

# ON THE USE OF NOTE ONSETS FOR IMPROVED LYRICS-TO-AUDIO ALIGNMENT IN TURKISH MAKAM MUSIC

Georgi Dzhabazov Ajay Srinivasamurthy  
Sertan Şentürk Xavier Serra

Music Technology Group, Universitat Pompeu Fabra, Barcelona, Spain

georgi.dzhabazov@upf.edu

## ABSTRACT

Lyrics-to-audio alignment aims to automatically match given lyrics and musical audio. In this work we extend a state of the art approach for lyrics-to-audio alignment with information about note onsets. In particular, we consider the fact that transition to next lyrics syllable usually implies transition to a new musical note. To this end we formulate rules that guide the transition between consecutive phonemes when a note onset is present. These rules are incorporated into the transition matrix of a variable-time hidden Markov model (VTHMM) phonetic recognizer based on MFCCs. An estimated melodic contour is input to an automatic note transcription algorithm, from which the note onsets are derived. The proposed approach is evaluated on 12 a cappella audio recordings of Turkish Makam music using a phrase-level accuracy measure. Evaluation of the alignment is also presented on a polyphonic version of the dataset in order to assess how degradation in the extracted onsets affects performance. Results show that the proposed model outperforms a baseline approach unaware of onset transition rules. To the best of our knowledge, this is the one of the first approaches tackling lyrics tracking, which combines timbral features with a melodic feature in the alignment process itself.

## 1. INTRODUCTION

Lyrics are one of the most important aspects of vocal music. When a performance is heard, most listeners will follow the lyrics of the main vocal line. The goal of automatic lyrics-to-audio alignment is to generate a temporal relationship between lyrics and recorded singing. In this particular work, the goal is to detect the start and end times of every phrase (1-4 words) from lyrics.

In recent years there has been a substantial amount of work on the extraction of pitch of predominant singing voice from polyphonic music [18]. Some algorithms have been tailored to the music characteristics of a particular

singing tradition [12]. This has paved the way to an increased accuracy of note transcription algorithms. One of the reasons for this is that a correctly detected melody contour is a fundamental precondition for note transcription. On the other hand, lyrics-to-audio alignment is a challenging task: to track the timbral characteristics of singing voice might not be straightforward [4]. An additional challenge is posed when accompanying instruments are present: their spectral peaks might overlap and occlude the spectral components of voice. Despite that, most work has focused on tracking the transitions from one phoneme to another only by timbral features [4, 14]. In fact, at a phoneme transition, in parallel to timbral change, a change of pitch or an articulation accent may be present, which contributes to the perception of a distinct vocal note onset. For example, a note onset occurs simultaneously with the first vowel in a syllable. This fact has been exploited successfully to enhance the naturalness of synthesized singing voice [21].

In this work we present a novel idea of how to extend a standard approach for lyrics-to-audio alignment by using automatically detected vocal note onsets as a complementary cue. We apply a state of the art note transcription method to obtain candidate note onsets. The proposed approach has been evaluated on time boundaries of short lyrics phrases on a cappella recordings from Turkish Makam music. An experiment on polyphonic audio reveals the potential of the approach for real-world applications.

## 2. RELATED WORK

### 2.1 Lyrics-to-audio alignment

The problem of lyrics-to-audio alignment has an inherent relation to the problem of text-to-speech alignment. For this reason most of current studies exploit an approach adopted from speech: building a model for each phoneme based on acoustic features [5, 14]. To model phoneme timbre usually mel frequency cepstral coefficients (MFCCs) are employed. A state of the art work following this approach [5] proposes a technique to adapt a phonetic recognizer trained on speech: the MFCC-based speech phoneme models are adapted to the specific acoustics of singing voice by means of Maximum Likelihood Linear Regression. Further, automatic segregation of the vocal line is performed, in order to reduce the spectral content of back-



© Georgi Dzhabazov, Ajay Srinivasamurthy, Sertan Şentürk, Xavier Serra. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Georgi Dzhabazov, Ajay Srinivasamurthy, Sertan Şentürk, Xavier Serra. "On the use of note onsets for improved lyrics-to-audio alignment in Turkish Makam music", 17th International Society for Music Information Retrieval Conference, 2016.

ground instruments. In general, in this approach authors consider only models of phonetic timbre and are thus focused on making them more robust as a mean to improve performance.

Few works for tracking lyrics combine timbral features with other melodic characteristics. For example in [7] a system for automatic score-following of singing voice combines melodic and lyrics information: observation probabilities of pitch templates and vowel templates are fused to improve alignment. In [13] lyrics-to-audio alignment has been aided on a coarser level by chord-to-audio alignment, assuming chord annotations are available in a paired chord-lyrics format. However, to our knowledge, no work so far has employed note onsets as additional cue to alignment.

### 2.2 Automatic note segmentation

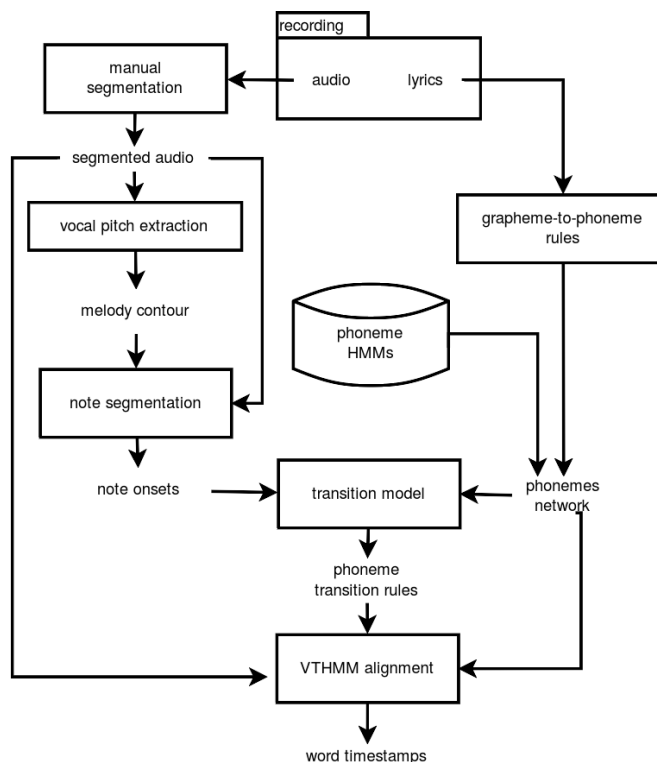
While the general problem of automatic music transcription has been long-investigated, automatic singing transcription has attracted the attention of MIR researchers only in recent years [6, 11, 15]. A fundamental part of singing transcription is automatic note segmentation. A probabilistic note event model, using a HMM trained on manual transcriptions is presented in [11]. The idea is that a note consists of different states representing its attack, sustain and decay phase. Then an onset is detected when the decoding path goes through an attack state of a new note.

A recent work on singing transcription with high onset accuracy has been developed for singing voice from the flamenco genre [12]. It consists of two stages: predominant vocal extraction and note transcription. As a primary step of the note transcription stage, notes are segmented by a set of onset detection functions based on pitch contour and volume characteristics, which take into account the peculiar for flamenco singing high degree of microtonal ornamentation.

## 3. PROPOSED APPROACH

A general overview of the proposed approach is presented in Figure 1. An audio recording and its lyrics are input. A variable time hidden Markov model (VTHMM), guided by phoneme transition rules, returns start and end timestamps of aligned words. For brevity in the rest of the paper our approach will be referred to as VTHMM.

First an audio recording is manually divided into segments corresponding to structural sections (e.g. verse, chorus) as indicated in a structural annotation, whereby instrumental-only sections are discarded. All further steps are performed on each audio segment. If we had used automatic segmentation instead, potential erroneous lyrics and features could have biased the comparison of a baseline system and VTHMM. As we focus on evaluating the effect of VTHMM, manual segmentation is preferred. In what follows each of the modules is described in details.



**Figure 1.** Overview of the modules of the proposed approach. One can see how phoneme transition rules are derived. Then together with the phonemes network and the features extracted from audio segments are input to the VTHMM alignment.

### 3.1 Vocal pitch extraction

To extract the melody contour of singing voice, we utilize a method that performs detection of vocal segments and in the same time pitch extraction for the detected segments [1]. It relies on the basic methodology of [19], but modifies the way in which the final melody contour is selected from a set of candidate contours, in order to reflect the specificities of Turkish Makam music: 1) It chooses a finer bin resolution of only 7.5 cents that approximately corresponds to the smallest noticeable change in Makam melodic scales. 2) Unlike the original methodology, it does not discard time intervals where the peaks of the pitch contours have relatively low magnitude. This accommodates time intervals at the end of the melodic phrases, where Makam singers might sing softer.

### 3.2 Note segmentation

In a next step, to obtain reliable estimate of singing note onsets, we adapt the automatic singing transcription method, developed for polyphonic flamenco recordings [12]. It has been designed to handle singing with high degree of vocal pitch ornamentation. We expect that this makes it suitable for material from Makam classical singing having heavily vibrato and melismas, too. We replace the original first stage predominant vocal extraction

method with the vocal pitch detection method described above.

The algorithm [12] considers two cases of onsets: interval onsets and steady pitch onsets. A Gaussian derivative filter detects interval onsets as long-term changes of the pitch contour, whereas steady-pitch onsets are inferred from pitch discontinuities. As in the current work phoneme transitions are modified only when onsets are present, we opt for increasing recall at the cost of losing precision. This is achieved by reducing the value of the parameter  $cF$ : the minimum output of the Gaussian filter. The extracted note onsets are converted into a binary onset activation at each frame  $\Delta n_t = (0, 1)$ . Recall rates of extracted note onsets are reported in Table 2.

### 3.3 Phoneme models

The formant frequencies of spoken phonemes can be induced from the spectral envelope of speech. To this end, we utilize the first 12 MFCCs and their delta to the previous time instant, extracted as described in [24]. For each phoneme a one-state HMM, for which a 9-mixture Gaussian distribution is fitted on the feature vector. The lyrics are expanded to phonemes based on grapheme-to-phoneme rules for Turkish [16, Table 1] and the trained HMMs are concatenated into a phonemes network. The phoneme set utilized has been developed for Turkish and is described in [16]. A HMM for silent pause  $sp$  is added at the end of each word, which is optional on decoding. This way it will appear in the detected sequence only if there is some non-vocal part or the singer makes a break for breathing.

### 3.4 Transition model

We utilize a transition matrix with time-dependent self-transition probabilities which falls in the general category of variable time HMM (VTHMM) [9]. For particular states, transitions are modified depending on the presence of time-adjacent note onset. Let  $t'$  be the timestamp of the closest to given time  $t$  onset  $\Delta n_{t'} = 1$ . Now the transition probability can be rewritten as

$$a_{ij}(t) = \begin{cases} a_{ij} - g(t, t')q, & R1 \text{ or } R3 \\ a_{ij} + g(t, t')q, & R2 \text{ or } R4 \end{cases} \quad (1)$$

$R1$  to  $R4$  stand for phoneme transition rules, which are applied in the phonemes network by picking the states  $i$  and  $j$  for two consecutive phonemes. The term  $q$  is a constant whereas  $g(t, t')$  is a weighting factor sampled from a normal distribution with its peak (mean) at  $t'$ :

$$g(t, t') = \begin{cases} f(t; t', \sigma^2) \sim \mathcal{N}(t', \sigma^2), & |t - t'| \leq \sigma \\ 0 & \text{else} \end{cases} \quad (2)$$

Since singing voice onsets are regions in time, they span over multiple consecutive frames. To reflect that fact,  $g(t, t')$  serves to smooth in time the influence of the discrete detected  $\Delta n_t$ , where  $\sigma$  has been selected to be 0.075

seconds. In this way an onset influences a region of 0.15 seconds - a threshold suggested for vocal onset detection evaluation by [6] and used in [12]. Furthermore, this allows to handle slight timestamp inaccuracies of the estimated note onsets.

#### 3.4.1 Phoneme transition rules

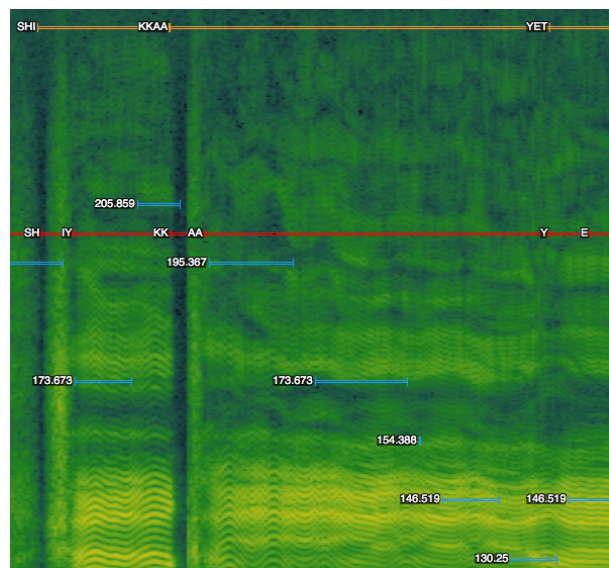
Let  $V$  denote a vowel,  $C$  denote a consonant and  $L$  denote a vowel, liquid (LL, M, NN) or the semivowel Y. Rules  $R1$  and  $R2$  represent inter-syllable transition, e.g. phoneme  $i$  is followed by phoneme  $j$  from the following syllable:

$$\begin{aligned} R1 : & \quad i = V \quad j = \neg L \\ R2 : & \quad i = C \quad j = L \end{aligned} \quad (3)$$

For example, for rule  $R2$  if a syllable ends in a consonant, a note onset imposes with high probability that a transition to the following syllable is done, provided that it starts with a vowel. Same rule applies if it starts with a liquid, according to the observation that pitch change takes place during a liquid preceding the vowel [21, timing of pitch change]. Rules  $R3$  and  $R4$  are for intra-syllabic phoneme patterns:

$$\begin{aligned} R3 : & \quad i = V \quad j = C \\ R4 : & \quad i = \neg L \quad j = V \end{aligned} \quad (4)$$

Essentially, if the current phoneme is vocal and the next is non-voiced (e.g.  $R1$ ,  $R3$ ), Eq. (1) discourages transition to next phoneme and encourages transition in the opposite cases. An example of  $R4$  can be seen for the syllable KK-AA in Figure 2 where the note onset triggers the change to the vowel AA, opposed, for example, to onset at Y for the syllable Y-E-T. Note that these rules assume



**Figure 2.** Ground truth annotation of syllables (in orange/top), phonemes (in red/middle) and notes (with blue/changing position). Audio excerpt corresponding to word şıkayet with syllables SH-IY, KK-AA and Y-E-T.

total #sections	#phrases per section	#words per phrase
75	2 to 5	1 to 4

**Table 1.** Phrase and section statistics about the dataset.

that a syllable has one vowel, which is the case for Turkish<sup>1</sup>. The optional silent phoneme *sp* is handled as a special case: transition probability from any phoneme to *sp* is derived according to intra-syllable rules, and the one from any phoneme skipping to the phoneme following *sp* is derived according to inter-syllable rules.

### 3.4.2 Alignment

The most likely state sequence is found by means of a forced alignment Viterbi decoding.

$$\delta_t(j) = \max_{i \in (j, j-1)} \delta_{t-1}(i) a_{ij}(t) b_j(O_t) \quad (5)$$

Here  $b_j(O_t)$  is the observation probability for state  $i$  for feature vector  $O_t$  and  $\delta_t(j)$  is the probability for the path with highest probability ending in state  $j$  at time  $t$  (complying with the notation of [17, III. B])<sup>2</sup>.

## 4. DATASET

The test dataset consists of 12 a cappella performances of 11 compositions with total duration of 19 minutes. The performances are drawn from *CompMusic* corpus of classical Turkish Makam repertoire with provided annotations of musical sections [23]. Solo vocal versions of the originals have been sung by professional singers, especially recorded for this study, due to the lack of appropriate a cappella material in this music tradition. A performance has been recorded in sync with the original recording, whereby instrumental sections are left as silence. This assures that the order, in which sections are performed, is kept the same. One of the contributions of this work is that we make available the annotated phrase boundaries<sup>3</sup>. A musical phrase spans 1 to 4 words depending on the duration of the words (as proposed in [10]). Table 1 presents statistics about phrases, while the total number of words in the dataset is 732.

Additionally, the singing voice for 6 recordings (with a total duration of 10 minutes) from the dataset has been annotated with MIDI notes complying to the musical score<sup>4</sup>. On annotation special care is taken to place the note onset on the time instant, at which the pitch becomes steady. Thus we avoid placing the onset on an unvoiced phoneme at the beginning of a syllable, which is assures rules R3 and R4 make sense (see Figure 2)<sup>5</sup>.

<sup>1</sup> Among one-vowel syllabic languages are also Japanese and to some extent Italian

<sup>2</sup> To encourage reproducibility of this research an efficient open-source implementation together with documentation is available at <https://github.com/georgid/AlignmentDuration/tree/noteOnsets>

<sup>3</sup> The audio and the annotations are available under a CC license at <http://compmusic.upf.edu/turkish-makam-acapella-sections-dataset>

<sup>4</sup> Creating the annotation is a time-consuming task, but we plan to annotate the whole dataset in the future

<sup>5</sup> Onset annotations are available at

## 4.1 Evaluation metric

Alignment is evaluated in terms of alignment accuracy as the percentage of duration of correctly aligned regions from total audio duration (see [5, Figure 9] for an example). A value of 100 means perfect matching of all phrase boundaries in the evaluated audio. Thus accuracy can be reported not only for an audio segment, but also on total for a recording, or as a total for all the recordings.

## 5. EXPERIMENTS

### 5.1 Experiment 1: alignment with oracle onsets

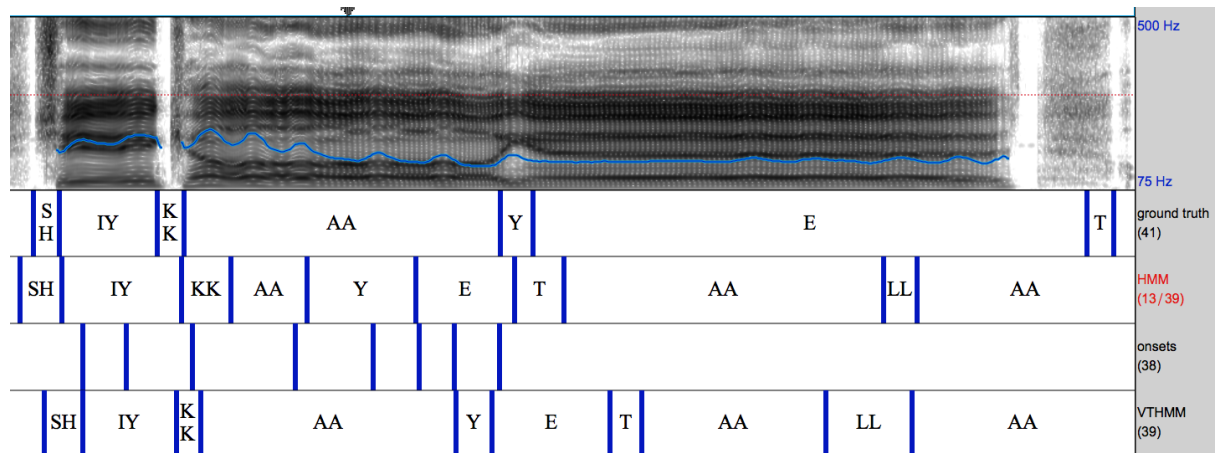
As a precursor to the following experiments, lyrics-to-audio alignment is run on these 6 recordings with manually annotated MIDI notes, which serve as an oracle for note onsets. This is done to test the general feasibility of the proposed model on the dataset, unbiased from errors in the note segmentation algorithm, and to set a glass-ceiling alignment accuracy. We have tested with different values of  $q$  from Eq. 1 achieving best accuracy of 83.5% at  $q = 0.23$ , which is used on all further reported experiments.

### 5.2 Experiment 2: recognition of phonemes

In general, the comparison to other lyrics alignment systems is not feasible, because there is no current work developed for Turkish language. However, to have an idea of how adequate the trained phoneme HMMs are, we have annotated phoneme boundaries for some excerpts of total length of 6 minutes. In [8] phonemes are recognized in a cappella singing with no lyrics given in advance. With phoneme MFCC-based HMMs - the same as our modeling setting - a phoneme recall rate of 44% is reported. Even though for forced alignment the recognition of phonemes is relatively easier, given that they are ordered in a sequence, we measured lower overall phoneme recall of 37%. This indicates that our phoneme models trained only on speech might not be the most optimal choice.

### 5.3 Experiment 3: comparison to a baseline

As a baseline we conduct alignment of the test dataset with unaffected phoneme transition probabilities, e.g. setting all  $\Delta n_t = 0$ , which resulted in alignment accuracy of 70.2%. Further, we measured the impact of the note segmentation module (introduced in Section 3.2), varying onset detection recall by changing the minimum output of the Gaussian filter (controlled by the parameter  $cF$ ). Table 2 summarizes the alignment accuracy with VTHMM depending on recall. On a cappella best improvement over the baseline is achieved at recall of 72.3% (at  $cF = 3.5$ ). This is somewhat lower than the best recall of 81-84% achieved for flamenco [12]. Setting recall higher than that degraded performance because there are too many false alarms, resulting in forcing false transitions.



**Figure 3.** Example of boundaries of phonemes for the word *şikayet* (SH-IY-KK-AA-Y-E-T): *on top*: spectrum and pitch; *then from top to bottom*: ground truth boundaries, phonemes detected with HMM, detected onsets, phonemes detected with VTHMM; (excerpt from the recording 'Kimseye etmem şikayet' by Bekir Unluater).

	$cF$	5	4.5	4.0	3.5	3.0
a cappella	<b>OR</b>	57.2	59.7	66.8	72.3	73.2
	<b>AA</b>	71.1	73.3	74.5	<b>75.7</b>	72.0
polyphonic	<b>OR</b>	52.8	58.2	65.9	66.2	68.4
	<b>AA</b>	61.2	63.3	<b>64.8</b>	64.6	60.3

**Table 2.** VTHMM performance on a cappella and polyphonic audio, depending on onset detection recall (OR). Alignment accuracy (AA) is reported as a total for all the recordings.

Figure 3 allows a glance at the level of detected phonemes: the baseline HMM switches to the following phoneme after some amount of time, similar for all phonemes. One reason for this might be that the waiting time in a state in HMMs with a fixed transition matrix cannot be randomly long [25]. In contrast, for VTHMM the presence of note onsets at vowels activates rules  $R1$  or  $R3$ , which allows waiting in the same state longer, as there are more onsets (for example AA from the word SH-IY-KK-AA-Y-E-T has five associated onsets). We chose to modify  $cF$  because setting it to lower values increases the recall of *interval onsets*: Often in our dataset several consecutive notes with different pitch correspond to the same vowel. In fact, it is characteristic of Turkish classical music that a single syllable may have a complex melodic progression spanning many notes (up to 12 in our dataset) [3]. However, for cases of vowels held long on same pitch, conceptually VTHMM is not capable of bringing any benefit. This is illustrated in Figure 3 by the prematurely detected end boundary of E from the word SH-IY-KK-AA-Y-E-T.

In addition to that, we examined alignment accuracy per recording (Figure 4). It can be observed that VTHMM performs consistently better than the baseline HMM (with some exceptions of where accuracy is close).

## 6. EXTENSION TO POLYPHONIC MATERIAL

To test the feasibility of the proposed approach on polyphonic material, the alignment is evaluated on the original versions of the recordings in the dataset. Typical for Turkish Makam is that vocal and accompanying instruments follow the same melodic contour in their corresponding registers, with slight melodic variations. However, the vocal line usually has melodic predominance. This special type of polyphonic musical interaction is termed heterophony [3]. In the dataset used in this study, a singer is accompanied by one to several string instruments.

We applied the vocal pitch extraction and note segmentation methods directly, since both are developed for singing voice in a setting that has heterophonic characteristics. However, instrumental spectral peaks deteriorate significantly the shape of the vocal spectrum. To attenuate the negative influence of instrumental spectrum, a vocal resynthesis step is necessary.

### 6.1 Vocal resynthesis

For the regions with predominant voice, the vocal content is resynthesized as separate vocal part. Resynthesis is conducted based on the harmonic model of [20]: Based on the extracted predominant pitch (see Section 3.1) and a set of peaks from the original spectrum, the harmonic partials of the predominant voice are selected and resynthesized. Then MFCCs are extracted from the resynthesized vocal part as if it were monophonic singing. This is a viable step, because the harmonic partials preserve well the overall spectral shape of the singing voice, including the formant frequencies, which encode the phoneme identities [22]<sup>6</sup>. More details and examples of the resynthesis can be found in previous work, which showed that the application of a harmonic model is suitable for aligning lyrics in Makam music [2]. A conceptually similar resynthesis

<sup>6</sup>The resynthesis allowed us to verify that vocals are intelligible despite some distortions from overlap with instrumental harmonic partials

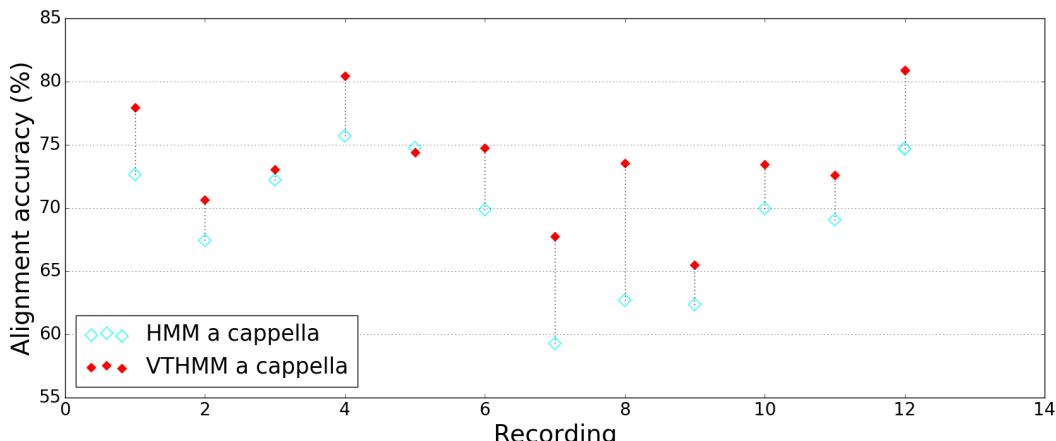


Figure 4. Comparison between results for VTHMM and baseline HMM on a cappella.

step is an established part also in current methods for alignment of lyrics in polyphonic Western pop music [5, 14].

6.2 Experiment 4: comparison of a cappella and polyphonic

The onset recall rates on polyphonic material after note segmentation are not much worse than a cappella as presented in Table 2. Even though the degree of degradation in onset detection is slight, degradation in alignment accuracy is significant. This can be attributed most probably to the fact that our MFCC-based models are not very discriminative and get confused by artifacts, induced from other instruments on resynthesis. However, applying VTHMM on polyphonic recordings still improves over the baseline (see Table 3). Note that the margin in accuracy between the baseline and the oracle glass ceiling is only about 6%, which is about twice much in the case of solo voice.

	HMM	VTHMM	oracle
a cappella	70.2	75.7	83.5
polyphonic	61.5	64.8	67.1

Table 3. Comparison of accuracy of baseline HMM, VTHMM and, VTHMM with oracle onsets. VTHMM shown are the best accuracies reported in Table 2. Alignment accuracy is reported on total for all recordings.

7. CONCLUSION

In this work we evaluated the behavior of a HMM-based phonetic recognizer for lyrics-to-audio alignment in two settings: with and without considering singing voice onsets as additional cue. Compared to existing work on lyrics alignment, this is, to our knowledge, the first attempt to include onsets of the vocal melody in the inference process. Updating transition probabilities according to onset-aware phoneme transition rules resulted in an improvement

of absolute 5.5 percent for aligning phrases of solo voice from Turkish Makam recordings. In particular, due to rules discouraging premature transition, the states of sustained vowels could have longer durations.

Alignment on same data with instrumental accompaniment brought also some small improvement over a baseline with no onset modeling. Having onset detection performing not substantially worse than a cappella indicates that improving the phoneme acoustic models in the future could probably lead to even more significant improvement.

A practical limitation of the current alignment system is the prerequisite for manual structural segmentation, which we plan to automate in the future.

**Acknowledgements** We thank Nadine Kroher for providing help with running the note segmentation module. This work is partly supported by the European Research Council under the European Union’s Seventh Framework Program, as part of the CompMusic project (ERC grant agreement 267583) and partly by the AGAUR research grant. We acknowledge as well financial support from the Spanish Ministry of Economy and Competitiveness, through the “María de Maeztu” Programme for Centres/Units of Excellence in R&D” (MDM-2015-0502)

8. REFERENCES

[1] Hasan Sercan Atlı, Burak Uyar, Sertan Şentürk, Barış Bozkurt, and Xavier Serra. Audio feature extraction for exploring Turkish makam music. In *3rd International Conference on Audio Technologies for Music and Media*, Ankara, Turkey, 2014. Bilkent University, Bilkent University.

[2] Georgi Dzhambazov, Sertan Şentürk, and Xavier Serra. Automatic lyrics-to-audio alignment in classical Turkish music. In *The 4th International Workshop on Folk Music Analysis*, pages 61–64, 2014.

- [3] Eric Bernard Ederer. *The Theory and Praxis of Makam in Classical Turkish Music 1910–2010*. University of California, Santa Barbara, 2011.
- [4] Hiromasa Fujihara and Masataka Goto. Lyrics-to-audio alignment and its application. *Dagstuhl Follow-Ups*, 3, 2012.
- [5] Hiromasa Fujihara, Masataka Goto, Jun Ogata, and Hiroshi G Okuno. Lyricsynchronizer: Automatic synchronization system between musical audio signals and lyrics. *IEEE Journal of Selected Topics in Signal Processing*, 5(6):1252–1261, 2011.
- [6] Emilia Gómez and Jordi Bonada. Towards computer-assisted flamenco transcription: An experimental comparison of automatic transcription algorithms as applied to a cappella singing. *Computer Music Journal*, 37(2):73–90, 2013.
- [7] Rong Gong, Philippe Cuvillier, Nicolas Obin, and Arshia Cont. Real-time audio-to-score alignment of singing voice based on melody and lyric information. In *Interspeech 2015*, Dresden, Germany, 06/09/2015 2015.
- [8] Jens Kofod Hansen. Recognition of phonemes in a-cappella recordings using temporal patterns and mel frequency cepstral coefficients. In *Proceedings of the 9th Sound and Music Computing Conference*, pages 494–499, Copenhagen, Denmark, 2012.
- [9] Michael T Johnson. Capacity and complexity of HMM duration modeling techniques. *Signal Processing Letters, IEEE*, 12(5):407–410, 2005.
- [10] M. Kemal Karasmanoğlu, Barış Bozkurt, Andre Holzapfel, and Nilgün Doğrusöz Dişiaçık. A symbolic dataset of Turkish makam music phrases. In *Fourth International Workshop on Folk Music Analysis (FMA2014)*, 2014.
- [11] Willie Krige, Theo Herbst, and Thomas Niesler. Explicit transition modelling for automatic singing transcription. *Journal of New Music Research*, 37(4):311–324, 2008.
- [12] Nadine Kroher and Emilia Gómez. Automatic transcription of flamenco singing from polyphonic music recordings. *IEEE/ACM Transactions on Audio, Speech and Language Processing*, 24(5):901–913, 2016.
- [13] Matthias Mauch, Hiromasa Fujihara, and Masataka Goto. Integrating additional chord information into hmm-based lyrics-to-audio alignment. *Audio, Speech, and Language Processing, IEEE Transactions on*, 20(1):200–210, 2012.
- [14] Annamaria Mesaros and Tuomas Virtanen. Automatic alignment of music audio and lyrics. In *Proceedings of the 11th Int. Conference on Digital Audio Effects (DAFx-08)*, 2008.
- [15] Emilio Molina, Lorenzo J Tardón, Ana M Barbancho, and Isabel Barbancho. Siph: Singing transcription based on hysteresis defined on the pitch-time curve. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(2):252–263, 2015.
- [16] Özgül Salor, Bryan L Pellom, Tolga Çiloğlu, and Mübeccel Demirekler. Turkish speech corpora and recognition tools developed by porting sonic: Towards multilingual speech recognition. *Computer Speech and Language*, 21(4):580 – 593, 2007.
- [17] Lawrence Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [18] Justin Salamon, Emilia Gómez, Dan Ellis, and Gaël Richard. Melody extraction from polyphonic music signals: Approaches, applications and challenges. *IEEE Signal Processing Magazine*, 31:118–134, 02/2014 2014.
- [19] Justin Salamon and Emilia Gómez. Melody extraction from polyphonic music signals using pitch contour characteristics. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(6):1759–1770, 2012.
- [20] Xavier Serra. A system for sound analysis/transformation/synthesis based on a deterministic plus stochastic decomposition. Technical report, 1989.
- [21] Johan Sundberg. The KTH synthesis of singing. *Advances in Cognitive Psychology*, 2(2-3):131–143, 2006.
- [22] Johan Sundberg and Thomas D Rossing. The science of singing voice. *the Journal of the Acoustical Society of America*, 87(1):462–463, 1990.
- [23] Burak Uyar, Hasan Sercan Atlı, Sertan Şentürk, Barış Bozkurt, and Xavier Serra. A corpus for computational research of Turkish makam music. In *1st International Digital Libraries for Musicology Workshop*, pages 57–63, London, United Kingdom, 2014.
- [24] Steve J Young. *The HTK hidden Markov model toolkit: Design and philosophy*. Citeseer, 1993.
- [25] Shun-Zheng Yu. Hidden semi-Markov models. *Artificial Intelligence*, 174(2):215–243, 2010.