

EFFECTIVE COVER SONG IDENTIFICATION BASED ON SKIPPING BIGRAMS

Xiaoshuo Xu, Xiaou Chen, Deshun Yang

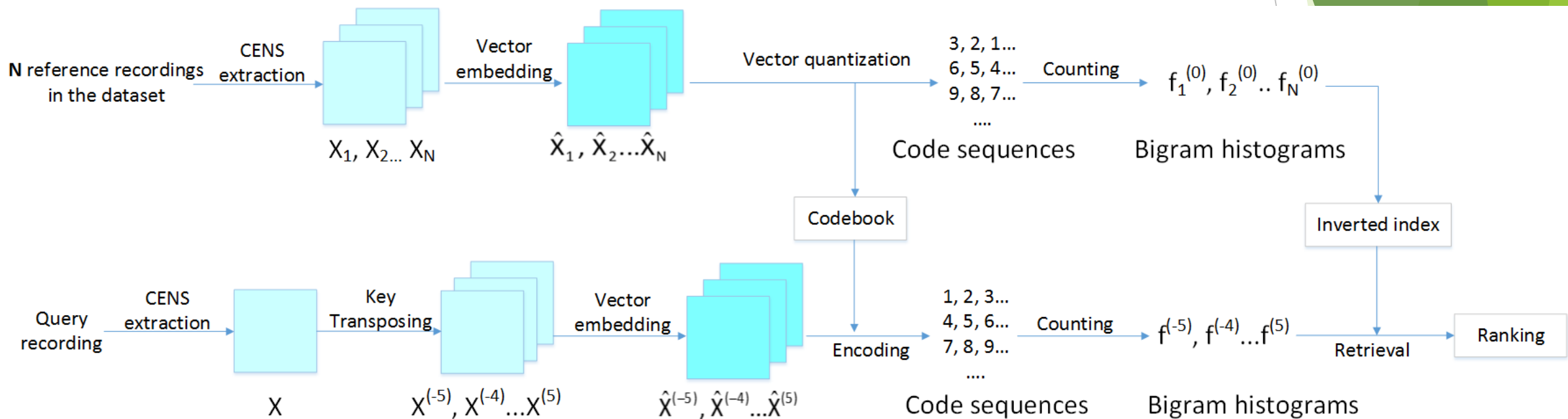
Institute of Compute Science and Technology, Peking University

128 Zhongguancun North Street, Haidian District, Beijing 100871, P.R.China

Outline

- ▶ What is cover song identification?
- ▶ Application: detect copyright infringement, music retrieval, etc.
- ▶ Challenge: Key transposition, structure and speed change
- ▶ Existing methods: Sequence alignment, Music representation
- ▶ Our approach
 - ▶ Represent music with skipping bigram histogram
 - ▶ Utilize inverted index to accelerate the calculation

Pipeline



Feature extraction

- ▶ Chroma Energy Normalized Statistics (CENS)

- ▶ Key transposition

- ▶ Given a CENS vector $\mathbf{x} = (x_0, x_1 \dots x_{11})^T$, the transposed vector is defined as follows:

$$\mathbf{x}^{(i)} = (x_{i\%12}, x_{(i+1)\%12} \dots x_{(i+11)\%12})^T$$

- ▶ Given a CENS sequence $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_M]$, the transposed sequence would be:

$$\mathbf{X}^{(i)} = [\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)} \dots \mathbf{x}_M^{(i)}]$$

- ▶ Vector Embedding

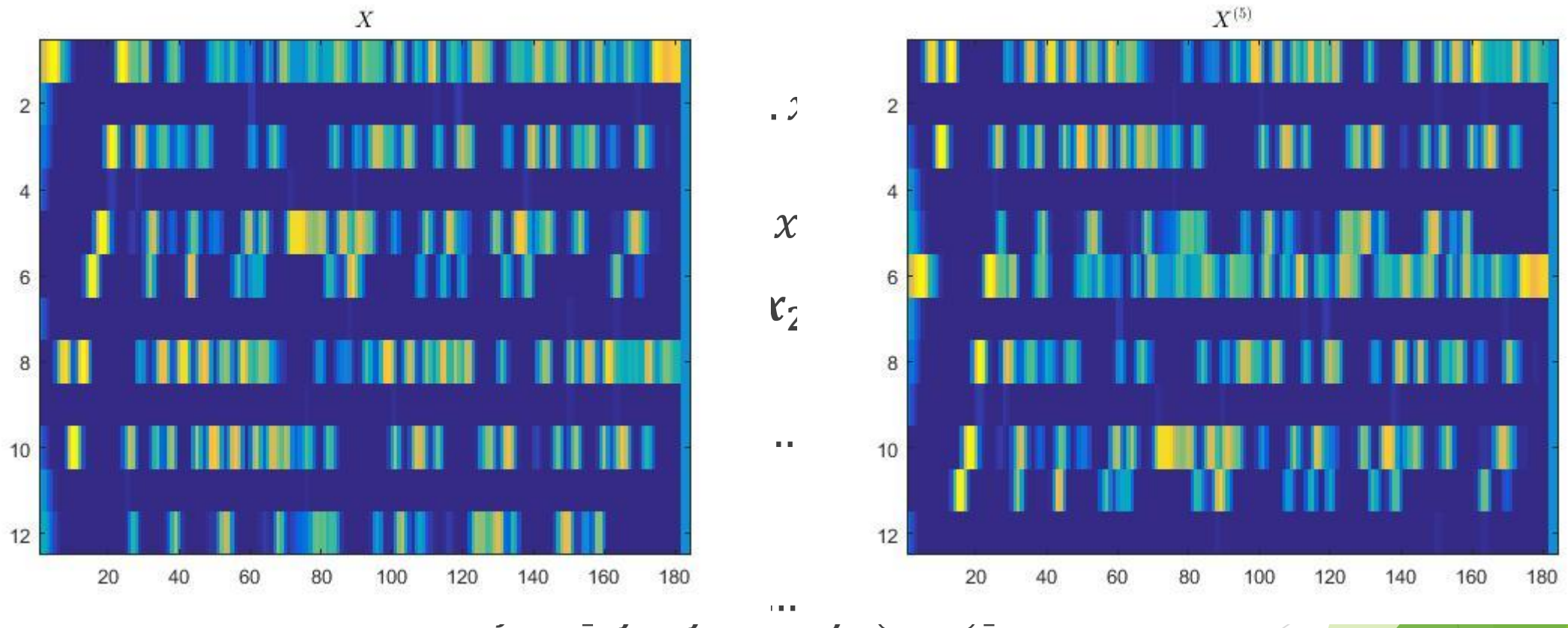
- ▶ Embedded vector: $\widehat{\mathbf{x}}_j = [\mathbf{x}_j^T, \mathbf{x}_{j-1}^T \dots \mathbf{x}_{j-(m-1)}^T], j = m, m+1 \dots M$

- ▶ Embedded sequence: $\widehat{\mathbf{X}} = [\widehat{\mathbf{x}}_m, \widehat{\mathbf{x}}_{m+1} \dots \widehat{\mathbf{x}}_M]$

- ▶ Transposed embedded sequence: $\widehat{\mathbf{X}}^{(i)}$

Feature extraction

- ▶ Chroma Energy Normalized Statistics (CENS)



- ▶ Embedded sequence: $\hat{X} = [\hat{x}_m, \hat{x}_{m+1} \dots \hat{x}_M]$
- ▶ Transposed embedded sequence: $\hat{X}^{(t)}$

Feature extraction

- ▶ Chroma Energy Normalized Statistics (CENS)

- ▶ Key transposition

- ▶ Given a CENS vector $\mathbf{x} = (x_0, x_1 \dots x_{11})^T$, the transposed vector is defined as follows:

$$\mathbf{x}^{(i)} = (x_{i\%12}, x_{(i+1)\%12} \dots x_{(i+11)\%12})^T$$

- ▶ Given a CENS sequence $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_M]$, the transposed sequence would be:

$$\mathbf{X}^{(i)} = [\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)} \dots \mathbf{x}_M^{(i)}]$$

- ▶ Vector Embedding

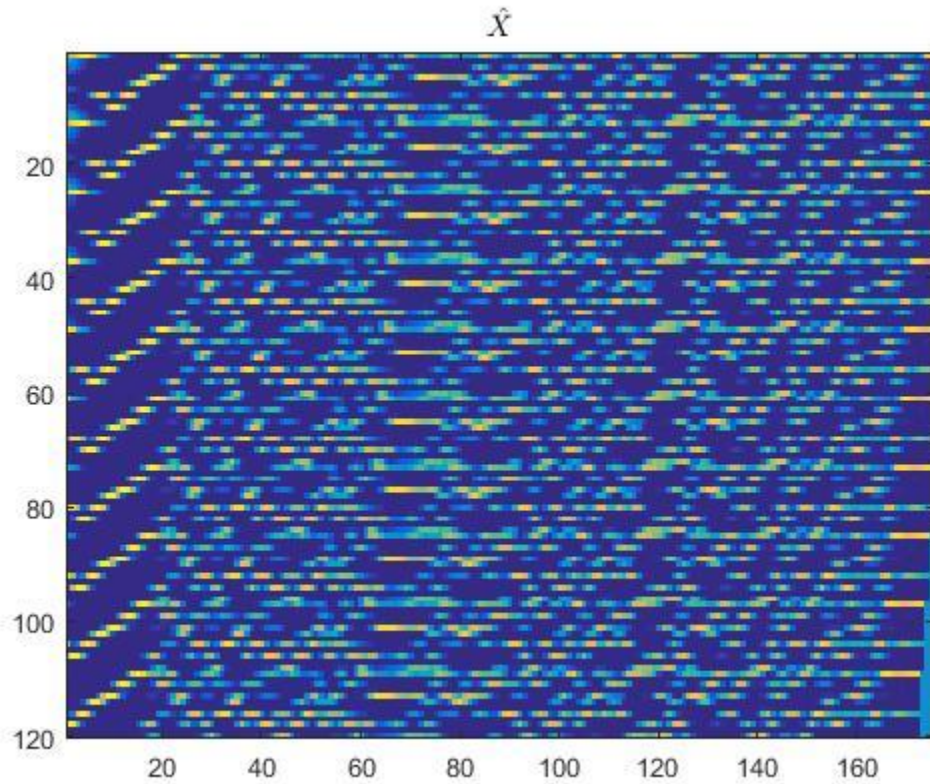
- ▶ Embedded vector: $\widehat{\mathbf{x}}_j = [\mathbf{x}_j^T, \mathbf{x}_{j-1}^T \dots \mathbf{x}_{j-(m-1)}^T], j = m, m+1 \dots M$

- ▶ Embedded sequence: $\widehat{\mathbf{X}} = [\widehat{\mathbf{x}}_m, \widehat{\mathbf{x}}_{m+1} \dots \widehat{\mathbf{x}}_M]$

- ▶ Transposed embedded sequence: $\widehat{\mathbf{X}}^{(i)}$

Feature extraction

► Chroma Energy Normalized Statistics (CENS)

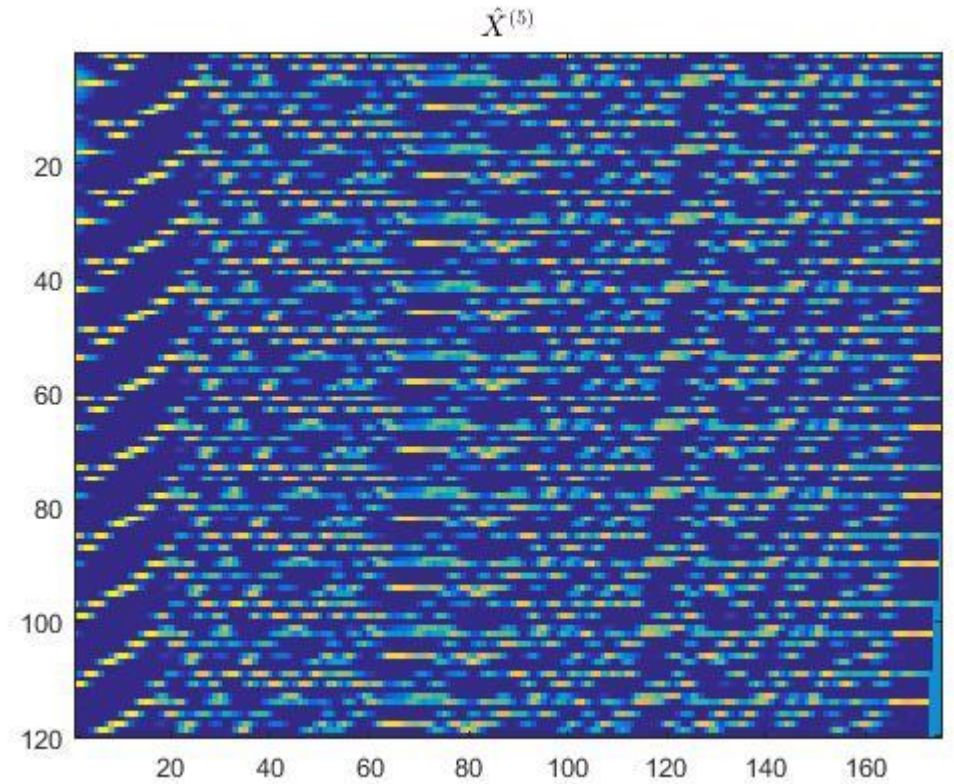


x_1

$x(i)$

$2 \cdot$

$\dots \lambda$



► Embedded vector: $x_j = [x_j, x_{j-1} \dots x_{j-(m-1)}], j = m, m+1, \dots, M$

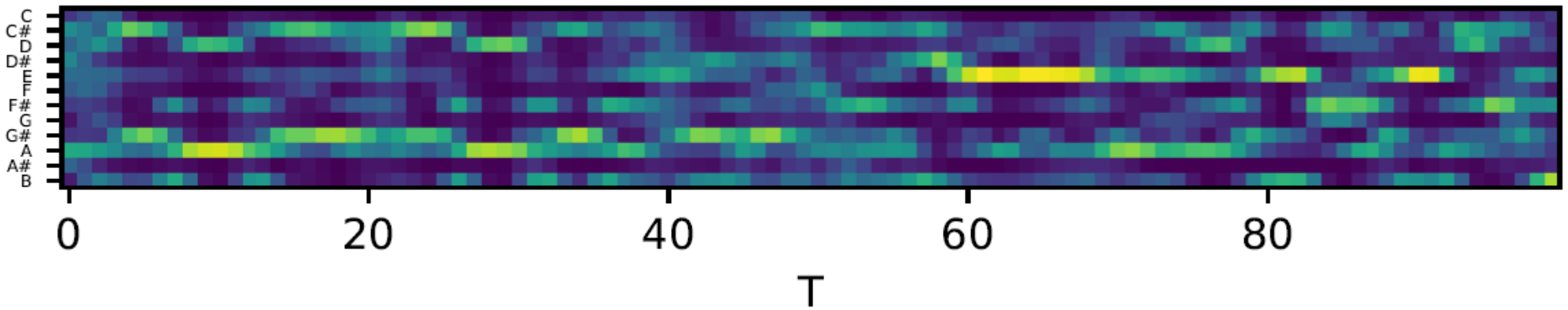
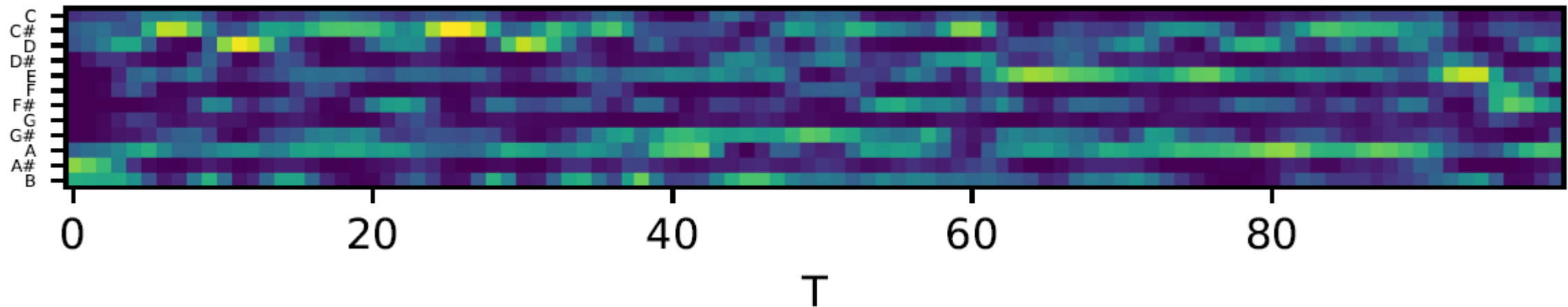
► Embedded sequence: $\hat{X} = [\hat{x}_m, \hat{x}_{m+1} \dots \hat{x}_M]$

► Transposed embedded sequence: $\hat{X}^{(i)}$

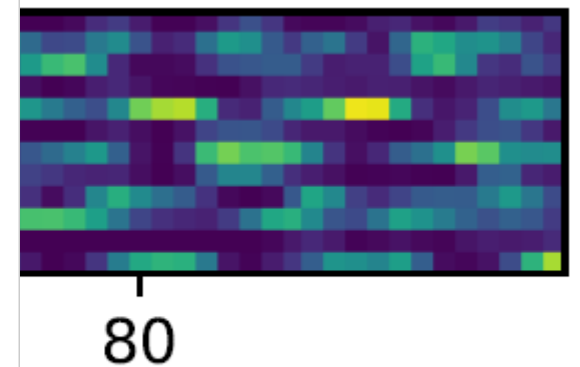
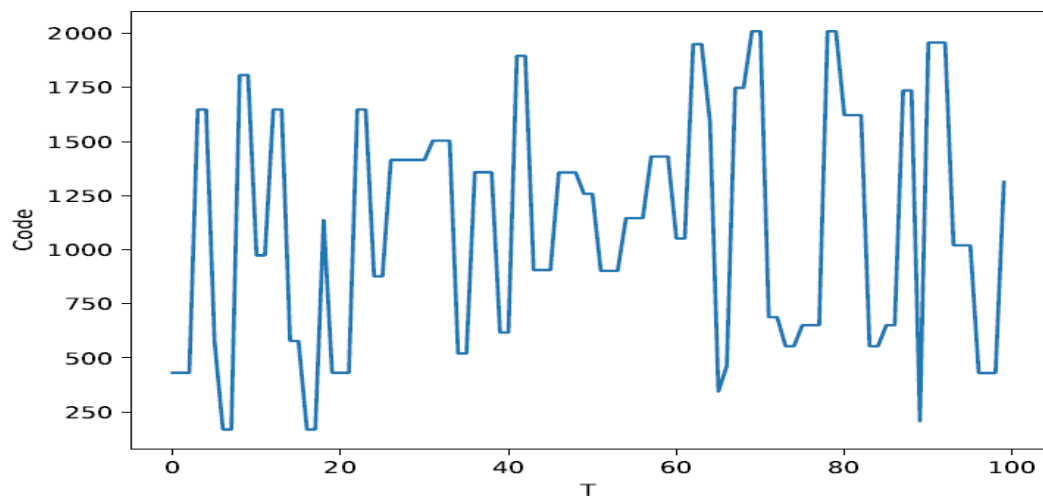
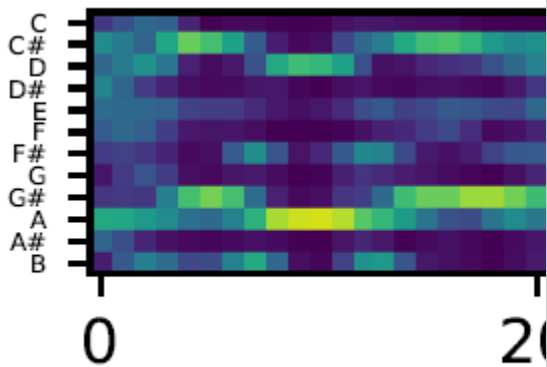
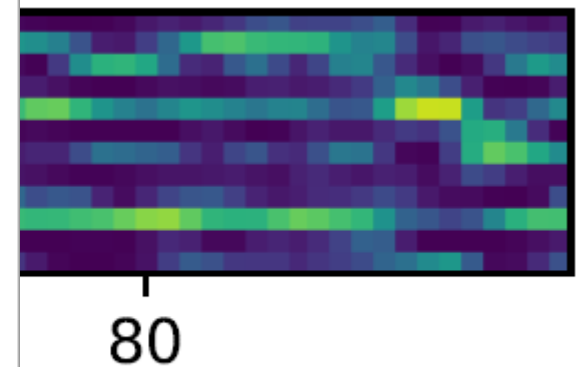
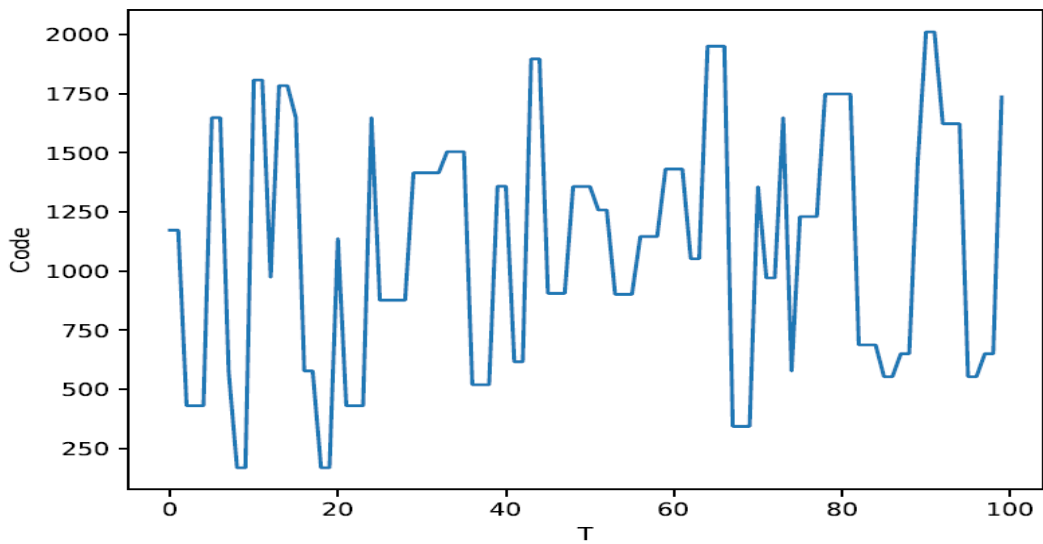
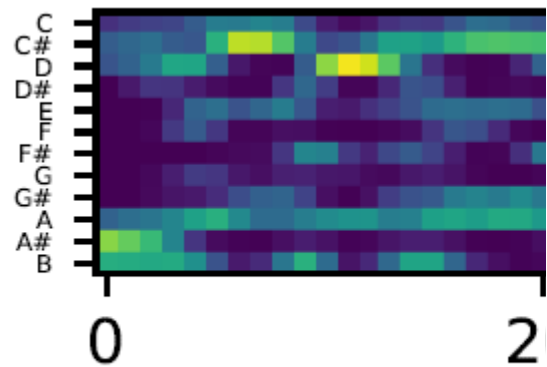
Vector quantization and encoding

- ▶ Vector quantization is used to cluster embedded vectors and a codebook is learnt for encoding.
- ▶ Reduce the impact of structural variations.
- ▶ Code sequences of cover songs reveal high similarity, while code sequences of different songs show little similarity.

Vector quantization and encoding



Vector quar



Bigram histogram and similarity

- ▶ Count the bigram histogram f
- ▶ The similarity between two songs is defined as:

$$S(u, v) = \max_i \sum_{a,b} \min\{f_u^{(i)}(a, b), f_v^{(0)}(a, b)\}$$

- ▶ Why use skipping bigram?
 - ▶ Consider the structural variations in cover songs
 - ▶ A simple example: consider two code sequences $\{1, 2, 3\}$ and $\{1, 3\}$, the similarity of bigram histogram is zero
 - ▶ Consider a gap s when constructing bigrams

Inverted index

- ▶ How to compute the similarity efficiently

$$S(u, v) = \max_i \sum_{a,b} \min\{f_u^{(i)}(a, b), f_v^{(0)}(a, b)\}$$

- ▶ A table is established to maintain the mapping from (a, b) to recording.
- ▶ Given a pair (a, b), we could get $\{(v, f_v^{(0)}(a, b)) \mid f_v^{(0)}(a, b) > 0\}$ quickly with the help of the table.

Retrieval

- ▶ Given a query u , code sequences are generated through embedding, transposition and encoding.
- ▶ Fixed i , for each bigram $(a, b) \in \{(a, b) \mid f_u^{(i)}(a, b) > 0\}$, we find $\{(v, f_v^{(0)}(a, b)) \mid f_v^{(0)}(a, b) > 0\}$ with the help of table.
- ▶ Enumerating $i \in \{-5, -4 \dots 5\}$, the algorithm computes the similarity between the query and the reference.

Experimental setting

▶ Evaluation metric

- ▶ Mean average precision (MAP)
- ▶ Precision at 10 (P@10)
- ▶ Mean rank of first correctly identified cover (MR1)

▶ Datasets

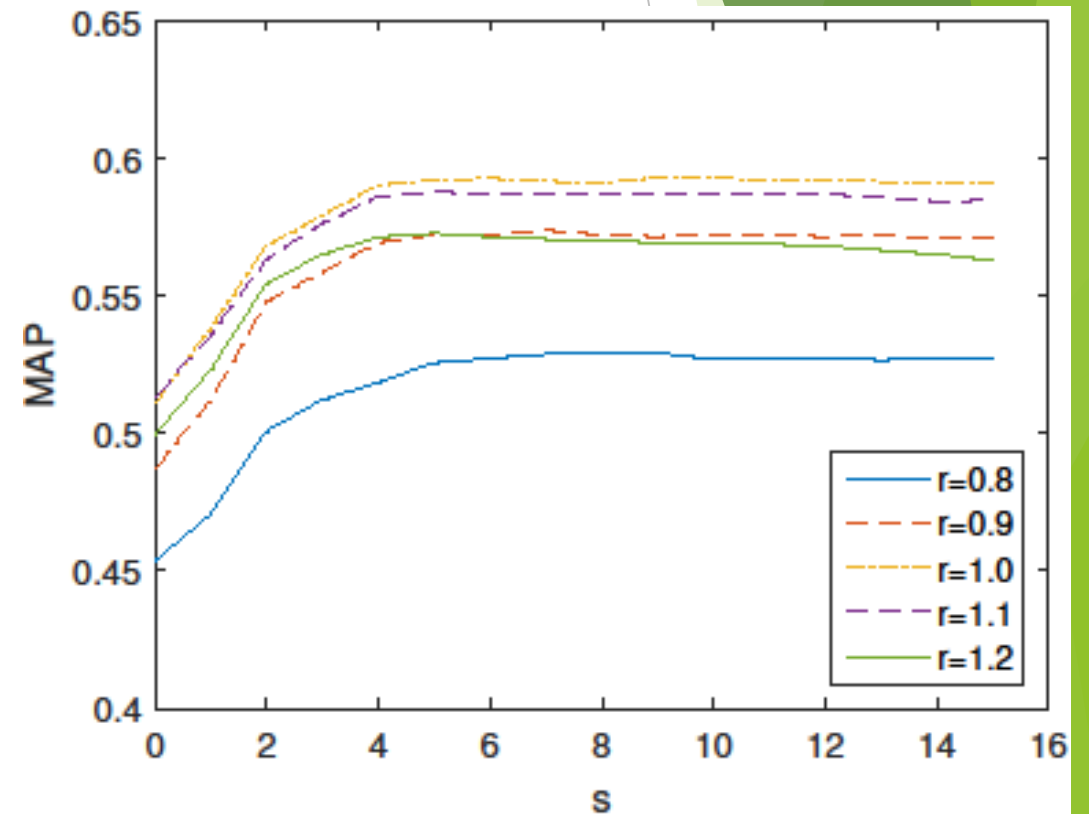
- ▶ Youtube350
- ▶ Music collection (MC)

Influence of hyperparameters

- ▶ Resample CENS sequences to simulate different speed
- ▶ Skipping bigrams help improve the precision

Influence of hyperparameters

- ▶ Resample CENS sequences to simulate different speed
- ▶ Skipping bigrams help improve the precision

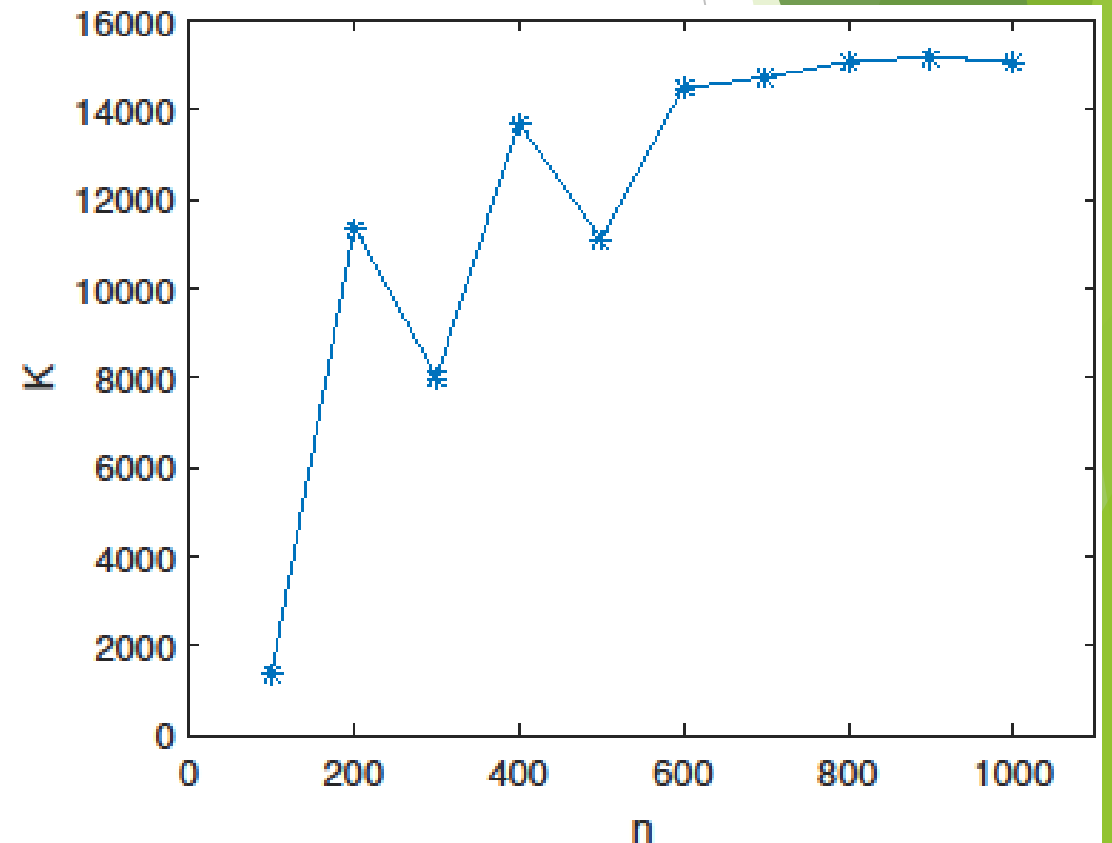


Influence of hyperparameters

- ▶ Explore how many codes are needed to ensure good performance
- ▶ Sub-linear relationship between N and K

Influence of hyperparameters

- ▶ Explore how many codes are needed to ensure good performance
- ▶ Sub-linear relationship between N and K



Comparison

- ▶ Highest P@10 and MR1 compared to state-of-the-art method
- ▶ Low time complexity

	MAP	P@10	MR1	Time/s	Complexity
DTW [19]	0.425	0.114	11.69	56.50	$O(NM^2)$
Silva et al. [19]	0.478	0.126	8.49	3.71	$O(NMS)$
Serra et al. [21]	0.525	0.132	9.43	2419.20	$O(NM^2)$
Silva et al. [18]	0.591	0.140	7.91	18.72	$O(NM \log M)$
Rafii CQT [22]	0.521	0.122	9.75	-	$O(NM^2)$
Rafii fingerprint [22]	0.648	0.145	8.27	-	$O(NM^2)$
Skipping bigrams	0.617	0.147	7.42	3.40	$O(M \log K)$

Conclusion & Future work

- ▶ Propose a skipping bigram model robust against structure and speed variations
- ▶ Design an inverted index for acceleration
- ▶ Achieve a high MAP with low time cost on a recent cover song dataset
- ▶ Adapt our approach to large-scale datasets

Thank you!