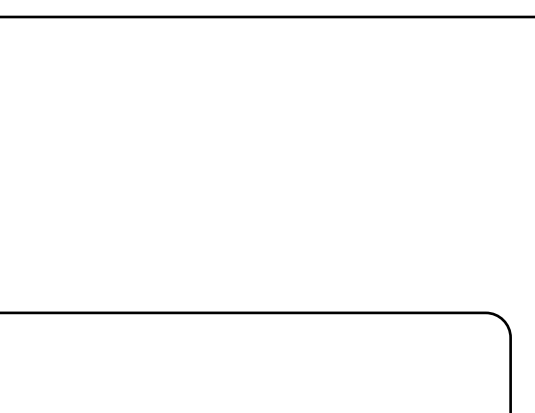
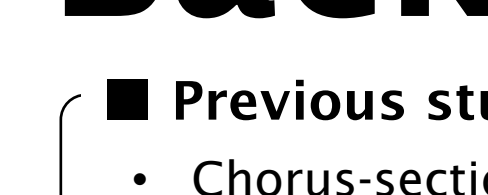


# A Chorus-Section Detection Method for Lyrics Text

Kento Watanabe and Masataka Goto

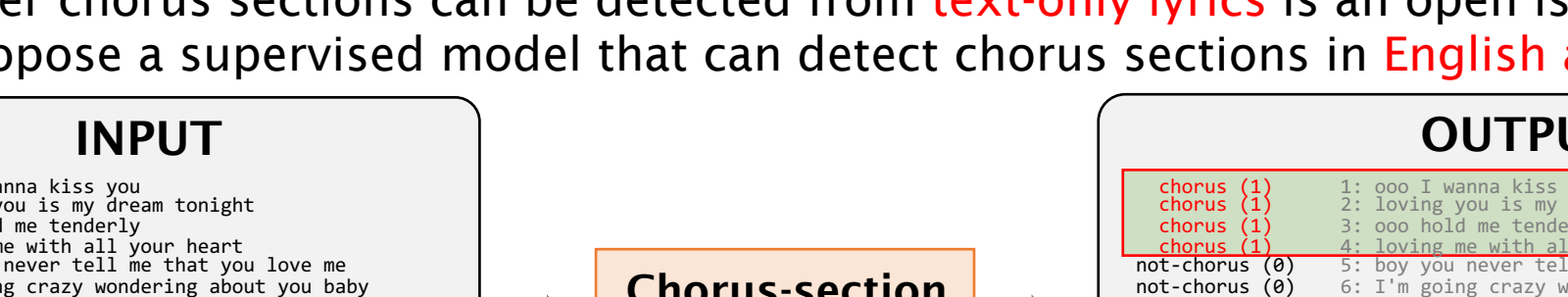
National Institute of Advanced Industrial Science and Technology (AIST)



## Background

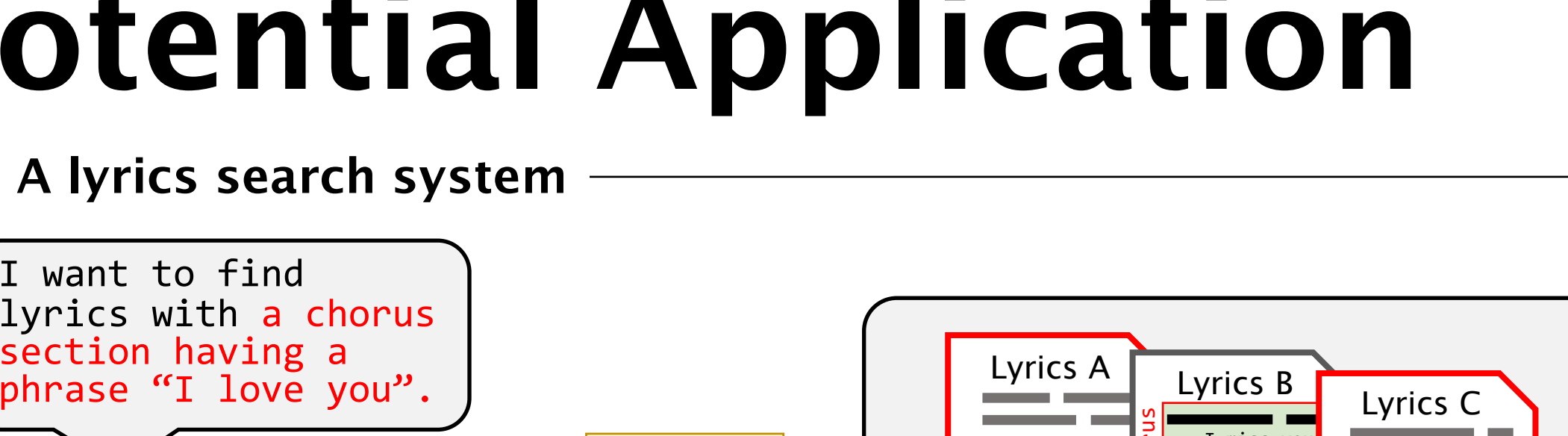
### Previous study

- Chorus-section detection using **audio signals**.



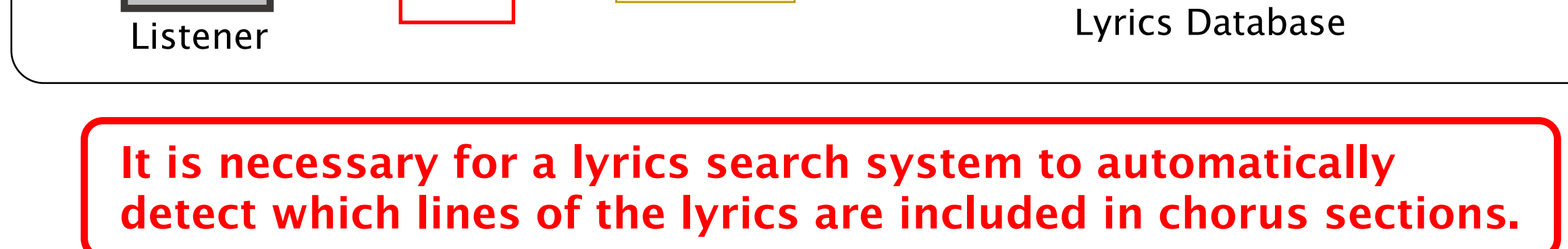
### This study

- Whether chorus sections can be detected from **text-only lyrics** is an open issue.
- We propose a supervised model that can detect chorus sections in **English and Japanese lyrics**.



## Potential Application

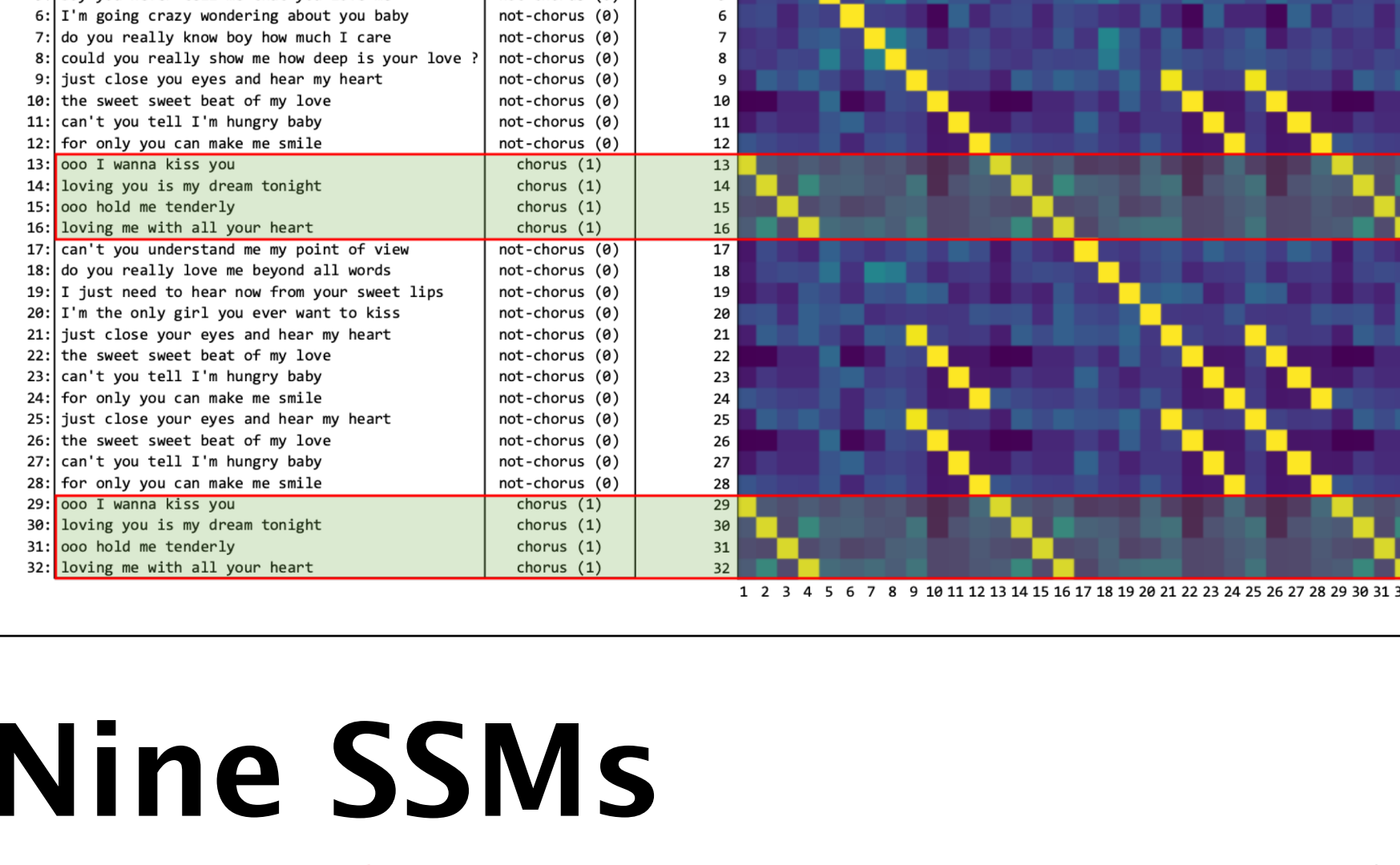
### A lyrics search system



It is necessary for a lyrics search system to automatically detect which lines of the lyrics are included in chorus sections.

## Key Idea

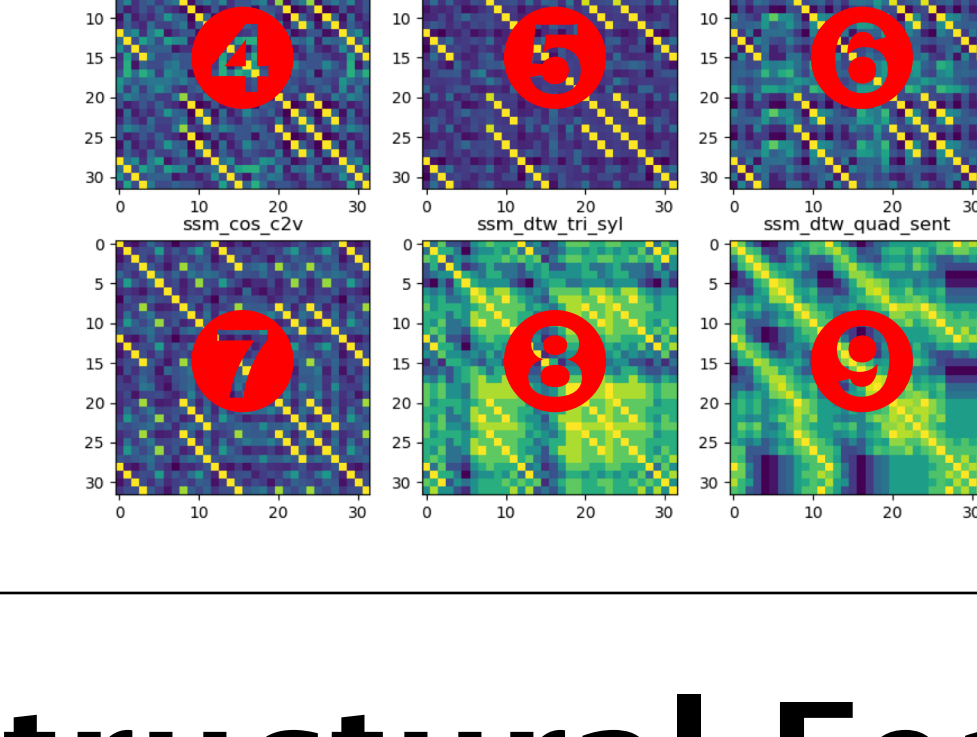
- We compute the **Self-similarity matrix (SSM)** from lyrics text.
- SSM representations are widely used in computational music structure analysis.



Repeated sections lead to high values in diagonals of the matrix, and those patterns are used to identify the structure.

## Nine SSMs

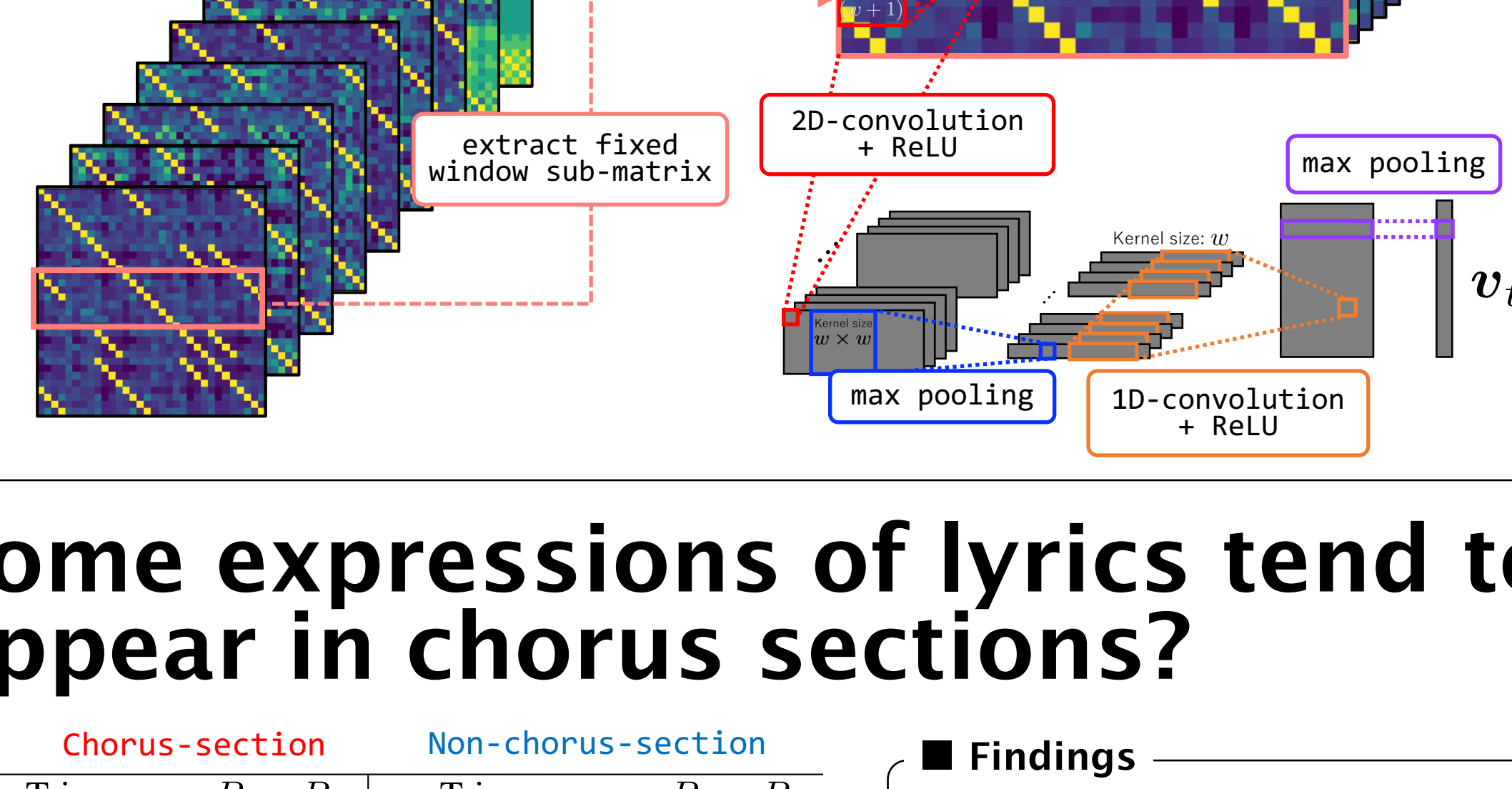
- The **design of the similarity measure** to compute each cell of the SSM is important.
- We propose to use the following nine variations of similarity measures.



- Edit distance**
  - String similarity
  - Head similarity
  - Tail similarity
  - Part-of-speech similarity
  - Phonetic similarity
- Cosine similarity**
  - Word vector similarity
  - Context vector similarity
- Dynamic time warping**
  - Word syllable count similarity
  - Lyric line syllable count similarity

## Structural Feature

To calculate feature vectors  $v_t$  from the above nine SSMs, we use a **CNN architecture\*** to detect textual macro structures from various patterns in SSMs regardless of their locations and relative sizes in SSMs.



## Some expressions of lyrics tend to appear in chorus sections?

Chorus-section		Non-gram	
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	0.05%	I's	0.03%
I'll	0.05%	I'd	0.02%
we're	0.04%	but I'	0.01%
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%

$P_c$  and  $P_n$  denote word tri-gram probabilities in the chorus and non-chorus sections, respectively.

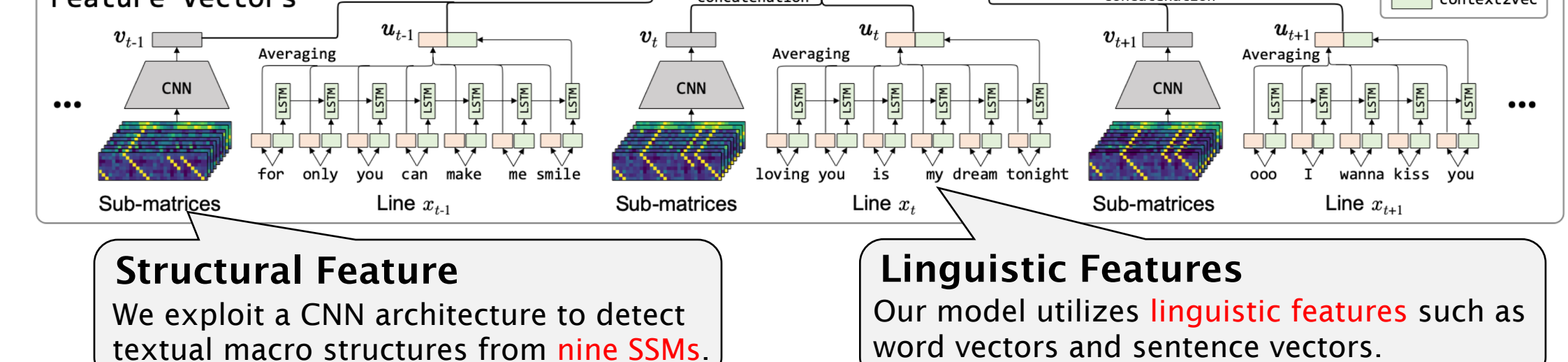
### Findings

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*).

It can be expected that some expressions will contribute to the chorus-section detection.

## Chorus-section Detection Model

- We propose a Bidirectional LSTM-based model using **two types of feature representations**.



**Structural Feature**  
We exploit a CNN architecture to detect textual macro structures from **nine SSMs**.

**Linguistic Features**  
Our model utilizes **linguistic features** such as word vectors and sentence vectors.

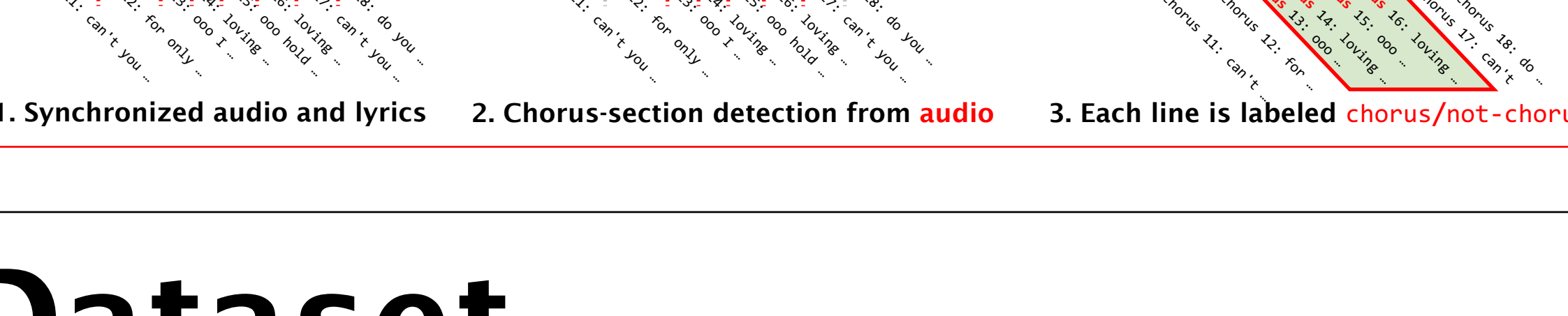
## Training Data

### Technical Problem for learning our model

No lyrics data with line-level chorus-section annotations are available.

### Key Idea: We generate training data with chorus-section annotation

- We prepared **100,772** pairs of musical audio signals and their corresponding manually time-aligned (temporally synchronized) lyrics.
- We detected chorus sections of every song automatically by **using its audio signals**.
- 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**.



## Dataset

### Dataset for training comparison

- EN\_auto: English lyrics data with the generated chorus-section annotations. (9,133 songs)
- JA\_auto: Japanese lyrics data with the generated chorus-section annotations. (91,459 songs)
- JA\_man: Japanese lyrics data with the gold chorus-section annotations. (1,103 songs)

### Dataset for tuning model parameters

- We annotated chorus-section label manually.
  - EN\_RWC: English lyrics data in RWC Music Database. (21 songs)
  - JA\_RWC: Japanese lyrics data in RWC Music Database. (79 songs)

### Dataset for testing

- We annotated chorus-section label manually.
  - EN\_test: English lyrics data. (118 songs)
  - JA\_test: Japanese lyrics data. (128 songs)

## Experimental Results (1)

### Question

To confirm the **effectiveness of our Bi-LSTM model** that can learn dependencies between adjacent lyric lines, we compared its performance with that of two baseline methods.

### Result

Method	Training and test data (# of training songs)	
	English (9,313 songs)	Japanese (91,459 songs)
Heuristic (extract repeated lines as chorus sections)	F-measure 57.8 %	F-measure 57.1 %
Multi-Layer Perceptron	F-measure 74.2 %	F-measure 80.6 %
<b>Bi-LSTM (Proposed)</b>	<b>F-measure 78.1 %</b>	<b>F-measure 83.4 %</b>

### Findings

- Methods based on supervised learning are better than a rule-based method.**
- The proposed method is the best for the chorus-section detection task.**

## Experimental Results (2)

### Question

To investigate the **effectiveness of structural and linguistic features**, we compared their use individually and in combination.

### Result

Feature	Training and test data (# of training songs)	
	English (9,313 songs)	Japanese (91,459 songs)
Structural feature	F-measure 77.9 %	F-measure 81.2 %
Linguistic feature	F-measure 57.4 %	F-measure 55.2 %
<b>Both</b>	<b>F-measure 78.1 %</b>	<b>F-measure 83.4 %</b>

### Findings

- The model with only the structural features greatly outperformed.**
- The additional use of linguistic features is helpful for detecting chorus sections.**

## Experimental Results (3)

### Question

We confirm that **our generated data is reliable enough for training purposes** by comparing the performance of the model trained on JA\_auto with that of the model trained on JA\_man.

### Result

	Training Data	F-measure (Japanese test data)
JA_auto: generated training data (91,459 songs)	F-measure	83.4 %
JA_man: human-annotated training data (1,103 songs)	F-measure	80.3 %

### Findings

- The model trained using generated data (91,459 songs) outperformed the model trained using human-annotated data (1,103 songs).**
- Even if generated annotations are not perfect, they are reliable enough for training.**

## Experimental Results (4)

### Question

- Can a model trained on a large amount of Japanese data detect English chorus sections?
- Can a model trained on both EN\_auto and JA\_auto perform better than one trained on only EN\_auto or JA\_auto?

★ Structural features based on the SSMs can be language independent because our SSMs simply represent patterns of repeating lyric lines, which could be universal in music. So we use the model without linguistic features in this experiment.

### Result

	Training Data	F-measure (English test data)
EN_auto: generated training English data (9,313 songs)	F-measure	77.9 %
JA_auto: generated training Japanese data (91,459 songs)	F-measure	80.3 %
<b>EJ_auto: EN_auto + JA_auto (100,772 songs)</b>	<b>F-measure</b>	<b>81.0 %</b>

### Findings

- The SSM-based model can detect chorus sections regardless of the language.**
- English and Japanese SSMs (i.e., patterns of repeating lyric lines) have similar structures.**
- Mixing different language data allows the model to learn the general structure of chorus sections and thereby perform better.**

## Conclusion

### Contributions

- We designed a variety of features to capture **structural and linguistic properties of chorus sections**.
- We proposed a sequence labeling model that can detect chorus sections in lyrics.
- We showed how to **generate a large training dataset** of lyrics with chorus-section annotations.
- We demonstrated that **our Bi-LSTM-based method outperforms** alternative baseline methods.
- 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**.