

2020

## Score Following with Hidden Tempo Using a Switching State-Space Model

Yucong Jiang

University of Richmond, [yjiang3@richmond.edu](mailto:yjiang3@richmond.edu)

Chris Raphael

Follow this and additional works at: <https://scholarship.richmond.edu/mathcs-faculty-publications>



Part of the [Computer Sciences Commons](#), [Data Science Commons](#), and the [Music Commons](#)

---

### Recommended Citation

Jiang, Yucong, and Christopher Raphael. "Score Following with Hidden Tempo Using a Switching State-Space Model." In *Proceedings of the 21st International Society for Music Information Retrieval (ISMIR) Conference*, 693–99, 2020.

This Poster is brought to you for free and open access by the Math and Computer Science at UR Scholarship Repository. It has been accepted for inclusion in Math and Computer Science Faculty Publications by an authorized administrator of UR Scholarship Repository. For more information, please contact [scholarshiprepository@richmond.edu](mailto:scholarshiprepository@richmond.edu).

# SCORE FOLLOWING WITH HIDDEN TEMPO USING A SWITCHING STATE-SPACE MODEL

Yucong Jiang (yjjiang3@richmond.edu)  
University of Richmond

Chris Raphael (craphael@indiana.edu)  
Indiana University Bloomington

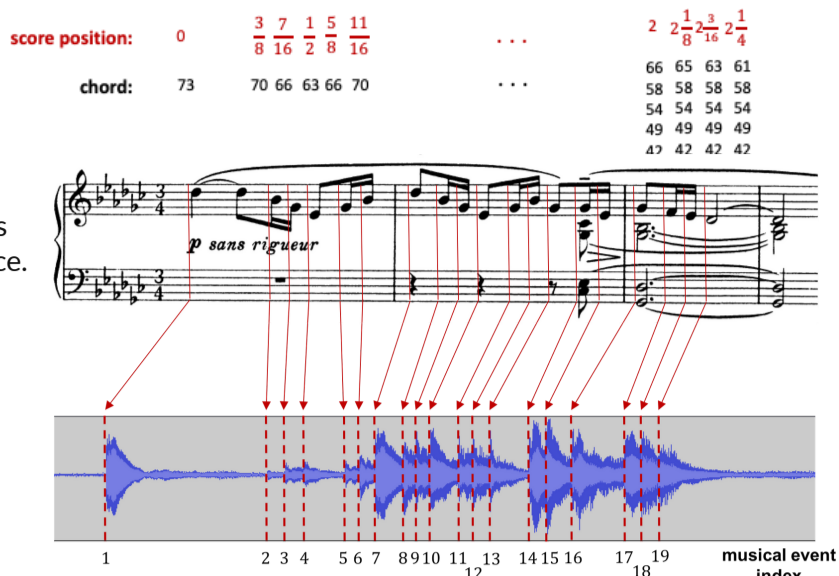
## 1 Introduction

### What is Score Following?

The score-following problem involves building a computer program that can trace musical events in a given musical score during a live performance.

### Why Score Following?

- Page turner
- Automatic accompaniment systems
- Virtual score composition
- Real-time audio enhancement/feedback



Monophonic music:



Polyphonic music:



## 2 Current Bottleneck

Existing score-following algorithms can still stumble on some challenging cases, especially when the data model is not reliable:

- Shared notes among neighboring chords
- Blurring effects caused by fast playing
- Pedaling

## 3 Research Aims

- To present a new method designed to improve the timing model — this aspect is especially meaningful in those challenging cases.
- To understand the nature of this problem better through empirical experiments.

- In the first diagram, the "time step" is the chord index.
- In the second diagram, the "time step" is the audio frame index.

## 4 The Model

### 1. Kalman Filter Model for Tempo

a linear dynamical system:

$$o_{k+1} = o_k + l_k t_k + \varepsilon_{k+1}$$

smooth tempo  $\rightarrow$   $t_{k+1} = t_k + \eta_{k+1}$

$$o_1 \sim N(\mu_{o,1}, \sigma_{o,1}^2)$$

$$t_1 \sim N(\mu_{t,1}, \sigma_{t,1}^2)$$

$$\varepsilon_k \sim N(0, \sigma_{\varepsilon,k}^2), k = 2, \dots, K$$

$$\eta_k \sim N(0, \sigma_{\eta,k}^2), k = 2, \dots, K$$

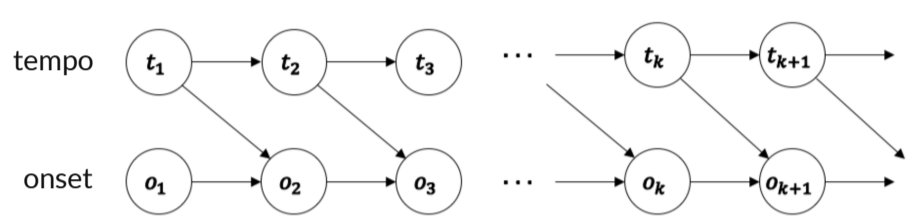
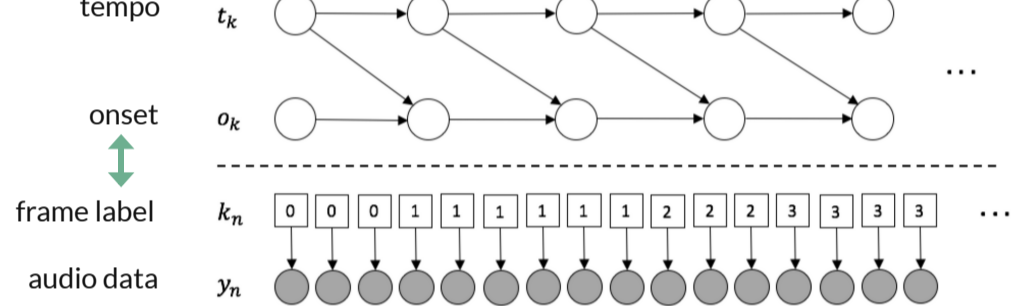
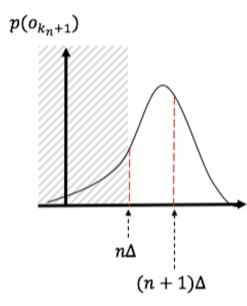


Figure 4.1: Linear dynamical system of the tempo and the onset.

### 2. Frame-wise Representation



$$k_n = \min \{k \in \{0, \dots, K\} : n\Delta < o_{k+1}\}$$

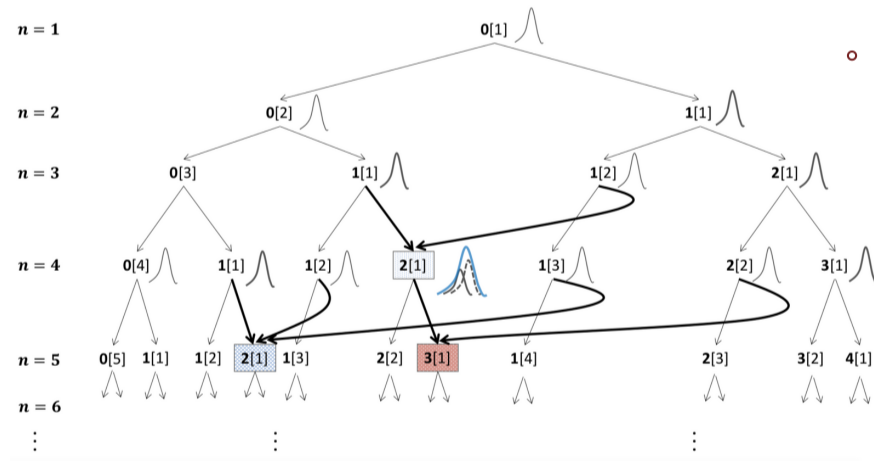
$$o_k \approx \Delta \cdot \min \{n \in \{1, \dots, N\} : k_n = k\}$$

## 5 Computation

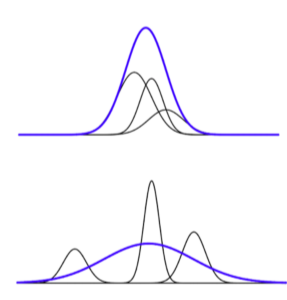
- [Chord Age]:** the number of frames a chord has lasted so far.
- Left Branch:** staying in the same chord.
- Right Branch:** moving on to the next chord.

The tree grows exponentially with time (or with n)!

### Approximation



- At each frame, merge nodes with the same label and age.
- Use a single Gaussian to approximate a mixture of Gaussians:



## 6 Experiments

- 15 solo piano pieces
- 50 excerpts
- Typical length: 40~90 seconds
- 48 minutes in total
- Sampling rate: 8k Hz
- Hop size: 16 ms
- Baseline: Music Plus One** (hidden Markov model)

### Evaluation Method

Frame-wise accuracy:

$$Acc_n = \sum_{\substack{k_n = \kappa_n \\ 1 \leq a_n \leq n - k_n + 1}} p(k_n, a_n | y_1^n)$$

correct hypotheses

$$Acc = \sum_n Acc_n / N$$

### Results

12 excerpts failed (accuracy < 40%) by either program:

- fatal error
- high uncertainty

|                   | baseline tempo tracking |       |
|-------------------|-------------------------|-------|
| # failed excerpts | 11                      | 9     |
| average accuracy  | 15.1%                   | 22.1% |

The proposed method is measurably better than the baseline:

| coefficient | value    |
|-------------|----------|
| $\bar{d}$   | -0.04    |
| $t$         | -2.7307  |
| df          | 37       |
| p-value     | 0.009623 |

|                  | baseline tempo tracking |       |
|------------------|-------------------------|-------|
| average accuracy | 65.0%                   | 69.1% |

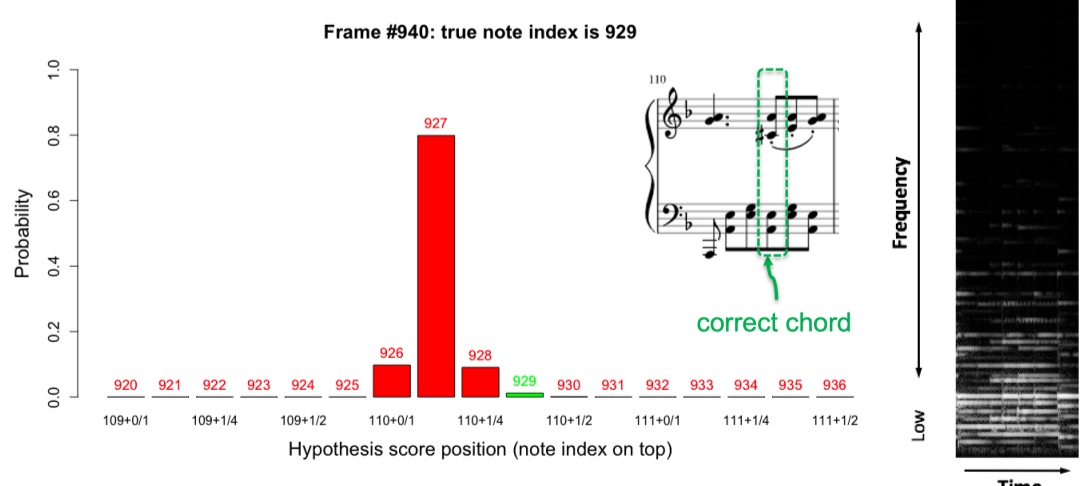
Table 5.10: Average accuracies of 38 excerpts.

### Case-by-case Analysis

An example of a typical scenario where the program was "confused" among neighboring chords, but was right about the general region:

Proposed method:

- is robust even with incorrect default tempo;
- recovers more easily from mistakes;
- however, can suffer at fast playing places.



Use a "movie" version of the barplot, which can show the hypothesis distribution one-by-one quickly through all frames of an excerpt, to inspect the nature of errors in these excerpts:

| Excerpt | Piece          | Baseline | Tempo | Observation   |
|---------|----------------|----------|-------|---|
| #23     | Liszt          | 25.7%    | 51.2% | Baseline got lost half-way through after a section of 14 repeated chords in a row.  |
| #43     | Beethoven_31   | 31.8%    | 46.5% | Incorrect default tempo.  |
| #46     | Beethoven_31   | 34.0%    | 62.7% | 14.4% higher accuracy among the other six excerpts (successfully followed).   |
| #15     | Chopin_ballade | 0.4%     | 27.2% | Baseline: completely lost near the beginning when the sound was blurring.   |
| #16     | Chopin_ballade | 0.3%     | 36.2% | Tempo: followed the region.   |
| #13     | Chopin_ballade | 11.8%    | 27.5% | Baseline: got lost starting around 1/3 through when the sound started to blur. Tempo: sometimes "confused," but always recovered. |
| #19     | Liszt          | 2.1%     | 18.6% | Both programs got lost near the beginning (repeated chords and patterns), but only the proposed method recovered.                 |
| #21     | Liszt          | 0.5%     | 13.4% |   |

## 7 Conclusion

We can speculate that treating the tempo as a variable helps the program adapt to unpredictable performance variations, and that modeling the tempo as smooth helps discriminate among hypotheses.

In conclusion, this paper presents an innovative new method for improved score-following, and suggests a promising direction for future research endeavors.