

# LEARNING INVARIANTS FOR POLYPHONIC INSTRUMENT RECOGNITION

Vincent Lostanlen, Carmine Emanuele Cella, and Stéphane Mallat  
École normale supérieure

## ABSTRACT

The abstract should be placed at the top left column and should contain about 150-200 words.

## 1. INTRODUCTION

Among the cognitive attributes of musical tones, pitch is distinguished by a combination of three properties. First, it is relative: ordering pitches from low to high gives rise to the notions of intervals and scales. Secondly, it is intensive: multiple pitches heard simultaneously produce a chord, not a single unified tone – contrary to loudness, which adds up with the number of sources. Thirdly, it is an invariant for source recognition: musical instruments have a single timbral identity over their whole range.

In the symbolic domain, melodies are represented by a piano-roll matrix  $\mathcal{P}$  over the continuous dimensions of time and pitch. In turn, other semantic attributes, such as instrument activations, are represented by a histogram  $\mathcal{H}$  over a finite number of categories. In this setting, music instrument recognition is the task of finding the active instruments at a given time point of an audio excerpt.

## 2. DEEP CONVOLUTIONAL NETWORKS

### 2.1 Time-frequency representation

We used the implementation from the librosa package [4] with  $Q = 12$  filters per octave, center frequencies ranging from 55 Hz to 14 kHz (8 octaves from A1 to A9), and a hop size of 23 ms. Furthermore, we applied nonlinear perceptual weighting of loudness in order to reduce the dynamic range between the fundamental partial and its upper harmonics. A 3-second sound excerpt  $x(t)$  is represented by a time-frequency matrix  $\mathbf{x}_1(t, k_1)$  of width  $T = 128$  samples and height  $K_1 = 96$  MIDI indices.

### 2.2 Architecture

First of all, we apply a family  $\mathbf{W}_2(\tau, \kappa_1, k_2)$  of  $K_2 = 50$  learned time-frequency convolutional operators, whose

supports are constrained to have width  $\Delta t$  and height  $\Delta k_1$ .

$$\mathbf{W}_2^{t, k_1} * \mathbf{x}_1 = \sum_{\substack{0 \leq \tau < \Delta t \\ 0 \leq \kappa_1 < \Delta k_1}} \mathbf{W}_2(\tau, \kappa_1, k_2) \mathbf{x}_1(t - \tau, k_1 - \kappa_1) \quad (1)$$

Furthermore, element-wise biases  $\mathbf{b}_2(k_2)$  are added to the convolutions, resulting in the tensor

$$\mathbf{y}_2(t, k_1, k_2) = \mathbf{b}_2(k_2) + (\mathbf{W}_2^{t, k_1} * \mathbf{x}_1)(t, k_1, k_2). \quad (2)$$

The second step is the application of a pointwise nonlinearity. We have chosen the *rectified linear unit* (ReLU) because of its popularity in computer vision and its computational efficiency.

$$\mathbf{y}_2^+(t, k_1, k_2) = \max(\mathbf{y}_2(t, k_1, k_2), 0) \quad (3)$$

To achieve invariance to translation as well as frequency transposition, we pool neighboring units in the time-frequency domain  $(t, k_1)$  over non-overlapping rectangles of width  $\Delta t$  and height  $\Delta k_1$ .

$$\mathbf{x}_2(t, k_1, k_2) = \max_{\substack{0 \leq \tau < \Delta t \\ 0 \leq \kappa_1 < \Delta k_1}} \left\{ \mathbf{y}_2^+(t - \tau, k_1 - \kappa_1, k_2) \right\} \quad (4)$$

We apply a family  $\mathbf{W}_3(\tau, \kappa_1, k_2, k_3)$  of  $K_3$  convolutional operators that perform a linear combination of time-frequency feature maps in  $\mathbf{x}_2$  along the channel variable  $k_2$ .

$$\mathbf{y}_3(t, k_1, k_3) = \sum_{k_2} \mathbf{b}_3(k_2, k_3) + \mathbf{W}_3(t, k_1, k_3) \overset{t, k_1}{*} \mathbf{x}_2(t, k_1, k_2) \quad (5)$$

After nonlinear rectification and max-pooling, the layer  $\mathbf{y}_3$  turns into a non-negative tensor  $\mathbf{x}_3(t, k_1, k_3)$ .

$$\mathbf{y}_4(k_4) = \sum_{t, k_1, k_3} \mathbf{W}_4(t, k_1, k_3, k_4) \mathbf{x}_3(t, k_1, k_3) \quad (6)$$

We apply nonlinear rectification, yielding  $\mathbf{x}_4(k_4) = \mathbf{y}_4^+(k_4)$ .  $\mathbf{y}_5(k_5) = \sum_{k_5} \mathbf{W}_5(k_4, k_5) \mathbf{x}_4(k_4)$ .

$$\mathbf{x}_5(k_5) = \frac{\exp \mathbf{y}_5(k_5)}{\|\exp \mathbf{y}_5\|_1} \quad (7)$$

The above ensures that the coefficients of  $\mathbf{x}_5$  are non-negative and sum to one, hence can be fit to a probability distribution. We define the categorical cross-entropy as

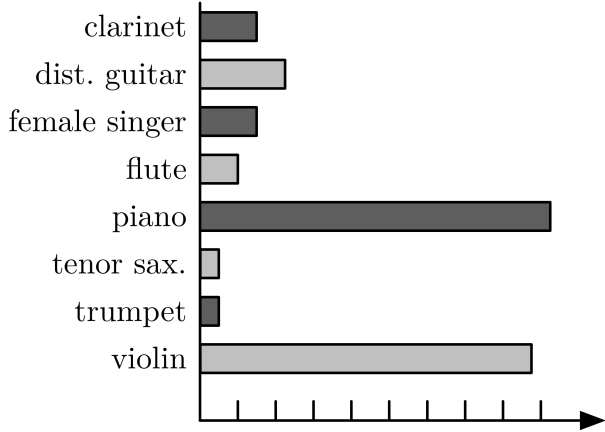
$$\mathcal{L}(\mathbf{x}_5, \mathcal{I}) = - \sum_{k_5 \in \mathcal{I}} \log \mathbf{x}_5(k_5). \quad (8)$$

The goal is to minimize the average loss  $\mathcal{L}(\mathbf{x}_5, \mathcal{I})$  for across all pairs  $(\mathbf{x}_5, \mathcal{I})$  in the training set.



© Vincent Lostanlen, Carmine Emanuele Cella, and Stéphane Mallat.

Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Vincent Lostanlen, Carmine Emanuele Cella, and Stéphane Mallat. “Learning invariants for polyphonic instrument recognition”, 16th International Society for Music Information Retrieval Conference, 2015.



**Figure 1:** Amount of training data per instrument in MedleyDB, in minutes.

### 2.3 Training

The network is trained on categorical cross-entropy over shuffled mini-batches of size 512 with uniform class distribution. The learning rate policy for each scalar weight in the network is *Adam* [3], a state-of-the-art online optimizer for gradient-based learning.

## 3. DEEP SUPERVISION OF MELODIC CONTOUR

### 3.1 Disentangling pitch from timbre

$$\mathcal{L}(\mathbf{x}_2, \mathcal{P}) = - \sum_{(t, k_1) \in \mathcal{P}} \log \sigma \left( \sum_{k_2} \mathbf{x}_2(t, k_1, k_2) \right) \quad (9)$$

### 3.2 Joint supervision

## 4. SINGLE-INSTRUMENT CLASSIFICATION

### 4.1 Experimental design

In order to train the proposed algorithms, we used MedleyDB v1.1. [1], a dataset of 122 multitracks annotated with instrument activations as well as melodic  $f_0$  curves when present. We extracted the monophonic stems corresponding to a selection of eight pitched instruments (see Figure 1). Stems with leaking instruments in the background were discarded. The evaluation set consists of 120 recordings of solo music collected by Joder et al. [2]. We discarded recordings with extended instrumental techniques, since they are under-represented in MedleyDB.

## 4.2 Results

## 5. POLYPHONIC CLASSIFICATION

### 5.1 Experimental design

### 5.2 Results

## 6. CONCLUSIONS

## 7. REFERENCES

- [1] Rachel Bittner, Justin Salamon, Mike Tierney, Matthias Mauch, Chris Cannam, and Juan Bello. Medleydb: a multitrack dataset for annotation-intensive mir research. *International Society for Music Information Retrieval Conference*, 2014.
- [2] Cyril Joder, Slim Essid, and Gaël Richard. Temporal integration for audio classification with application to musical instrument classification. *IEEE Transactions on Audio, Speech and Language Processing*, 17(1):174–186, 2009.
- [3] Diederik P. Kingma and Jimmy Lei Ba. Adam: a Method for Stochastic Optimization. *International Conference on Learning Representations*, pages 1–13, 2015.
- [4] Brian McFee, Matt McVicar, Colin Raffel, Dawen Liang, Oriol Nieto, Eric Battenberg, Josh Moore, Dan Ellis, Ryuichi Yamamoto, Rachel Bittner, Douglas Repetto, Petr Viktorin, Joo Felipe Santos, and Adrian Holovaty. librosa: 0.4.1. zenodo. 10.5281/zenodo.18369, October 2015.