# COMBINING SCORE AND FILTER BASED MODELS TO PREDICT TEMPO FLUCTUATIONS IN EXPRESSIVE MUSIC PERFORMANCES

**Florian Krebs**
Department of Computational Perception
Johannes Kepler University, Linz, Austria
`http://www.cp.jku.at/people/krebs`

**Maarten Grachten**
Department of Computational Perception
Johannes Kepler University, Linz, Austria
`http://www.cp.jku.at/people/grachten`

## ABSTRACT

Tempo variations in classical music are an important means of artistic expression. Fluctuations in tempo can be large and sudden, making applications like automated score following a challenging task. Some of the fluctuations may be predicted from (tempo annotations in) the score, but prediction based only on the score is unlikely to capture the internal coherence of a performance. On the other hand, filtering approaches to tempo prediction (like the Kalman filter) are suited to track gradual changes in tempo, but do not anticipate sudden changes. To combine the advantages of both approaches, we propose a method that incorporates score based tempo predictions into a Kalman filter model of performance tempo. We show that the combined model performs better than the filter model alone.

## 1. INTRODUCTION AND RELATED WORK

Interpreting the world around us is easier when we know what to expect. As in other modes of perception, this appears to be true for music listening as well. Musicians, especially in classical music, use tempo and timing of individual notes as a means of expression, thereby transmitting musically important information such as emotion [1], and a structural interpretation of the piece [2]. These cues are helpful to human listeners, who are familiar with conventions regarding the expressive performance of music.

However, computer systems designed to 'listen to music' typically have no access to expressive performance conventions, and interpreting expressively performed music has proved to be challenging for such systems. A scenario where this is particularly relevant is automatic online tracking of music performances with respect to a score (known as *automatic score following*). This problem has been addressed by several authors in different contexts, like real-time accompaniment of a soloist by a recorded orchestra [3,4], live visualizations of performances and automatic page turning [5]. One of the main challenges in this problem is to accommodate for the fluctuations of tempo in the performance that is being tracked. When fluctuations in tempo are not known in advance, they can lead to ambiguities in the interpretation of the music, that can easily

cause score following systems to lose track of the performance, in the sense of forming an incorrect hypothesis of the position in the score that corresponds to the current position in the performance.

The question of estimating the tempo of a performance at a particular score position, can basically be answered from two complementary perspectives. On the one hand, one can predict performance tempo based solely on score information. This can be regarded as an a priori hypothesis of performance tempo, that does not take into account any information about an actual performance of the piece. On the other hand, one can treat the question as a time series prediction problem, in which the tempo prediction at a given position in the score is based primarily on the observed performance at previous score positions.

Arzt and Widmer have shown that even a simple tempo model based on an weighted average filter improves the automatic score following performance [6]. Furthermore, Bayesian methods for time series modeling have been applied to tempo tracking in various works: Hidden Markov Models [4,7], Kalman filter models [8,9] and particle filters [10,11].

One problem of filtering approaches is that by definition they make predictions using only information from the past. They succeed in modeling slow trends in the data but will fail to model faster changes as sudden jumps. However, such jumps appear frequently in expressive music performances. Some of these jumps are means of articulation induced by the performer, but some are also intended by the composer and annotated in the score (e.g., *ritenuto*: a sudden decrease of the tempo) and hence can be predicted by a score model.

To combine the advantages of filtering approaches with the benefits of a score based approach, we propose a method that incorporates score based tempo predictions into a Kalman filter model of performance tempo. The score based model for performance tempo is an adaption of the model for expressive dynamics proposed in [12]. The Kalman filter is able to exploit prior knowledge about the variance of performance tempo and of the deviations of performed note onsets with respect to the local performance tempo. Additionally, the score model is used to "correct" the filter model whenever it has some additional information from the score.

We show that a combination of score and filter based models performs better than filter models alone, both by large scale evaluation on a corpus of classical piano per-

formances, and by an illustrative example.

The outline of the paper is as follows: In section 2, we describe a baseline filtering model (subsection 2.1), the Kalman filter (subsection 2.2), the linear basis model for score based tempo prediction (subsection 2.3), and the extended Kalman filter tempo model that allows the integration of score based tempo predictions (subsection 2.4). In section 3, we describe the data and experimental setting used to evaluate the model. Finally, we discuss the experimental results (section 4) and present conclusions (section 5).

## 2. METHOD

We represent the tempo using the log beat period $\Delta_k^{\log}$ at beat $k$ using

$$\Delta_k^{\log} = \log_2(\tau_{k+1} - \tau_k) \qquad (1)$$

where $\tau_k$ is the beat time of the $k - th$ beat (see 3.1.1 for how to compute the beat times). The logarithm was also proposed in [8] and is closer to the human perception of tempo: Relative tempo fluctuations become independent of the absolute tempo, e.g., a doubling of the tempo always corresponds to a subtraction of 1 in the log beat period domain.

### 2.1 Baseline approach

To compare the Kalman filter model to a baseline we use a (very) simple tempo model similar to the one proposed in [6]. The beat period of the current note $\Delta_k^{\log}$ is predicted as the weighted average of the beat periods of the $n$ recent notes:

$$\hat{\Delta}_k^{\log} = \frac{\sum_{i=k-n}^{k-1} (\Delta_i^{log} \times i)}{\sum_{i=k-n}^{k-1} i} \qquad (2)$$

### 2.2 Kalman filter

The tempo prediction problem can be formulated as an inference problem in a (continuous) state space model. We are interested in the one-step predictive distribution $p(\hat{\mathbf{s}}_t | o_{1:t-1})$, where $\hat{\mathbf{s}}_t$ is the state prediction at time $t$, and $\mathbf{o}_{1..t-1}$ are the observations from time 1 to $t-1$. Kalman filters [13] provide a very efficient way of recursively computing $p(\hat{\mathbf{s}}_t | o_{1:t-1})$ in linear dynamical systems with Gaussian noise. The restrictions to linear dynamical systems can be loosened by using a linear approximation of the model, so that also non-linear models can be described by an (extended) Kalman filter (see [14] for a comprehensive tutorial).

A Kalman filter model for tempo tracking when no score information is present was first proposed in [8] and is reviewed here briefly: Considering a beat sequence indexed by $k = 1..K$ we denote the state of the system at beat $k$ as $\mathbf{s}_k = (\tau_k, \Delta_k^{\log})^T$, where $\tau_k$ is the beat time of the $k - th$ beat and $\Delta_k^{\log}$ is its log beat period. The states of the model evolve as follows:

$$\tau_k = \tau_{k-1} + 2^{\Delta_{k-1}^{\log}} + v_k^\tau \qquad (3)$$

where $v_k^\tau$ is Gaussian noise with zero mean and variance $Q^{(1)}$.

Fluctuations of the beat period $\Delta_k^{\log}$ are assumed to be Gaussian distributed:

$$\Delta_k^{\log} = \Delta_{k-1}^{\log} + v_k^\Delta \qquad (4)$$

where $v_k^\Delta$ is Gaussian noise with zero mean and variance $Q^{(2)}$. The covariance matrix $\mathbf{Q} = diag(Q^{(1)}, Q^{(2)})$ describes the amount of tempo fluctuation we expect. Hence, the state prediction due to the transition model is written as

$$\hat{\mathbf{s}}_k = f(\mathbf{s}_{k-1}) + \mathbf{v}_k \qquad (5)$$

where $f$ is a non-linear function defined by equation 3 and 4 and $\mathbf{v}_k \sim \mathcal{N}(0, \mathbf{Q})$.

Our observation $o_k$ consists of the beat time $\tau_k$. These beat times incorporate fluctuations like expressive elements and timing errors and hence are modeled as a Gaussian distribution with mean $\tau_k$ and variance $R$:

$$o_k \sim \mathcal{N}(\tau_k, R) \qquad (6)$$

Having observed $o_k$, $\mathbf{s}_k$ is computed as a weighted average of the prediction $\hat{\mathbf{s}}_k$ and the observation $o_k$:

$$\mathbf{s}_k = \hat{\mathbf{s}}_k + \mathbf{K}(o_k - \hat{s}(1)_k) \qquad (7)$$

where $\mathbf{K}$ is called the *Kalman gain*. For details on how to compute $\mathbf{K}$ we refer the reader to [13, 14]

### 2.3 Linear basis modeling

Linear basis modeling refers to a technique introduced in [12], to model expressive parameters in a musical performance as a linear sum of *basis functions*, which describe features of the musical score. Although the approach in [12] was described for loudness, it can be applied without significant modifications to other expressive parameters, such as tempo, articulation, and chord spread. In this paper, we will focus exclusively on modeling tempo fluctuations, quantified inversely as the log beat period.

The intuition behind the linear basis model is that any musical knowledge about the score that could be relevant for shaping expressive tempo is encoded a priori in the basis functions. Basis functions can be as simple as indicator functions for notated accents, or fermata in the score, but they may be arbitrarily complex. A basis function may for example express phrase shapes that are heuristically computed from the score (as in [12]).

The parameters of the model are the weights for summing the basis functions to approximate the expressive tempo curve. They can be easily estimated from data, as described below. The estimated parameters can be used for musicological research (e.g. to quantify differences between performers [15], or the relevance of different musical aspects for expressive tempo), but also to predict expressive tempo for new scores.

Given a musical score, represented as a list of $N$ notes $\mathbf{x} = (x_1, \cdots, x_N)$, and a set of $L$ predefined basis functions $\boldsymbol{\varphi} = (\varphi_1, \cdots, \varphi_L)$, the sequence of log beat periods $\boldsymbol{\Delta}^{\log}$ is modeled as a weighted sum of the basis functions:

$$\boldsymbol{\Delta}^{\log} = f(\mathbf{x}, \mathbf{w}) = \boldsymbol{\varphi}(\mathbf{x})\mathbf{w} + \boldsymbol{\epsilon} \qquad (8)$$
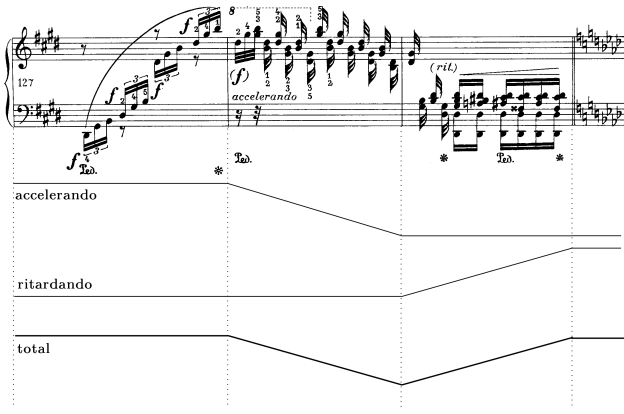
**Figure 1**. Two basis functions and their weighted sum for modelling the log beat period in Chopin Nocturne Op. 1 (bars 127-129)

where we use the notation $\varphi(\mathbf{x})$ to denote the $N \times M$ matrix with element $\varphi_{i,m} = \varphi_m(x_i)$, $\mathbf{w}$ is a vector of $M$ weights and $\epsilon$ is an error vector. Note that $L \leq M$. This is because basis functions can be instantiated multiple times. For example, if a basis function $\varphi_l$ is associated to a score annotation *ritardando* (see subsection 2.3.3), $\varphi(\mathbf{x})$ will contain an instantiation of $\varphi_l$ for each *ritardando* annotation that occurs in $\mathbf{x}$ (shifted in time to be aligned with the annotation). It is important to note that even if $\mathbf{w}$ does not contain exactly $L$ elements, each element of $\mathbf{w}$ can be associated with one of the $L$ unique basis functions $\varphi_l$. We will use the notation $[\mathbf{w}]_l$ to denote the set of weights in $\mathbf{w}$ that are associated with unique basis function $\varphi_l$.

### 2.3.1 Learning basis function weights from data

Given a training set $D = (\ (\mathbf{x}_1, \mathbf{\Delta}_1^{\log}), \cdots, (\mathbf{x}_P, \mathbf{\Delta}_P^{\log})\ )$ of $P$ piece/performance pairs, we solve equation (8) for $\mathbf{w}$, for each piece/performance pair $1 \leq p \leq P$ individually:

$$\hat{\mathbf{w}}_p = \text{argmin}_{\mathbf{w}} \left\| \mathbf{\Delta}_p^{\log} - \varphi(\mathbf{x}_p)\mathbf{w} \right\| \qquad (9)$$

where $\|\cdot\|$ denotes the $\ell^2$-norm.

Then, we compute the estimated weight $\hat{w}_l$ for each unique basis function $\varphi_l$ simply as the mean of all weights associated to $\varphi_l$, as computed from the training data:

$$\hat{w}_l = \begin{cases} \overline{W}_l & \text{if } |W_l| > 0 \\ 0 & \text{otherwise} \end{cases} \qquad (10)$$

that is, $\hat{w}_l$ is the average over the set of weights $W_l$, where:

$$W_l = \bigcup_{p=1}^{P} [\hat{\mathbf{w}}_p]_l \qquad (11)$$

### 2.3.2 Prediction

For a new score $\mathbf{x}$, the sequence of log beat periods is then estimated as:

$$\hat{\mathbf{\Delta}}^{\log} = f(\mathbf{x}, \hat{\mathbf{w}}) = \varphi(\mathbf{x})\hat{\mathbf{w}} \qquad (12)$$

The weight vector $\hat{\mathbf{w}}$ is constructed by first constructing the matrix $\varphi(\mathbf{x})$, and then for each of the $L$ basis functions in $\varphi(\mathbf{x})$, selecting the associated weight $\hat{w}_l$.

### 2.3.3 Basis functions used

In the current context we define basis functions that represent tempo annotations in the score. The basis functions are intended to represent the prototypical effect of the annotation on the tempo curve. We distinguish three prototypes. Firstly, there are annotations that indicate a constant tempo (e.g. *lento*, *allegro*, or *a tempo*[1]). Such annotations are represented by a step function, with value one over the range where the annotation is active, and zero elsewhere. Secondly there are dynamic tempo annotations, prescribing a decrease (e.g. *ritardando*), or increase (e.g. *accelerando*) in tempo, respectively. They are represented by a ramp function over the range where the annotation is active, and zero elsewhere. A simple example of two such basis functions is shown in figure 1. Lastly, there are annotations with only an instantaneous effect, such as *fermata*. Such annotations are represented by an impulse function with value one at the time of the annotation, and zero elsewhere.

It is important to note that although the prototype of the tempo curve is fixed by the basis function, the basis functions do not impose a specific direction, such as a tempo increase in case of a *accelerando*, and a decrease in case of a *ritardando*. Both annotations are represented by the same ramp function, and only the sign of the weights (which are learned from real performances), determines if the effect of, say a *ritardando* annotation, will be an increase or rather a decrease in tempo.

## 2.4 Combined model

To combine score and filter information, we first compute the first-order temporal difference of the linear and logarithmic score model predictions (Eq. 12):

$$u_k^{(1)} = 2^{\hat{\Delta}_k^{\log}} - 2^{\hat{\Delta}_{k-1}^{\log}} \qquad (13)$$

$$u_k^{(2)} = \hat{\Delta}_k^{\log} - \hat{\Delta}_{k-1}^{\log} \qquad (14)$$

Using $u_k^{(1)}$ and $u_k^{(2)}$ as control input, the prediction of the combined model becomes:

$$\hat{\mathbf{s}}_k = f(\mathbf{s}_{k-1}) + b \begin{bmatrix} u_k^{(1)} \\ u_k^{(2)} \end{bmatrix} \qquad (15)$$

where $b$ is a control input parameter.

## 3. EXPERIMENTS

To evaluate the proposed tempo model we use it to predict the next ($log$) beat period given all previous beats of real performances.

---

[1] Referential annotations such as *a tempo*, or *tempo primo* are handled by combining the both referent and referee annotations into a single basis function

## 3.1 Data set

For the evaluation we use the Magaloff corpus [16] a data set that comprises live performances of virtually the complete Chopin piano works, as played by the Russian- Georgian pianist Nikita Magaloff (1912-1992). The music was performed in a series of concerts in Vienna, Austria, in 1989, on a Bösendorfer SE computer-controlled grand piano [17] that recorded the performances onto a computer hard disk. Symbolic scores were obtained from scanned sheet music using *optical music recognition* (OMR). Performances were aligned to the score automatically, and were corrected manually.

Tempo markings in the score are available as far as they have been recognized in the OMR process. Unfortunately, many tempo markings are not correctly recognized, and recognized markings are not always positioned correctly. [2] Even if this is expected to limit the potential benefit of the score model, the current annotations are sufficient to demonstrate the benefit.

The data set consists of 131 pieces, adding up to approximately 9 hours of music, and 59917 beats.

### 3.1.1 Preprocessing

To make the data set usable for the evaluation several preprocessing steps are undertaken: Firstly, performance onsets, that happen at the same time in the score, are replaced by their arithmetic mean. Secondly, we deal with the fact that the performance data differs from the actual score in the number of notes: *Deletion* of notes happens when the performer skips notes and *insertion* happens when he plays additional notes that are not present in the score. To account for the errors, which would emerge when computing the inter-onset intervals, the missing sections are cut out manually from the performance (insertion case) or the score (deletion case). Thirdly, the beat times are obtained by interpolation of the unique performance note onsets and finally, all predicted and target sequences are normalized to a total duration of 1 to be independent of the mean tempo.

### 3.1.2 Parameter estimation

The noise covariance $\mathbf{Q}$ and variance R of the Kalman filter have been optimized using a grid search over the whole data set using leave-one-out cross validation. We have found that parameters optimized using an Expectation - Maximization (EM) algorithm do not produce better tempo predictions (as measured by the log beat period). Sometimes the target sequence even has a lower likelihood under our model assumptions than a predicted sequence.

The window size of the weighted average filtered was also optimized by grid search. [3]

The weights of the linear basis model have been trained for each piece individually, also using a leave-one-out approach on the complete data set, according to the method described in subsection 2.3.1.

## 3.2 Evaluation measures

As in [12], we use two quantities to quantify how well the model is able to capture tempo variations of the performances: $r$ is the Pearson product-moment correlation coefficient, denoting how strong the observed log beat periods and the log beat periods of the fitted model linearly depend on each other. The quantity $R^2$ is the coefficient of determination and is a measure for how much of the tempo variance is accounted for by the model.

## 4. RESULTS AND DISCUSSION

Table 1 lists the results for all algorithms on the complete data set described in section 3.1. A one-tailed paired t-test on the $R^2$ measures indicates that the Kalman filter significantly outperforms the weighted average filter ($t(130) = 2.942, p = 0.005$). Integration of the proposed score model improves the performance even more and has been found to perform significantly better than the weighted average filter ($t(130) = 4.019, p = 0.00005$) and also significantly better than the Kalman filter alone ($t(130) = 2.271, p = 0.025$) . As shown in figure 2 the weighted average filter and the Kalman filter capture the main tempo fluctuations but are always one step too late. This is because in the prediction step it always uses the tempo that was computed in the last update step. In contrast, information from the score is available off-line and can tell the filter models in time when sudden changes happen. This is shown in figure 2, where a *ritenuto* causes the performer to slow down between beats 222 and 227 and then, after an *in tempo*, returns to the initial tempo at beat 228. The same happens between beats 271 and 278 (*ritenuto* from beat 271 to 274 and *in tempo* from beat 275 to 278) ).

We find that the Kalman filter works best when the tempo variance parameter $Q^{(2)}$ is very low [4] in comparison to the other noise variance parameters, resulting in smooth beat period state sequences $\mathbf{\Delta}^{\log}$. Due to large tempo fluctuations it is convenient to have a stable tempo estimate that captures the slow fluctuations only. Fast fluctuations are modeled by the beat time state $\tau$ with higher variance $Q^{(1)}$. We have also incorporated additional tempo states with different transition variances as proposed in [10], but have not found it beneficial for the predictions.

Note that although the explained proportion of variance may seem rather low, this is not unexpected for models of such complicated human behaviors as music performance. It reflects that we still know relatively little of the factors that contribute to musical expression. Also, it is likely that not all model assumptions are met. For one thing, the score model assumes that annotated tempo directives will always be followed by the performer (to an extent determined by the training data). This assumption is unlikely to be true in all cases, for example when the performer used an other edition of the score, possibly containing different annotations than those used by the model.

The fact that a rather simple score model improves the

---

[2] For a subset of the data (around 5%, including Op. 62, No. 2, in figure 2), tempo annotations have been corrected manually.

[3] We obtained the best results using n=6 beats

[4] For good performance, typical values of $Q^{(1)}$ ranged from 8 to 12, typical values of $Q^{(2)}$ ranged from 0 to $10^{-4}$, typical values of $R$ ranged from 26 to 38 and typical values of $b$ ranged from 0.035 to 0.045.
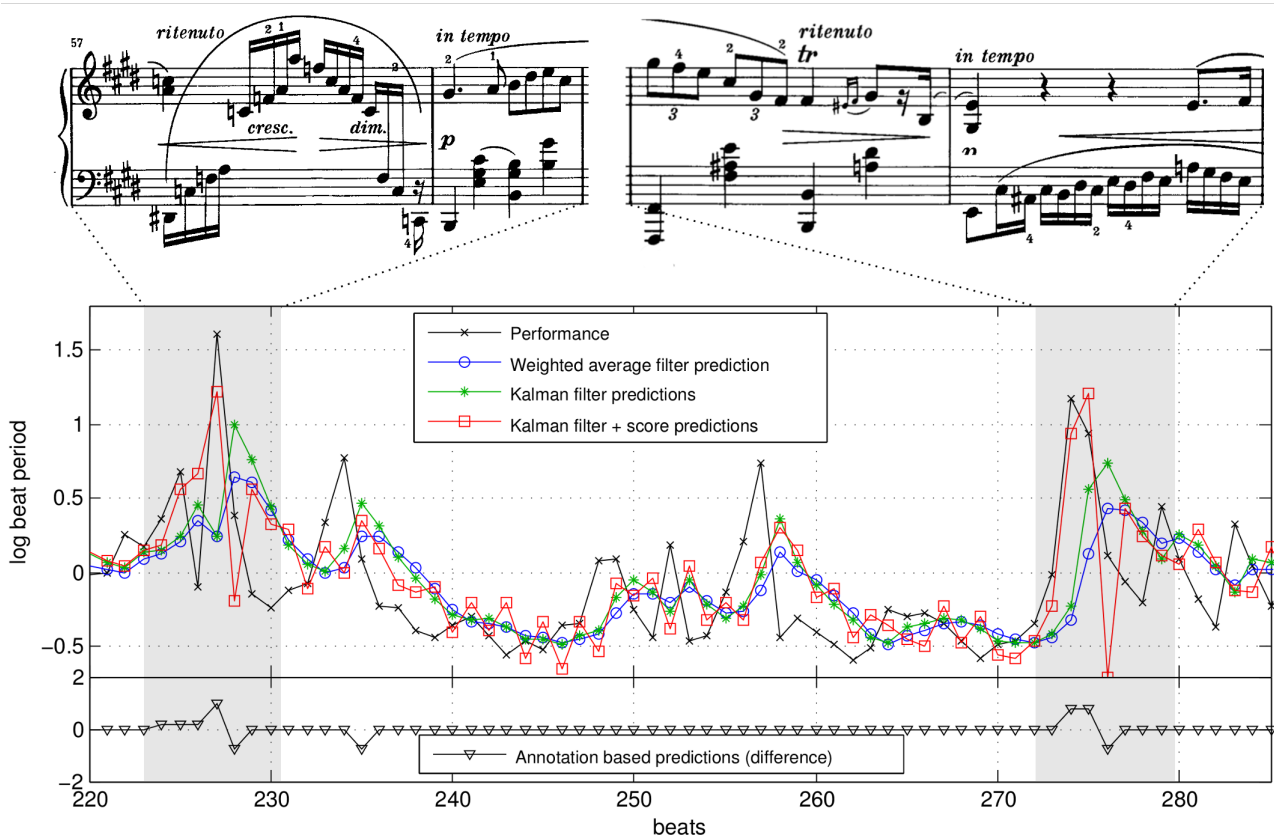
**Figure 2**. Tempo predictions for Chopin Nocturne Op. 62, No. 2 (bars 57–71)

| Algorithm | $R^2$ | | $r$ | |
|---|---|---|---|---|
| | avg. | std. | avg. | std. |
| WAF | 0.217 | 0.232 | 0.499 | 0.190 |
| EKF | 0.224 | 0.242 | 0.520 | 0.183 |
| EKF+LBM | **0.242** | 0.229 | **0.537** | 0.165 |

**Table 1**. Accuracy of different tempo prediction models on the Magaloff corpus (WAF: Weighted average filter; EKF: Extended Kalman Filter; LBM: Linear basis model)

filtering approach in spite of overly strong assumptions, and partially incomplete tempo annotations (see subsection 3.1), shows that the score is a robust source of information for modeling expressive tempo in music performances.

## 5. CONCLUSIONS AND FUTURE WORK

Applications like score following benefit from models that can predict the tempo fluctuations of a performance. Current (filter based) tempo models only take into account notes that were played in the past and hence can not anticipate any sudden tempo changes. This main drawback of the filtering approaches can be alleviated by incorporating information from the score. In this paper we present a method to improve tempo models by combining them with a scored based model of tempo. We show that the Kalman filtering approach outperforms the weighted average filter and its performance further increases when used in combination with the proposed score based model.

Note that also other sources of tempo knowledge - like tempo curves learned from previous performances of the same piece [6] - could be integrated into the model using the proposed procedure.

In future work we would like to jointly estimate parameters of the filter and the score model using a probabilistic formulation of the score annotation model. Besides, we plan to investigate other combination methods, in addition to the one presented in subsection 2.4. In cases where annotations over longer sections occur (e.g., *accelerando*, *ritardando*) it could be effective to increase the Kalman gain in order to put more trust in the observations than in the transition model. In this manner, the filter is informed that some relevant changes are going to happen and it should react more sensible to incoming observations.

Finally, we plan to integrate the new tempo model into a score following application to prove its usability in practice.

### Acknowledgments

# 6. REFERENCES

[1] A. Gabrielsson and P. N. Juslin, "Emotional expression in music performance: Between the performer's intention and the listener's experience," *Psychology of Music*, vol. 24, no. 1, 1996.

[2] E. F. Clarke, "Generative principles in music," in *Generative Processes in Music: The Psychology of Performance, Improvisation, and Composition*, J. Sloboda, Ed. Oxford University Press, 1988.

[3] R. Dannenberg, "An on-line algorithm for real-time accompaniment," in *Proceedings of the 1984 International Computer Music Conference.* International Computer Music Association, 1984.

[4] C. Raphael, "A bayesian network for real-time musical accompaniment," in *Advances in Neural Information Processing Systems, NIPS 14.* MIT Press, 2001.

[5] A. Arzt, G. Widmer, and S. Dixon, "Automatic page turning for musicians via real-time machine listening," in *In Proc. of the 18th European Conference on Artificial Intelligence (ECAI*, 2008.

[6] A. Arzt and G. Widmer, "Simple tempo models for real-time music tracking," in *Proc. of the Sound and Music Computing Conference (SMC), Barcelona, Spain*, 2010.

[7] A. Cont, "A coupled duration-focused architecture for real-time music-to-score alignment," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 32, no. 6, pp. 974–987, 2010.

[8] A. Cemgil, B. Kappen, P. Desain, and H. Honing, "On tempo tracking: Tempogram representation and kalman filtering," *Journal of New Music Research*, vol. 29, no. 4, pp. 259–273, 2000.

[9] Y. Shiu and C. Kuo, "Musical beat tracking via kalman filtering and noisy measurements selection," in *Circuits and Systems, 2008. ISCAS 2008. IEEE International Symposium on.* IEEE, 2008, pp. 3250–3253.

[10] A. Cemgil and B. Kappen, "Monte carlo methods for tempo tracking and rhythm quantization," *Journal of Artificial Intelligence Research*, vol. 18, no. 1, pp. 45–81, 2003.

[11] S. Hainsworth and M. Macleod, "Particle filtering applied to musical tempo tracking," *EURASIP Journal on Applied Signal Processing*, vol. 2004, pp. 2385–2395, 2004.

[12] M. Grachten and G. Widmer, "Explaining expressive dynamics as a mixture of basis functions," in *Proceedings of the Eighth Sound and Music Computing Conference (SMC)*, Padua, Italy, 2011.

[13] R. Kalman *et al.*, "A new approach to linear filtering and prediction problems," *Journal of basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.

[14] G. Terejanu, "Extended kalman filter tutorial," http://users.ices.utexas.edu/~terejanu/files/tutorialEKF.pdf, University at Buffalo: Buffalo, NY, USA, Tech. Rep., 2003, accessed on 21 May 2012.

[15] M. Grachten and G. Widmer, "A method to determine the contribution of annotated performance directives in music performances," in *Proceedings of the International Symposium of Performance Science*, Toronto, Canada, 2011.

[16] S. Flossmann, W. Goebl, M. Grachten, B. Niedermayer, and G. Widmer, "The Magaloff Project: An Interim Report," *Journal of New Music Research*, vol. 39, no. 4, pp. 369–377, 2010.

[17] R. A. Moog and T. L. Rhea, "Evolution of the Keyboard Interface: The Bösendorfer 290 SE Recording Piano and the Moog Multiply-Touch-Sensitive Keyboards," *Computer Music Journal*, vol. 14, no. 2, pp. 52–60, 1990.