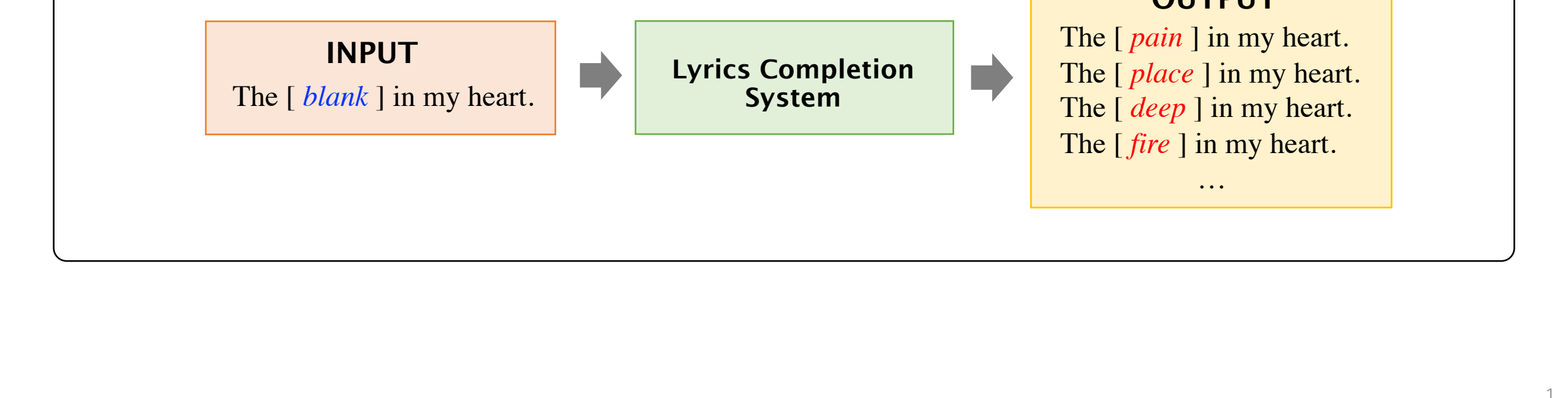


Lyrics Completion Task

We propose a lyrics completion task to recommend **candidate words** for a blank in a given sentence.



What Words Should be Suggested?

Suggesting words that are (1) **atypical** but (2) **suitable for the musical audio signals**.

- Atypicality** is important in creative lyrics. To make lyrics attractive, lyric writers usually consider both typical and atypical phrases. However, previous research on automatic lyrics generation has used **language models** that **predict highly frequent (i.e., typical) words**.
 ⓧ The *pain* in my heart. Many songs use this phrase...
 Ⓢ The *agony* in my heart. This is unique!
LSTM-LM, GPT2, GPT3, Transformer...
- Lyrics depend on the moods of music. (e.g., "love" is often used in ballad songs, "kill" is often used in metal songs).
 Ballad song: ... love you ...
 Metal song: ... kill you ...

Previous Studies and Our Study

◆ Previous study of general word completion task

- Generating **highly frequent (i.e., typical) words**.

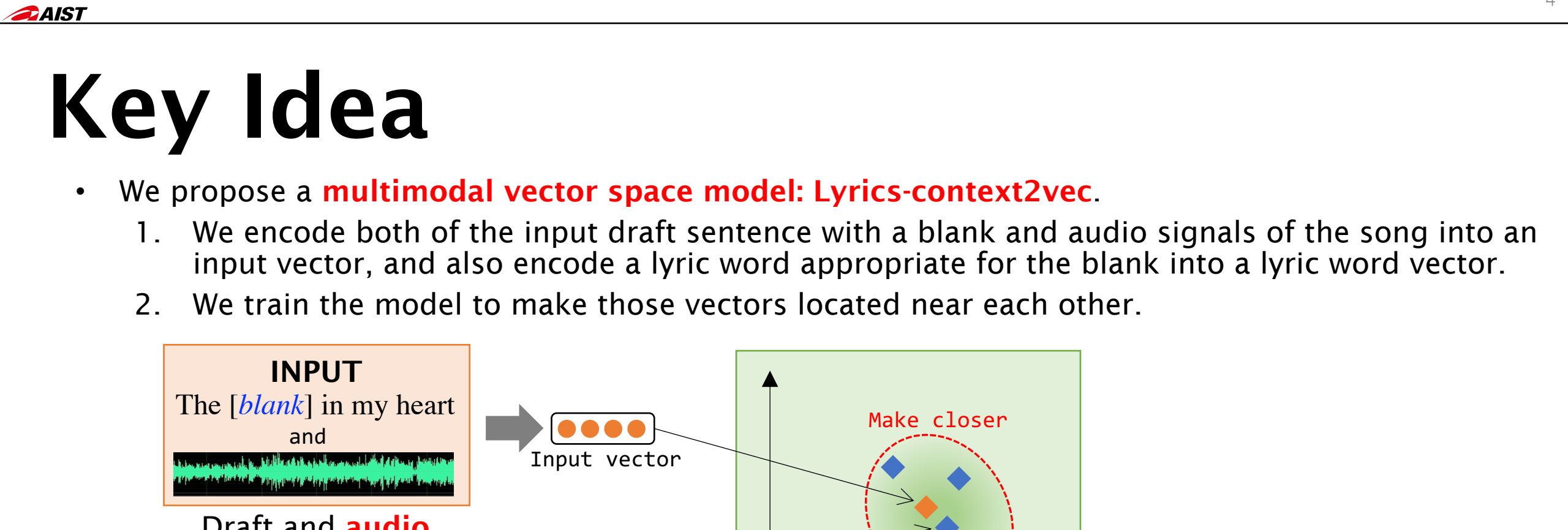


◆ This study

- Suggesting words that are (1) **atypical** but (2) **suitable for musical audio signals**.



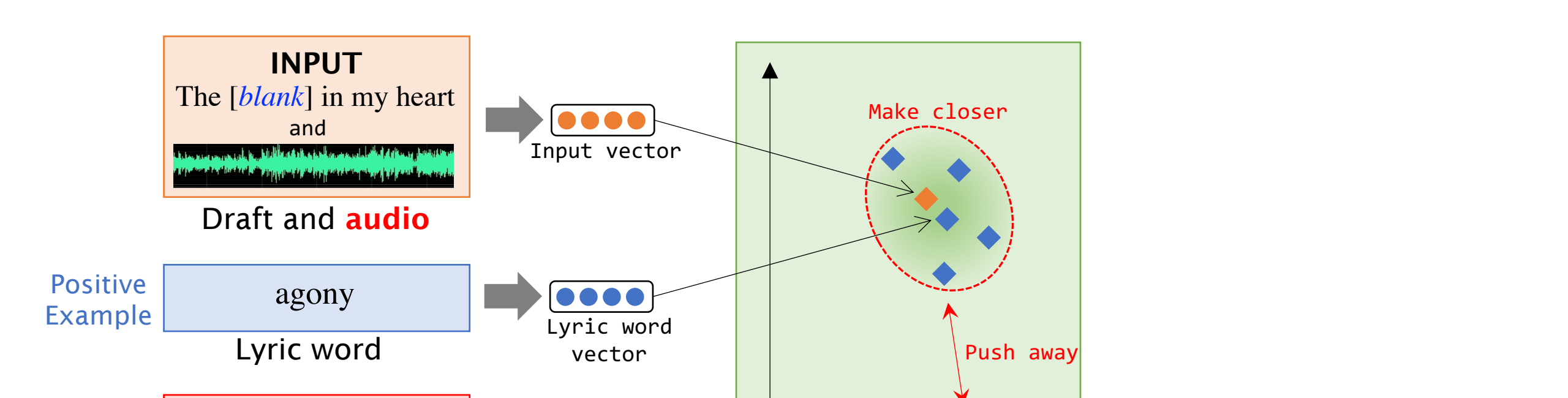
Demo System



Key Idea

- We propose a **multimodal vector space model: Lyrics-context2vec**.

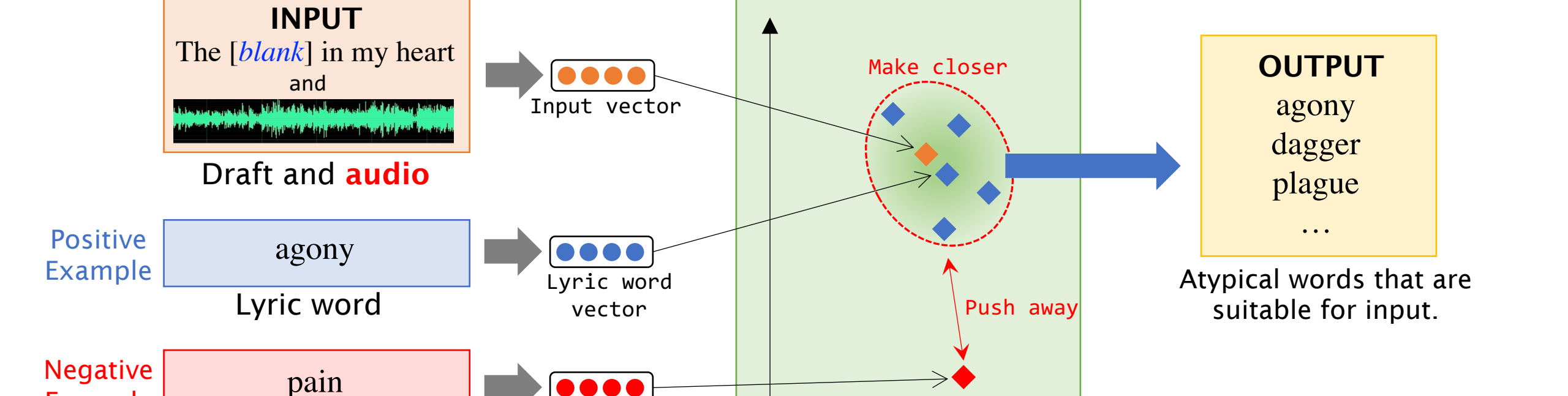
- We encode both of the input draft sentence with a blank and audio signals of the song into an input vector, and also encode a lyric word appropriate for the blank into a lyric word vector.
- We train the model to make those located near each other.



Key Idea

- We propose a **multimodal vector space model: Lyrics-context2vec**.

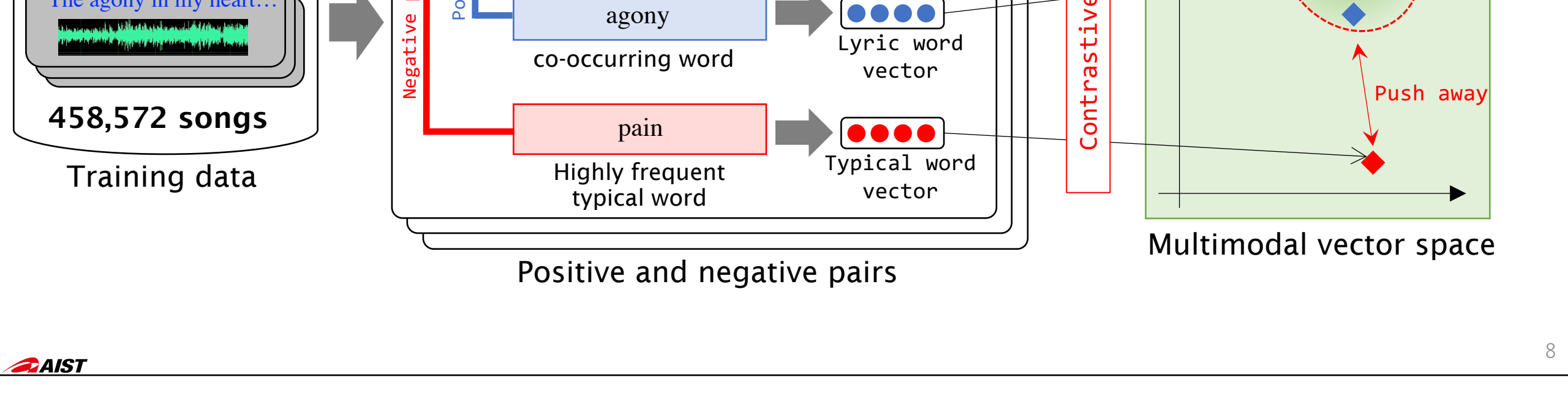
- A highly frequent word is encoded into a typical word vector.
- We use it as a negative example to make it located far away from the input vector.



Key Idea

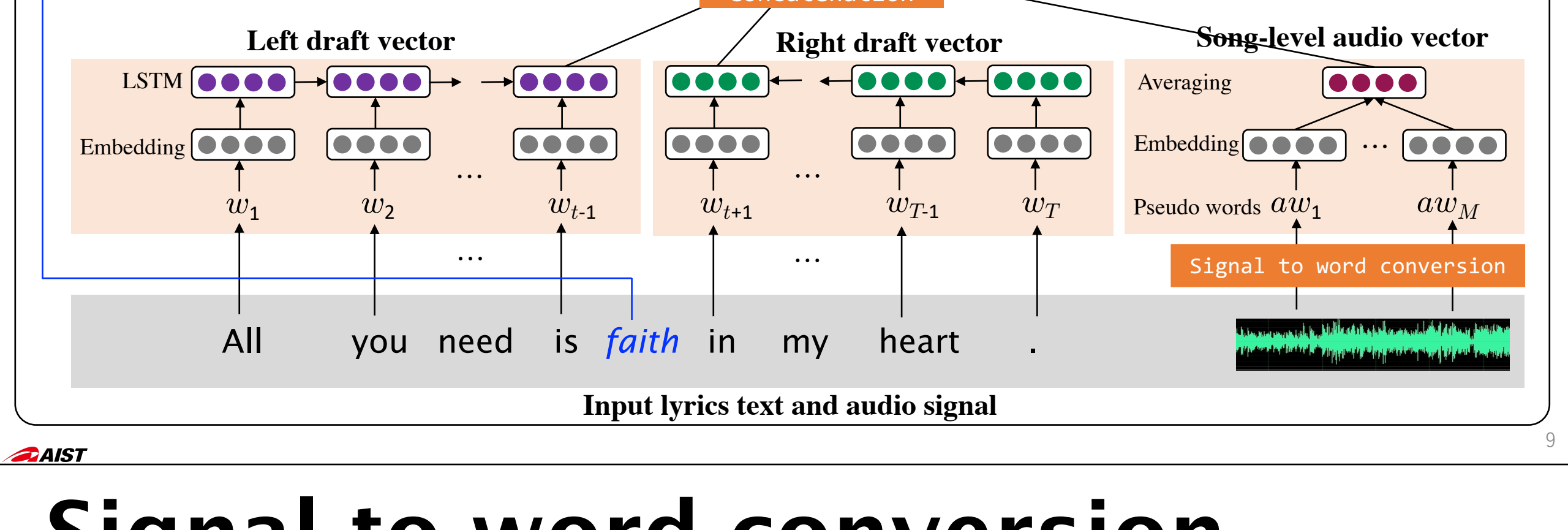
- We propose a **multimodal vector space model: Lyrics-context2vec**.

- Only atypical words are suggested since they are close to the input vector of a draft sentence and audio signals.



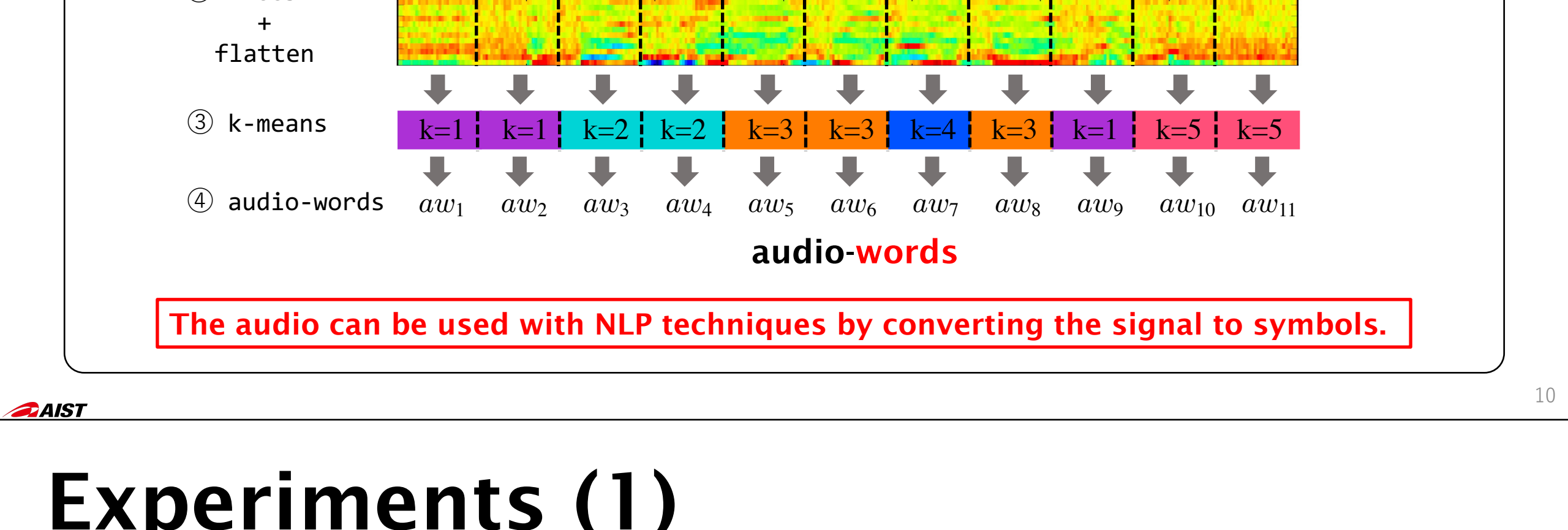
Training Lyrics-context2vec

We used contrastive learning by extracting positive and negative pairs from a large-scale dataset.



Model Construction

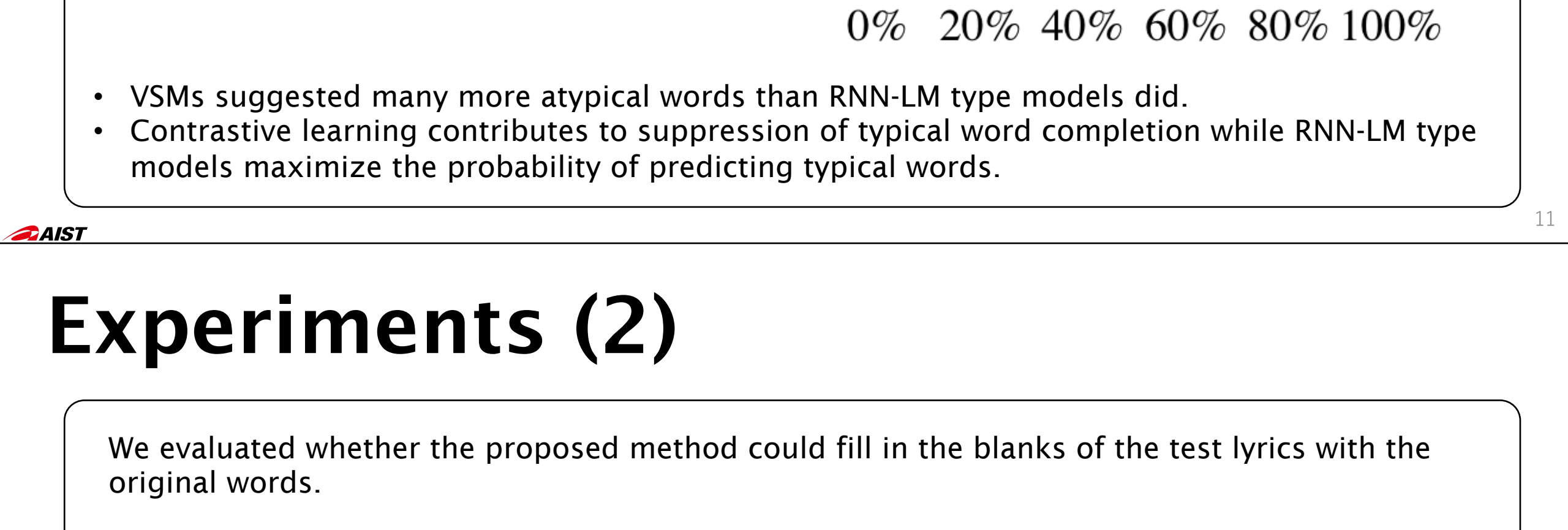
◆ Lyrics-context2vec



Signal to word conversion

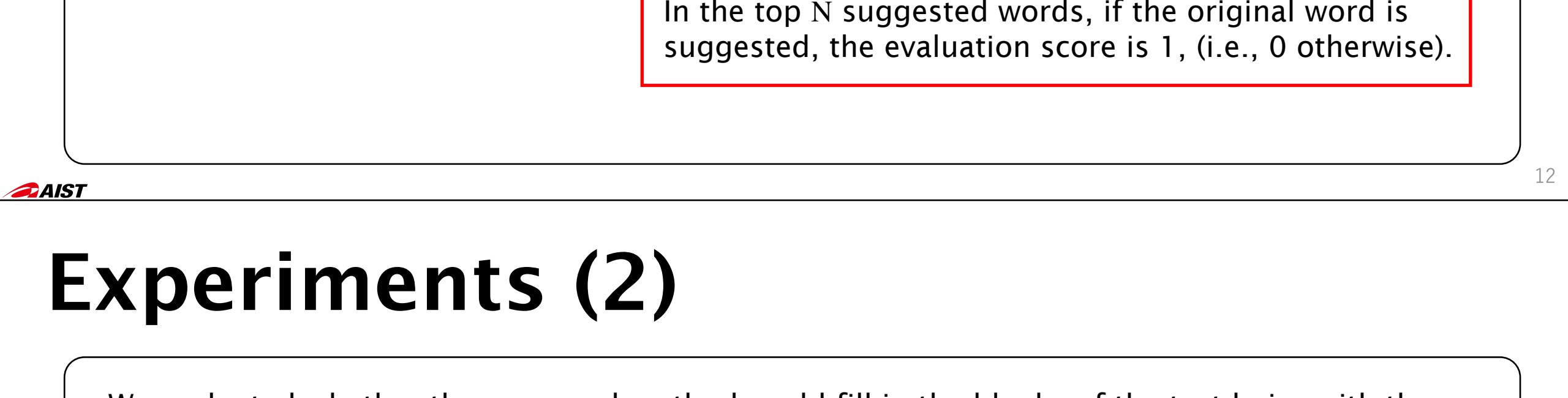
◆ Bag of Audio-words

To embed audio signals, we use a discrete symbol called an **audio-word**.



Experiments (1)

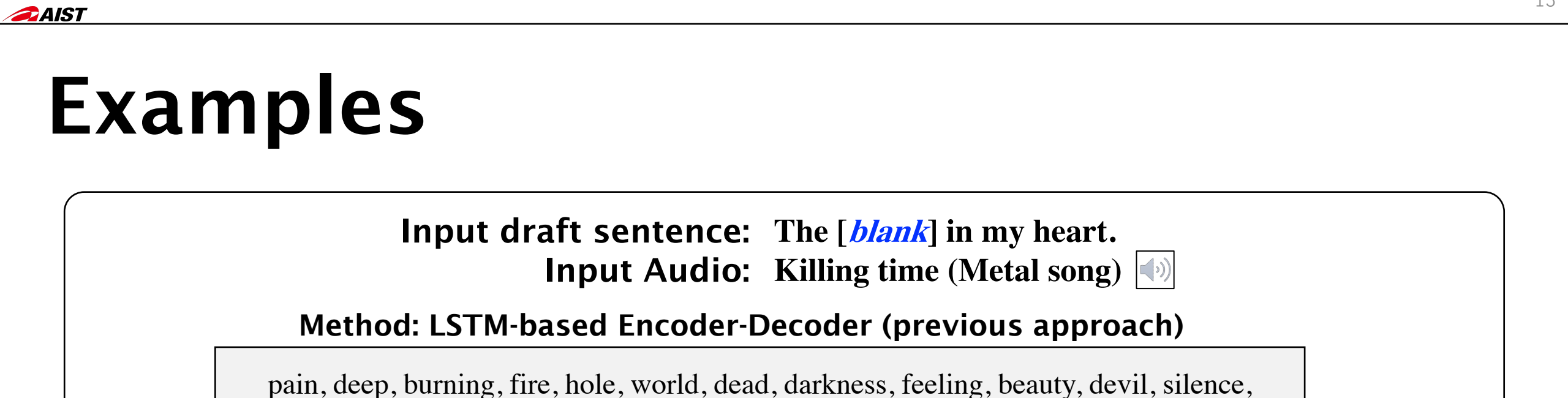
- We calculated how many of the top 20 words suggested by each model were **typical** or **atypical**.
- We calculated the document frequency of words and assumed that **the top 10% of them are typical words** (i.e., the remaining 90% are atypical words).



- VSMs suggested many more atypical words than RNN-LM type models did.
- Contrastive learning contributes to suppression of typical word completion while RNN-LM type models maximize the probability of predicting typical words.

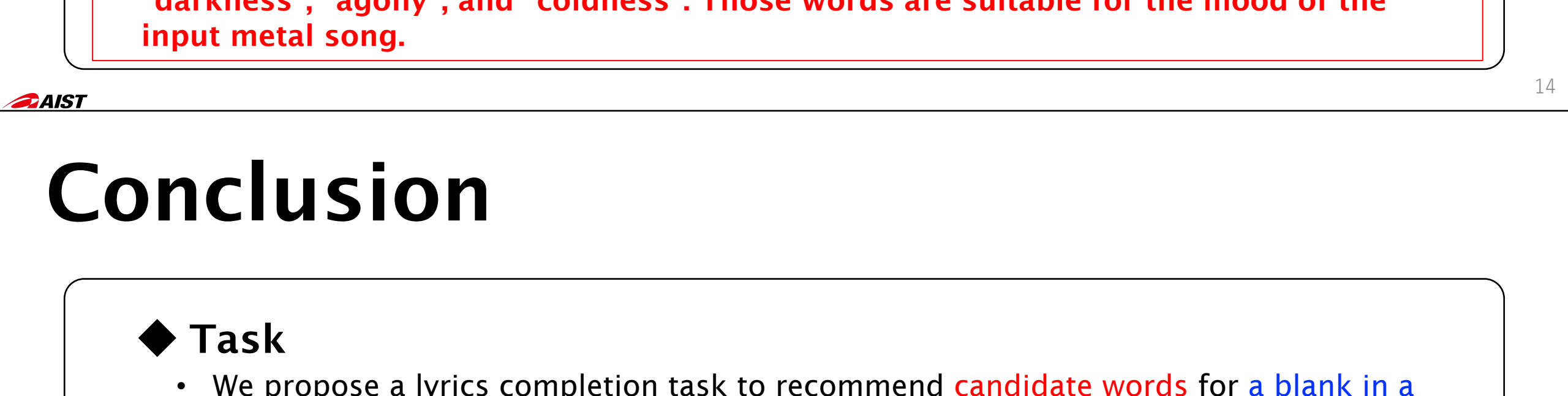
Experiments (2)

We evaluated whether the proposed method could fill in the blanks of the test lyrics with the original words.



Experiments (2)

We evaluated whether the proposed method could fill in the blanks of the test lyrics with the original words.



- Lyrics-context2vec predicted atypical words suitable for music audio better than all other models.**
- Our model captures both the atypicality and the relationship between a music audio and words simultaneously.**

Examples

Input draft sentence: The [blank] in my heart.
 Input Audio: Killing time (Metal song)

- Our lyrics-context2vec can successfully suggest some rare words in bold fonts, such as "dagger" and "agony".
- Lyrics-context2vec successfully suggested explicit and negative words, such as "darkness", "agony", and "coldness". Those words are suitable for the mood of the input metal song.

Conclusion

◆ Task

- We propose a lyrics completion task to recommend **candidate words** for a blank in a given sentence.
- Our task aims to suggest words that are (1) **atypical** and (2) **suitable for the musical audio signal**.

◆ Method

- We proposed **lyrics-context2vec**, a **multimodal vector space model** that suggests atypical but appropriate words for the given music audio and draft sentence.
- Input vector and output word vector located near each other.
- Vectors of highly frequent word located far away from the input vector.

◆ Findings

- Our contrastive learning strategy contributes to suggesting atypical words.
- Embedding audio signals contributes to suggesting words suitable for the mood of the provided music audio.