

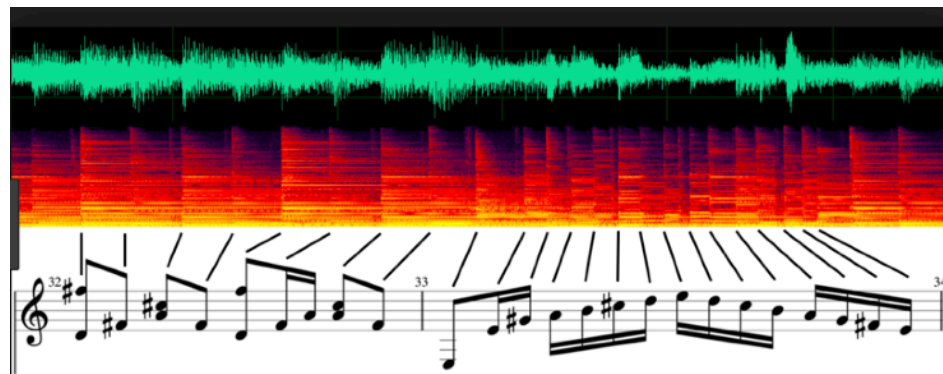
## Introduction

### – Goal :

- obtain a transcription of the content of an audio signal
  - which notes ? which chords ? which instruments ?
  - equivalent to Automatic Speech Recognition
- difficulty :
  - all instruments are super-imposed in the audio signal
- sub-tasks :
  - only transcribe the notes /  $f_0$  (in Hz)
  - only transcribe the dominant melody

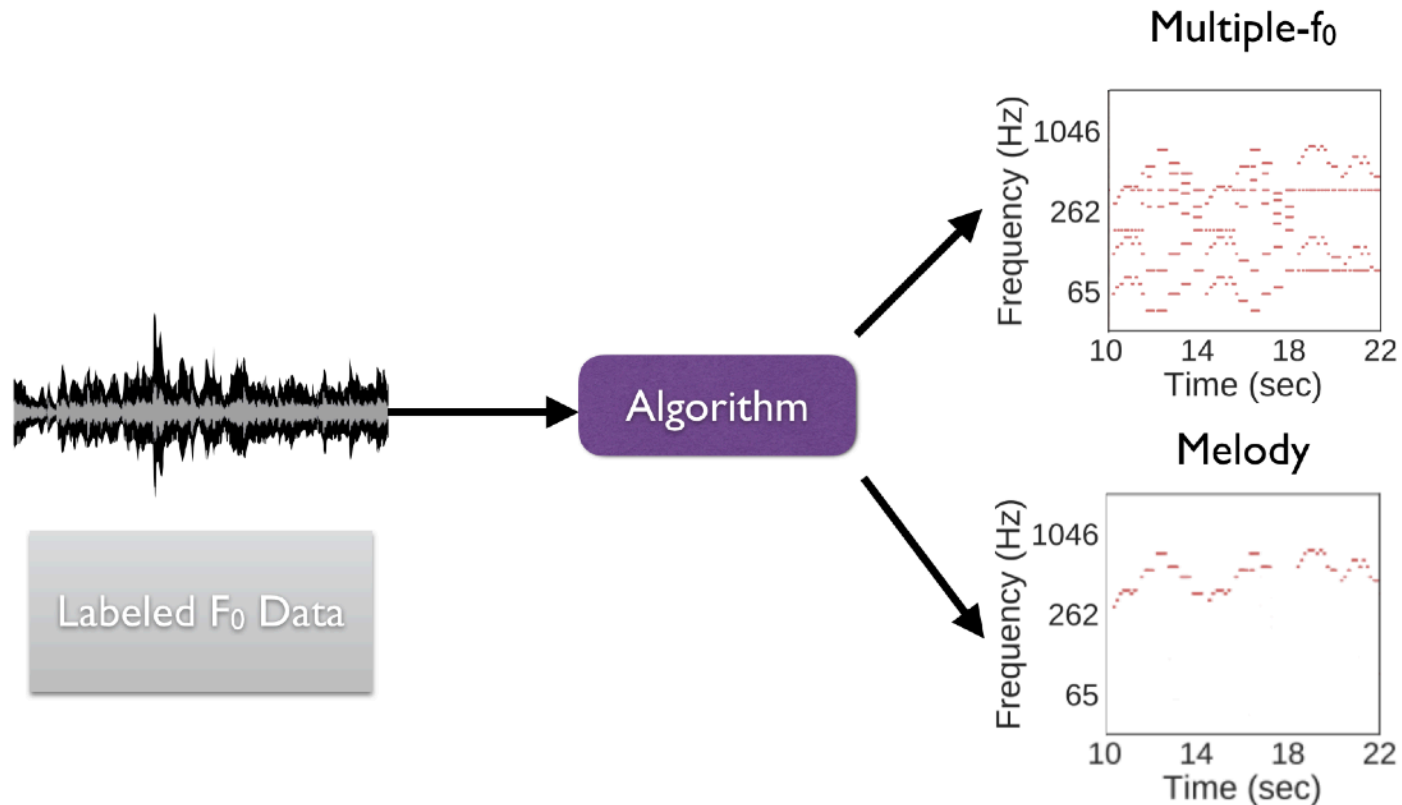
### – How can we do that ?

- many methods proposed in the past : sinusoidal model, NMF, PLCA, ... ,
- large improvement since 2017 using ... deep learning

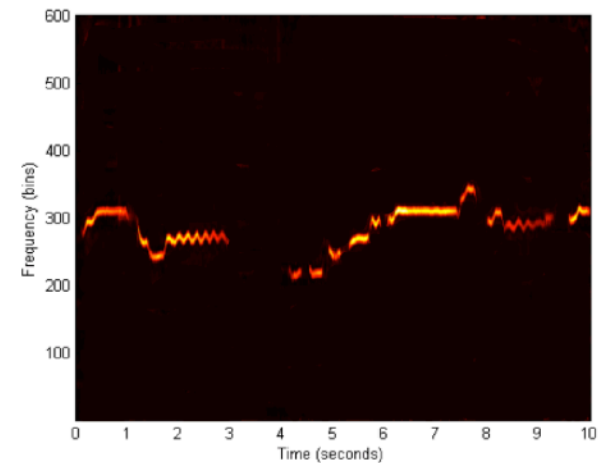
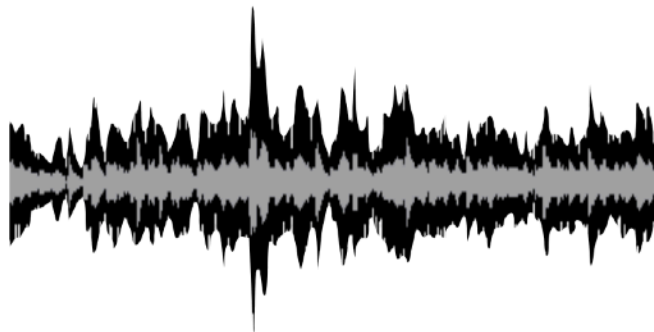


# 2017 → Music Multiple- $f_0$ Estimation using Harmonic-CQT

## Introduction

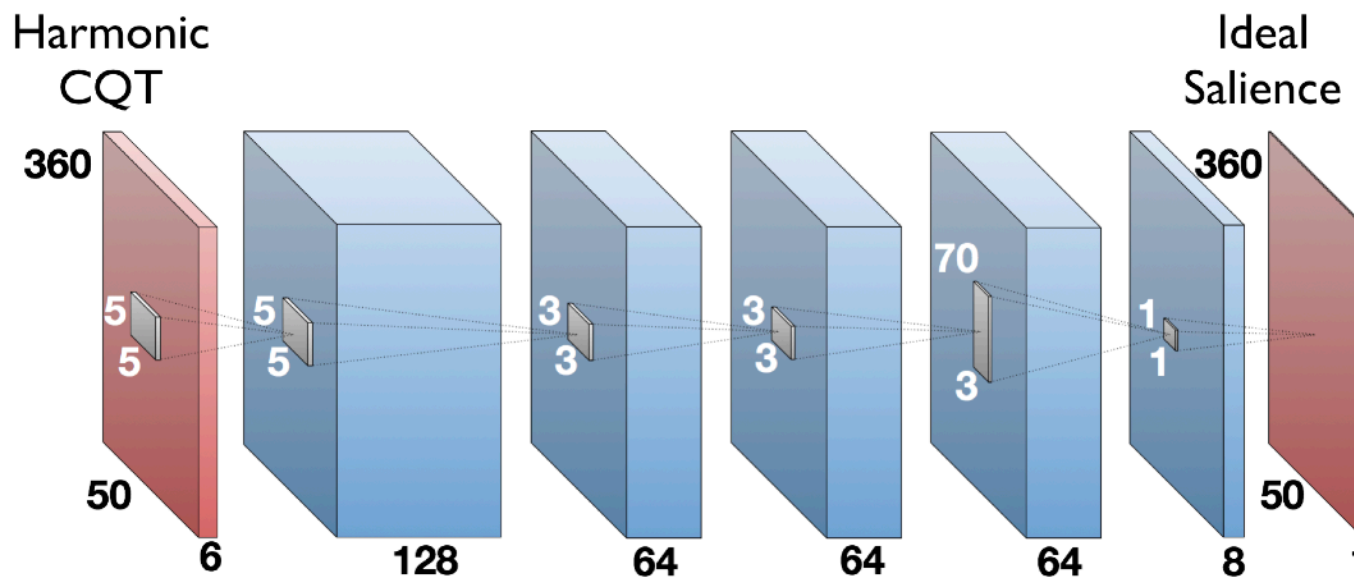


## Learning a "Saliency" Representation



## Using Convolutional Neural Network

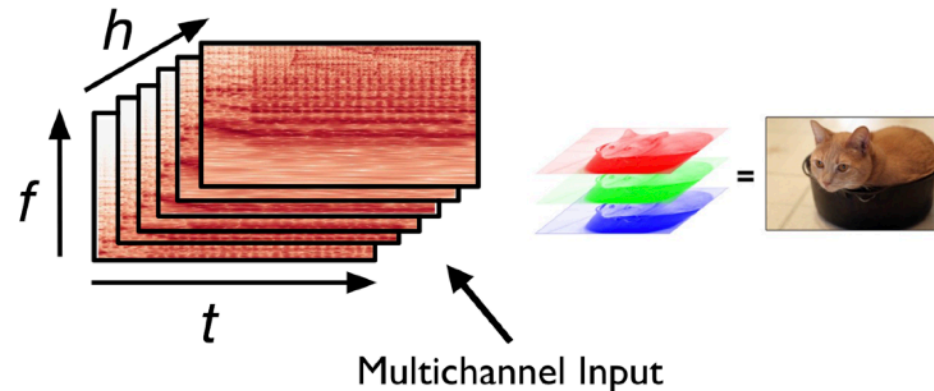
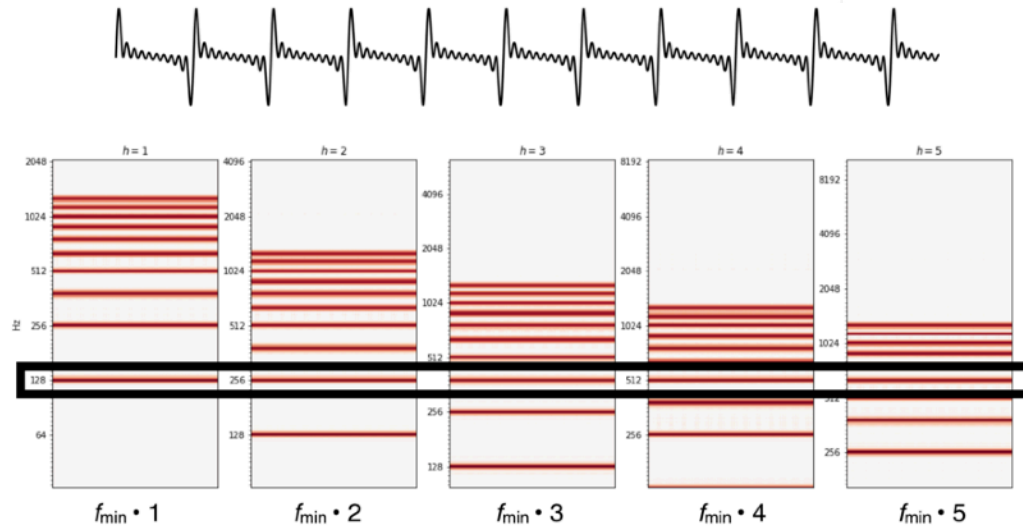
- **Inputs :**
  - magnitude spectrogram **with depth** (Harmonic CQT)
- **Outputs :**
  - image of the same size of input but with only saliency information



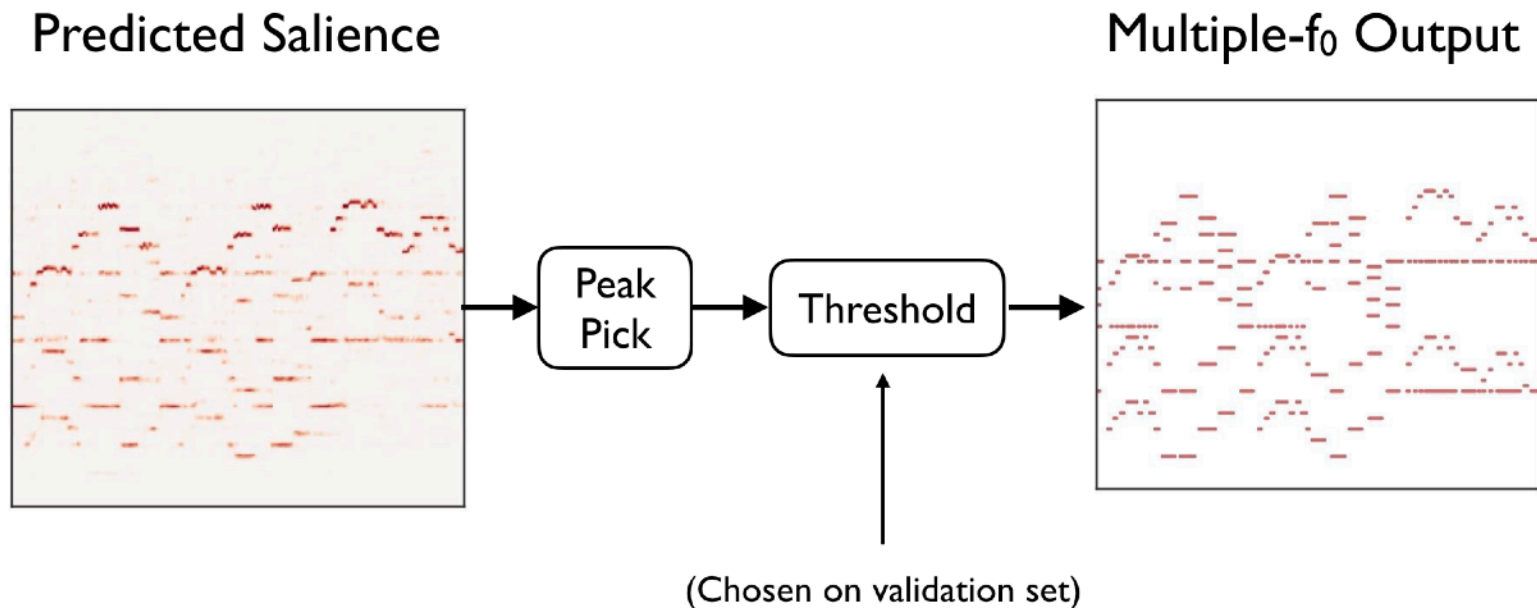
$$L(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

# 2017 → Music Multiple- $f_0$ Estimation using Harmonic-CQT

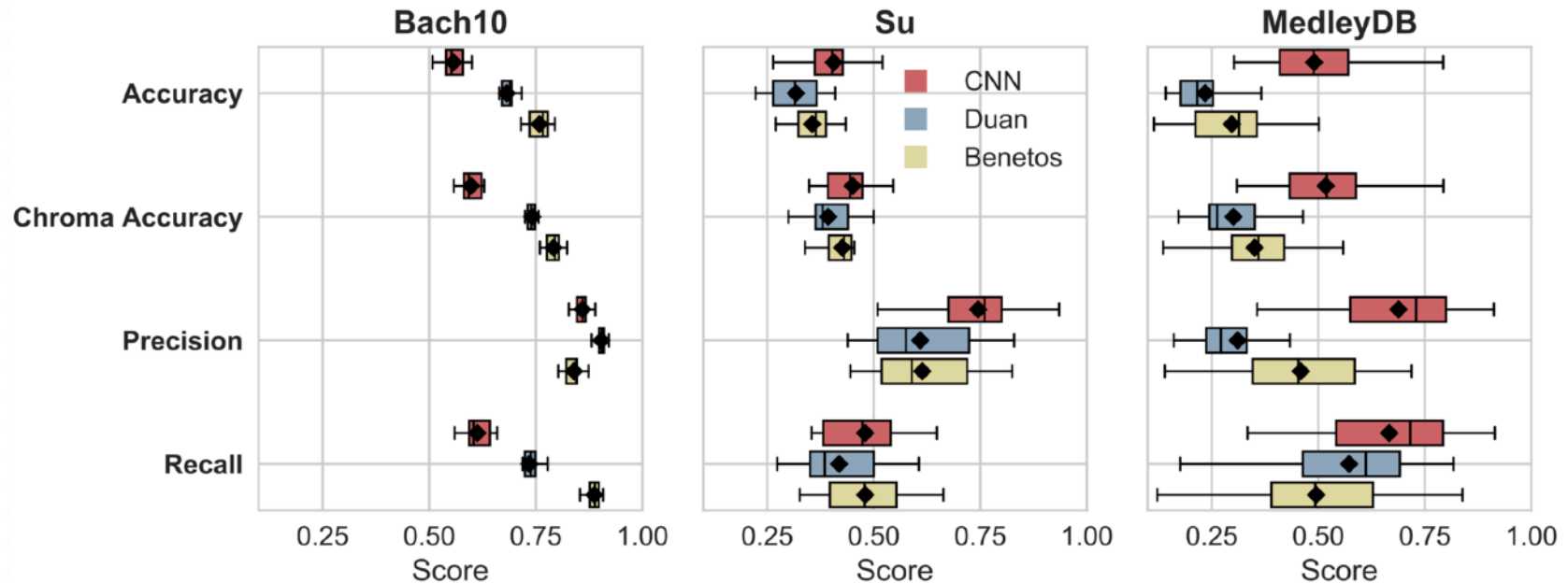
## Input Representation: a Harmonic CQT



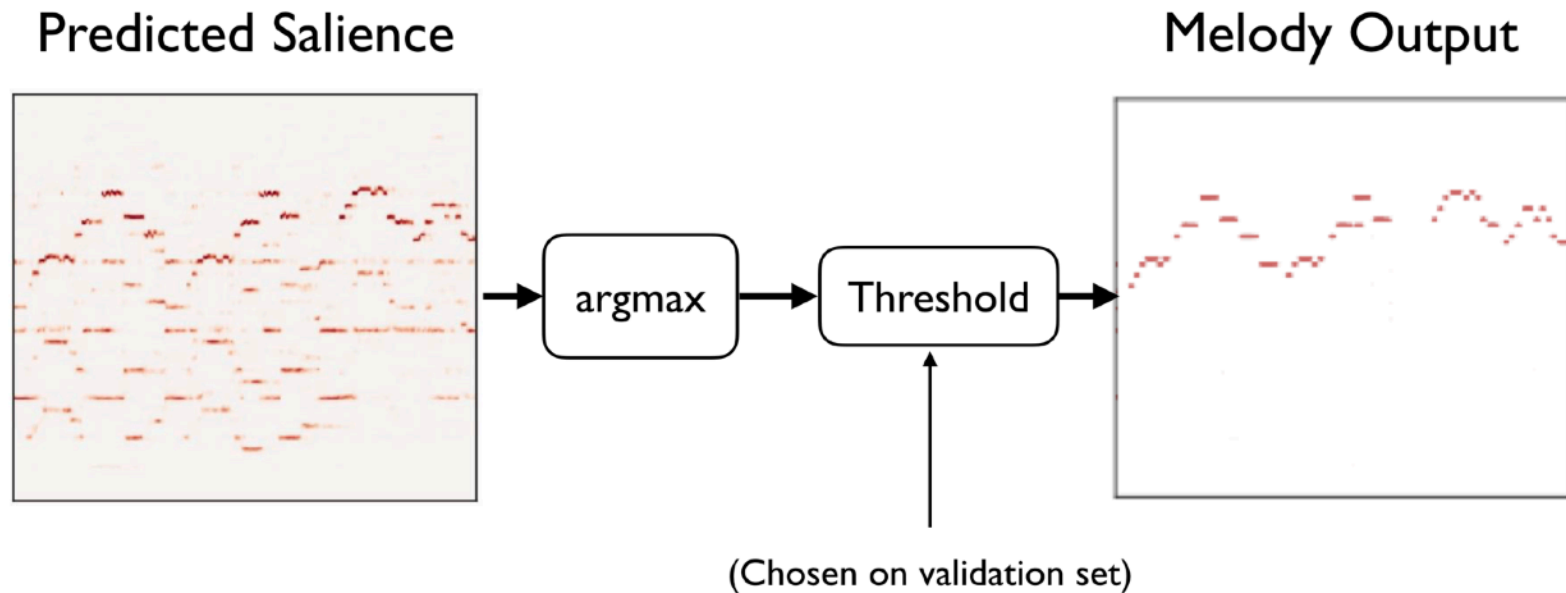
## Multiple- $f_0$ Estimation



## Multiple- $f_0$ Estimation: Results

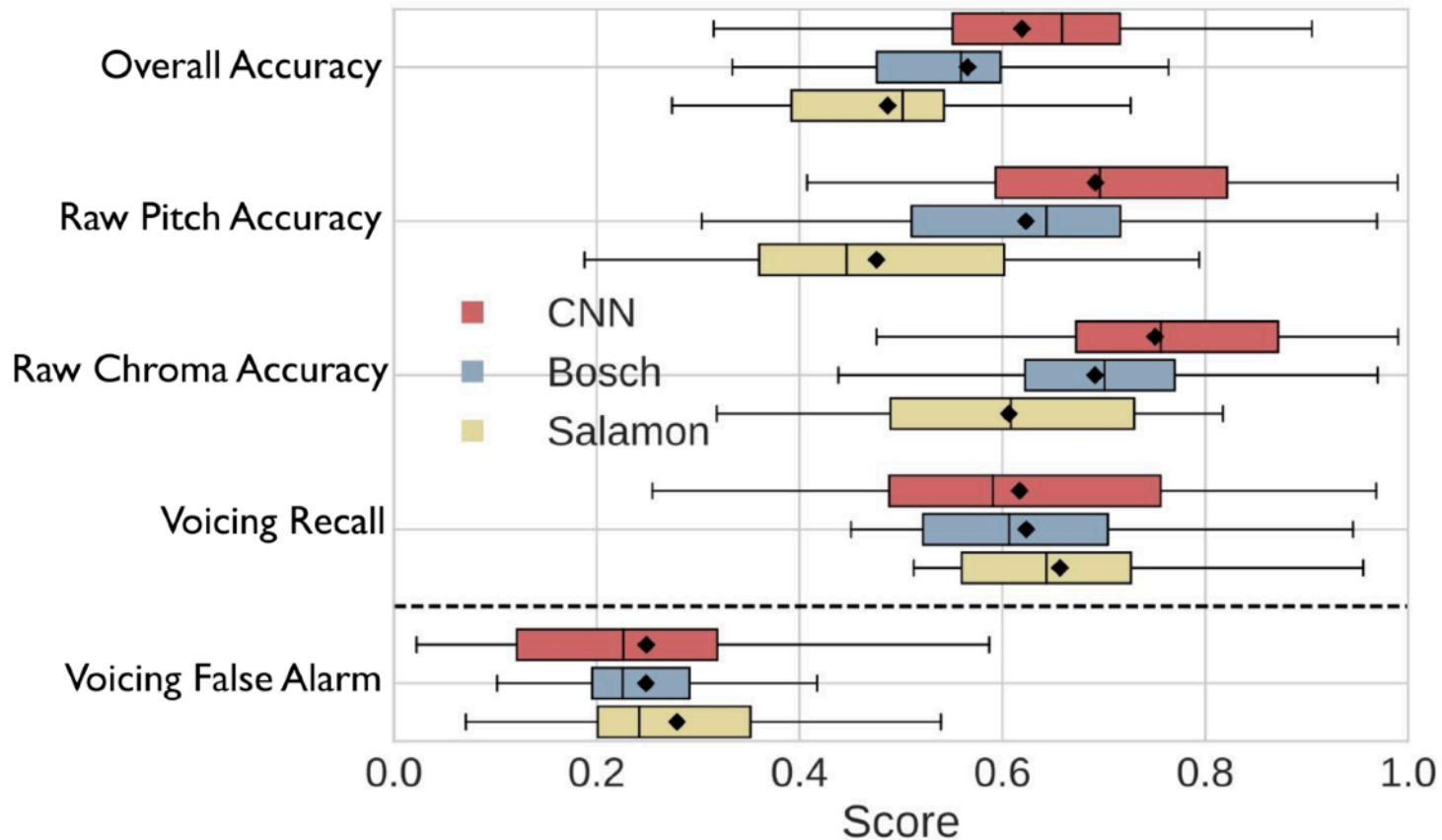


## Dominant Melody Estimation





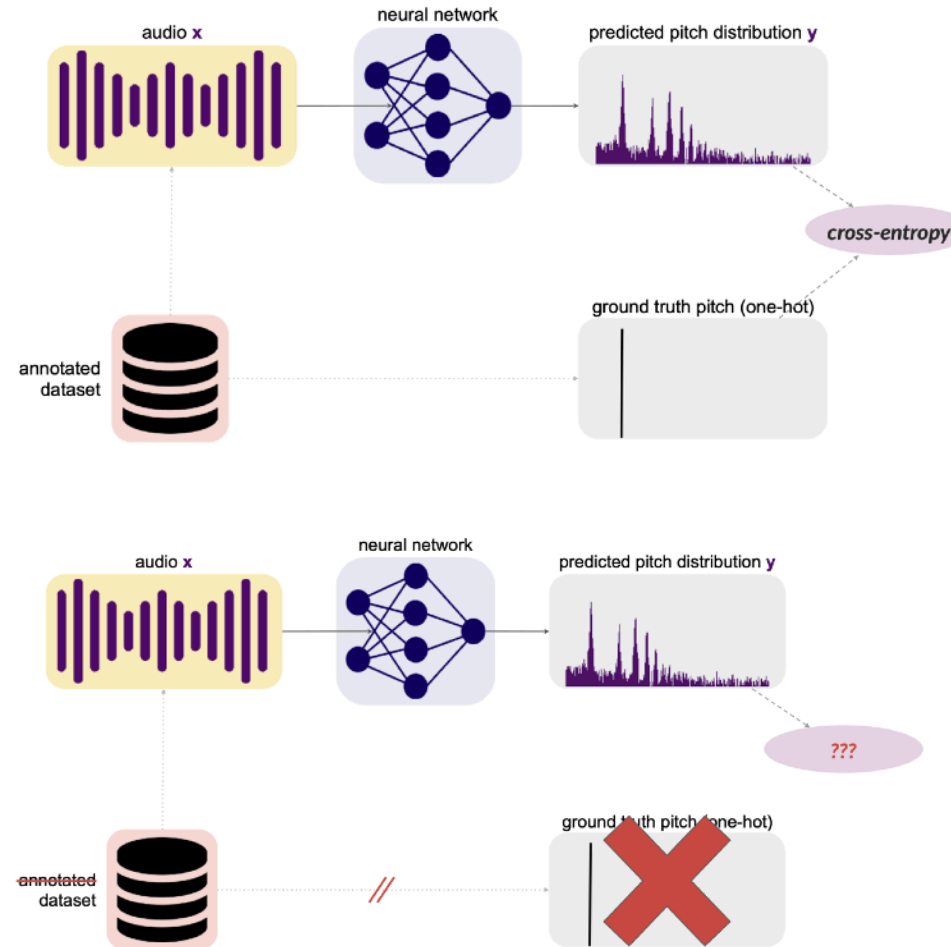
## Dominant Melody Estimation: Results



# 2023 → PESTO: Self-Supervised Pitch Estimation: Equi-variance

## Introduction

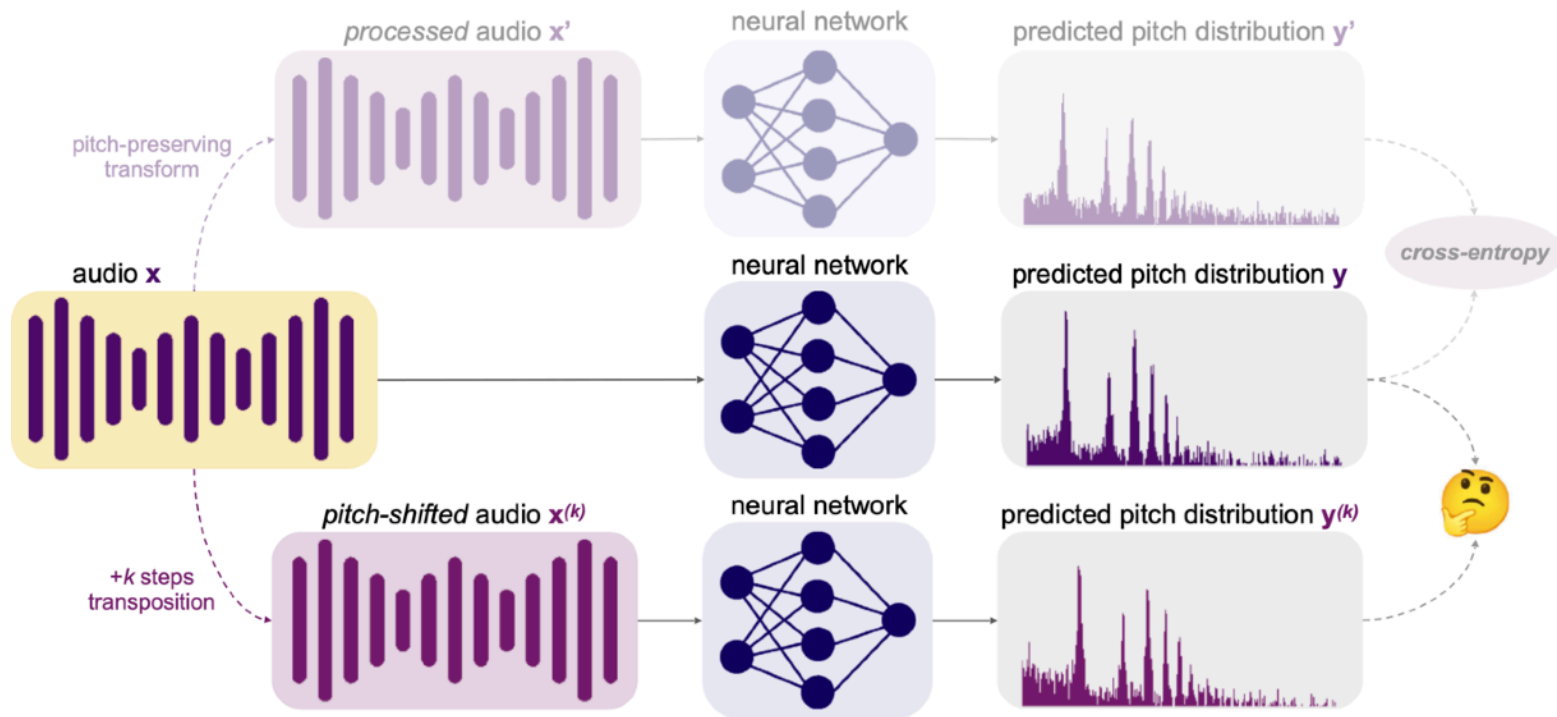
- **Traditional** way of training a neural network
  - **Supervised** training (compare prediction to ground-truth)
- **Our goal:**
  - train a neural network for pitch estimation **without** ground-truth annotations
  - = **Self-Supervised-Learning** (SSL)



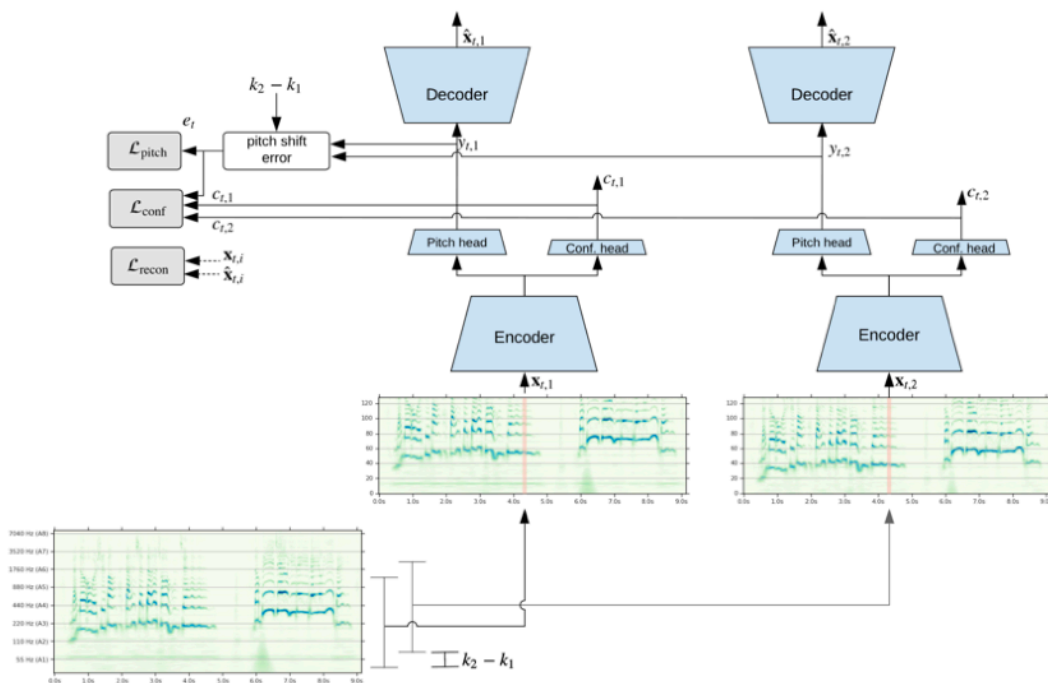
## Self-Supervised-Learning

### – Siamese Network, SimCLR,

- given two views  $x, x'$  (augmented versions) of the same data  $x$  trains a neural network such that the corresponding outputs  $y, y'$  are similar  $\Rightarrow$  feature learning
- can we make the output equivariant instead of invariant ?



# 2019 → SPICE: Self-supervised Pitch Estimation



equivariance loss

$$e_t = |(y_{t,1} - y_{t,2}) - \sigma(k_{t,1} - k_{t,2})|$$

with

$$\mathcal{L}_{pitch} = \frac{1}{T} \sum_i h(e_i)$$

$$h(x) = x^2/2 \text{ if } |x| \leq \tau$$

$$h(x) = \tau^2/2 + \tau(|x| - \tau) \text{ if } |x| > \tau$$

$$\sigma = \frac{1}{B[\log_2(f_{max}/f_{min})]}$$

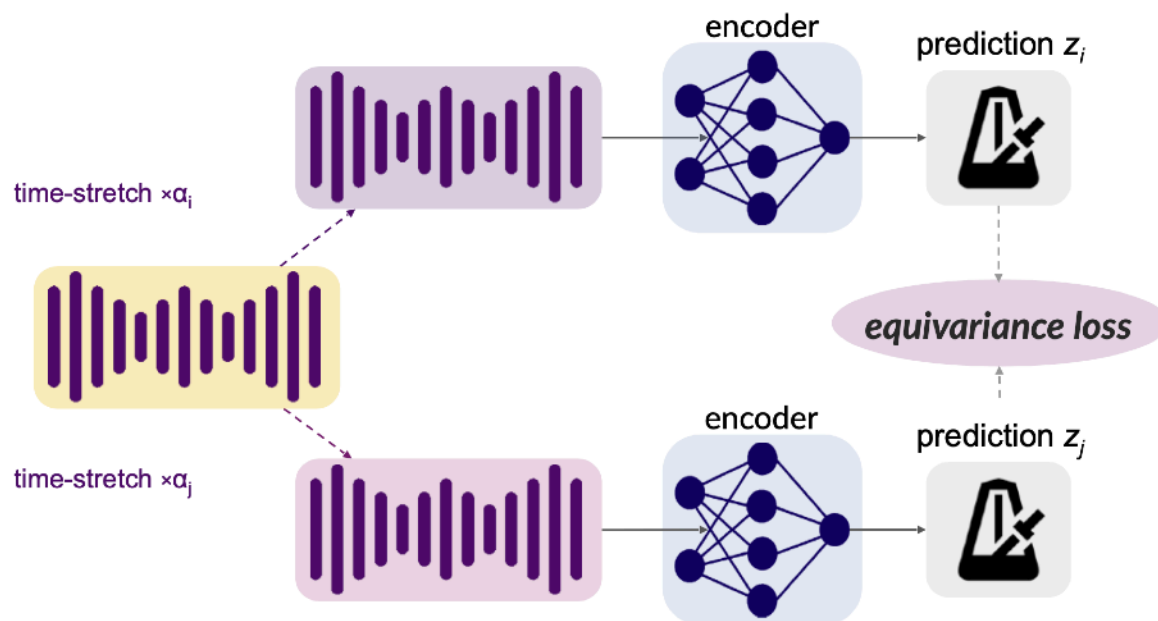
reconstruction loss

$$\mathcal{L}_{recon} = \frac{1}{T} \sum_i ||x_{t,i} - \hat{x}_{t,i}||_2^2$$

total loss

$$\mathcal{L} = w_{pitch}\mathcal{L}_{pitch} + w_{recon}\mathcal{L}_{recon}$$

# 2022 → Equivariant Self-Supervision for Tempo Estimation

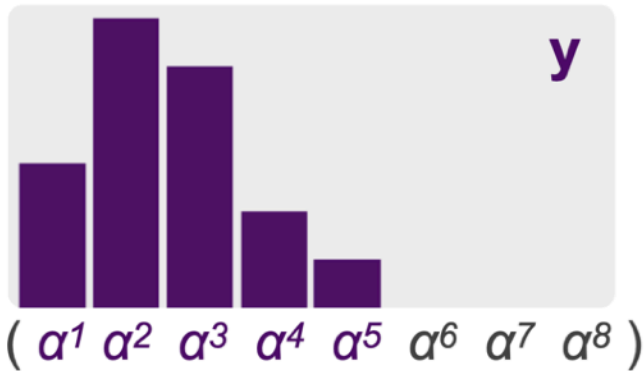


$$\mathcal{L} = \left| \frac{z_i}{z_j} - \frac{\alpha_i}{\alpha_j} \right|$$

## Self-Supervised-Learning

### – Our proposal: Class-based equivariance loss

- scalar product between the softmax output by a geometric series  $\alpha$



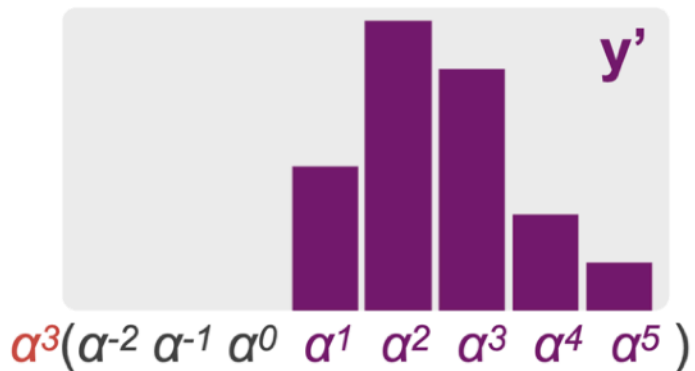
- Define  $\mathbf{a} = (\alpha, \alpha^2, \dots, \alpha^d)^\top, \alpha > 0$
- Compute the **scalar** products  $\mathbf{a}^\top \mathbf{y}$  and  $\mathbf{a}^\top \mathbf{y}'$
- If  $\mathbf{y}$  and  $\mathbf{y}'$  are equal up to a shift of  $k$ , then

$$\mathbf{a}^\top \mathbf{y}' = \alpha^k \mathbf{a}^\top \mathbf{y}$$

- **Equivariance** loss:

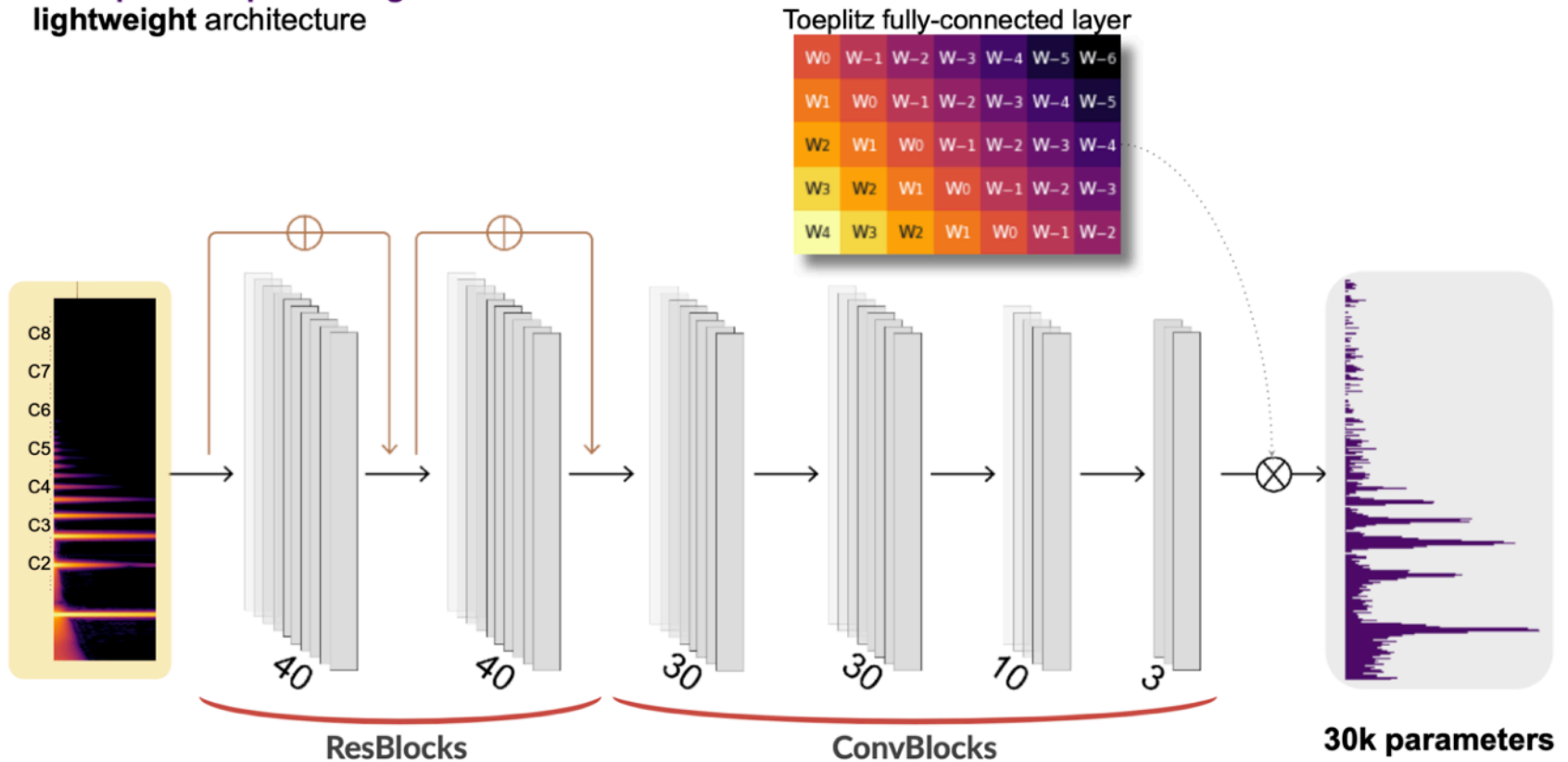
$$\mathcal{L}_{\text{equiv}}(\mathbf{y}, \mathbf{y}', k) = \left\| \frac{\mathbf{a}^\top \mathbf{y}'}{\mathbf{a}^\top \mathbf{y}} - \alpha^k \right\|$$

- translation of  $k$  between  $\mathbf{y}$  and  $\mathbf{y}'$   $\Rightarrow$  equivariance loss = 0



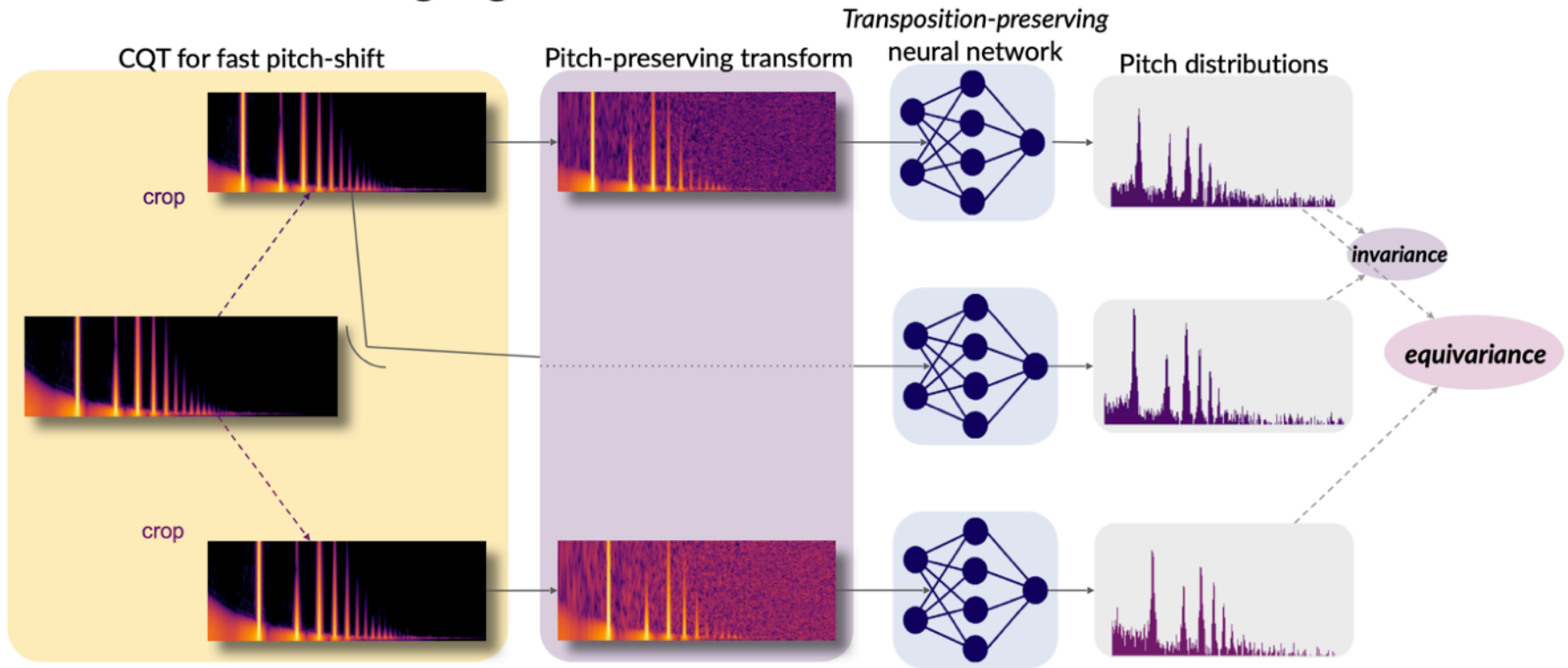
## Lightweight architecture

### Transposition-preserving lightweight architecture



# 2023 → PESTO: Self-Supervised Pitch Estimation: Equi-variance

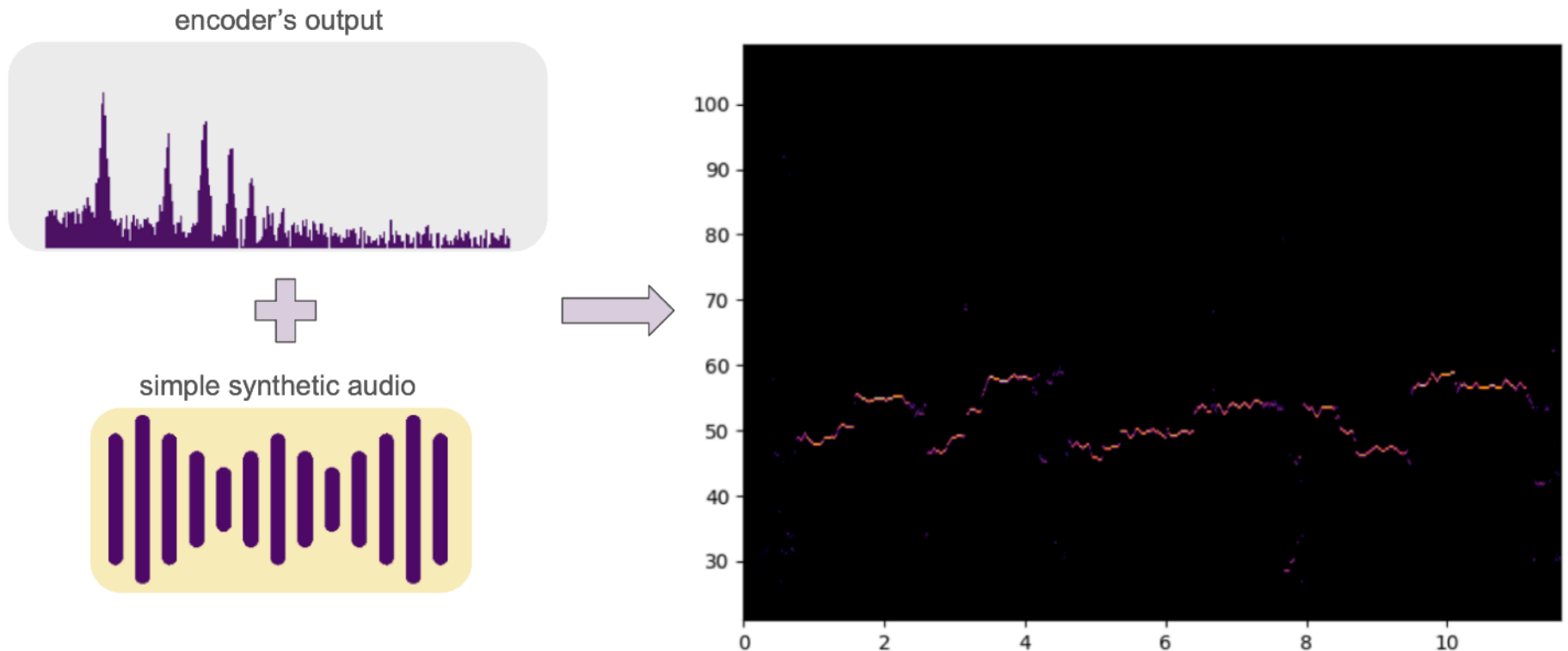
## Training data





# 2023 → PESTO: Self-Supervised Pitch Estimation: Equi-variance

## From relative pitch to absolute pitch



## Evaluation

Model	# params	Trained on	Raw Pitch Accuracy	
			<i>MIR-1K</i>	<i>MDB-stem-synth</i>
SPICE [19]	2.38M	private data	90.6%	89.1%
DDSP-inv [45]	-	<i>MIR-1K / MDB-stem-synth</i>	91.8%	88.5%
PESTO (ours)	28.9k	<i>MIR-1K</i>	<b>96.1%</b>	94.6%
PESTO (ours)	28.9k	<i>MDB-stem-synth</i>	93.5%	<b>95.5%</b>
CREPE [16]	22.2M	many (supervised)	<b>97.8%</b>	<b>96.7%</b>

- SOTA in self-supervised pitch estimation
- Strong generalization performances
- Extremely lightweight model