

# Simulating Human Tonal Expectation using Recurrent Neural Networks

Carlos Cancino-Chacón<sup>1,2</sup>, Maarten Grachten<sup>2</sup>, and Kat Agres<sup>3</sup>

<sup>1</sup>Austrian Research Institute for Artificial Intelligence

<sup>2</sup>Johannes Kepler University of Linz

<sup>3</sup>Institute of High Performance Computing, A\*STAR, Singapore

Technical Report; Version 0.01

## Abstract

This document includes some of the technical details for reproducing the work presented in

**From Bach to the Beatles: The simulation of human tonal expectation using ecologically-trained predictive models,**

presented at the 18th International Society for Music Information Retrieval Conference (ISMIR 2017) in Suzhou, China.

## 1 Data preprocessing

An audio signal can be represented using a Constant-Q Transform (CQT) [Brown, 1991], which is a discrete frequency domain representation of audio.

The  $t$ -th slice of a CQT spectrogram is a non-negative vector  $\mathbf{x}_t \in \mathbb{R}^{334}$  that represents frequencies between 27.5 and 16744.04 Hz with a resolution of 36 frequency bins per octave. The audio in all recordings has a sample rate of 44.1 kHz. We normalize each slice so that all of its components lie between 0 and 1 as

$$\hat{\mathbf{x}}_t \leftarrow \frac{\tanh(\mathbf{x}_t)}{\tanh(\max(\mathbf{x}_t) + \varepsilon)}, \quad (1)$$

where  $\max(\mathbf{x}_t)$  represents the value of the largest component of  $\mathbf{x}_t$ ,  $\tanh(\cdot)$  is the element-wise hyperbolic tangent function and  $\varepsilon$  is a machine epsilon for numerical stability. In this work, we set  $\varepsilon = 10^{-10}$ . Slightly abusing notation, in the following we will assume that all CQT slices have been normalized, and therefore simply refer to the normalized CQT slices  $\hat{\mathbf{x}}$  as  $\mathbf{x}$ .

## 2 Recurrent Neural Networks

### 2.1 Architectures

An RNN is a neural architecture that allows for modeling dynamical systems [Graves, 2013]. Let  $\mathbf{x}_1, \dots, \mathbf{x}_t$  be a sequence of  $N$ -dimensional (normalized) input vectors and

$\mathbf{y}_1, \dots, \mathbf{y}_t$  be its corresponding sequence of outputs. An RNN provides a natural way to model  $\mathbf{x}_{t+1}$ , the next event in the sequence, by using the outputs of the network to parametrize a predictive distribution given by

$$p(x_{t+1,i} | \mathbf{x}_t, \dots, \mathbf{x}_1) = y_{t,i} \quad (2)$$

where  $x_{t+1,i}$  and  $y_{t,i}$  are the  $i$ -th component of  $\mathbf{x}_{t+1}$  and  $\mathbf{y}_t$  respectively.

The basic component of an RNN is the *recurrent layer*, whose activation at time  $t$  depends on both the input at time  $t$  and its activation at time  $t-1$ . Although theoretically very powerful, in practice RNNs with *vanilla* recurrent layers are known to have problems learning long term dependencies due to a number of problems, including vanishing and exploding gradients [Pascanu et al., 2013]. Other recurrent layers such as LSTM layers [Hochreiter and Schmidhuber, 1997] and gated recurrent units (GRUs)[Chung et al., 2014] try to address some of these problems by introducing special structures within the layer, such as purpose-built memory cells and gates to better store information. More recently, recurrent layers with multiplicative integration (MI-RNNs) [Wu et al., 2016] have been shown to extend the expressivity of traditional additive RNNs by changing the way the information from different sources is aggregated within the layer while introducing just a small number of extra parameters.

The predictions of the network can be written as

$$\mathbf{y}_t = \sigma(\mathbf{W}_y \mathbf{h}_t(\mathbf{x}_t | \theta_h) + \mathbf{b}_y), \quad (3)$$

where  $\mathbf{h}_t \in \mathbb{R}^H$  is the output of the hidden layer at time  $t$ ,  $\theta_h$  are the parameters of the layer,  $\mathbf{W}_y \in \mathbb{R}^{334 \times H}$  is a matrix of weights connecting the recurrent layer to the output  $\mathbf{b}_y \in \mathbb{R}^{334}$  is a bias vector and  $\sigma(\cdot)$  is the element-wise sigmoid function.

A detailed description of the recurrent layers used in this paper can be found in Appendix A.

## 2.2 Training

Given a set of  $T$  training sequences  $\mathbf{X} = \{\{\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_{T_i-1}^{(k)}\} | 1 \leq k \leq T\}$ , where  $T_k$  is the length of the  $k$ -th sequence, and their respective targets  $\mathbf{T} = \{\{\mathbf{x}_2^{(k)}, \dots, \mathbf{x}_{T_i}^{(k)}\} | 1 \leq k \leq T\}$ , i.e. a shifted version of the input sequences, the parameters of the network are learned by minimizing the mean cross entropy (MCE) [Bishop, 2006]

$$MCE = -\frac{1}{TN} \sum_{k=1}^T \frac{1}{T_k} \sum_{t=1}^{T_k} \sum_{i=1}^N CE_{kti}, \quad (4)$$

where  $CE_{kti}$  is the cross entropy of the  $i$ -th feature at time-step  $t$  for the  $k$ -th sequence given by

$$CE_{kti} = t_{kti} \log(p_{kti}) + (1 - t_{kti}) \log(1 - p_{kti}), \quad (5)$$

and  $p_{kti} = p(x_{t,i}^{(k)} | \mathbf{x}_{t-1}^{(k)}, \dots, \mathbf{x}_1^{(k)})$  and  $t_{kti} = x_{t,i}^{(k)}$ .

### 2.2.1 Biasing learning towards predicting change

A crucial question when applying discrete time recurrent models to a continuous stream of data such as audio is how to choose the rate of discrete time steps with

respect to the absolute time of the data. This choice depends on the approximate rate or temporal density of relevant events in the data—in our case the notes that make up the musical material. Ideally, we would like the discrete time steps to be small enough to capture the occurrence of even the shortest notes individually, but if the discrete time step is chosen much smaller than the median event rate, this leads to strong correlations between data at consecutive time steps. A result of this is that training models to predict the data at time-step  $t + 1$  teaches them to strongly expect the data at  $t + 1$  to be approximately equal to the data at  $t$ . Choosing a larger discrete step size for the model alleviates this problem, but has the disadvantage that the data the model sees at a particular time may actually be an average over consecutive events that happened within that larger step.

We slightly revise the training objective of the models as a remedy to this unfortunate trade-off. This revised objective biases the models to care more about correctly predicting the data at  $t + 1$  when the change from  $t$  to  $t + 1$  is large (e.g. the start of a new note) than when it is small (e.g. a transition without any starting or ending note events). This allows us to use a relatively small step size without causing the models to trivially learn to expect the data to stay constant between consecutive time steps.

More specifically, we modify the original cross-entropy objective  $CE_{kti}$  by multiplying it with a time-varying weight  $w_{kt}$  as follows:

$$\tilde{CE}_{kti} \leftarrow w_{kt} CE_{kti}, \quad (6)$$

where  $w_{kt}$  is given by

$$w_{kt} = \begin{cases} 1 & \text{if } \sum_i^N |x_{t+1,i}^{(k)} - x_{t,i}^{(k)}| > \varepsilon \\ \beta & \text{otherwise} \end{cases} \quad (7)$$

where  $\varepsilon \in \mathbb{R}$  acts as a threshold distinguishing small and large change transitions, and  $\beta \in \mathbb{R}$  controls the relative influence of prediction errors on the training in the case of small change transitions<sup>1</sup>. Based on an informal inspection of the model predictions in a grid search on  $\beta$  and  $\varepsilon$ , we choose  $\beta = 10^{-3}$ , and  $\varepsilon$  such that

$$P_{training} \left( \sum_i^N |x_{t+1,i} - x_{t,i}| \leq \varepsilon \right) = 0.505 \quad (8)$$

where  $P_{training}(X)$  denotes the empirical probability of event  $X$  under the training data.

## 2.3 Hyperparameters

The models are trained using RMSProp [Tieleman and Hinton, 2012], a variant of stochastic gradient descent that adaptively updates the step-size using a moving average of the magnitude of the gradients. The initial learning rate is set to  $10^{-3}$ . The gradients are computed using truncated back propagation through time, where computation of the gradients is truncated after 100 steps and are clipped at 1. Each training batch consists of 20 sequences of 100 CQT slices. Each sequence is selected randomly out of the training data. Thus, an epoch of training corresponds to the

---

<sup>1</sup>We empirically found a binary distinction between small and large change transitions to be more effective than a gradual weighting scheme

model seeing roughly the same number of time steps as in the whole fold. Early stopping is used after 100 epochs without any improvement in the test set.

### 3 Probe-tone experiments

Let  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$  be an input musical stimulus, also referred to as context, and  $\mathbf{T} = \{\tau_1, \dots, \tau_{12}\}$  the set of probe-tones each corresponding to one of the 12 pitch classes. In order to quantitatively assess how well a probe-tone  $\tau$  fits the musical stimulus, we compare  $\mathbf{y}^*$ , the next prediction of the RNN given the input stimulus  $\mathbf{X}$ , and the probe-tone using the Kullback-Leibler (KL) divergence, given by

$$D_{KL}(\tau || \mathbf{y}^*) = \sum_{i=1}^N \tau_i \log \left( \frac{\tau_i}{y_i} \right), \quad (9)$$

where  $\tau_i$  and  $y_i$  are the  $i$ -th components of  $\tau$  and  $\mathbf{y}^*$  respectively. In this case, we can see that the KL divergence is small when the expectations of the model are similar to the probe-tones.

We can construct the probe-tone profile for a stimulus in key  $key$ , denoted by  $\mathbf{P}_{key}^{model} \in \mathbb{R}^{12}$ , as

$$\mathbf{P}_{key}^{model} = \text{minmax} \begin{pmatrix} -D_{KL}(\tau_1 || \mathbf{y}_{key}^*) \\ \vdots \\ -D_{KL}(\tau_{12} || \mathbf{y}_{key}^*) \end{pmatrix}, \quad (10)$$

where  $\text{minmax}(\cdot)$  is a normalization function given by

$$\text{minmax}(\mathbf{x}) = \frac{\mathbf{x} - \min(\mathbf{x})\mathbf{1}}{\max(\mathbf{x}) - \min(\mathbf{x})}, \quad (11)$$

where  $\min(\mathbf{x})$  and  $\max(\mathbf{x})$  are the minimal and maximal components of the input  $\mathbf{x}$  and  $\mathbf{1} \in \mathbb{R}^{12}$  is a vector of ones.

The average profile for a stimulus is computed as

$$\mathbf{P}^{model} = \frac{1}{12} \sum_{key} \mathbf{P}_{key}^{model}. \quad (12)$$

We also normalize the human profiles reported by Krumhansl and Kessler by using the minmax function. We use Pearson's correlation coefficient to compare the profiles of the model to those of humans

$$r = \frac{\sum_{i=1}^{12} (P_i^{model} - \bar{P}^{model})(P_i^{human} - \bar{P}^{human})}{\sqrt{\sum_{i=1}^{12} (P_i^{model} - \bar{P}^{model})^2} \sqrt{\sum_{i=1}^{12} (P_i^{human} - \bar{P}^{human})^2}}, \quad (13)$$

where  $P_i$  represents the  $i$ -th component of a profile  $P$ , and  $\bar{P} = \frac{1}{12} \sum_i P_i$ .

## References

[Bishop, 2006] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer Verlag, Microsoft Research Ltd.

- [Brown, 1991] Brown, J. C. (1991). Calculation of a constant Q spectral transform. *The Journal of the Acoustical Society of America*, 89(1):425–434.
- [Chung et al., 2014] Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.
- [Graves, 2013] Graves, A. (2013). Generating Sequences With Recurrent Neural Networks. *arXiv*, 1308:850.
- [Hochreiter and Schmidhuber, 1997] Hochreiter, S. and Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780.
- [Pascanu et al., 2013] Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1–9, Atlanta, Georgia, USA.
- [Tieleman and Hinton, 2012] Tieleman, T. and Hinton, G. (2012). Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. In *COURSERA Neural Networks for Machine Learning*.
- [Wu et al., 2016] Wu, Y., Zhang, S., Zhang, Y., Bengio, Y., and Salakhutdinov, R. (2016). On Multiplicative Integration with Recurrent Neural Networks. In *30th Conference on Neural Information Processing Systems (NIPS 2016)*, Barcelona, Spain.

## A Recurrent layers

### A.1 LSTM with Multiplicative Integration

Let  $\mathbf{x}_t \in \mathbb{R}^N$  be an input vector representing the CQT at time  $t$ . The activation of the MI-LSTM layer is given by

$$\mathbf{z}_t = \tanh(\beta_{z,1} \odot \mathbf{U}_z \mathbf{h}_{t-1} + \beta_{z,2} \odot \mathbf{W}_z \mathbf{x}_t + \alpha_z \odot \mathbf{W}_z \mathbf{x}_t \odot \mathbf{U}_z \mathbf{h}_{t-1} + \mathbf{b}_z) \quad (14)$$

$$\mathbf{i}_t = \sigma(\beta_{i,1} \odot \mathbf{U}_i \mathbf{h}_{t-1} + \beta_{i,2} \odot \mathbf{W}_i \mathbf{x}_t + \alpha_i \odot \mathbf{W}_i \mathbf{x}_t \odot \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i) \quad (15)$$

$$\mathbf{f}_t = \sigma(\beta_{f,1} \odot \mathbf{U}_f \mathbf{h}_{t-1} + \beta_{f,2} \odot \mathbf{W}_f \mathbf{x}_t + \alpha_f \odot \mathbf{W}_f \mathbf{x}_t \odot \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f) \quad (16)$$

$$\mathbf{o}_t = \sigma(\beta_{o,1} \odot \mathbf{U}_o \mathbf{h}_{t-1} + \beta_{o,2} \odot \mathbf{W}_o \mathbf{x}_t + \alpha_o \odot \mathbf{W}_o \mathbf{x}_t \odot \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o) \quad (17)$$

$$\mathbf{c}_t = \mathbf{i}_t \odot \mathbf{z}_t + \mathbf{f}_t \odot \mathbf{c}_{t-1} \quad (18)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t), \quad (19)$$

where  $\mathbf{z}_t, \mathbf{h}_t \in \mathbb{R}^{N_h}$  are the *block input* and *output* (i.e. the activation of the MI-LSTM layer), respectively,  $\mathbf{c}_t \in \mathbb{R}^{N_h}$  is the *cell state* and  $\mathbf{i}_t, \mathbf{f}_t, \mathbf{o}_t \in \mathbb{R}^{N_h}$  are the *input*, *forget* and *output gates*, respectively. The parameters of an MI-LSTM layer are  $\{\mathbf{W}_k \mid k \in \{z, i, f, o\}\}$ , input-to-hidden weight matrices in  $\mathbb{R}^{N_h \times N}$ ,  $\{\mathbf{U}_k \mid k \in \{z, i, f, o\}\}$ , hidden-to-hidden weight matrices in  $\mathbb{R}^{N_h \times N_h}$ ,  $\{\mathbf{b}_k \mid k \in \{z, i, f, o\}\}$ , bias vectors in  $\mathbb{R}^{N_h}$  and  $\{\alpha_k, \beta_{k,1}, \beta_{k,2} \mid k \in \{z, i, f, o\}\}$ , gating bias vectors for multiplicative integration in  $\mathbb{R}^{N_h}$ .  $\sigma(\cdot)$  and  $\tanh(\cdot)$  are the elementwise sigmoid and hyperbolic tangent functions, respectively.