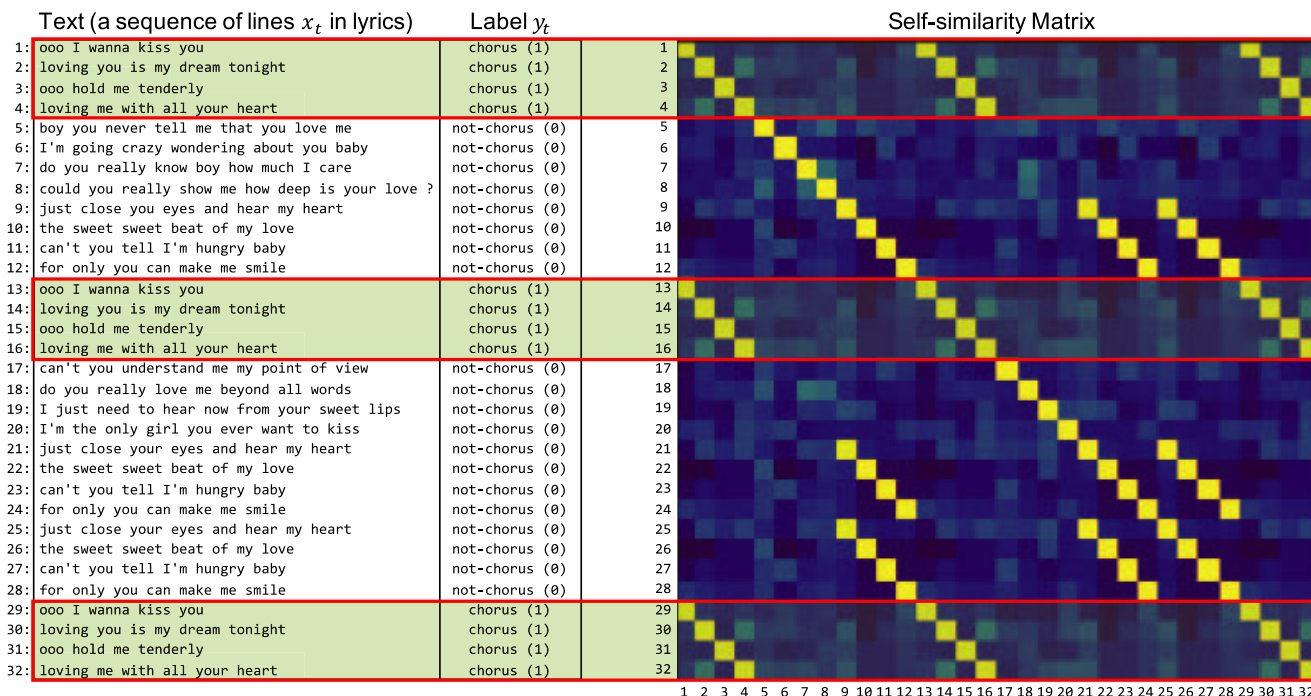


Fig. 1 Example of lyrics with chorus-section annotations and corresponding self-similarity matrix in which each cell represents the similarity between two lyric lines. These lyrics are from “How Deep Is Your Love?” (RWC-MDB-P-2001 No. 81 in the RWC Music Database [46]).

matically generated dataset performs better than the model trained with a smaller manually annotated dataset and that the model trained on Japanese lyrics can detect chorus sections surprisingly well in English lyrics.

2. Lyrics Chorus-Section Detection Task

The left side of Fig. 1 shows an example of lyrics with chorus-section annotations (labels). The lyrics of a song are a sequence of lyric lines, each line having a sentence or phrase. In this example there are three highlighted chorus sections that have exactly the same four lines, though in other songs, lyrics of chorus sections are repeated with some modifications. To maximize the applicability, as shown in this example, we assume that the input text of lyrics does not have any section boundaries. Even though some lyrics contain empty lines at those boundaries, those lines are deleted in advance. We also assume that the input text does not have explicit chorus labels such as “(chorus)” at the beginnings of chorus sections. Even though some lyrics contain those labels, they are deleted as well. When lyrics contain a repetition label such as “(* repeat)”, it is manually replaced with the corresponding lyric lines.

We formulate this chorus-section detection task as a sequence labeling problem: predicting the chorus or not-chorus status for each lyric line. Let X_s be the lyrics of a song s composed of T lines of text: $X_s = \{x_1, \dots, x_t, \dots, x_T\}$. Each lyric line x_t has a binary label y_t . If $y_t = 1$, x_t is in a chorus section. If $y_t = 0$, x_t is not in a chorus section. Y_s denotes a sequence of labels correspond-

ing to X_s : $Y_s = \{y_1, \dots, y_t, \dots, y_T\}$. In the training step, the model learns the conditional probability $P(Y_s|X_s)$. In the validation/testing step, the trained model has to predict labels Y_s for given lyric lines X_s .

Chorus sections cannot be detected by simply extracting repeated lines since those lines often correspond to non-chorus sections. For example, lyric lines 9–12 and 21–24 in Fig. 1 are exactly repeated, but those lines are not in chorus sections. It is also difficult to manually define a set of rules to find various chorus sections. We therefore prepare various features that could be useful for machine learning to deal with various types of chorus sections.

3. Computational Modeling of Chorus Sections in Lyrics

We propose a neural-network-based model for sequence labeling by using structural features that are self-similarity matrix (SSM) representations. SSM representations are widely used in computational music structure analysis, but we use different representations for lyrics. In addition to structural features, our model utilizes two kinds of linguistic features widely used in natural language processing (NLP): (1) word vectors and sentence vectors calculated from word2vec [43] and context2vec [44], and (2) sentence vectors calculated from BERT [45].

In the following sections, we first describe nine SSMs for capturing patterns of repeating lyric lines and explain how to encode the SSMs for neural networks (Sect. 3.1). We then describe the linguistic features obtained by vectorizing

the semantic/syntactic information of lines using word2vec, context2vec, and BERT (Sect. 3.2). Finally, we describe a neural-network-based sequence labeling model with these structural and linguistic features (Sect. 3.3).

3.1 Structural Features

Most previous work on music structure analysis for audio signals [24]–[42] identifies repeated musical sections by using a SSM like that shown in Fig. 1. Repeated sections lead to high values in diagonals of the matrix, and those patterns are used to identify the structure. To capture repeated lyric lines that often appear in chorus sections, we also compute the SSM from lyrics text, but the design of the similarity measure to compute each cell of the SSM is important. We propose to use the following nine variations of similarity measures sim_m , where m denotes the variation. Some of the similarities are based on previous studies [19], [20].

1. **String similarity** (sim_{str}): a normalized Levenshtein edit distance [47] between the characters of two lyric lines.
2. **Head similarity** (sim_{head}): a normalized Levenshtein edit distance between the characters of the first two words of two lyric lines.
3. **Tail similarity** (sim_{tail}): a normalized Levenshtein edit distance between the characters of the last two words of two lyric lines.
4. **Phonetic similarity** (sim_{phone}): To capture rhymes in the lyrics, we calculate a normalized Levenshtein edit distance between the phonetic transcriptions of two lyric lines. We use the CMU pronunciation dictionary[†] to extract the phonetic transcription. For example, the phonetic transcription of “I love you” is [AY1, L, AH1, V, Y, UW1].
5. **Part-of-speech similarity** (sim_{pos}): To capture similarities in grammatical structure, we calculate a normalized Levenshtein edit distance between the part-of-speech (POS) sequences of two lines. We use the default POS tagger in the NLTK package [48].
6. **Word vector similarity** (sim_{w2v}): To capture the semantic similarity between two lyric lines, we simply average vectors of the words of each lyric line by using pre-trained word2vec [43] and compute their cosine similarity. This “bag of words” representation does not differentiate “dog bites person” from “person bites dog”.
7. **Context vector similarity** (sim_{c2v}): To consider the word order, we vectorize the lyric lines using pre-trained context2vec [44], an extension of word2vec, which encodes a sequence of words by using Long Short-Term Memory (LSTM) networks [49]. We then compute their cosine similarity to obtain sim_{c2v} .
8. **Word syllable count similarity** (sim_{syw}): Since repeated phrases sometimes have the same number of syllables even if their words are different, we use a

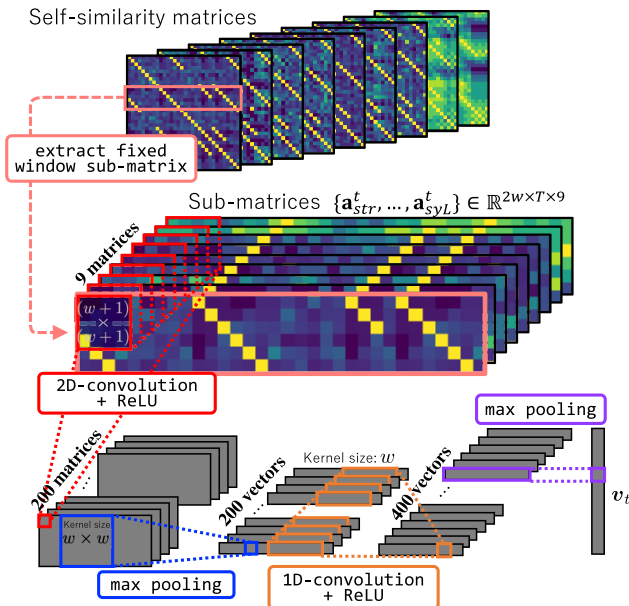


Fig. 2 Convolutional neural network for SSMs.

sequence of word syllable counts on each lyric line. For example, the word syllable counts of the two lyric lines “Sometimes you lost yourself away” and “Every-time you just close your eyes”^{††} are {2, 1, 1, 2, 1} and {2, 1, 1, 1, 1}, respectively. When successive lyric lines have similar syllable count sequences, they are likely to correspond to the repetition of sections. We use dynamic time warping (DTW) [50] to calculate the similarity between syllable count sequences.

9. **Lyric Line syllable count similarity** (sim_{syL}): We can also use the total syllable count of all words in each lyric line. For example, in all the chorus sections shown in Fig. 1, the total syllable count of the first lyric line is 6 and that of the second line is 8. We calculate the similarity of such total syllable counts of each pair of lyric lines by using the following procedure. (1) We extract a window of four lyric lines $L_t = \{x_t, x_{t+1}, x_{t+2}, x_{t+3}\}$ and shift it over the entire lyrics of a song. (2) The similarity between the lyric lines x_t and $x_{t'}$ is calculated by DTW of L_t and $L_{t'}$.

We thus calculated nine SSMs $\mathbf{A}_m \in \mathbb{R}^{T \times T}$, where each cell is a sim_m explained above. Then, to calculate feature vectors from the above nine SSMs, we exploit a convolutional neural network (CNN) architecture to detect textual macro structures from various patterns in SSMs regardless of their locations and relative sizes in SSMs. Except for network parameters, this CNN architecture is the same as that of Fell et al. [19], as we share the same motivation: to extract translation, scaling and rotation invariant features from the input image (in our case, nine SSMs).

Figure 2 illustrates the CNN structure. After calculating the nine SSMs, we extract fixed-size elongated-rectangle

[†]<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

^{††}These lyrics are taken from the RWC Music Database (RWC-MDB-P-2001 No. 92) [46].

sub-matrices centered on the target lyric line x_t : $\mathbf{a}_m^t = \mathbf{A}_m[t-w+1, \dots, t+w; 1, \dots, T] \in \mathbb{R}^{2w \times T}$, where w is a fixed window size. The input of the CNN is nine sub-matrices $\{\mathbf{a}_{str}^t, \dots, \mathbf{a}_{syL}^t\} \in \mathbb{R}^{2w \times T \times 9}$, where the number of channels corresponds to the number of SSMs (i.e., 9). The kernel size of the first 2D-convolutional layer is $(w+1) \times (w+1)$ so that each feature can capture a prospective chorus section. This first 2D-convolutional layer has 200 kernels, so this convolution produces 200 matrices from the nine sub-matrices (bottom left of Fig. 2). Each of the resulting 200 matrices is downsampled to a vector by max-pooling with $w \times w$ kernel size, so this produces 200 vectors (bottom center of Fig. 2). We then apply the 1D-convolutional layer with a kernel size of w to the 200 vectors (200 channels). This 1D-convolutional layer has 400 kernels, so the convolution produces 400 vectors from the 200 vectors (bottom right of Fig. 2). The final max-pooling layer reduces the dimension of each of the resulting 400 vectors to one dimension, resulting in the final 400-dimensional vector \mathbf{v}_t .

In this network, all convolutional layers employ the ReLU function. We perform the above procedure independently for each lyric line x_t and obtain the CNN-based feature vector \mathbf{v}_t . We call this vector \mathbf{v}_t the structural feature sim_{all} .

3.2 Linguistic Features

Some expressions tend to appear in chorus sections. To quantify this tendency, we calculate the difference between word tri-gram probabilities in the chorus and non-chorus sections. Table 1 shows the word tri-grams that frequently appear in both of the sections. Here, P_c and P_n denote word tri-gram probabilities in the chorus and non-chorus sections, respectively. As shown in this table, we found that phrases about the future (e.g., “I’ll” and “Let’s”) tend to appear in chorus sections more often than do phrases about the past (e.g., “have been” and “didn’t”). To exploit such tendencies, we compute two kinds of linguistic features:

- word2vec/context2vec-based linguistic feature ($ling_{ave+seq}$):** For each lyric line x_t , we first calculate the average of word vectors obtained using pre-trained word2vec [43], skipping out-of-vocabulary

Table 1 Frequent word tri-grams in chorus and non-chorus sections. An apostrophe is regarded as a word.

Tri-gram	$P_c - P_n$	Tri-gram	$P_n - P_c$
I'm	0.12%	there's	0.04%
don't	0.11%	I've	0.03%
oh oh oh	0.05%	's a	0.03%
I'll	0.05%	I'd	0.02%
we're	0.04%	but I'	0.02%
you're	0.04%	's not	0.01%
'll be	0.04%	what's	0.01%
I don'	0.04%	na na na	0.01%
Let's	0.03%	yeah yeah yeah	0.01%
you got ta	0.03%	've been	0.01%
I can'	0.03%	't take	0.01%
can't	0.03%	didn't	0.01%

words. Since the word order cannot be modeled by word2vec, we then use pre-trained context2vec [44] that puts a sequence of the word vectors based on word2vec into the LSTM to obtain a sentence vector. We finally concatenate the averaged word2vec-based word vector with the context2vec-based sentence vector to obtain the concatenated vector $\mathbf{u}_t^{ave+seq}$. We call this vector $\mathbf{u}_t^{ave+seq}$ the linguistic feature $ling_{ave+seq}$.

- BERT-based linguistic feature ($ling_{BERT}$):** Since BERT has been reported to improve performance in various NLP tasks [45], we calculate the BERT-based feature vector \mathbf{u}_t^{BERT} in addition to $\mathbf{u}_t^{ave+seq}$ so that we can compare them. We first feed each lyric line x_t to the pre-trained BERT model. Among the output vectors of BERT, we then obtain the vector resulting from the position of the [CLS] token and use it as the sentence embedding vector \mathbf{u}_t^{BERT} . We call this vector \mathbf{u}_t^{BERT} the linguistic feature $ling_{BERT}$.

As $ling_{BERT}$ is expected to be better than $ling_{ave+seq}$, we compared their performances in our experiments.

3.3 Neural-Network-based Sequence Labeling Model

To solve the sequence labeling problem, we use the standard Bidirectional Long Short-Term Memory (Bi-LSTM) networks [51] to compute the conditional probability $P(Y_s|X_s)$. The neural network structure is illustrated in Fig. 3.

The input to the Bi-LSTM layer at each time step t (lyric line x_t) is a concatenation of two different types of feature vectors: (1) the structural feature vector \mathbf{v}_t encoded from the nine variations of SSMs in Sect. 3.1 and (2) the linguistic feature vector $\mathbf{u}_t^{ave+seq}$ or \mathbf{u}_t^{BERT} encoded in Sect. 3.2. Formally, the conditional probability $P(Y_s|X_s)$ is calculated by using a softmax function:

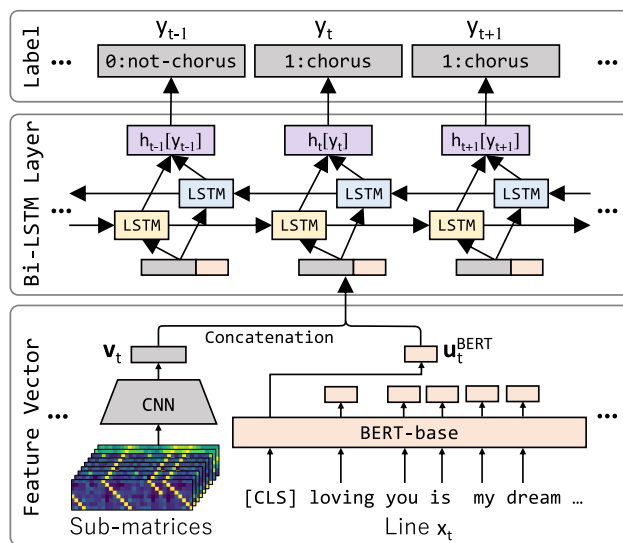


Fig. 3 Neural-network-based sequence labeling model for chorus-section detection. The BERT-based feature vector \mathbf{u}_t^{BERT} is used as the linguistic feature in this figure.

$$P(Y_s|X_s) = \frac{\exp(\text{Score}(X_s, Y_s))}{\sum_{Y'_s} \exp(\text{Score}(X_s, Y'_s))}. \quad (1)$$

The Score() is defined as

$$\text{Score}(X_s, Y_s) = \sum_t \text{BN}(h_t[y_t]), \quad (2)$$

where $h_t[y_t]$ is the output of the Bi-LSTM for each time step t and $\text{BN}()$ denotes batch normalization [52]. In the model training step, we use a binary cross-entropy loss.

4. Experiment

Inspired by audio-based chorus-section detection [40], we evaluated the proposed method by using the F-measure (F) that is a harmonic mean of precision (P) and recall (R), $F = (2 \cdot R \cdot P)/(R + P)$, where

$$P = \frac{\# \text{ of lyric lines in correctly detected chorus sections}}{\# \text{ of lyric lines in detected chorus sections}},$$

$$R = \frac{\# \text{ of lyric lines in correctly detected chorus sections}}{\# \text{ of lyric lines in correct (annotated) chorus sections}}.$$

We also used the pair-wise F-measure (p - F), normalized conditional entropy F-measure (n - F) and V-measure (V) that are provided by the Python module `mir_eval` and commonly used to evaluate computational music structure analysis [53].

4.1 Methods Compared

To confirm the effectiveness of our Bi-LSTM method based on the Bi-LSTM model that can learn dependencies between adjacent lyric lines, we compared its performance with those of two baseline methods:

1. **Heuristic:** We implemented the heuristic that “if lines at the end of the lyrics are repeated with small modifications, all those repeated lines are chorus sections” by the following procedure: (i) From the SSM that is the average of the nine SSMs, we extracted diagonals whose cells had values higher than a threshold λ , which was tuned on a development set to be 0.62. (ii) From the extracted diagonals, we selected the shortest diagonal among diagonals placed at the bottom of the SSM (e.g., the diagonal starting at the cell `SSM[29; 1]` in Fig. 1). (iii) Successive lines corresponding to the rows where the selected diagonal was located (e.g., lyric lines 29–32 in Fig. 1) were assigned the label **chorus**. (iv) Other successive lines that were similar to the chorus lines (e.g., lyric lines 1–4 and 13–16 in Fig. 1) were also assigned **chorus** labels.
2. **Multi-Layer Perceptron (MLP):** Similar to the Bi-LSTM method, but with the Bi-LSTM model replaced by a standard MLP model. This method ignores transitions between adjacent lyric lines and predicts y_t from x_t only.

We chose the number of kernels for the first and second CNNs to be 200 and 400, respectively. We used $w = 3$

for the window size. In the MLP and Bi-LSTM methods, we chose the dimension of the hidden state to be 600. The word2vec [43] and context2vec [44] were pre-trained on lyrics and were not updated in the model training step of our method. The dimension of their output vectors was 300. We used pre-trained BERT models [45] that are publicly available[†]. The dimension of the BERT-based feature vector is 768. We used AdamW for parameter optimization [54]. The initial learning rate was 0.001 with an exponential decay. We used a mini-batch size of 64. Training was run for 100 epochs, and the model used for testing was the one that achieved the best F-measure on the development set.

4.2 Dataset

To train our computational model that predicts whether the label of each lyric line is **chorus** or **not-chorus**, we needed a large amount of lyrics data with line-level chorus-section annotations like those illustrated in Fig. 1. Since there was no dataset for this, we generated a large amount of such lyrics data by the following procedure:

1. We prepared 100,772 pairs of musical audio signals and their corresponding manually time-aligned (temporally synchronized) lyrics^{††}. To avoid unreliable lyrics, we made sure that all lyrics had more than eight lines and less than 120 lines.
2. We detected chorus sections of every song automatically by using its audio signals. In our experiments, we used the RefraiD method [40] to obtain the start and end times of each chorus section, but other methods could also be used.
3. If the start time of a lyric line was within any chorus section detected in audio signals, that line was labeled **chorus**; otherwise, it was labeled **not-chorus**.

Of course, not all generated annotations were correct, but by using over 100,000 training data, the model could be robustly trained without being influenced by errors or outliers. The generated training data consisted of 9,313 English and 91,459 Japanese songs, and we called them `EN_auto` and `JA_auto`, respectively.^{†††}

Furthermore, we manually annotated three sets of lyrics data with more reliable line-level chorus-section annotations for three different purposes:

- (a) **For training comparison:** We annotated 1,103

[†]We used <https://huggingface.co/bert-base-uncased> for English and <https://huggingface.co/cl-tohoku/bert-base-japanese-v2> for Japanese.

^{††}In our experiments, English and Japanese lyrics text as well as the start time of every lyric line were provided by a lyrics distribution company. Automatic lyrics-to-audio synchronization [8]–[18] could also be used to estimate such start times.

^{†††}The main genres are Rock (33%), Pop (25%), and Alternative (12%) for `EN_auto` and are J-Pop (53%), Rock (20%), and Anime (9%) for `JA_auto`.

Table 3 Experimental result: Importance of using both structural and linguistic features.

Feature	Training data / Testing data							
	EN_auto / EN_test				JA_auto / JA_test			
	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>
<i>sim_{all}</i>	77.9	76.1	48.6	45.5	81.2	82.7	63.6	59.6
<i>ling_{ave+seq}</i>	57.4	59.9	16.5	6.9	55.2	61.8	22.1	16.7
<i>ling_{BERT}</i>	61.7	61.1	19.9	12.7	58.3	66.6	29.3	23.3
<i>sim_{all} + ling_{ave+seq}</i>	78.1	77.7	50.8	47.3	83.4	83.5	64.9	61.4
<i>sim_{all} + ling_{BERT}</i>	79.7	78.2	51.7	47.8	85.4	83.3	65.3	62.7

Table 2 Experimental result: Comparison of different methods (the unit is %).

Method	Training data / Testing data							
	EN_auto / EN_test				JA_auto / JA_test			
	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>
Heuristic	57.8	73.8	43.0	35.8	57.1	73.2	43.6	36.3
MLP	76.4	72.7	43.6	39.4	79.8	83.0	62.3	58.9
Bi-LSTM	79.7	78.2	51.7	47.8	85.4	83.3	65.3	62.7

Japanese lyrics and called them JA_man[†]. By comparing the performance of the model trained on JA_auto with that of the model trained on JA_man, we could confirm that our generated data is reliable enough for training purposes.

- (b) **For tuning model parameters:** We annotated the lyrics of 21 English and 79 Japanese songs from RWC-MDB-P-2001 and called them EN_RWC and JA_RWC, respectively. These were used to tune model parameters.
- (c) **For testing:** We annotated the lyrics of 118 other English songs and 128 other Japanese songs and called them EN_test and JA_test, respectively^{††}. These were used to test the chorus-section detection methods.

4.3 Comparison of Different Methods

Table 2 summarizes the evaluated performances of Heuristic, MLP, and the proposed Bi-LSTM. Note that in this section, MLP and Bi-LSTM are trained using the structural feature *sim_{all}* and the linguistic feature *ling_{BERT}*. We found that MLP and Bi-LSTM performed better than Heuristic. This indicates that methods based on supervised learning are better than a rule-based method. We also found that Bi-LSTM was better than MLP and thus confirmed the importance of learning dependencies between adjacent lines.

Since we concluded from these results that the proposed Bi-LSTM is the best for the chorus-section detection task, in the subsequent experiments reported here we used only Bi-LSTM.

[†]To investigate the accuracy of the automatic annotation method we used for generating EN_auto and JA_auto, we applied the same method to the songs (audio signals and corresponding manually time-aligned lyrics) in JA_man. The accuracy of the generated annotations was $F = 68.0\%$, so the automatic annotation method seems to work fairly well.

^{††}The chorus and not-chorus labels were annotated only on the lyrics. No audio signal is available for these test data.

4.4 Importance of Using Both Structural and Linguistic Features

To investigate the effectiveness of structural and linguistic features, we compared their use individually and in combination. Table 3 summarizes the results.

The top three entries in Table 3 show that the models trained solely on the structural feature *sim_{all}* greatly outperformed the models trained solely on the linguistic feature *ling_{ave+seq}* or *ling_{BERT}*. This result confirms that capturing the repetitive structure is effective in detecting chorus sections not only in audio signals but also in lyrics text. We also confirmed that the use of the BERT-based linguistic feature *ling_{BERT}* outperformed the use of the word2vec/context2vec-based linguistic feature *ling_{ave+seq}*, as expected.

The bottom two entries in Table 3 show that combining the structural feature *sim_{all}* with the linguistic feature *ling_{ave+seq}* or *ling_{BERT}* improved performance. This not only confirms the importance of using SSMS, as had been shown for the audio-based detection of chorus sections, but also confirms that the additional use of linguistic features is helpful for detecting chorus sections, which has not been shown before.

Table 3 also shows that the combination of *sim_{all} + ling_{BERT}* outperformed the combination of *sim_{all} + ling_{ave+seq}*, as expected. However, the performance difference between those combinations (the bottom two entries) was much smaller than the performance difference between *ling_{BERT}* and *ling_{ave+seq}* (the second and third entries from the top). This confirms that although BERT is superior to word2vec/doc2vec as a linguistic feature, the combination of structural and linguistic features reduces its superiority since the structural feature is much more effective for the lyric chorus-section detection task.

Since we concluded from these results that *ling_{BERT}* is superior to *ling_{ave+seq}* whether it is used alone or in combination, in the subsequent experiments reported here we used only *ling_{BERT}*.

4.5 Reliability of Generated Annotations

To investigate whether JA_auto containing some annotation errors is reliable enough for training purposes, we compared the model trained using JA_auto with the model trained using JA_man that does not contain annotation errors. Table 4

Table 4 Experimental result: Reliability of automatically generated annotations.

Training data	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>
JA_auto (91,459 songs)	85.4	83.3	65.3	62.7
JA_man (1,103 songs)	81.1	77.8	54.7	52.0

Table 5 Experimental result: Can the Japanese model detect English chorus sections?

Training data	Testing data	<i>F</i>	<i>p-F</i>	<i>n-F</i>	<i>V</i>
EN_auto (9,313 songs)	EN_test	77.9	76.1	48.6	45.5
JA_auto (91,459 songs)	EN_test	80.3	80.6	58.1	54.4
EJ_auto (100,772 songs)	EN_test	81.0	82.3	60.7	57.4

clearly shows that the model trained using JA_auto, automatically generated data the amount of which can be large, outperformed the model trained using JA_man, manually annotated data, the amount of which is usually very limited because of the laborious manual effort its creation requires. The result also confirms that even if annotations generated automatically are not perfect, they are reliable enough for training the model.

4.6 Training Data Size and Language Dependency

Tables 2 and 3 also show that the performances for English lyrics were worse than those for Japanese lyrics. Since the amount of Japanese training data was about 10 times that of English training data, we think that the amount of training data greatly affects the performance of the proposed model. We are thus interested in answering the question ‘‘Can a model trained on a large amount of Japanese data detect English chorus sections?’’ In fact, although linguistic features are language dependent and the process of computing SSMs is also language dependent, structural features based on the resulting SSMs can be language independent because our SSMs simply represent patterns of repeating lyric lines, which could be universal in music. Therefore, in this section, we compare models solely trained on the structural features, without using the linguistic features.

As shown in the upper half of Table 5, which shows results obtained without using linguistic features, we found that the structural-feature-based model trained on Japanese data JA_auto succeeded in detecting English chorus sections in EN_test and its performance was better than that of the model trained on the smaller dataset EN_auto. This result indicates that the SSM-based model trained on a large amount of data can detect chorus sections regardless of the language of the test set. Moreover, this result is further evidence that Japanese and English SSMs (i.e., patterns of repeating lyric lines) have similar structures.

Obviously, the above result raises another question: ‘‘Can a model trained on both EN_auto and JA_auto perform better than one trained on only EN_auto or JA_auto?’’ To answer this question, we created training data EJ_auto by including both EN_auto and JA_auto and constructed yet another structural-feature-based model with EJ_auto. As shown in the lower half of Table 5, we found that the model

trained on both languages performed better than the model trained on only one.

These results confirm that chorus sections can be detected by a model trained on data in another language, that patterns of repeating lyric lines are language-independent and that mixing different language data allows the model to learn the general structure of chorus sections and thereby perform better. This could have an impact on low-resource languages because large-scale training data can be created by mixing other available language resources.

5. Related Work

Previous work in the community of music information retrieval has addressed musical structure analysis and chorus-section detection based on repeated patterns in musical audio signals [24]–[42]. Studies in the chorus-section detection for audio signals typically used SSMs to capture repeated structures, and we share this motivation. Our approach differs from those audio-based approaches in that it exploits multiple lyrics-based SSMs and linguistic features within chorus sections.

On the other hand, recent work in the community of natural language processing has tackled lyrics segmentation and summarization tasks by exploiting SSMs. Fell et al. [19] and Watanabe et al. [20] proposed a neural network model and logistic regression model for segmenting paragraphs (sections) without labeling them by using SSMs as features. Those tasks, however, are essentially different from detecting all chorus sections that are the most representative sections in lyrics text. Addressing a task similar to chorus-section detection, Fell et al. [55] proposed a method of summarizing lyrics by combining general document summarization methods with audio thumbnailing methods. They focus on extracting individual informative lines as a summary from lyrics text, not redundant repeated lines. On the other hand, the focus of our paper is to detect chorus sections whose successive lines are often repeated in lyrics text.

6. Conclusion

This paper has addressed the novel task of detecting chorus sections in English and Japanese lyrics. We proposed a neural-network-based sequence labeling model that learns structural (i.e., phrase-repetition) and linguistic features to detect lyric lines of chorus sections. We also generated over 100,000 training data with chorus-section annotations. No previous work has ever conducted chorus-section detection for text-only lyrics with this much data.

The contributions of this study are summarized as follows: (1) We designed a variety of features to capture structural and linguistic properties of chorus sections. (2) We proposed a sequence labeling model that can detect chorus sections in lyrics. (3) We showed how to generate a large training dataset of lyrics with chorus-section annotations. (4) We demonstrated that our Bi-LSTM-based method outperforms alternative baseline methods. (5) We thoroughly

investigated this detection task and the nature of chorus sections of lyrics from different perspectives such as the importance of features, the amount of training data, and language dependency.

We plan to extend our method to detect other sections, such as verse and bridge sections. Future work will also develop music applications using our method, such as those discussed in Sect. 1.

Acknowledgments

The authors appreciate SyncPower Corporation for providing lyrics data. This work was supported in part by JST ACCEL Grant Number JPMJAC1602, JST CREST Grant Number JPMJCR20D4, and JSPS KAKENHI Grant Number JP20K19878.

References

- [1] K. Watanabe and M. Goto, "A chorus-section detection method for lyrics text," Proceedings of the 21st International Society for Music Information Retrieval Conference (ISMIR 2020), pp.351–359, 2020.
- [2] K. Watanabe and M. Goto, "Query-by-Blending: A music exploration system blending latent vector representations of lyric word, song audio, and artist," Proceedings of the 20th International Society for Music Information Retrieval Conference (ISMIR 2019), pp.144–151, 2019.
- [3] K. Tsukuda, K. Ishida, and M. Goto, "Lyric Jumper: A lyrics-based music exploratory web service by modeling lyrics generative process," Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR 2017), pp.544–551, 2017.
- [4] S. Sasaki, K. Yoshii, T. Nakano, M. Goto, and S. Morishima, "LyricsRadar: A lyrics retrieval system based on latent topics of lyrics," Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014), pp.585–590, 2014.
- [5] A. Tsaptsinos, "Lyrics-based music genre classification using a hierarchical attention network," Proceedings of the 18th International Society for Music Information Retrieval Conference (ISMIR 2017), pp.694–701, 2017.
- [6] R. Mayer and A. Rauber, "Music genre classification by ensembles of audio and lyrics features," Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011), pp.675–680, 2011.
- [7] R. Mayer, R. Neumayer, and A. Rauber, "Rhyme and style features for musical genre classification by song lyrics," Proceedings of the 9th International Conference on Music Information Retrieval (ISMIR 2008), pp.337–342, 2008.
- [8] B. Sharma, C. Gupta, H. Li, and Y. Wang, "Automatic lyrics-to-audio alignment on polyphonic music using singing-adapted acoustic models," Proceedings of the 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2019), pp.396–400, 2019.
- [9] D. Stoller, S. Durand, and S. Ewert, "End-to-end lyrics alignment for polyphonic music using an audio-to-character recognition model," Proceedings of the 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2019), pp.181–185, 2019.
- [10] C. Gupta, E. Yilmaz, and H. Li, "Acoustic modeling for automatic lyrics-to-audio alignment," Proceedings of the 20th Annual Conference of the International Speech Communication Association (Interspeech 2019), pp.2040–2044, 2019.
- [11] C. Gupta, R. Tong, H. Li, and Y. Wang, "Semi-supervised lyrics and solo-singing alignment," Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR 2018), pp.600–607, 2018.
- [12] S. Chang and K. Lee, "Lyrics-to-audio alignment by unsupervised discovery of repetitive patterns in vowel acoustics," IEEE Access, vol.5, pp.16635–16648, 2017.
- [13] S.W. Lee and J. Scott, "Word level lyrics-audio synchronization using separated vocals," Proceedings of the 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2017), pp.646–650, 2017.
- [14] Y.-R. Chien, H.-M. Wang, and S.-K. Jeng, "Alignment of lyrics with accompanied singing audio based on acoustic-phonetic vowel likelihood modeling," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol.24, no.11, pp.1998–2008, 2016.
- [15] G.B. Dzhambazov and X. Serra, "Modeling of phoneme durations for alignment between polyphonic audio and lyrics," Proceedings of the 12th Sound and Music Computing Conference (SMC 2015), pp.281–286, 2015.
- [16] M. Mauch, H. Fujihara, and M. Goto, "Integrating additional chord information into HMM-based lyrics-to-audio alignment," IEEE Transactions on Audio, Speech, and Language Processing, vol.20, no.1, pp.200–210, 2012.
- [17] H. Fujihara, M. Goto, J. Ogata, and H.G. Okuno, "LyricSynchronizer: Automatic synchronization system between musical audio signals and lyrics," IEEE Journal of Selected Topics in Signal Processing, vol.5, no.6, pp.1252–1261, 2011.
- [18] M.-Y. Kan, Y. Wang, D. Iskandar, T.L. Nwe, and A. Shenoy, "LyricAlly: Automatic synchronization of textual lyrics to acoustic music signals," IEEE Trans. Speech Audio Process., vol.16, no.2, pp.338–349, 2008.
- [19] M. Fell, Y. Nechaev, E. Cabrio, and F. Gandon, "Lyrics segmentation: Textual macrostructure detection using convolutions," Proceedings of the 27th International Conference on Computational Linguistics (COLING 2018), pp.2044–2054, 2018.
- [20] K. Watanabe, Y. Matsubayashi, N. Orita, N. Okazaki, K. Inui, S. Fukayama, T. Nakano, J.B.L. Smith, and M. Goto, "Modeling discourse segments in lyrics using repeated patterns," Proceedings of the 26th International Conference on Computational Linguistics (COLING 2016), pp.1959–1969, 2016.
- [21] A. Baratè, L.A. Ludovico, and E. Santucci, "A semantics-driven approach to lyrics segmentation," Proceedings of the 8th International Workshop on Semantic and Social Media Adaptation and Personalization (SMAP 2013), pp.73–79, 2013.
- [22] J.P.G. Mahedero, Á. Martínez, P. Cano, M. Koppenberger, and F. Gouyon, "Natural language processing of lyrics," Proceedings of the 13th ACM International Conference on Multimedia (ACM Multimedia), pp.475–478, 2005.
- [23] K. Watanabe, Y. Matsubayashi, K. Inui, S. Fukayama, T. Nakano, and M. Goto, "Modeling storylines in lyrics," IEICE Transactions on Information and Systems, vol.E101.D, no.4, pp.1167–1179, 2018.
- [24] G. Shibata, R. Nishikimi, E. Nakamura, and K. Yoshii, "Statistical music structure analysis based on a homogeneity-, repetitiveness-, and regularity-aware hierarchical hidden semi-markov model," Proceedings of the 20th International Society for Music Information Retrieval Conference (ISMIR 2019), pp.268–275, 2019.
- [25] A. Maezawa, "Music boundary detection based on a hybrid deep model of novelty, homogeneity, repetition and duration," Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2018), pp.206–210, 2018.
- [26] G. Sargent, F. Bimbot, and E. Vincent, "Estimating the structural segmentation of popular music pieces under regularity constraints," IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol.25, no.2, pp.344–358, 2017.
- [27] T. Cheng, J.B.L. Smith, and M. Goto, "Music structure boundary detection and labelling by a deconvolution of path-enhanced self-similarity matrix," Proceedings of the 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2018), pp.106–110, 2018.

- [28] J.B.L. Smith and M. Goto, "Using priors to improve estimates of music structure," Proceedings of the 17th International Society for Music Information Retrieval Conference (ISMIR 2016), pp.554–560, 2016.
- [29] T. Grill and J. Schlüter, "Music boundary detection using neural networks on combined features and two-level annotations," Proceedings of the 16th International Society for Music Information Retrieval Conference (ISMIR 2015), pp.531–537, 2015.
- [30] B. McFee and D. Ellis, "Analyzing song structure with spectral clustering," Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014), pp.405–410, 2014.
- [31] G. Peeters and V. Bisot, "Improving music structure segmentation using lag-priors," Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014), pp.337–342, 2014.
- [32] H. Grohganz, M. Clausen, N. Jiang, and M. Müller, "Converting path structures into block structures using eigenvalue decompositions of self-similarity matrices," Proceedings of the 14th International Society for Music Information Retrieval Conference (ISMIR 2013), pp.209–214, 2013.
- [33] O. Nieto and T. Jehan, "Convex non-negative matrix factorization for automatic music structure identification," Proceedings of the 2013 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2013), pp.236–240, 2013.
- [34] F. Kaiser and G. Peeters, "A simple fusion method of state and sequence segmentation for music structure discovery," Proceedings of the 14th International Society for Music Information Retrieval Conference (ISMIR 2013), pp.257–262, 2013.
- [35] J. Serrà, M. Müller, P. Grosche, and J.L. Arcos, "Unsupervised detection of music boundaries by time series structure features," Proceedings of the 26th AAAI Conference on Artificial Intelligence, pp.1613–1619, 2012.
- [36] M. Müller, P. Grosche, and N. Jiang, "A segment-based fitness measure for capturing repetitive structures of music recordings," Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR 2011), pp.615–620, 2011.
- [37] J. Paulus, M. Müller, and A. Klapuri, "State of the art report: Audio-based music structure analysis," Proceedings of the 11th International Society for Music Information Retrieval Conference (ISMIR 2010), pp.625–636, 2010.
- [38] J. Paulus and A. Klapuri, "Music structure analysis using a probabilistic fitness measure and a greedy search algorithm," IEEE Transactions on Audio, Speech, and Language Processing, vol.17, no.6, pp.1159–1170, 2009.
- [39] M. Müller and S. Ewert, "Joint structure analysis with applications to music annotation and synchronization," Proceedings of the 9th International Conference on Music Information Retrieval (ISMIR 2008), pp.389–394, 2008.
- [40] M. Goto, "A chorus section detection method for musical audio signals and its application to a music listening station," IEEE Transactions on Audio, Speech, and Language Processing, vol.14, no.5, pp.1783–1794, 2006.
- [41] M. Cooper and J. Foote, "Summarizing popular music via structural similarity analysis," Proceedings of the 2003 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA 2003), pp.127–130, 2003.
- [42] J. Foote, "Automatic audio segmentation using a measure of audio novelty," Proceedings of the 2000 IEEE International Conference on Multimedia and Expo (IEEE ICME 2000), pp.452–455, 2000.
- [43] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," Proceedings of the 26th Annual Conference on Neural Information Processing Systems (NeurIPS 2013), pp.3111–3119, 2013.
- [44] O. Melamud, J. Goldberger, and I. Dagan, "context2vec: Learning generic context embedding with bidirectional LSTM," Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, pp.51–61, 2016.
- [45] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019), pp.4171–4186, 2019.
- [46] M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka, "RWC Music Database: Popular, classical and jazz music databases," Proceedings of the 3rd International Conference on Music Information Retrieval (ISMIR 2002), pp.287–288, 2002.
- [47] L. Yujian and L. Bo, "A normalized Levenshtein distance metric," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol.29, no.6, pp.1091–1095, 2007.
- [48] S. Bird, "NLTK: the natural language toolkit," Proceedings of the 21st International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics (ACL 2006), pp.69–72, 2006.
- [49] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol.9, no.8, pp.1735–1780, 1997.
- [50] D.J. Berndt and J. Clifford, "Using dynamic time warping to find patterns in time series," Proceedings of Workshop on Knowledge Discovery in Databases, pp.359–370, 1994.
- [51] A. Graves, A. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," Proceedings of the 2013 IEEE International Conference on Acoustics, Speech and Signal Processing (IEEE ICASSP 2013), pp.6645–6649, 2013.
- [52] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," Proceedings of the 32nd International Conference on Machine Learning (ICML 2015), pp.448–456, 2015.
- [53] C. Raffel, B. McFee, E.J. Humphrey, J. Salamon, O. Nieto, D. Liang, and D.P.W. Ellis, "mir-eval: A transparent implementation of common MIR metrics," Proceedings of the 15th International Society for Music Information Retrieval Conference (ISMIR 2014), pp.367–372, 2014.
- [54] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," Proceedings of the 7th International Conference on Learning Representations (ICLR 2019), 2019.
- [55] M. Fell, E. Cabrio, F. Gandon, and A. Giboin, "Song lyrics summarization inspired by audio thumbnailing," Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019), pp.328–337, 2019.



Kento Watanabe received the B.E. degree in 2013 and Ph.D. degree in 2018, both from Tohoku University. He is currently a Researcher at the National Institute of Advanced Industrial Science and Technology (AIST), Japan. His research interests include lyrics information processing, natural language processing, machine learning, and human computer iteration.



Masataka Goto received the Doctor of Engineering degree from Waseda University in 1998. He is currently a Prime Senior Researcher at the National Institute of Advanced Industrial Science and Technology (AIST), Japan. Over the past 30 years he has published more than 300 papers in refereed journals and international conference proceedings and has received 64 awards, including several best paper awards, best presentation awards, the Tenth Japan Academy Medal, and the Tenth JSPS

PRIZE. He has served as a committee member of over 120 scientific societies and conferences, including as the General Chair of ISMIR 2009 and 2014. As the research director, he began the OngaACCEL project in 2016 and the RecMus project in 2021, which are five-year JST-funded research projects (ACCEL and CREST) related to music technologies.