

SEARCHING LYRICAL PHRASES IN A-CAPELLA TURKISH MAKAM RECORDINGS

First author

Affiliation1

author1@ismir.edu

Second author

Retain these fake authors in

submission to preserve the formatting

Third author

Affiliation3

author3@ismir.edu

ABSTRACT

Search by lyrics, the problem of locating the exact occurrences of a phrase from lyrics in musical audio, is a recently emerging research topic. Unlike key-phrases in speech, lyrical key-phrases have durations that bear important relation to other musical aspects like the structure of a composition. In this work we propose an approach that address the differences of syllable durations, specific for singing. First a phrase is expanded to MFCC-based phoneme models, trained on speech. Then, we apply dynamic time warping between the phrase and audio to estimate candidate audio segments in the given audio recording. Next, the retrieved audio segments are ranked by means of a novel score-informed hierarchical hidden Markov model, in which durations of the syllables within a phrase are explicitly modeled. The proposed approach is evaluated on 12 a-capella audio recordings of Turkish makam music. Relying on standard speech phonetic models, we arrive at an f-measure, comparable to a state-of-the-art work for the related task of keyword-spotting in singing. To the best of our knowledge, this is the first work tackling the problem of search by lyrical key-phrases. We expect that it can serve as a baseline for further research on singing material with similar musical characteristics.

1. INTRODUCTION

Searching by lyrics is the problem of locating the exact occurrences of a key-phrase from textual lyrics in musical signal. It has inherent relation to the problem of keyword spotting in speech, which long-investigated research task. In keyword spotting a user is interested to find at which time position in speech a relevant keyword (presenting a topic of interest) is spoken.

Most of the work on searching for keyword/phrases in singing has lended concepts and ideas from its equivalent problem for speech. For spoken utterances phonemes have relatively similar duration across speakers. Unlike that, in singing durations of phoneme (especially vowels) have

higher variation [7]. When being sung, vowels are prolonged according to musical note values. Therefore adopting an approach from speech recognition might lack some singing-specific semantics, including among others durations of sung syllables. Furthermore, in comparison to keyword spotting, we believe that key-phrase detection has higher potential to be integrated with other relevant MIR-applications, because lyrical key-phrases bear semantics in the context of the musical idiom: a line form lyrics is correlated, for example, to musical structure. For most types of music a section-long lyrical phrase is a feature that represents the corresponding structural section (e.g. chorus) in a unique way. Therefore correctly retrieved audio segments for, for example, the first lyrics line for a chorus can serve as a structure discovery tool.

In this work we investigate searching by lyrics in the case when a query represents an entire section or phrase from the textual lyrics of a particular composition. Unlike keyword-spotting or query-by-humming where a hit would be a document from an entire collection, in our case a hit is the occurrence of a phrase, being retrieved only from the performances of the given composition. In this respect the problem setting is more similar to linking melodic patterns from score to musical audio [11], rather than to keyword-spotting. We assume that the musical score with lyrics is present for the composition of interest. The proposed approach has been tested on a small dataset of a-cappella performances from a repertoire of Turkish makam music. For a given performance, the composition is known in advance, but no information about the structure is given. Characteristic for makam music is that, in a performance there might be reordering or repetitions of score sections.

2. RELATED WORK

2.1 Keyword spotting in singing

A recent work proved that keyword spotting in singing voice is a hard problem even when singing material is acapella (from pop songs in English) [7]. The authors adopt an approach from spotting keywords for speech, using a compound hidden Markov model (HMM) with keyword and filler model. Keywords are automatically extracted from a collection of lyrics. The best classifier (multi-layer perceptron) yielded an f-measure of 44%, averaged over top 50% of keywords. Notably the achieved results on singing material are not very different from results on spoken utterances of same keywords and data.



© First author, Second author, Third author.

Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** First author, Second author, Third author. "Searching Lyrical Phrases in A-capella Turkish Makam Recordings", 16th International Society for Music Information Retrieval Conference, 2015.

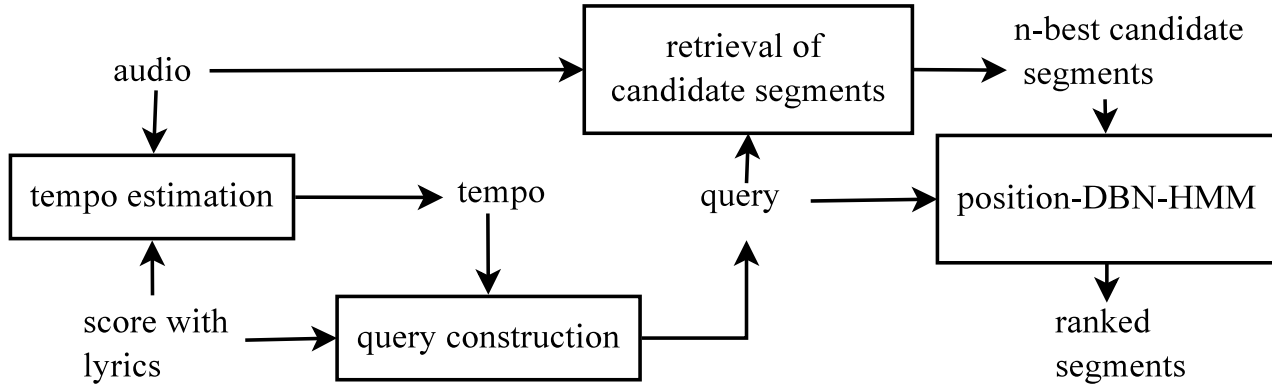


Figure 1. System overview: A key-phrase query is constructed in two variants: in the first stage candidate segments from audio are retrieved. In the second stage, the query is modeled by a DBN-HMM aware of the position in music score. The DBN-HMM decodes and ranks candidate segments

One of the few attempts to go beyond keywords is the work of [3]. Their goal is to automatically link phrases that appear in the lyrics of one song to the same phrase in another song. To this end, a keyword-filler model for detecting characteristic phrases (of 2-3 words) in sung audio. The method has been evaluated on polyphonic audio from Japanese pop. A non-HMM-based approach has been chosen in [2]. The authors propose subsequence dynamic time warping (DTW) to find a subsequence of features from a target recording similar to an example utterance of a keyword. [The approach has been evaluated with f-measure of 39%.](#)

In summary, performance of state-of-the-art work in keyword and key-phrase spotting for singing is not sufficiently good for practical applications. A probable reason for this is that hitherto approaches do not take into account the duration of syllables, which, as stated above, is an important factor that distinguishes speech from singing. In addition to that, syllable durations have been shown to be a strong reinforcing cue for the related task of automatically synchronizing lyrics and singing voice [1].

2.2 Position-aware DBN-HMMs

The modeling in most of the above mentioned approaches relies on classical HMMs with single hidden variable. A drawback of HMMs is that their capability to model exact state durations is restricted. The wait time in a state is implicitly set to a geometric distribution (derived from the self-transition likelihood).

One alternative to tackle durations can be seen in dynamic Bayesian networks (DBN), which allow modeling of any relations among phenomena in terms of probabilistic dependencies [10]. DBNs can be used to represent HMMs with more than one latent variables. Thus makes them attractive to model interdependent musical aspects. Such an aspect is for example the position in musical score. In [4] it has been shown that the dependence of score position on structural sections makes it possible to link musical performances to score. [These approaches have in common the strategy to have a variable that points to a musical aspect with sequential nature. In this paper, for brevity we](#)

[will refer to HMMs, which use DBNs to describe their hidden states as DBN-HMMs. Inspired by concepts from these works, we model the position of a musical section by a DBN-HMM, in a way that integrates syllable durations and the features, capturing phoneme acoustics.](#)

3. SYSTEM OVERVIEW

Figure 1 presents an overview of the proposed approach. We suggest a two-stage retrieval approach: On the first stage a subsequence DTW retrieves a set of candidate audio segments that are acoustically similar to a simple query with no duration information. On the second stage a the phonemes from the query and syllable durations are modeled by a novel DBN-HMM aware of the position in music score (in short position-DBN-HMM). We run a Viterbi decoding on each candidate segment separately. Viterbi assures that only one (the most optimal) path is detected for the candidate audio segment. Any query-to-audio fullpath match is considered as a hit and all hit results are ranked according to the weights derived from the Viterbi decoded path. In what follows each of the two stages is described in details, preceded by remarks on how a query keyphrase is handled.

3.1 Tempo Estimation

Often a performance is not played at the tempo indicated in the score. We follow the method of [12] to extract tempo by matching pitch values extracted from performance and pitch-values for the first section from the score. The angle of the most salient diagonal matching line segment in the distance matrix reveals the actual tempo T in for the recording.

3.2 Query construction

A selected lyrical phrase serves as a query twice: first a *simple query* for retrieval of candidate segments and then a *duration-informed query* for the decoding with position-DBN-HMM.

3.2.1 Acoustic features

For each of the 38 Turkish phonemes (and for a silent pause model) a 3-state HMM is trained from a corpus of Turkish speech. The 3 states represent respectively the beginning, middle and ending acoustic state of a phoneme. The acoustic properties (most importantly the formant frequencies) of spoken phonemes can be induced by the spectral envelope of speech. To this end, we utilize the first 12 MFCCs and their delta to the previous time instant. For each state a 9-mixture Gaussian distribution is fitted on the feature vector.

3.2.2 Simple query

For the first step no score-position information is utilized: lyrics is merely expanded to its constituent phoneme models. Let $\lambda_n \in \Lambda$ be a model state at position n in the query, where Λ is a set of all 3×38 states for the 38 Turkish phonemes. No transition probabilities are taken into account.

3.2.3 Duration-informed query

The idea is to associate each particular score position p_n in a musical section s to a phoneme state λ_n by exploiting the note-to-syllable mappings, present in sheet music.

Let ν be the duration of a short note from score than can be at most one-syllable-long. In this work we opted for 64-th note $\nu = \frac{1}{64}$. For each syllable a reference duration (in units of ν) is derived by aggregating values of its associated musical notes. Then the reference duration is spread among its constituent phonemes in a rule-based manner, resulting in reference durations R_p for all phonemes.¹

To query a particular performance, R_p are multiplied by τ : the number of time frames per ν according to some reference musical tempo T :

$$\tau = \frac{1}{h} \frac{T}{60} \nu$$

where T is in *bpm* and h is hop size of a frame in seconds ($h = 0.01$ in this work). Now let $D(s_n) = \sum_p R_p$ be the total duration for a section s_n . More formally, we define a mapping

$$f(p_n, s_n) \rightarrow \lambda_n \quad (1)$$

that determines the true state of a phoneme sung at time n from the position p_n ($p_n \in \{1, 2, \dots, D(s_n)\}$) within a section s_n .

4. RETRIEVAL OF CANDIDATE SEGMENTS

Subsequence-DTW proved to be effective when the feature series of an audio query are subsequence of features of a target audio [2]. In our case a query of phoneme models Λ with length M can be seen as subsequence of the series of MFCC features Y with length N , extracted from the whole recording. To this end we define a distance metric for an

audio frame y_m and model state λ_n as a function of the posterior probability.

$$d(m, n) = -\log P(y_m | \lambda_n) \quad (2)$$

where for phoneme state model λ_n

$$P(y_m | \lambda_n) = \sum_{c=1}^9 w_{c, \lambda_n} \cdot N(y_m; \mu_{c, \lambda_n}, \Sigma_{c, \lambda_n}) \quad (3)$$

with N being the Gaussian distribution from a 9-component mixture with weights w_{c, λ_n} . Based on the distance metric 2 a distance matrix $D^{N \times M}$ is constructed. Let a warping path Ω be a sequence of L points $(\omega_1, \dots, \omega_L)$, $l \in [1, L]$ and $\omega_l = (m, n)$ refers to an entry $d(m, n)$ in D .

4.1 Path computation

To generate a warping path Ω we select step sizes $\omega_l - \omega_{l-1} \in \{(1, 1), (1, 0), (1, 2)\}$ corresponding respectively to diagonal, horizontal and skip step according to the notation of [9]. A horizontal step means staying in the same phoneme in next audio frame. The step size $(0, 1)$ is disallowed because each frame has to map to exactly one phonetic model. To counteract the preference for the diagonal and the skip step $(1, 2)$ local weights w_d and w_s are introduced (as suggested in [9]). We set rather high values (empirically found $w_d = 6.5$ and $w_s = 11$) to assure that a path will not exit from long vocals prematurely.

A list of candidate paths $(\Omega_1^*, \dots, \Omega_r^*)$ is computed by iteratively detecting the current path with maximum score. After having detected a path Ω^* with final position in frame n^* a small region of 5% of M : $(n^* - 5\%M, n^* + 5\%M)$ is blacklisted from further iterations. This is an idea inspired by the machine learning principle of simulated annealing and assures that the iterative procedure will not get stuck in a current maximum and which might lead to detecting many paths from its vicinity. We pick $r = 30$ (a relatively high number compared to 20 in [Dittmar]), because our goal is to cover as many relevant candidates as possible.

4.2 Candidate segment selection

We define an overlapped audio segment to be any segment from the target audio, for which any frame y_m belongs to more than one path Ω . Let us call the set of paths for an overlapped audio segment an overlapped path set (see figure 2). Then as candidate segment is considered the audio segment span from the initial timestamp of the leftmost path until the final timestamp of the rightmost path. We assume that a concentration of found paths within a segment indicates at least one true hit.

5. POSITION-DBN-HMM

In this section, we present a novel probabilistic model for modeling a lyrical phrase. The main idea of the model is to incorporate the phonetic identities of lyrics and the syllable durations, available from musical score, into a coherent unit. The dependence of the observed acoustic features (that capture the phonetic identity) on musical velocity and score position are presented as DBN in figure 3.

¹ In this work a simple rule is applied: consonants are assigned a fixed duration and the rest is assigned to the vowel.

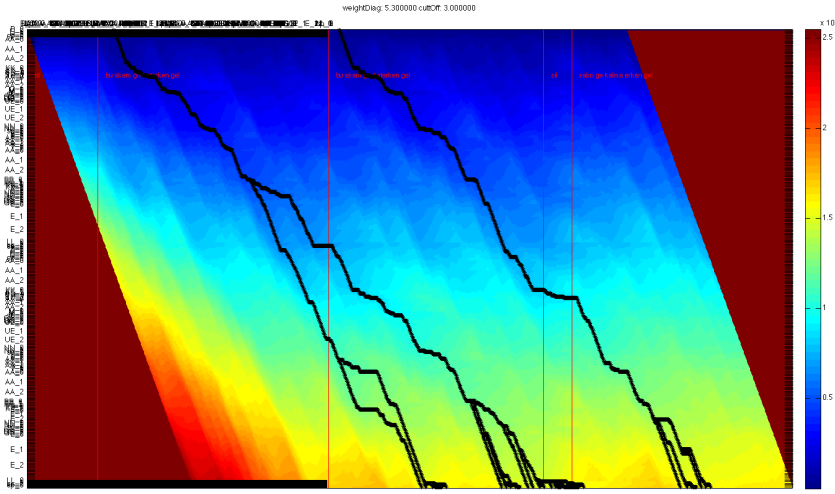


Figure 2. Example candidate audio segment of many overlapped paths (two main paths can be seen)

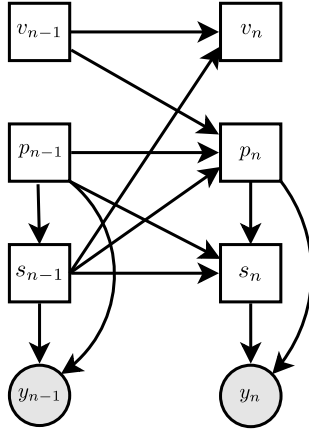


Figure 3. Representation of the hidden layers of the proposed model as a dynamic Bayesian network. Hidden variables (not shaded) are v - velocity, p - score position and s - section. The observed feature vector y is not shaded. Squares and circles denote respectively continuous and discrete variables

The derivation of model dependencies and parameters is the main contribution of this work.

5.1 Hidden variables

1. Position p_n from musical score for a section ($p_n \in \{1, \dots, D(s_n)\}$). $D(s_n = Q)$ is the total duration for a section s_n as defined in section 3.2.3 Note that $D(s_n)$ for a given section is different for two performances with different tempo T .
2. Velocity $v_n \in \{1, 2, \dots, V\}$. The number of units, which the position variable has to jump over at next time frame. Staying in state $v_n = 2$, for example, means that the actual tempo is steady and around 2 times faster than the reference one.

3. Structural section $s_n \in \{Q, F\}$ where Q is the queried section and F is a filler section. A filler section represents any non-key-phrase audio regions, and practically allows with equal probability being in any phoneme state (see section 5.3)

For the experiments reported in this paper, we chose $V = 5$ and set $D(s_n = F) = V$. This assures that even in fastest tempo there is option of entering the filler section. In deriving $D(s_n)$ instead of T we use $2T$ to allow handling local tempo fluctuations to a twice-slower tempo when $V = 1$.

The proposed model is different than the model proposed in [6] in two aspects:

- $D(s_n = Q)$ is not fixed but depends on the section of interest and the detected tempo T of performance
- a section (a pattern in the original model) is not fixed for a recording, but can alternate between states.

Since all the hidden variables are discrete, one can reduce this model to a regular HMM by merging all variables into a single 'meta-variable' x_n :

$$x_n = [v_n, p_n, s_n] \quad (4)$$

Note that the state space becomes the Cartesian product of the individual variables.

5.2 Transition model

Due to the conditional independence relations presented in figure 3, the transition model reduces to

$$P(x_n|x_{n-1}) = \frac{P(v_n|v_{n-1}, s_{n-1}) \times P(p_n|v_{n-1}, p_{n-1}, s_{n-1}) \times P(s_n|p_{n-1}, s_{n-1}, p_n)}{P(v_n|v_{n-1}, s_{n-1}) \times P(p_n|v_{n-1}, p_{n-1}, s_{n-1}) \times P(s_n|p_{n-1}, s_{n-1}, p_n)} \quad (5)$$

5.2.1 Velocity transition

$$p(v_n|v_{n-1}) = \begin{cases} \phi/2, & v_n = v_{n-1} \pm 1 \\ 1 - \phi, & v_n = v_{n-1} \\ 0, & \text{else} \end{cases} \quad (6)$$

where ϕ is a constant probability of change in velocity.

5.2.2 Position transition

The score position is defined deterministically according to:

$$p_n = (p_{n-1} + v_{n-1} - 1) \mod D(s_{n-1}) + 1 \quad (7)$$

where the modulus operator resets the position to be in a beginning of a new section after it exceeds the duration of previous section $D(s_{n-1})$

5.2.3 Section transition

$$P(s_n|p_{n-1}, s_{n-1}, p_n) = \begin{cases} P(s_n|s_{n-1}), & p_n \leq p_{n-1} \\ 1, & p_n > p_{n-1} \ \& \ s_n = s_{n-1} \end{cases} \quad (8)$$

A lack of increase in the position is an indicator that a new section should be started. $P(s_n|s_{n-1})$ is set according to a transition matrix $A = \{a_{ij}\}$ where $i \in \{Q, F\}$ and self transitions a_{QQ} and a_{FF} for query and filler section respectively can be set to reflect the expected structure of the target audio signal. In this work we set $a_{QQ} = 0$, since we expect that a query might be decoded at most once in a candidate audio segment. The value $a_{FF} = 0.9$ is determined empirically.

5.3 Observation model

The probability of the observed feature vector in position p_n from section s_n is computed for the model state $\lambda_n = f(p_n, s_n)$, where the mapping function 1 was introduced in section 3.2. A similar mapping function has been proposed for the first time in the DBN-HMM in [4].

Then

$$P(y_n|p_n, s_n = Q) = P(y_n|\lambda_n)$$

which reduces to applying the distribution defined in 3.

In case of the filler section we allow with equal probability any phoneme state.

$$P(y_n|p_n, s_n = F) = \max_{\lambda \in \Lambda} P(y_n|\lambda)$$

Note that position p_n plays a role only in tracking the total section duration $D(s_n = F)$.

5.4 Inference

An exact inference of the 'meta-variable' x can be performed by means of the Viterbi algorithm. A key-phrase is detected whenever the Viterbi path $\bar{\Omega}$ passes through a section $s_n = Q$ and the last position $D(s_n = Q)$ is reached.

statistic	value
#section queries	32
average cardinality \bar{C}_q	3.2
#words per section	5-14
#sections per recording	6-16
#phonemes per section	26-63

Table 1. Statistics about queries (lyrics sections with unique lyrics) in the test dataset

5.5 Paths ranking

Then the detected paths $\{\bar{\Omega}\}$ for all performances of the composition of interest, are ranked according to a weight: a weight measures the acoustic matching between the segment and the query. A weight for a path with length L is computed from the average likelihood:

$$w = \frac{1}{L} \sum_{l=1}^L P(\omega_l) \quad (9)$$

where $P(\omega_l) = P(y_m|\lambda_n)$ is the observation probability of a point $\omega_l = (m, n)$ from the Viterbi path.

6. DATASET

The **anonymous** test dataset consists of 12 a-cappella performances of 9 compositions with total duration of 18:40 minutes, drawn from classical Turkish Makam repertoire. It has been sung by semiprofessional singers and recorded especially for this study. Scores are provided in the machine-readable *sympTr* format [5], which contain marks of section divisions. A performance has been recorded in-sync with the original recording, whereby instrumental sections are left as silence. This assures that the sequence in which sections are performed is kept the same as in the original performance.

We consider as a query q each section from the scores, which has unique lyrics. In a given recording the boundary timestamps of an occurrence of each query have been annotated. Let C_q be the total number of relevant occurrences (cardinality) of a query q . Table 1 presents the average cardinality \bar{C}_q and other relevant statistics about sections. The low values of \bar{C}_q is mainly due to the small number of performances per composition.

7. EVALUATION

7.1 Evaluation metrics

Having a ranked list of occurrences of each lyrical query, the search-by-lyrics can be interpreted as a ranked retrieval problem, in which the users are interested in checking only the top K relevant results [8]. This allows to reject irrelevant results by considering only top K results in the evaluation metric. This strategy is appropriate when a query has low average cardinality (in our dataset $\bar{C}_q = 3.2$). Therefore we report results in terms of mean average precision (MAP). Let the relevance of ranked results for a query q be

$[r_q(1), \dots, r_q(n_q)]$ where n_q is the number of retrieved occurrences. Note that a detected audio segment is either hit or not, making $r_q(k) \in \{0, 1\}$. The results are ranked by the weights defined in equation 9. Now the recall and precision of a query at rank index K are defined respectively as:

$$R_q(K) = \frac{1}{C_q} \sum_{k=1}^K r_q(k), \quad P_q(K) = \frac{1}{K} \sum_{k=1}^K r_q(k) \quad (10)$$

and the average precision as

$$\bar{P}_q = \frac{1}{C_q} \sum_{k=1}^{n_q} r_q(k) P_q(k)$$

For a each of the 32 score sections \bar{P}_q is computed, ranking together the hits from all recordings of the composition, from which the section is taken. MAP is the average over all 32 \bar{P}_q .

7.2 Experiments

To assess the benefit of the proposed modeling of positions, we conduct a comparison of the performance of the complete system and a baseline version without the position-DBN-HMM. For the baseline, as result set we consider the audio segments corresponding to the list of candidate paths ($\Omega_1^*, \dots, \Omega_r^*$) derived after subsequence-DTW (see section 4.1). Note that DTW-paths were sorted according to their total distance of distance matrix D , which is derived from observation probability. Thus, the strategies for computing weights for both DTW and Viterbi are based on acoustic similarity (and thus are consistent with each other). We report results at different values for K in table 2.

The results confirm the expectation that the performance of subsequence DTW alone is inferior. Retrieving relevant candidate paths in the set S seemed to be very dependent on the weights w_d and w_s for the diagonal and skip steps. We noted that adapting weights for a recording according to detected tempo T might be beneficial, but did not conduct related experiments in this work. The optimal values ($w_d = 6.5$ and $w_s = 11$) in fact guaranteed good coverage of relevant segments in the slowest estimated tempo T .

The MAP for the complete system has a peak when considering first 3 hits from the ranked result list. This is expected because $\bar{C}_q = 3.2$. In general MAPs after decoding of candidate segment with position-modeling, is substantially better than the baseline, which suggests that modeling syllable durations is beneficial. Another reason might be that the position-DBN-HMM models tempo and is thus insusceptible to the difference between the tempo indicated in the score and the real performance tempo.

7.3 Comparison to related work

Unfortunately, no direct comparison to previous work on keyword spotting [7] is possible because these rely on speech models for languages different than ours.

Furthermore, the evaluation setting in none of the work is comparable to ours. In [7] a result is considered true

K	1	2	3	4	6	8
<i>DTW</i>	28.3	23.7	38.3	36.2	37	39
<i>DBN-HMM</i>	48.5	59.3	64.6	61	41.3	42.4

Table 2. MAPs (in percent) for ranked result segments for two system variants: baseline with subsequence-DTW and complete with position-DBN-HMM. Results for $K \geq 8$ are omitted because the maximum cardinality for a query is 8.

K	1	2	3	4	6	8
<i>DBN-HMM</i>	32.4	34.5	49.6	63.2	58.2	59.3

Table 3. F-measure (in percent) for the position-DBN-HMM for ranked results segments

positive if a keyword is detected at any position in an expected audio clip. The authors argue that since a clip spans one line of lyrics (only 1 to 10 words) this is sufficiently exact, whereas we are interested in detecting the exact timestamps of a key-phrase. In addition to that, their longest query has 8 phonemes, which is much less than the average in our setting.

For the only publication (to our knowledge) on longer key-phrases [3], the accuracy of the key-phrase spotting module is not reported, but instead only the percentage of the links correctly connecting key-phrases from a song to another song. On creating a link for a given key-phrase only the candidate section with highest score for a song has been considered, which might ignore any other true positives.

For the sake of comparison to any future work, we report results in terms of f-measure in table 3. These results makes us suggest that, although the evaluation setting of [7] is different, the length of the keyphrases is crucial for higher retrieval rates, because in all cases when $K \geq \bar{C}_q$, our system surpasses their f-measure of 44%.

8. CONCLUSION

In this study we have investigated an important problem that has started to attract attention of researchers only recently. We tackle the linking between audio and structural sections from the perspective of lyrics: a proposed a method for searching in musical audio the occurrences of a characteristic section-long lyrical phrase. We presented a novel DBN-based HMM for tracking in a score-informed way sung phonemes and their durations. Evaluation on a-cappella material from Turkish makam music shows that the search with the proposed model brings substantial improvement compared to a baseline system, unaware of syllable-duration information.

We plan to focus future work on applying the proposed model to the case of polyphonic singing. We expect further, that this work can serve as a baseline for further research on singing material with similar musical characteristics.

We want to point as well that, the proposed score-informed scheme is applicable not necessarily only when musical

scores are available. Scores can be replaced by any format, from which duration information can be inferred: for example annotated melodic contour or singer-created indications along the lyrics.

9. REFERENCES

- [1] On the use of lyrical duration from musical score for automatic lyrics-to-audio alignment.
- [2] Christian Dittmar, Pedro Mercado, Holger Grossmann, and Estefania Cano. Towards lyrics spotting in the syncglobal project. In *Cognitive Information Processing (CIP), 2012 3rd International Workshop on*, pages 1–6. IEEE, 2012.
- [3] Hiromasa Fujihara, Masataka Goto, and Jun Ogata. Hyperlinking lyrics: A method for creating hyperlinks between phrases in song lyrics. In *ISMIR*, pages 281–286, 2008.
- [4] Andre Holzapfel, Umut Simsekli, Sertan Sentürk, and Ali Taylan Cemgil. Section-level modeling of musical audio for linking performances to scores in turkish makam music.
- [5] M Kemal Karaosmanoğlu. A turkish makam music symbolic database for music information retrieval: Symbtr. *Proc. Int. Society for Music Information Retrieval (ISMIR)*, 2012.
- [6] Florian Krebs, Sebastian Böck, and Gerhard Widmer. Rhythmic pattern modeling for beat and downbeat tracking in musical audio. In *ISMIR*, pages 227–232, 2013.
- [7] Anna M Kruspe and IDMT Fraunhofer. Keyword spotting in a-capella singing.
- [8] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008.
- [9] Meinard Müller. *Information retrieval for music and motion*, volume 2. Springer, 2007.
- [10] Kevin Patrick Murphy. *Dynamic bayesian networks: representation, inference and learning*. PhD thesis, University of California, 2002.
- [11] Senturk. Composition identification.
- [12] Sertan Şentürk, Sankalp Gulati, and Xavier Serra. Score informed tonic identification for makam music of turkey. In *Proceedings of 14th International Society for Music Information Retrieval Conference (ISMIR)*, pages 175–180, Curitiba, Brazil, 04/11/2013 2013.