



# A Model You Can Hear: Audio Identification with Playable Prototypes

Romain Loiseau<sup>1, 2</sup>

Baptiste Bouvier<sup>3</sup>

Yann Teytaut<sup>3</sup>

Elliot Vincent<sup>1, 4</sup>

Mathieu Aubry<sup>1</sup>

Loïc Landrieu<sup>2</sup>

<sup>1</sup>LIGM, ENPC

<sup>2</sup>LASTIG, IGN

<sup>3</sup>STMS, IRCAM

<sup>4</sup>Inria and DIENS



## Motivation

### Context:

- Recent methods often rely on representations in high-dimensional abstract space. Those methods perform well but are difficult to interpret.

### Contributions:

- We adapt the **transformation-invariant clustering** paradigm for audio in both supervised and unsupervised settings.
- We provide an audio identification model based on **prototypical sounds** that can be heard directly.
- Our model reaches **state-of-the-art** results for **audio classification and clustering** tasks while remaining easily **interpretable**.

## Method

For each **input audio clip**  $x$  of **label**  $y$  characterized by its **log-spectrogram** with  $T$  time steps and  $F$  frequency bins, we define the following reconstruction model and losses for the  $k$ -th **prototype**. **Deep Transformation-Invariant Prototyping**:

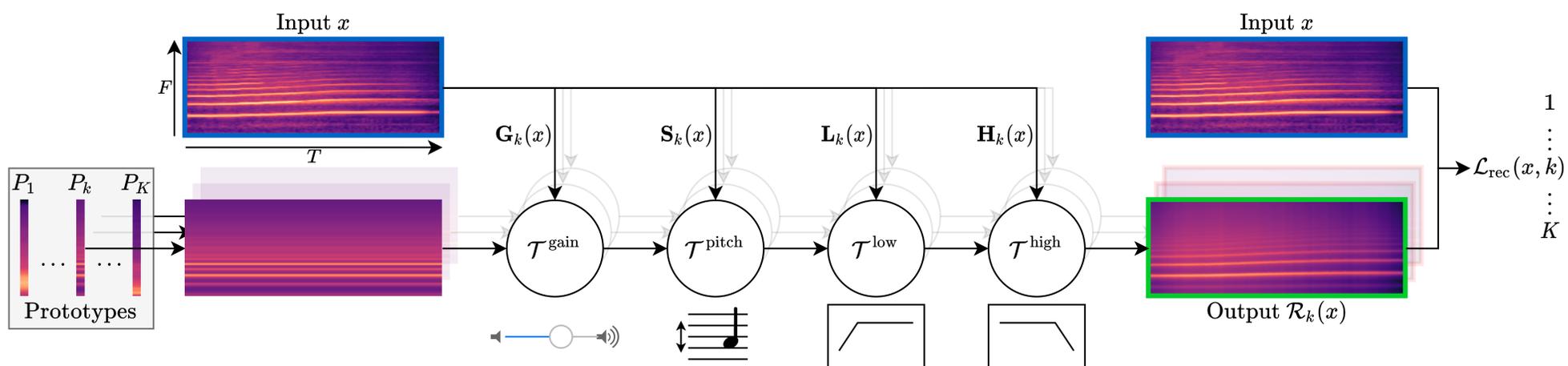
$$\mathcal{R}_k(x)[t] = \mathcal{T}_{\mathbf{H}_k(x)[t]}^{\text{high}} \circ \mathcal{T}_{\mathbf{L}_k(x)[t]}^{\text{low}} \circ \mathcal{T}_{\mathbf{S}_k(x)[t]}^{\text{pitch}} \circ \mathcal{T}_{\mathbf{G}_k(x)[t]}^{\text{gain}}(P_k)$$

### Supervising Reconstruction and Cross-Entropy:

$$\mathcal{L}_{\text{rec}}(x, k) = \frac{1}{T} \sum_{t=1}^T \|x[t] - \mathcal{R}_k(x)[t]\|^2$$

$$\mathcal{L}_{\text{ce}}(x, y) = -\log \left( \exp(-\beta \mathcal{L}_{\text{rec}}(x, y)) / \sum_{k=1}^K \exp(-\beta \mathcal{L}_{\text{rec}}(x, k)) \right)$$

## Method Overview



Given an **input sound**, we predict for each prototype a *gain*, a *pitch* shift, as well as *low* and *high frequency filters* at each timestamp to generate the **output**. Prototypes and transformations are learned jointly using a reconstruction loss in either a supervised or unsupervised setting.

## Audio Identification Results

	OA	AA
<b>SOL [1,2]</b>		
Autoencoder + K-means	28.7	12.3
† APNet [7] + K-means	<b>37.3</b>	<b>18.2</b>
Ours w/o supervision	34.5	15.4
<b>LibriSpeech [6]</b>		
Autoencoder + K-means	11.0	11.1
† APNet [7] + K-means	36.3	36.4
Ours w/o supervision	<b>48.6</b>	<b>49.5</b>

### Clustering Results. Clustering performances on the test sets.

†: Note that APNet requires labels at training time.

$$\mathcal{L}_{\text{clustering}}(x) = \min_{k=1}^K \mathcal{L}_{\text{rec}}(x, k)$$

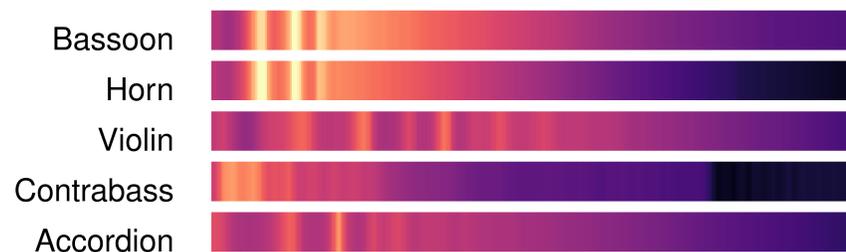
	OA	AA	$\mathcal{L}_{\text{rec}}$
<b>SOL [1,2]</b>			
Direct Classification	97.8	94.8	—
APNet [7]	95.3	91.3	<b>0.1</b>
Ours w supervision	<b>99.3</b>	<b>95.8</b>	2.6
<b>LibriSpeech [6]</b>			
Direct Classification	99.4	99.5	—
APNet [7]	97.8	97.8	<b>0.2</b>
Ours w supervision	<b>99.9</b>	<b>99.9</b>	2.6

### Classification Results. Accuracy and reconstruction error.

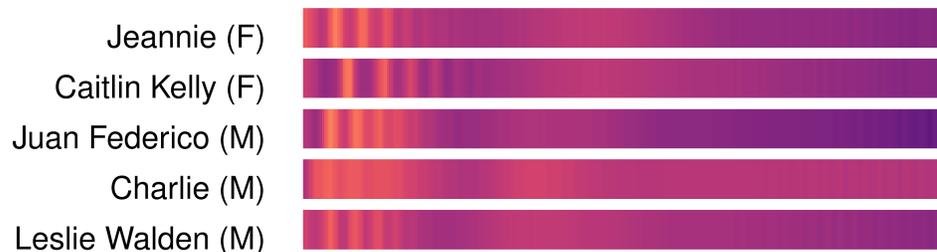
$$\mathcal{L}_{\text{classif}}(x, y) = \mathcal{L}_{\text{rec}}(x, y) + \lambda_{\text{ce}} \mathcal{L}_{\text{ce}}(x, y)$$

## Meaningful and Interpretable Prototypes

- Prototypes learn characteristics of their assigned class that go beyond simple harmonic components.
- The model captures elements of the **timbre**, such as their spectral envelopes under various conditions: notes, techniques, words, etc.
- Such results pave the way for **instrumental or vocal timbre transfers** as natural applications.



**Prototypes learned** on the instrument dataset SOL [1,2]



**Prototypes learned** on the speech dataset LibriSpeech [6]

## Bibliography

[1] Guillaume Ballet *et al.* Journées d'Informatique Musicale, 1999 [2] Carmine Emanuele Cella *et al.* ICMC, 2020 [3] Romain Loiseau *et al.* 3DV, 2021 [4] Vincent Lostanlen *et al.* DLFM, 2018 [5] Tom Monnier *et al.* NeurIPS, 2020 [6] Vassil Panayotov *et al.* ICASSP, 2015 [7] Pablo Zinemanas *et al.* Electronics, 2021