

# An Efficient Frequent Melody Indexing Method to Improve the Performance of Query-By-Humming Systems

Jinhee You; Sanghyun Park<sup>1</sup>

*Department of Computer Science, Yonsei University, Korea*

## Abstract

In recent years, the need to efficiently store and quickly retrieve large amounts of musical information has increased. Currently, text-based approaches are most widely used in the field of music information retrieval. However, users cannot use text-based approaches when they cannot remember the keywords associated with the desired music. In addition, since text-based approaches only search for the music that exactly matches a user query, they are not appropriate for either content-based searching or similarity searching. In this paper, we design and implement a Query-By-Humming (QBH) system, which can retrieve the melodies similar to users' humming. To make this QBH system reliable and efficient, the following three methods were proposed. First, we convert the melodies to be indexed into the corresponding strings, in order to increase search speed. The conversion method is designed to tolerate the errors involved in humming. Second, we extract significant melodies, which have a high chance of being queried, from music and then build a couple of indexes from them. For this task, we propose reliable methods are proposed for extracting melodies that occur frequently and for melodies that begin after a long rest. Third, we propose a three-step index searching method for minimizing database access and obtaining query answers as early as possible. Through the extensive experiments with a real-world data set, it was verified that this system has noticeable improvements over the N-gram approach, which is presently one of the most typical methods for QBH systems.

Keywords: Multimedia database; Music information retrieval; Query-by-humming system; Content-based indexing method

---

<sup>1</sup> Correspondence to: Sanghyun Park, Department of Computer Science, Yonsei University, Korea, [sanghyun@cs.yonsei.ac.kr](mailto:sanghyun@cs.yonsei.ac.kr).

## Jinhee You and Sanghyun Park

### 1. Introduction

As the use of multimedia information is becoming more common, users require new methods to retrieve multimedia information more easily and more quickly than the current methods. In order to satisfy these demands, many studies have been conducted. Currently, the most common approach for finding multimedia content is text-based retrieval. The text-based retrieval method finds multimedia content using keywords that are previously annotated. However, users cannot use the text-based method when they do not remember the keywords associated with desired multimedia content. Also, since the text-based method only searches for the multimedia content that exactly matches a user's query, it is not appropriate for similarity searching. In particular, the text-based retrieval method is most common in the field of music information retrieval. Using the text-based method, users find the desired music by submitting keywords such as the music title, the album title, and the name of the singer. However, the text-based method is not appropriate for music information retrieval when users do not remember relevant keywords or they want to find melodies that are similar to the desired music.

There have been many studies [1, 2, 3] which employ the content-based retrieval method in order to solve the problems aforementioned. Unlike the text-based method, the content-based method finds the desired information by directly comparing multimedia content. In addition, since the content-based method computes the similarity between multimedia content during the retrieval process, it is capable of performing similarity searching.

Query-By-Humming (QBH) [4] is a study of the content-based method in the area of music information retrieval. The approach of QBH is to find the music that contains melodies that are similar to the queries hummed via microphone. Since QBH systems are capable of retrieving the desired music even when users' humming is inaccurate, they are useful for people who are not talented singers or who have a speech impediment. However, QBH systems also have drawbacks. That is, if a database contains a large amount of music, then the operation to find the music similar to user's humming requires complicated computation and long search times. This is because QBH systems have to sequentially read music from a database, and then measure the similarity between each one of its melodies and the user's humming by calculating the distance with representative numeric data.

In this paper, we proposed an efficient QBH system. In order to reduce the excessive cost for calculating the distance between numeric melodies, we convert the numeric melodies into the corresponding strings and then compare the strings using some of existing string matching algorithms. Also, to guarantee the accuracy of the search results, we design the conversion method to tolerate the errors involved in humming. In addition, to reduce the search time, we extract significant melodies, which have a high chance of being queried, from music and then build a couple of indexes from them. The definition of significant melodies is divided into two types in this paper: 1) a melody that appears frequently in a single piece of music, and 2) a rest-unit melody that begins after a long rest within the music. Furthermore, we propose reliable methods for extracting melodies that occur frequently and for melodies that begin after a long rest.

One of the representative indexing methods for QBH systems is N-gram [5, 6, 7], where N consecutive notes are represented as a single symbol. The N-gram method has an advantage of simple indexing. However, since the N-gram method extracts N-grams from every possible position in the music, the number of N-grams to be stored is directly proportional to the length of the piece of music. Therefore, indexes are likely to become large when a music database is big. In addition, since the errors contained in humming are not considered when converting N-grams into corresponding symbols, the search result becomes less accurate as a users' humming contains more errors. In this paper, we try to overcome a number of the limitations of the N-gram method. Through experimental comparison with the N-gram method, we verify that our system has noticeable improvements in terms of the speed and accuracy of the search operations.

As shown in Figure 1, the proposed QBH systems have two major steps including system construction and query processing. In system construction, 1) feature data is extracted from the MIDI music, 2) feature data is stored in a feature database, 3) feature data of numeric type is converted to corresponding symbols of character type, and 4) significant string melodies are extracted and indexed. In query processing, 1) the user's humming is preprocessed and, as a result, the corresponding query string is obtained, and 2) three-step index searches are performed in order to minimize database access and obtain query answers as early as possible.

## Jinhee You and Sanghyun Park

The rest of the paper is organized as follows. Section 2 describes typical feature data and indexing methods commonly used in the area of content-based music retrieval. Section 3 and Section 4 explain in detail the system construction and query processing steps of the proposed QBH system, respectively. Section 5 shows and analyzes the results of experimental comparison with the N-gram method. Finally, concluding remarks are given in Section 6.

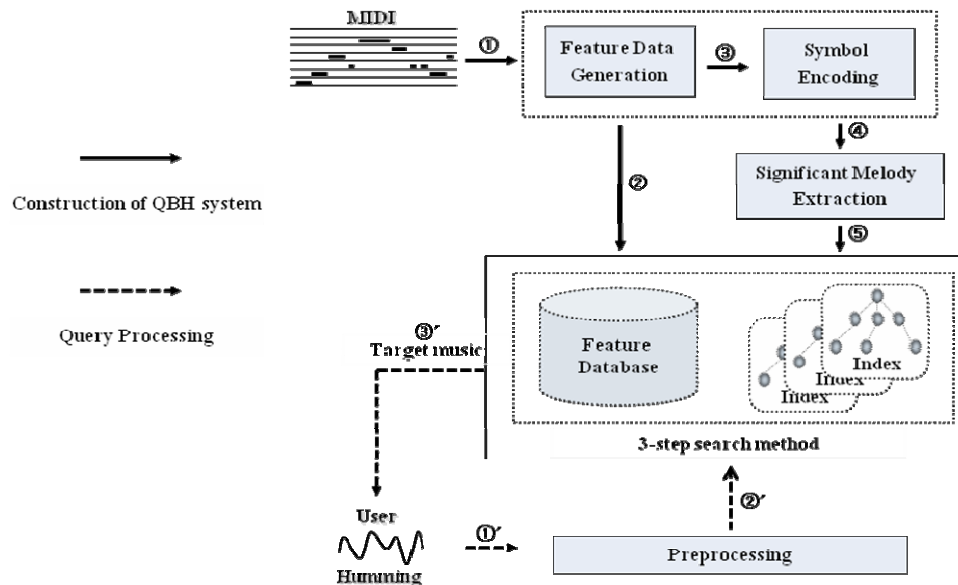


Figure 1. Construction and query processing of the proposed QBH system

## 2. Related Work

Recently, studies on finding feature data suitable for content-based music retrieval have been actively conducted. Feature data that is most appropriate for a specific content-based music retrieval system depends on the type of query, which is either an accurate query such as an original song or an inaccurate query such as humming. In particular, the interest in finding feature data that is suited to QBH systems has been growing [4, 8]. UDS strings [9] are one of the typical feature data sets for QBH systems. UDS strings express the relationship between pitches as one of three letters: 'U', 'D', and 'S'. If the pitch of a note is higher than that of its previous note, it is denoted as 'U'. If the pitch of a note is lower than that of its previous note, it is denoted as 'D'. And, if the pitch of a note is the same as its previous note, it is denoted as 'S'. UDS strings have a limitation in accuracy because they represent only pitch contours while ignoring the durations of notes. Therefore, the accuracy of this type of search results is very low when the user's humming contains many errors. There is a QBH system that uses tempo as its feature data [10]. However, since tempo data ignores the pitch and duration information, this type of system becomes inaccurate when the user's humming is not accurate.

Kline *et al.* [7] proved that QBH systems that use feature data that combine the pitch interval and the duration ratio of adjacent notes produce the most accurate search results. This is because both the pitch and note duration is taken into account, and it expresses them using accurate numeric values rather than rough melody contours. Based on this study, it was decided to use pitch interval and the duration ratio of adjacent notes as the feature data in the proposed QBH system.

## Jinhee You and Sanghyun Park

Another study on content-based music retrieval is to devise an efficient indexing method. The motivation for this study is that, when there are no appropriate index structures in a large music database, the operation to retrieve target music requires sequential access to the database and a consequently enormous number of comparisons and long search times. One of the most typical indexing methods for content-based music retrieval is N-gram [11, 12], where N adjacent notes are extracted at every possible position, which is just like sliding windows with N notes, and represented as a single symbol. Since the number of N-grams to be stored is proportional to the length of the music, indexes are likely to become large when a music database is big. In addition, since the errors contained in humming are not considered when converting N-grams into corresponding symbols, the search result becomes less accurate as the users' humming contains more errors. There are other studies that are based on spatial access methods. These studies treat the pitch and duration of a note as a point in two-dimensional space. One of the representative methods in this category uses R-tree [13]. Its operation to traverse down an R-tree is efficient. However, since there can be many notes that are mapped to the closely located points, its selectivity may not be good.

Most indexing methods for content-based music retrieval build the index for a specific part of a melody [15, 16, 17]. Assuming that users easily remember the beginning part of the music, Andreas [15] proposed to index only the initial part of the music. However, the proposed method does not guarantee good search performance when users query with the part of the music other than the beginning. Liu *et al.* [16] proposed to extract and index what is known as 'frequent melody'. However, since their method extracts frequent melodies without regular units, the extraction process is very complicated. In addition, since their method creates many candidates for frequent melody, their index becomes too large when the database is big. In contrast to the methods mentioned above, there is a different type of method [17], which does not index the music in databases, but rather indexes frequent queries using UDS strings. This has the advantage of obtaining short search times when the users query the melody that many other users have queried before. However, this has the drawback of accessing the entire database when some melodies are first queried. Moreover, UDS strings lead to low accuracy.

### 3. Construction of the QBH System

#### 3.1. Feature Extraction and Symbol Encoding

The first task in constructing the QBH system is to extract the melodies from the music. The melody is composed of a series of notes, which have the two basic components of pitch and duration. Therefore, when MIDI music is given, these two components are extracted from each note of the music. The pitch of a note is obtained directly from its Octave that is previously defined as the value of pitch in MIDI format. The duration of the notes and rests are determined using the length of a quarter note (i.e., time base). Note that, since music may have a different time base, we must scale the time bases of the music to an identical value. After extracting these two components from each note, we must calculate the boundary of each bar using the time signature, the time base, and the duration of all the notes and rests.

Next, we calculate the feature data (i.e., pitch interval and duration ratio), in order to express the melody precisely. Given two adjacent notes ( $P_i, D_i$ ) and ( $P_{i+1}, D_{i+1}$ ), their feature data are calculated using the expressions in Figure 2. The feature data between two adjacent notes takes the form of two-dimensional vectors and it is stored in the feature database as the information of the former note.

A rest between adjacent notes is treated as shown in Figure 3. Most QBH system studies have ignored the treatment of rests, but users are likely to hum passages that include rests. Unlike notes, rests have only duration values. Therefore, we represent the pitch of the notes as '-'. Feature data of a rest is calculated using either one of two ways according to the position. Consider the case where a rest follows a note. Since it is assumed that the pitch of the note does not change during the rest, the pitch interval is 0.00. The duration ratio is calculated using the expression shown in Figure 3 with the duration values of the former note and the latter rest. The calculated

## Jinhee You and Sanghyun Park

feature data is stored as the information of the former note. Unfortunately, the same feature data can be generated if two notes of the same pitch are adjacent to each other. However, because it is not easy for users to distinguish these two situations and hum them differently, they are not handled separately. Next, consider the case where a note follows a rest. In this case, the rest cannot be expressed as a melody (i.e., dumb melody). Therefore, the feature data of the rest is not calculated. As mentioned earlier, all of the extracted feature data are stored in the feature database.

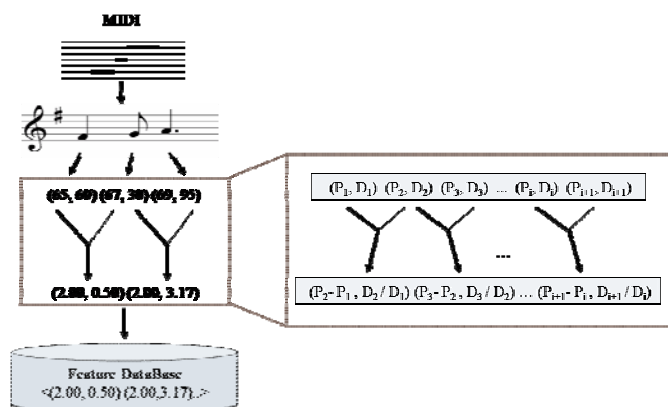


Figure 2. Extraction of feature data

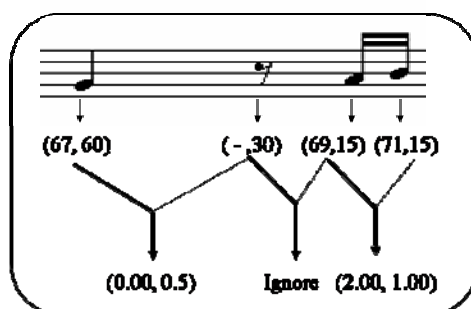


Figure 3. Treatment of rests

To calculate the similarity between the query and each melody in the feature database, many methods use the Euclidean distance. However, the operation to calculate the Euclidean distance has a high CPU cost. Therefore, the performance of the methods based on the Euclidean distance deteriorates as the size of database increases. To resolve this problem, we take the approach of converting the melodies, which are a series of two-dimensional numeric vectors, into corresponding strings of characters, and then employ efficient string matching algorithms to calculate the similarity of the melodies. Figure 4 shows the proposed Character Encoding Standard (CES) with which two-dimensional numeric vectors of pitch interval and duration ratio are changed into corresponding characters. More precisely, the two-dimensional space that has the pitch interval on the X axis and the duration ratio on the Y axis is divided into multiple cells of possibly different sizes and then each cell is assigned a unique character. Under this conversion scheme, it is possible for more than one vector located closely in the two-dimensional space to be changed to the same character. This has the effect of tolerating the errors included in humming. Figure 5 shows the result of encoding a series of two-dimensional vectors into a corresponding string using the example of the CES in Figure 4. To achieve high accuracy, the CES must be chosen optimally for the

## Jinhee You and Sanghyun Park

underlying music database. As explained in Section 5.3, we take an experimental approach where we generate a set of CES suitable for the underlying music database and choose the one that has the highest accuracy.

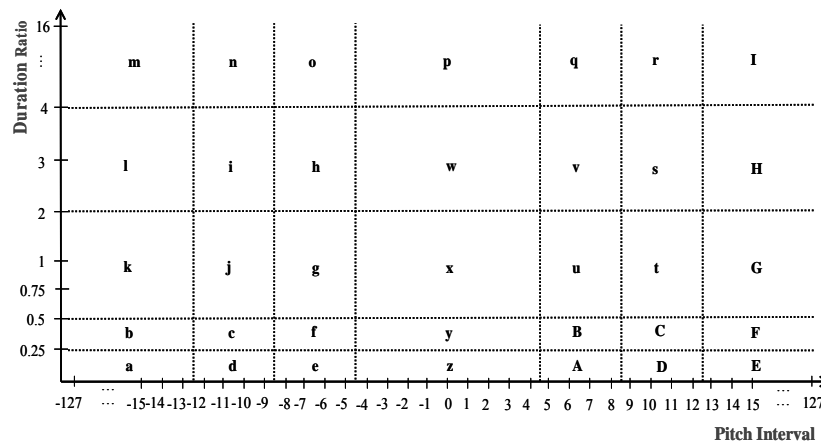


Figure 4. Example of CES for changing two-dimensional vectors of pitch interval and duration ratio to corresponding characters

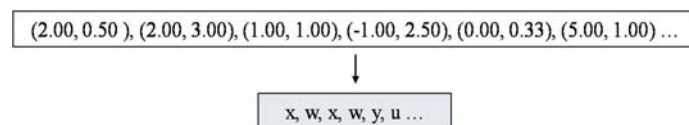


Figure 5. Result of encoding a series of two-dimensional vectors into a corresponding string using the example of CES in Figure 4

### 3.2. Significant Melody Extraction

Among the many melodies that comprise a single piece of music, the significant melody is the one that is able to represent the music. In this paper, significant melodies, which have a high probability of being queried by users, are divided into two types: frequent melodies that occur repetitively within the music and rest-unit melodies that begin after a long rest. Users typically remember frequent melodies and rest-unit melodies more easily than the other melodies. Therefore, we extract both types of melodies from the music and indexes are built from them. Since the process of extracting rest-unit melodies is straightforward, the focus here will be on the process of extracting frequent melodies.

Given a piece of music, we first carry out a string pattern analysis, treating each bar as a single processing unit. That is, initially a set of strings is extracted, each of which represents a bar, and then execute a string matching algorithm such as the KMP (Knuth-Morris-Pratt) algorithm [17] to count the number of occurrences of each string. Thus, the following data is obtained for each string: 1) a list of bar positions (*BarId\_list*), 2) the number of occurrences of the string (*C*), and the descending order of value *C* (*R*). As an example, Table 1 is from the song titled 'After for a long time', which is a popular Korean ballad, after executing the proposed string pattern analysis.

After executing the string pattern analysis, we create the groups of bars. As Figure 6 indicates, each bar whose rank value (i.e., *R*) is 1 produces a separate group and becomes the key bar (KB) of its own group. Therefore, the initial number of groups is the same as the number of bars whose rank values are 1. Next, to construct frequent melodies, each group expands in both forward and backward directions by joining the bars that are adjacent to the group and that have a high *C* value. This joining step requires two parameters, bar selection threshold (BST) and

## Jinhee You and Sanghyun Park

frequent melody length threshold (FMLT). The BST is for allowing only frequent bars to be joined into the groups. That is, in the case that a bar has a C value less than the bar selection threshold, it cannot be joined with the group even if it is adjacent to the group. The FMLT is for controlling the length of frequent melodies. Whenever a bar is joined with the group being expanded, the difference from the R value of the bar to the R value of the key bar of the group is calculated and then added to the variable BSum (i.e., bar sum), whose initial value is 0. If the joining of a bar to a group makes the value of BSum exceed the frequent melody length threshold, the bar is not allowed to join the group.

Bar String	List of Bar Positions	Count (C)	Rank (R)
<j, x, k, j, i, x>	18, 22, 35, 39	4	1
<k, k, j, x, g>	19, 36	2	2
<j, x, j, k, k, k>	20, 37	2	2
<k, j, j, w, g>	21, 38	2	2
...	...	...	...
<j, x, g, k, E, x, i, x>	17	1	3

Table 1. Data obtained from the song titled ‘After for a long time’, a popular Korean ballad, after executing the proposed string pattern analysis

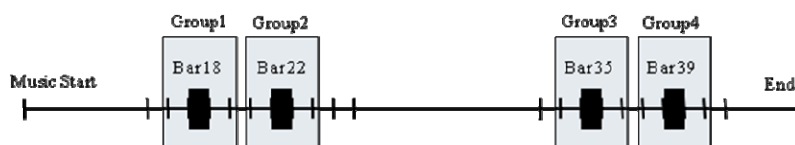


Figure 6. A set of initial groups constructed from the bars of Table 1 whose rank values are 1

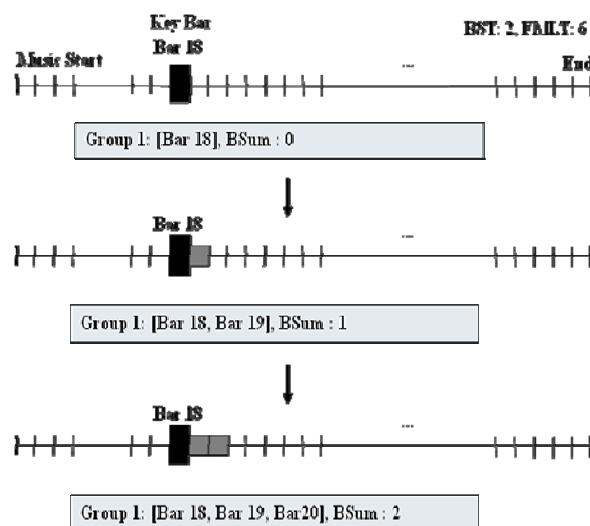


Figure 7. First few stages of creating the frequent melody around bar 18

## Jinhee You and Sanghyun Park

Based on Table 1, Figure 7 shows the first few stages of creating the frequent melody around bar 18. This example assumes the values of BST and FMLT are 2 and 6, respectively. The first bar with the greatest value of C, which is bar 18, becomes the key bar of the first group. Next, both sides of the key bar are considered, which are bar 17 and bar 19, one by one. Since the C value of bar 17 is less than the value of BST (i.e.  $1 < 2$ ), bar 17 cannot be joined with the first group and the extension to the left direction cannot be executed further. Consider the joining of bar 19 to the first group. The C value of bar 19 is not smaller than the value of BST (i.e.  $2 \geq 2$ ). The difference from the R value of bar 19 to the R value of the key bar (i.e., bar 18) of the first group is 1 (i.e.,  $2-1=1$ ) and the current value of BSum is 0. If this difference is added to the current value of BSum, 1 will be obtained, which does not exceed the value of FMLT. Since bar 19 satisfies both parameters, BST and FMLT, bar 19 is joined to the first group, the value of BSum is updated to 1, and the extension to the right direction (i.e., bar 20) is continued. Generally, such an extension process ends when progress cannot proceed to either the left or right directions. At this moment, the extent of the group becomes a frequent melody.

As shown in Figure 8, it is possible for a frequent melody to overlap with other frequent melodies. In this case, rather than treating each one of overlapped frequent melodies separately, they are merged into a single one in order to reduce the storage space. Eventually, as shown in Figure 9, the two frequent melodies from the example of Table 1 are obtained.

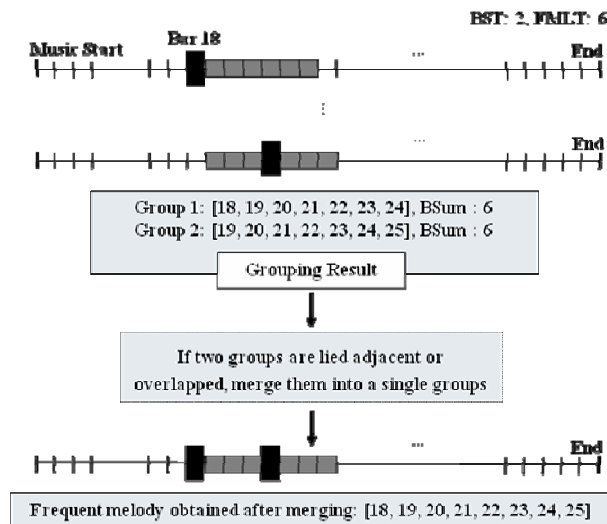


Figure 8. Handling of frequent melodies overlapped each other

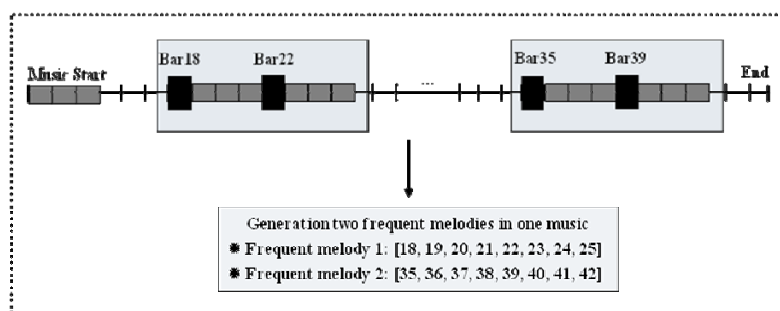


Figure 9. Frequent melodies obtained finally from the example of Table 1



## Jinhee You and Sanghyun Park

### 3.3. Indexing

To speed up the retrieval of the desired music, three suffix trees [18] are built, as shown in Figure 10. The idea of suffix tree comes from trie [19]. A trie is an index structure used for indexing sets of keywords of varying sizes. A suffix trie [18] is a trie whose set of keywords comprises the suffixes of a string. Nodes with a single outgoing edge can be collapsed, yielding the structure known as the suffix tree.

First, for each song stored in database, we retrieve the corresponding encoded string, extract every suffix from the string, and then insert them into the suffix tree named *all music suffix tree index* or *1<sup>st</sup> index*. The main purpose of this index is to minimize the direct access to the database that may occur when users query melodies other than frequent melodies or melodies starting after a long rest. Next, for each frequent melody discovered in the music database, every suffix is extracted and inserted into the suffix tree named *frequent melody suffix tree index* or *2<sup>nd</sup> index*. If users query some part of the frequent melodies, the desired music can be quickly retrieved using this index. Finally, for each rest-unit melody that begins after a long rest, every suffix is extracted and inserted into the suffix tree named *rest-unit melody suffix tree index* or *3<sup>rd</sup> index*. If users query some part of the rest-unit melodies, the desired music can be quickly retrieved using this index.

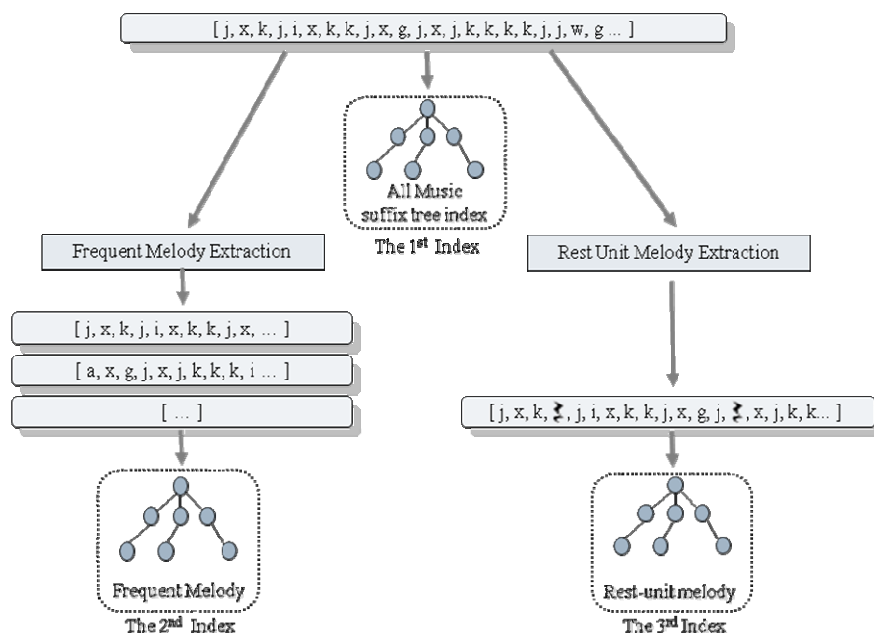


Figure 10. Three suffix trees for speeding up the retrieval of the desired music

## 4. Query Processing

Queries hummed by users are processed with two stages, preprocessing and searching. In preprocessing, a query is transformed into a series of two-dimensional feature vectors including the pitch interval and duration ratio, and then into the corresponding string using the same CES employed at the time of index construction.

As shown in Figure 11, there are three steps in the searching stage. In the first step, the frequent melody and rest-unit melody suffix tree indexes are searched in parallel for the string that corresponds to the submitted query. If the query looks for some part of the frequent melodies or rest-unit melodies, the target music will be found during the first searching step. Otherwise, the second step involves searching all music suffix tree indexes for the

**Jinhee You and Sanghyun Park**

query string. If the target music is not found during the second searching step, then the third searching step is utilized, in which the Euclidean distance is used as a similarity measure, where the feature database is searched for the series of two-dimensional feature vectors that correspond to the submitted query. Remember that the goal of this query processing approach is to minimize the direct access to the database. This goal can be attained by processing as many queries as possible within the first or second searching step.

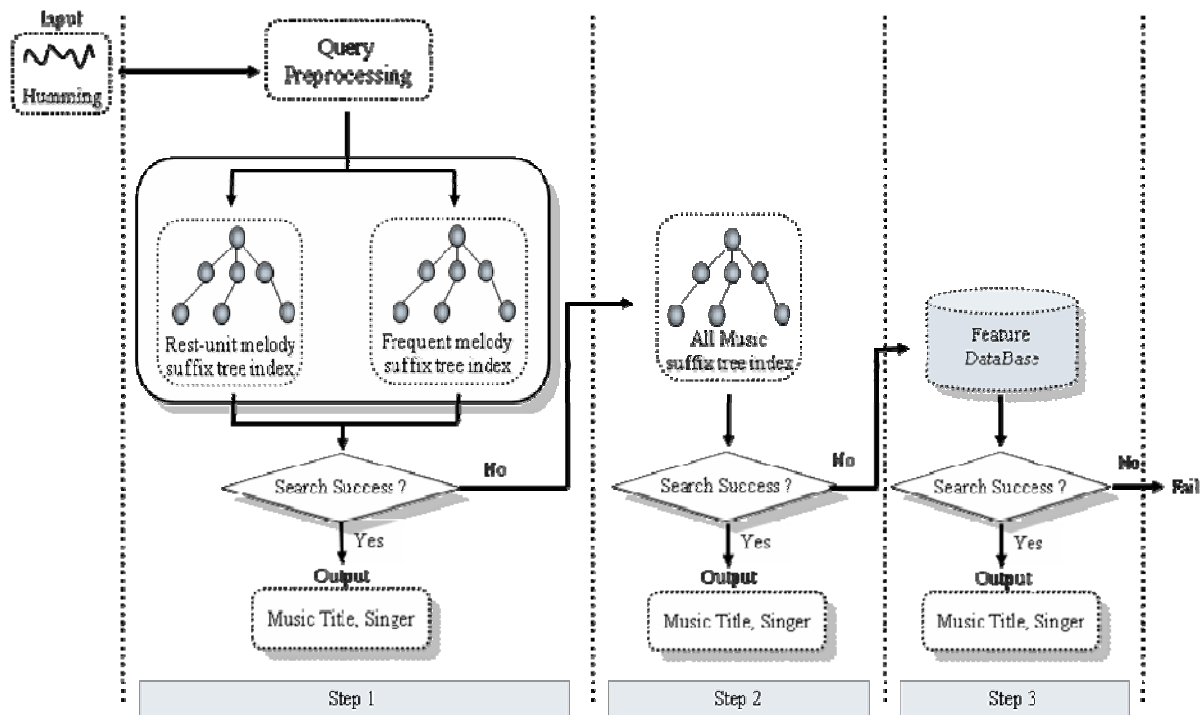


Figure 11. Query processing with two stages, preprocessing and searching

## 5. Performance Evaluation

In this section, we evaluate the proposed QBH system through experiments with real data. We collected 3000 Korean songs to construct a music database and, as shown in Table 2 and Table 3, let users hum the songs they wanted to retrieve. As shown in Table 2, for each  $i$  from 1 to 8, 30 queries were collected whose Euclidean distances to the corresponding original melodies in the music database are  $2*i$ . Also, as shown in Table 3, for each  $i$  from 1 to 6, 30 queries were collected that were hummed for  $5*i$  seconds and calculated the average number of notes contained in the 30 queries.

When the QBH system was designed, we assumed that the users would query frequent melodies or rest-unit melodies most often. To verify this assumption, we collected 700 queries of real humming from 30 people and then analyzed the queries. Figure 12 indicates that 76% of the queries contained frequent melodies, 38% of the queries contained rest-unit melodies, and 40% of the queries were the beginning parts of music. Note that some melodies can be included in more than one melody type. For example, it is possible for a frequent melody to start at the beginning of music or after a long rest. This analysis gives strong evidence that the assumption is realistic.

Jinhee You and Sanghyun Park

Euclidean distance between queries and their corresponding original melodies	Number of queries
2	30
4	30
6	30
8	30
10	30
12	30
14	30
16	30

Table 2. Number of queries whose Euclidean distances to the corresponding original melodies in the music database are  $2*i$  for each  $i$  from 1 to 8

Query time (sec.)	Average number of notes in queries	Number of queries
5	7	30
10	13	30
15	20	30
20	26	30
25	30	30
30	33	30

Table 3. For each  $i$  from 1 to 6, the number of queries that were hummed for  $5*i$  seconds and their average number of notes

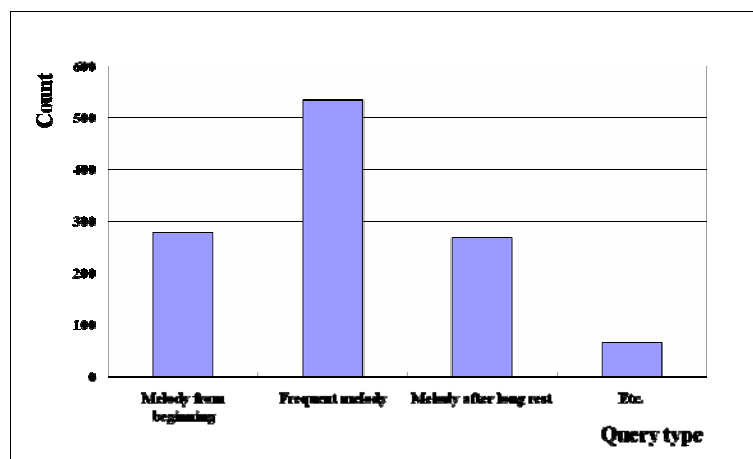


Figure 12. Number of queries contained in frequent melodies, rest-unit melodies, and melodies starting at the beginning of music, respectively

### 5.1. Determination of an Optimal Value of Bar Selection Threshold (BST)

Remember that we used two parameters, BST and FMLT, to control the length of the frequent melodies when the QBH system was constructed. In this subsection, we determine the optimal value of the BST by examining the pattern changes in the index size, the search time, and the accuracy of the QBH system according to varying values of BST. The optimal value of FMLT is determined in the next subsection.

**Jinhee You and Sanghyun Park**

To analyze the pattern change in the index size, we observed the total number of characters stored in the frequent melody suffix tree index while changing the BST value from 1 to 4. For this experiment, we set the value of FMLT to 6, generated frequent melodies, and selected from the frequent melodies 30 queries of 15 seconds. The reason of setting the value of FMLT to 6 is that it makes frequent melodies longer than 15 seconds. Figure 13 shows the result of this experiment. The frequent melody suffix tree index is largest when the value of BST is 1, but it becomes smaller as the value of BST increases. The result indicates that frequent melodies that are stored in the index get longer as the value of BST increases.

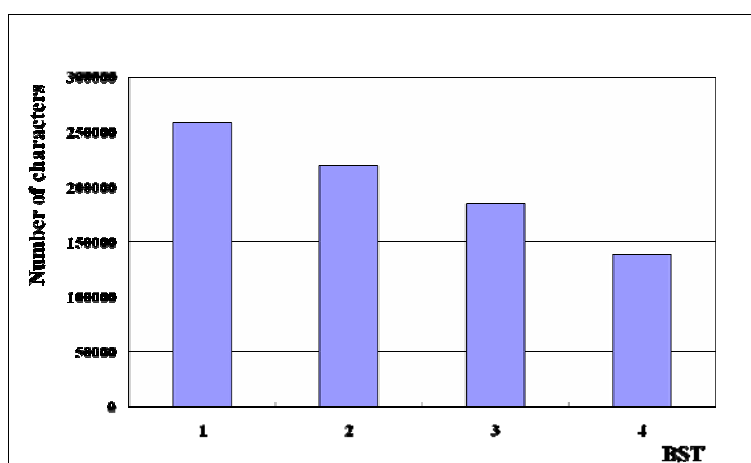


Figure 13. Total number of characters stored in the frequent melody suffix tree index when the value of BST is 1, 2, 3, and 4, respectively (FMLT=6, 3000 songs)

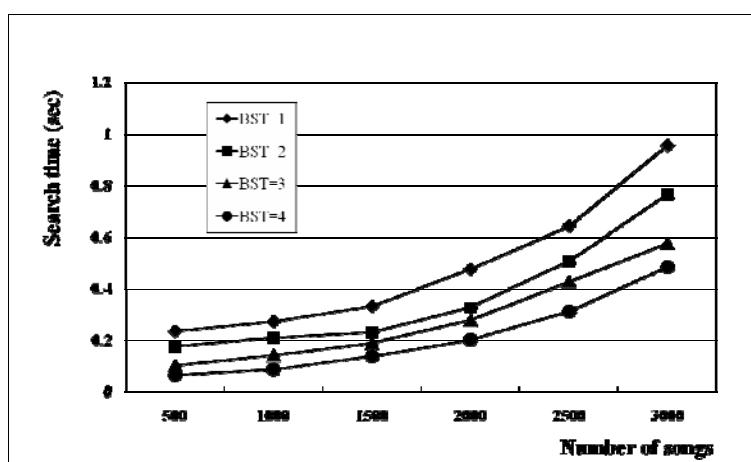


Figure 14. Search time of the frequent melody suffix tree index when the value of BST is 1, 2, 3, and 4, respectively (FMLT=6, query length=15 seconds)

Next, we examined the pattern change in the amount of search time within the frequent melody suffix tree index, while increasing the value of BST from 1 to 4. Figure 14 shows the result of this experiment. The X-axis is the number of songs stored in the database and the Y-axis is the search time of the frequent melody suffix tree index. Regardless of the amount of music stored in the database, the search time increases conspicuously as the

**Jinhee You and Sanghyun Park**

value of BST gets larger. This result can be interpreted as the growth of the BST value causes the increment of the index size and subsequently the increase in search time.

Lastly, we examined the pattern change accuracy (i.e., recall and precision) of the frequent melody suffix tree index, while increasing the value of BST from 1 to 4. Figure 15 and Figure 16 show the result of this experiment. Both recall and precision get lower as the BST increases. More specifically, both recall and precision get lower slightly when the value of BST changes from 1 to 2, but they decrease abruptly when the value of BST changes to 3 and then to 4. There are two possible reasons for this phenomenon. The first reason is that a significant number of frequent melodies obtained when the value of BST is 3 or 4 are shorter than the hummed queries. In this case, it may not be possible to obtain the target music by searching only the frequent melody suffix tree index. The second reason is that some of the frequent melodies may be unreliable when the value of BST is too large. For example, when the value of BST is 4, users may not get the target music if they enquire the melody that appears three times in the music. Considering the trade-offs among the index size, search time, and accuracy obtained by the above experiments, we can conclude that the optimal value of BST is 2.

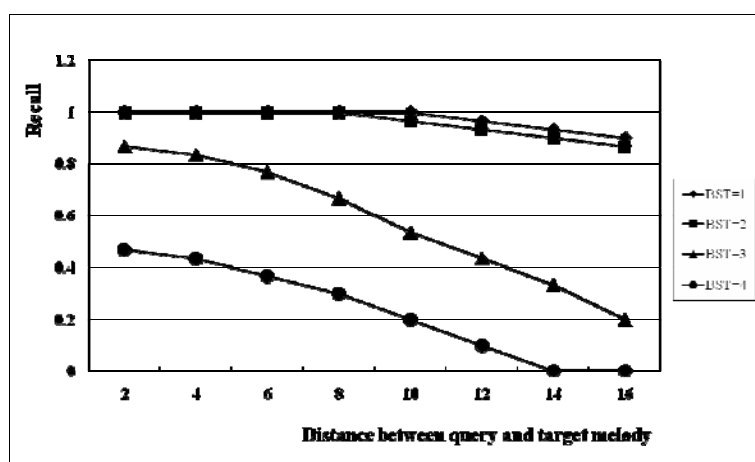


Figure 15. Recall of the frequent melody suffix tree index when the value of BST is 1, 2, 3, and 4, respectively (FMLT=6, 3000 songs, query length=15 seconds)

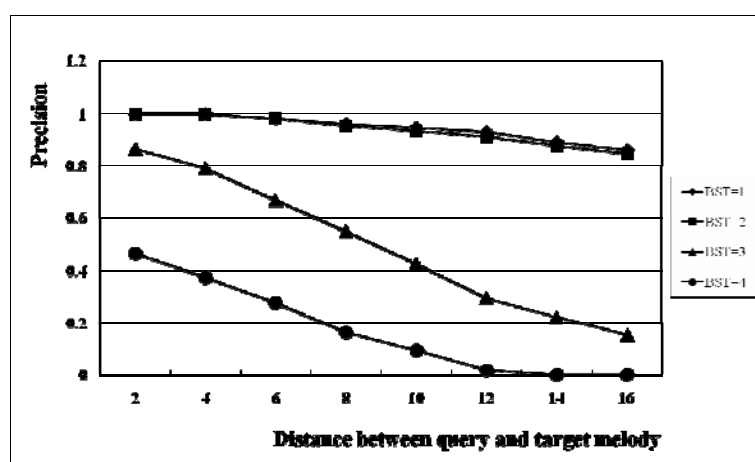


Figure 16. Precision of the frequent melody suffix tree index when the value of BST is 1, 2, 3, and 4, respectively (FMLT=6, 3000 songs, query length=15 seconds)

**Jinhee You and Sanghyun Park****5.2. Determination of an Optimal Value of Frequent Melody Length Threshold (FMLT)**

In this subsection, we determined the optimal value of FMLT by examining the pattern changes in the index size, search time, and accuracy of the QBH system according to varying FMLT. Just like the experiment for determining the optimal value of BST, in this experiment, we constructed a music database with 3000 songs and then extracted queries hummed for 15 seconds. Using the result of the previous experiment, we set the value of BST to 2 in this experiment. Figure 17 shows the total number of characters stored in the frequent melody suffix tree index when the value of FMLT is 2, 4, 6, 8, and 10, respectively. The figure indicates that the size of the index grows almost linearly as the value of FMLT increases. Figure 18 shows the search time of the frequent melody suffix tree index for each FMLT. Here, the X-axis is the number of songs stored in the database and the Y-axis is the search time of the frequent melody suffix tree index. The figure indicates that, with any number of songs stored in the database, the search time of the index increases as the value of FMLT grows.

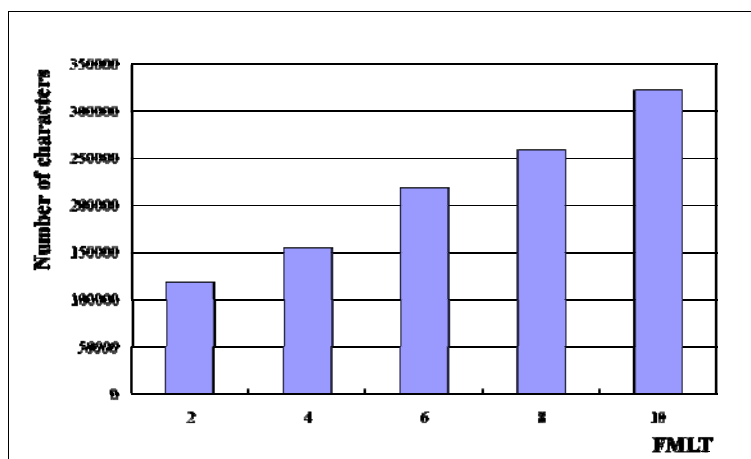


Figure 17. Total number of characters stored in the frequent melody suffix tree index when the value of FMLT is 2, 4, 6, 8, and 10, respectively (BST=2, 3000 songs)

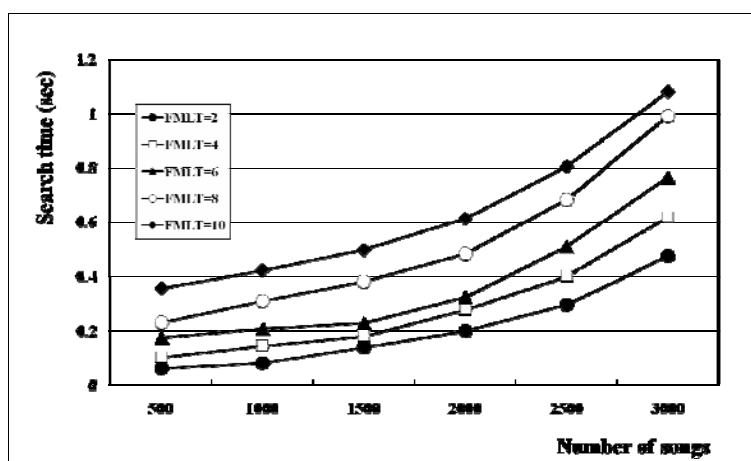


Figure 18. Search time of the frequent melody suffix tree index with various values of FMLT and various numbers of songs stored in the music database (BST=2, query length=15 seconds)

## Jinhee You and Sanghyun Park

Next, we examined the pattern change in accuracy (i.e., recall and precision) of the frequent melody suffix tree index, while increasing the value of FMLT from 2 to 10 by increments of 2. Figure 19 and Figure 20 show the result of this experiment. In these two figures, the X-axis is the distance between queries and their corresponding target music and the Y-axis is the recall and precision, respectively. For any value of FMLT, both recall and precision decrease as the distance between the queries and their corresponding target music increases. Also, for any distances between queries and their corresponding target music, both recall and precision increase as the value of FMLT increases. This result indicates that a significant number of frequent melodies obtained when the value of FMLT is small are shorter than the hummed queries.

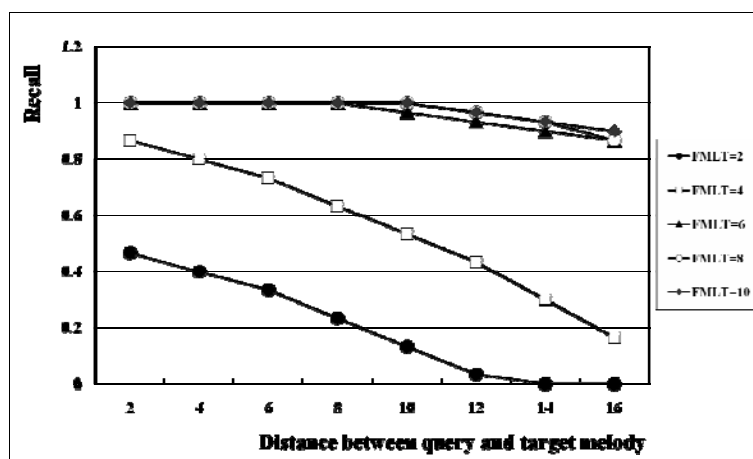


Figure 19. Recall of the frequent melody suffix tree index with various values of FMLT and various distances between queries and their corresponding target music (BST=2, 3000 songs, query length=15 seconds)

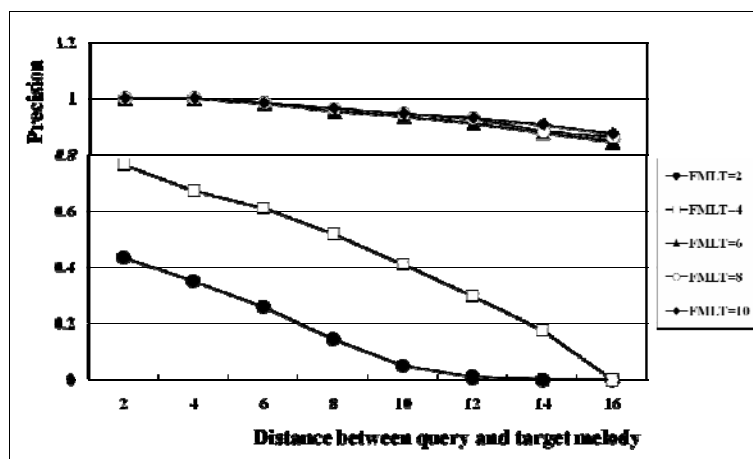


Figure 20. Precision of the frequent melody suffix tree index with various values of FMLT and various distances between queries and their corresponding target music (BST=2, 3000 songs, query length=15 seconds)

Considering the trade-offs among the index size, search time, and accuracy obtained by the above experiments, it is concluded that the optimal FMLT is 6, which leads to the creation of the index that is fairly fast, relatively small, and shows accuracy that is comparable to the indexes with a larger value of FMLT (i.e., 8 or 10).

## Jinhee You and Sanghyun Park

Note that the value of FMLT must be determined considering the typical length of queries, in order to avoid the situation where frequent melodies to be stored in index are shorter than the queries.

### 5.3. Selection of an Optimal Character Encoding Standard

The accuracy of the QBH system is highly affected by the CES employed when converting two-dimensional vectors of feature data to their corresponding characters. Therefore, in this section, the experimental approach to determine the optimal standard for character encoding is explained. For this experiment, as shown in Figure 21, we generated four types of CES by dividing each axis of the two-dimensional space with non-overlapping ranges of possibly different sizes. Then the QBH system was constructed with each one of the four CESs and evaluated its accuracy. According to the result of previous experiments, we set the value of the BST and FMLT to 2 and 6, respectively, in this experiment.

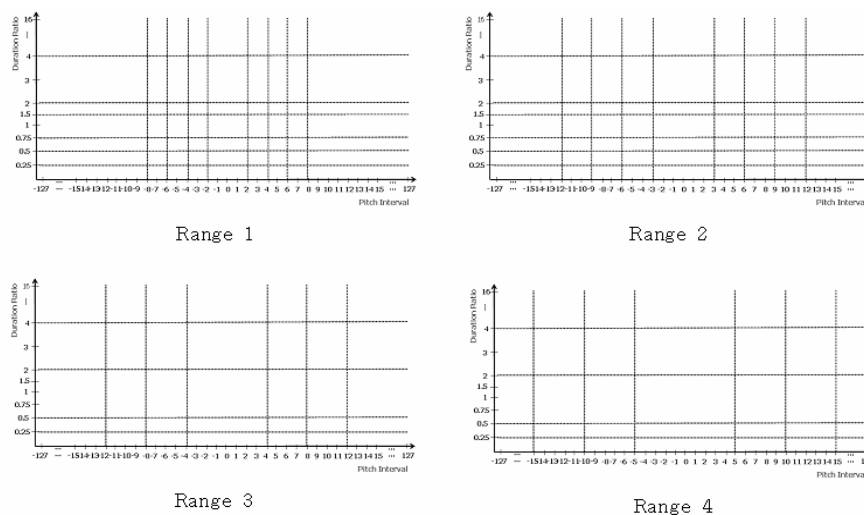


Figure 21. Four types of character encoding standard

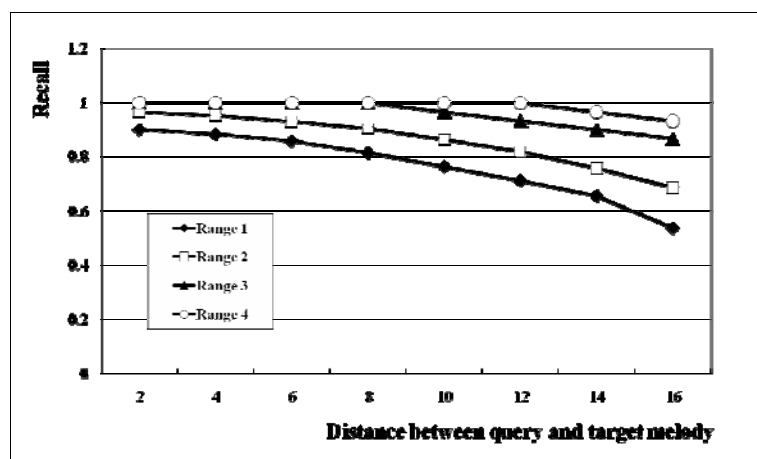


Figure 22. Recall of four character encoding standards with increasing distance between the queries and their corresponding target music (BST=2, FMLT=6, 3000 songs, query length=15 seconds)



## Jinhee You and Sanghyun Park

Figure 22 compares the recall of four character encoding standards including 'Range 1', 'Range 2', 'Range 3', and 'Range 4', while changing the distances between the queries and their corresponding target music. The recall of all four standards decreases gradually as the distances between the queries and their target music increase. This figure also indicates that the recall of Range 4 is higher than the other three standards. This result can be interpreted as follows: 1) the cells in the middle area, where lots of feature data are positioned, of Range 4 are bigger than those of other standards, 2) more feature data in the middle area are assigned the same character, and therefore, 3) even though errors are contained in humming, it is more likely for users to obtain the target music. On the contrary, the recall of Range 1 is lower than the other three methods. The reason for this is that the cells in the middle area of Range 1 are smaller than those of other standards and therefore the accuracy of the query result is very sensitive to the errors contained in humming.

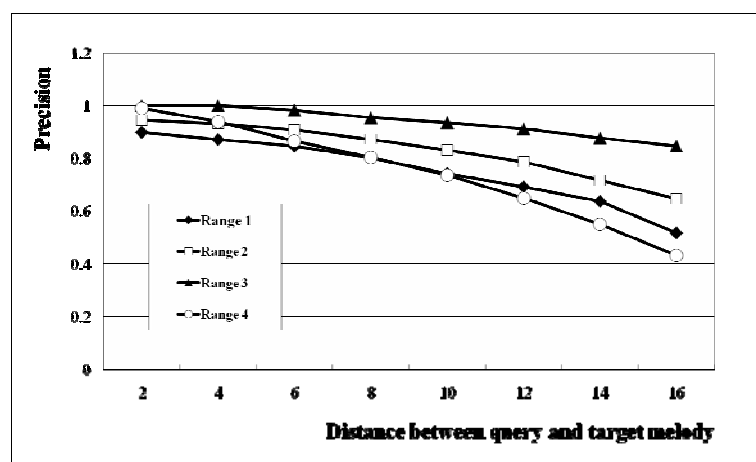


Figure 23. Precision of four character encoding standards with increasing distance between the queries and their corresponding target music (BST=2, FMLT=6, 3000 songs, query length=15 seconds)

Figure 23 compares the precision of four character encoding standards, while changing the distances between the queries and their corresponding target music. The precision of all four standards falls noticeably as the distances between the queries and their target music increase. This is because wrong answers are retrieved more often as the distances between the queries and their target music increase. Note that the precision of Range 4 drops most rapidly because the cells in its middle area are largest. Considering the trade-offs between the pattern changes in recall and precision, Range 3 was selected as the optimal standard of character encoding.

### 5.4. Comparison with the N-gram Method

To evaluate the performance improvements of the proposed QBH system, we compared the search time and accuracy of the proposed system with those of the N-gram method. As shown in Table 4, the methods involved in this comparison are 1) a **naive DB scan** that searches only the feature database without looking up any indexes, 2) the **2-step search method** that first searches all of the music suffix tree index and then, if necessary, searches the feature database, 3) the **3-step search method** that first searches the frequent melody suffix tree index and rest-unit melody suffix tree index in parallel and then, if necessary, searches all of the music suffix tree index and lastly, if necessary, searches the feature database, and 4) the **N-gram method**.

For this experiment, a music database was constructed with up to 3000 songs and collected 300 queries that last for up to 15 seconds, whose distances to the corresponding target music are 4. The standard of character encoding was Range 4 in Figure 21, and value of BST and FMLT were 2 and 6, respectively.

Jinhee You and Sanghyun Park

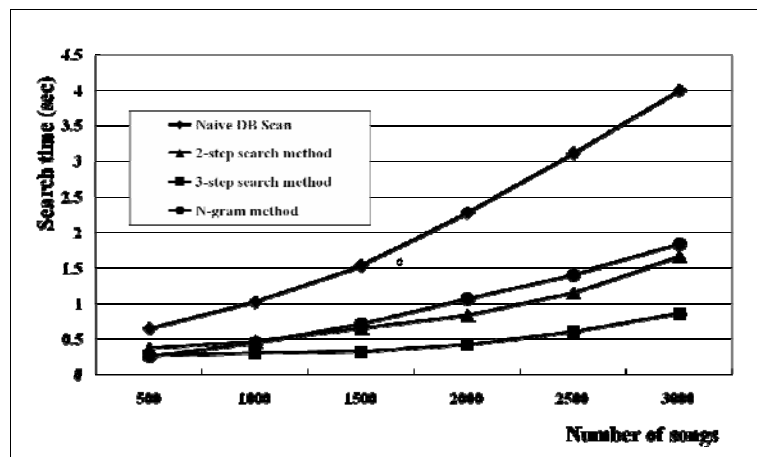


Figure 24. Search time of four methods with increasing number of songs contained in music database (BST=2, FMLT=6, query length=15 seconds)

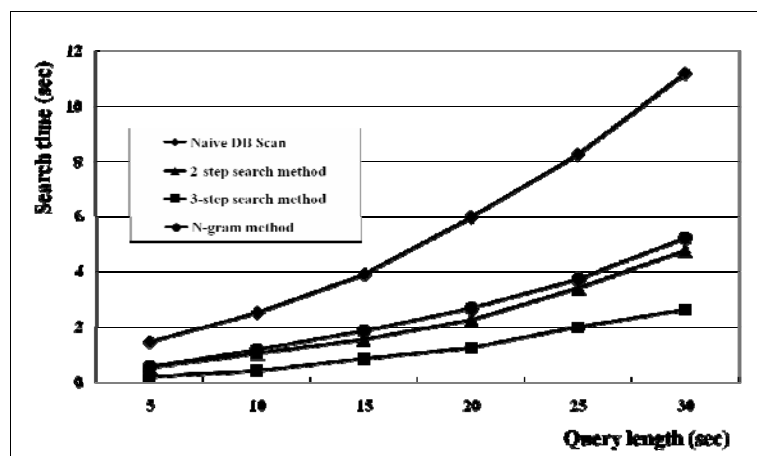


Figure 25. Search time of four methods with increasing length of queries from 5 seconds to 30 seconds (BST=2, FMLT=6, 3000 songs)

Figure 24 compares the search time of the four methods with increasing number of songs, from 500 to 3000 by increments of 500, contained in the music database. It is clear that the search time of all of the methods increases as the size of database increases. However, the **3-step search method** consistently takes the least amount of search time irrespective of the database size. The reason for this result is that a large portion of the queries are processed by just searching the frequent melody suffix tree index or the rest-unit melody suffix tree index. Since the **N-gram method** has to sequentially examine the strings stored in the music database, it took about twice as long as the **3-step search method** on average.

Figure 25 compares the search time of the four methods with increasing length of queries from 5 seconds to 30 seconds. Just like the result of the previous experiment, the search time of all of the methods increases as the length of the queries increases and the **3-step search method** takes the least amount of search time with any query length.

## Jinhee You and Sanghyun Park

Next, we compared the recall and precision of the four methods while increasing the distance between the queries and their corresponding target music. Figure 26 and Figure 27 show the result. As clearly shown in these two figures, the **3-step search method** produces a more accurate result compared to the **N-gram method** with any distance between the queries and their corresponding target music.

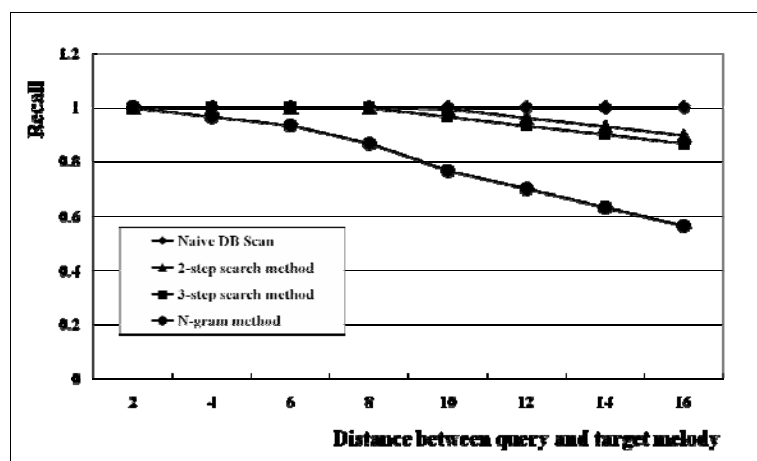


Figure 26. Recall of four methods with increasing distance between queries and their corresponding target music (BST=2, FMLT=6, 3000 songs, query length=15 seconds)

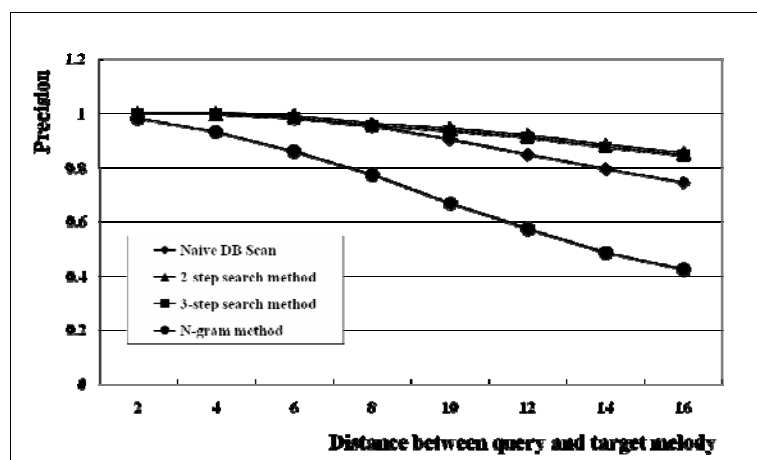


Figure 27. Precision of four methods with increasing distance between queries and their corresponding target music (BST=2, FMLT=6, 3000 songs, query length=15 seconds)

## 6. Conclusion

Query-By-Humming (QBH) is one of the content-based method studies in the area of music information retrieval. Since QBH systems are capable of retrieving the desired music even when the users' humming is inaccurate, they are useful for people who are not talented singers or who have a speech impediment.

## Jinhee You and Sanghyun Park

In this paper, we proposed an efficient Query-By-Humming (QBH) system. Considering the characteristics of humming, the feature data included the pitch interval and the duration ratio between adjacent notes. Since the operation for calculating the distance of melodies of numeric type is expensive, we convert the melodies of numeric type into the corresponding strings with the consideration of errors involved in humming. The strings were then compared using existing string matching algorithms. For an additional improvement in search speed, significant melodies were extracted, which have a high chance of being queried, from music and then build a couple of indexes from them. The definition of significant melodies is divided into two types in this paper: 1) the most frequent melody that appears in a single piece of music, and 2) the rest-unit melody that begins after a long rest within the music. In addition, it was suggested that a three-step index search method be used for minimizing database access and obtaining query answers as early as possible. Through experiments with real-world data, it was verified that users are more likely to look for the desired music by humming some parts of frequent melodies or rest-unit melodies and the proposed QBH system is much faster and more accurate than the N-gram method.

As a further study, it is planned to devise a systematic tool to find an optimal standard for character encoding and optimal values of two parameters including the bar selection threshold (BST) and the frequent melody length threshold (FMLT). Compression of the suffix trees employed in the proposed QBH system is another direction of further research.

## 7. Acknowledgements

This research is accomplished as the result of the promotion project for culture contents technology research center supported by Korea Culture & Content Agency (KOCCA).

## 8. References

- [1] A. Uitdenbogard, J. Zobel, Melodic Matching techniques for Large Music Databases, *Proceedings of the 7th ACM International Conference on Multimedia '99* (ACM, Orlando, 1999) 57-66.
- [2] P. Rolland, G. Raskinis, and J. Ganascia, Musical Content-Based Retrieval: An Overview of the Melodiscov Approach and System, *Proceedings of the 7th ACM International Conference on Multimedia '99* (ACM, Orlando, 1999) 81-84.
- [3] A. Lippincott, Issues in content-based music information retrieval, *Journal of information science* (Bowker-Saur, London, 2002) 137-142.
- [4] A. Ghias, J. Logan, and D. Chamberlin, Query By Humming, *Proceedings of the 3rd ACM International Conference on Multimedia '95* (ACM, San Francisco, 1995) 231-236.
- [5] S. Pauws, CubyHum: A Fully Operational Query by Humming System, *Proceedings of the 3rd International Conference on Music Information Retrieval* (2002) 187-196.
- [6] S. Doraisamy, S. Ruger, Robust Polyphonic Music Retrieval with N-grams, *Journal of Intelligent Information System* (2002) 53-70.
- [7] R. L. Kline, E. P. Glinert, Approximate matching algorithms for music information retrieval using vocal input, *Proceedings of the 11th ACM International Conference on Multimedia 2003* (ACM, Berkeley, 2003) 130-139.
- [8] R. B. Dannenberg, W. P. Birmingham, G. Tzanetakis, C. Meek, N. Hu, and B. Pardo, The Musart Testbed for Query-By-Humming Evaluation, *Proceedings of the 4th International Conference on Music Information Retrieval* (2003) 41-50.

---

**Jinhee You and Sanghyun Park**

- [9] N. Kosugi, Y. Nishihara, T. Sakata, M. Yamamuro, and K. Kushima, A Practical Query-By-Humming System for a Large Music Database, *Proceedings of the 8th ACM International Conference on Multimedia* (ACM, Los Angeles, 2000) 333-342.
- [10] E. Scheirer, Tempo and Beat Analysis of Acoustic Musical Signals, *Journal of the Acoustic Society of America* (1998) 588-601.
- [11] A. Pienimäki, Indexing Music Databases Using Automatic Extraction of Frequent Phrases, *Proceedings of the 3rd International Conference on Music Information Retrieval* (2002) 25-30.
- [12] J. Downie, Evaluating a Simple Approach to Music Information Retrieval: Conceiving Melodic N-Grams as Text, *Thesis, University of Western Ontario* (1999).
- [13] A. Guttman, R-trees: A Dynamic Index Structure for Spatial Searching, *Proceedings of the ACM SIGMOD* (1984) 47-57.
- [14] K. Andreas, Themefinder: A Web-based Melodic Search Tool, Melodic Similarity: Concepts, Procedures and Applications, *Computing in Musicology* (1998) 231-236.
- [15] C. Liu, J. Hsu and A. L. P. Chen, Efficient Theme and Non-trivial Repeating Pattern Discovering in Music Databases, *Proceedings of the 15th International Conference on Data Engineering '99* (1999) 14-21.
- [16] S. Rho, S. Lee and E. Hwang, An Efficient Audio Retrieval Scheme in Large Audio Databases, *Journal of SIGDB* (2003) 46-57.
- [17] D. E. Knuth, J. H. Morris and V. B. Pratt, Fast Pattern Matching In Strings, *SIAM J. Comput.*, (1977) 6:323-350.
- [18] G. A. Stephen, String Searching Algorithms, *World Scientific Publishing Co.*, (1994).
- [19] E. H. Fredkin, Trie memory, *Communications of the ACM*, (1960) 490-500.