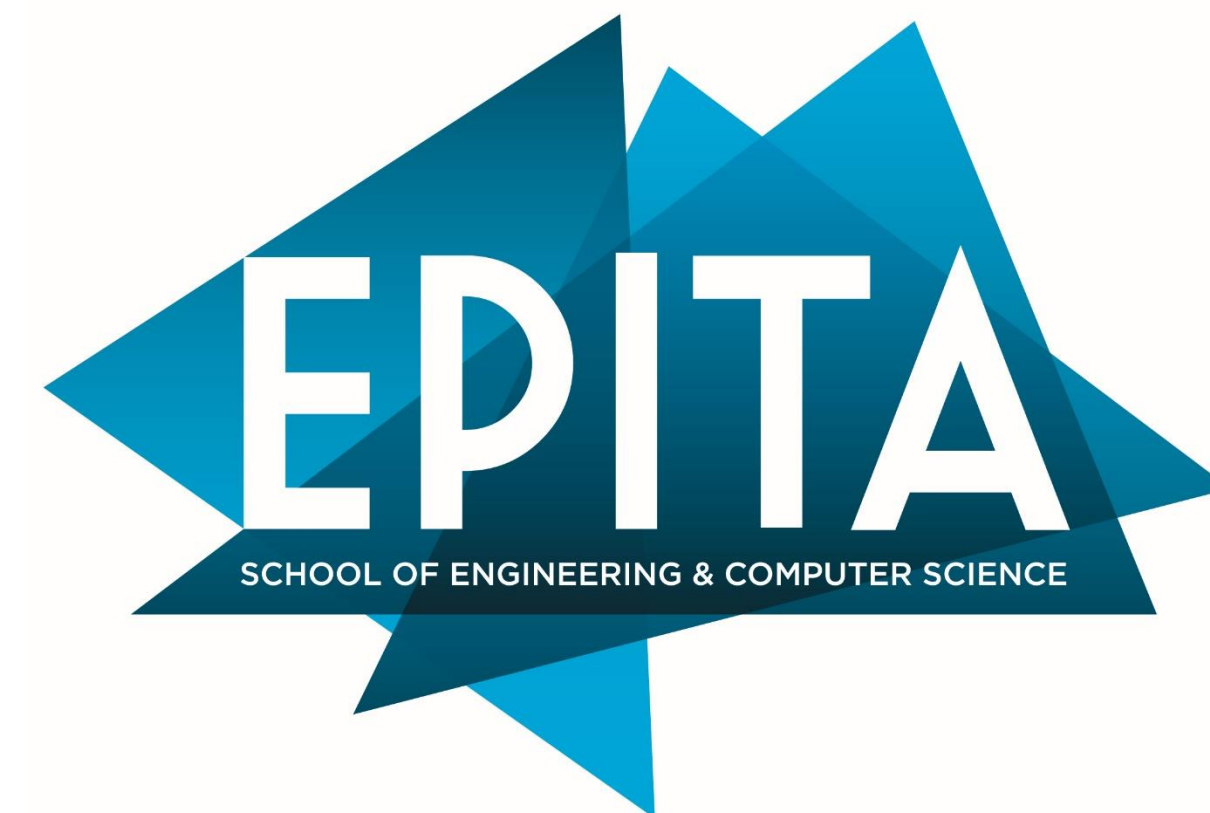




Arab Republic of Egypt
Ministry of Communications
and Information Technology

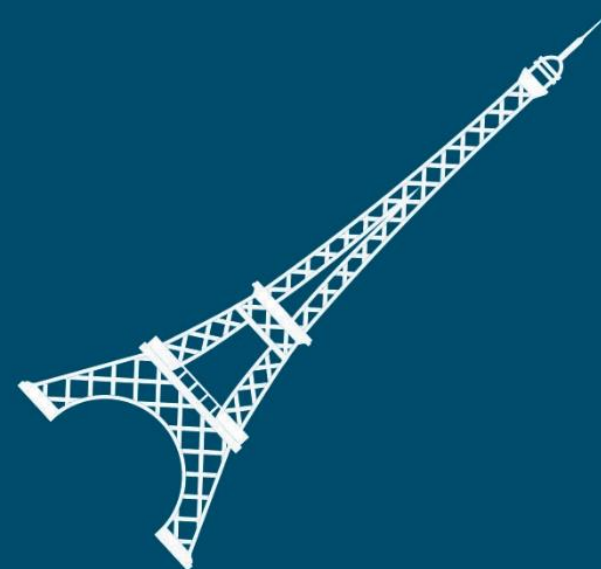


Information
Technology
Institute

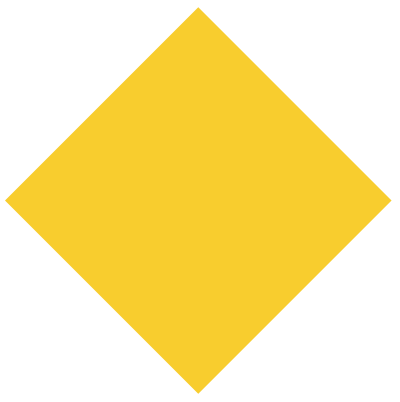


A.I. IN AUDIO & SIGNAL PROCESSING

Session 2: Deep learning for audio and speech processing



SESSION 2: DEEP LEARNING FOR AUDIO & SPEECH PROCESSING



Quick Summary

1. Approaches for feature learning

- a) Input représentations
- b) Filters shape
- c) Signal models
- d) Generative models
- e) How to visualize learnt représentations ?

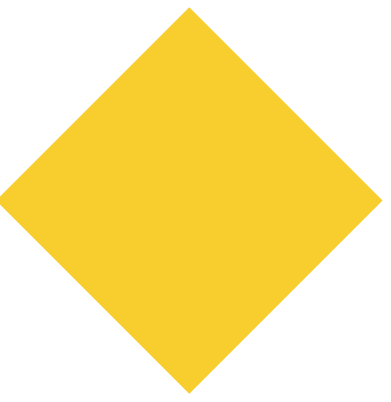
2. New learning paradigms

- a) Classification
- b) Auto-encoders & variational auto-encoder
- c) Metric learning
- d) Semi-supervised learning

DEEP LEARNING FOR AUDIO AND SPEECH PROCESSING.

Approaches for feature
learning

APPROACHES FOR FEATURE LEARNING



2D representation (time-frequency)

Representation

- spectrogram (STFT magnitude)
- Mel-gram
- Constant-Q-transform

Basic idea

Considering time-frequency representation as a 2-D image as an input of a Conv2D network

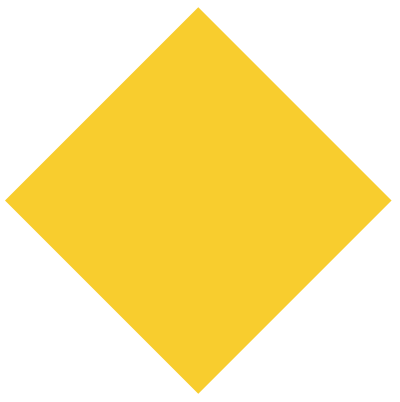
Problem

A time-frequency representation is not an image

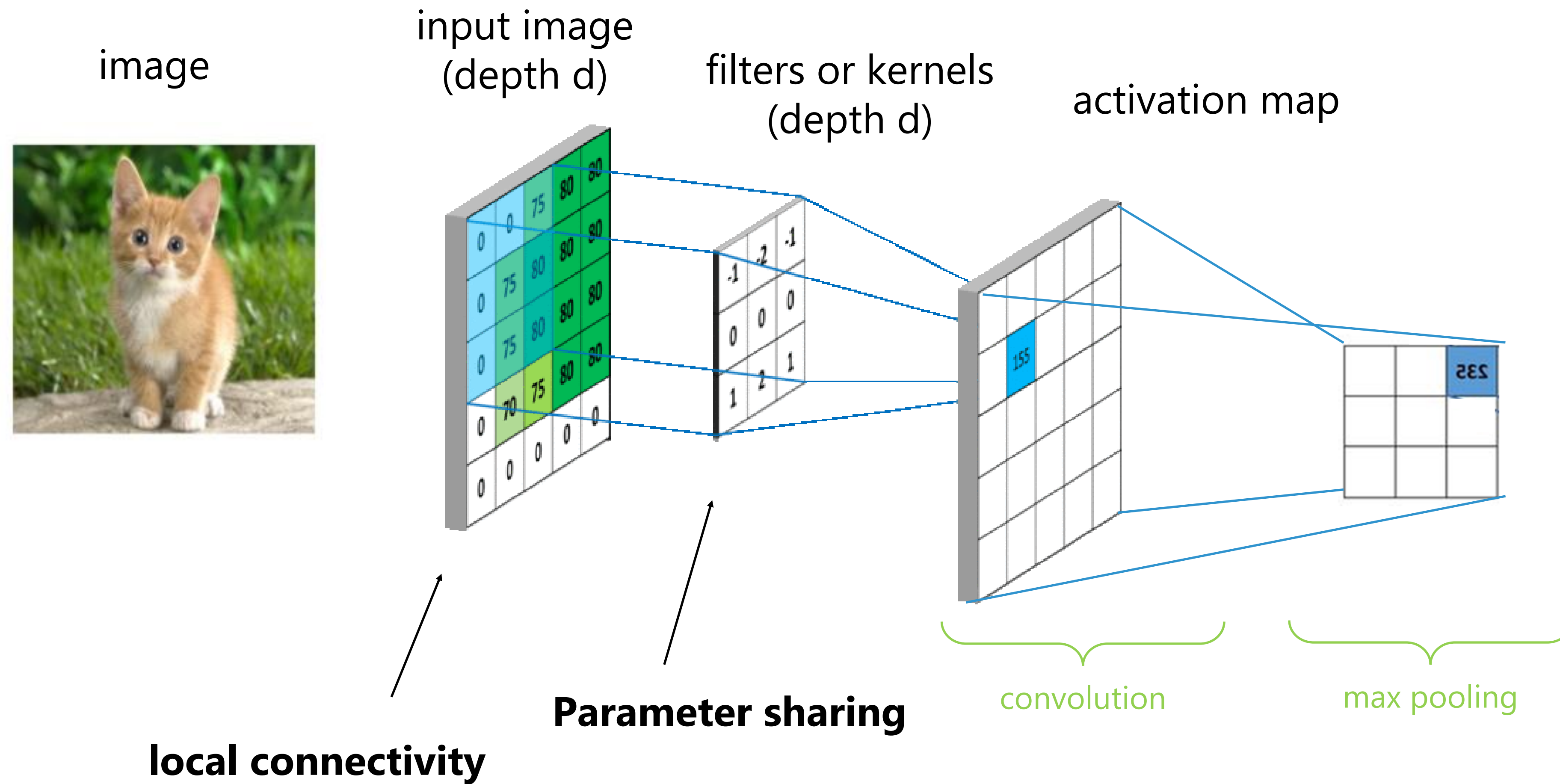
Research direction

Choice of 1st Conv2D layer filters

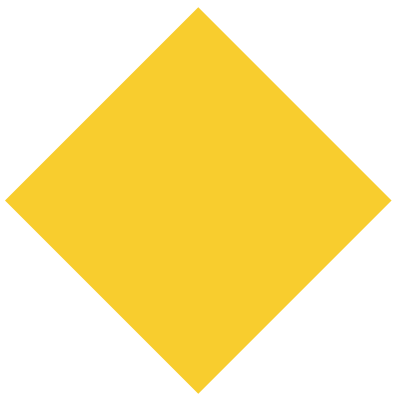
APPROACHES FOR FEATURE LEARNING



Concept of Conventional Neural Network (CNN or ConvNet)



APPROACHES FOR FEATURE LEARNING



2D representation (time-frequency)

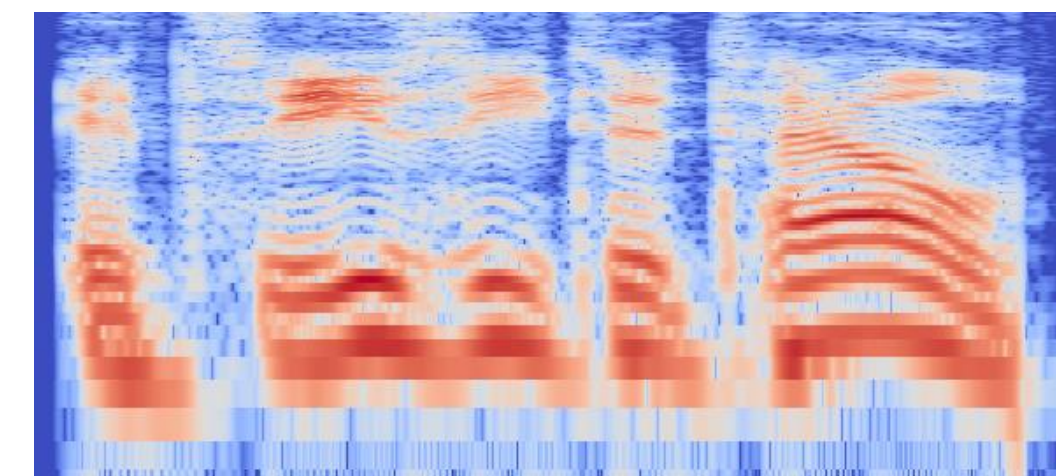
Natural images

- the 2 axis represent the same concept
- Whatever the position of an element is, it represents the same thing
 - spatial invariance
 - weight sharing in the 2 dimensions
- close pixels are often strongly correlated
- close and similar pixel usually belong to a same object

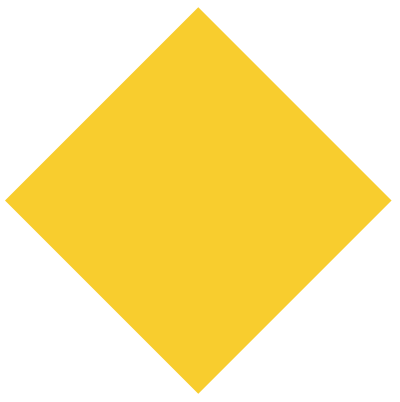


Harmonic sounds

- The 2 axis represent completely different things
- Properties of a sound event:
 - same signification whenever it is played
 - usually, different signification depending on its frequency
 - no invariance in frequency (even in log-scale)
- Frequencies of a same source are not distributed locally on spectrogram (sparsity)



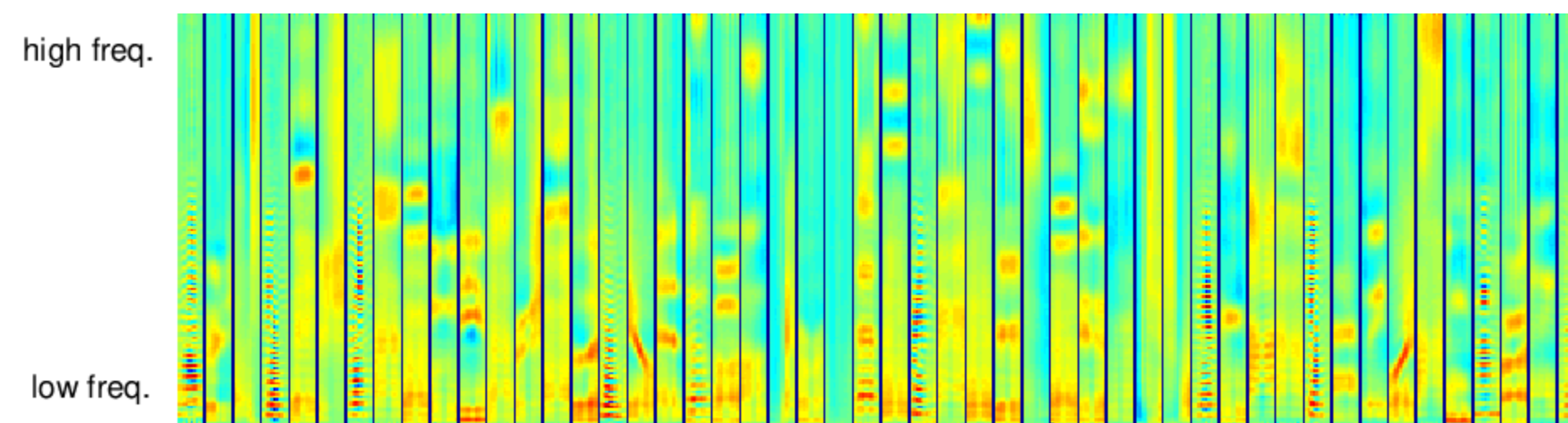
APPROACHES FOR FEATURE LEARNING



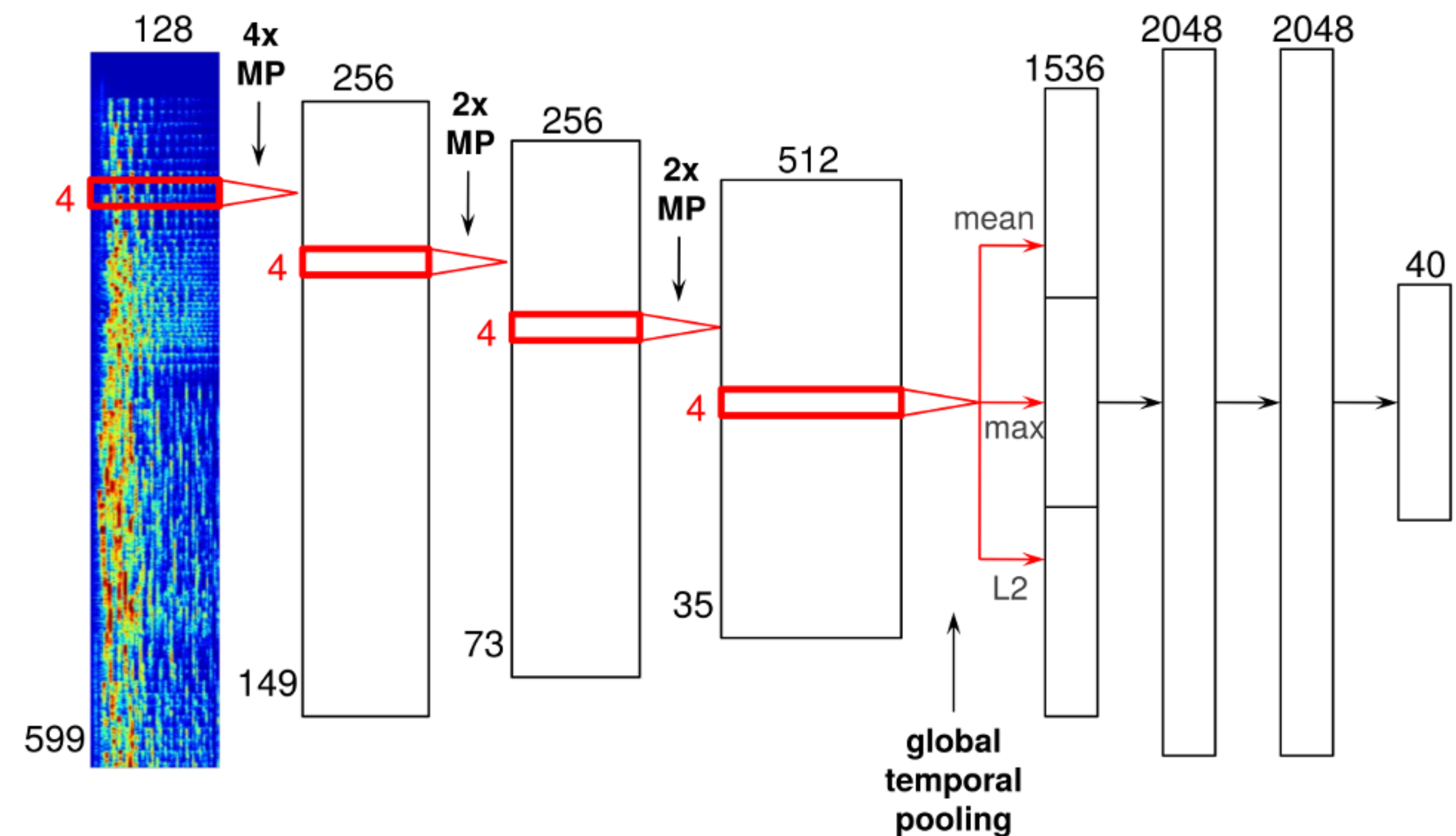
An adapted shape of 2D filters for 1st layer

1) Filter covering the entire frequency bandwidth

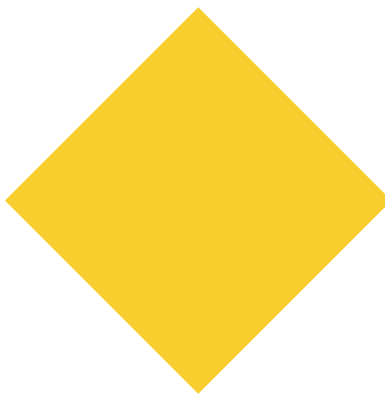
→ convolution only over time



Each column represents a « temporal receptive field » of a 1st layer basis in the spectrogram space



APPROACHES FOR FEATURE LEARNING

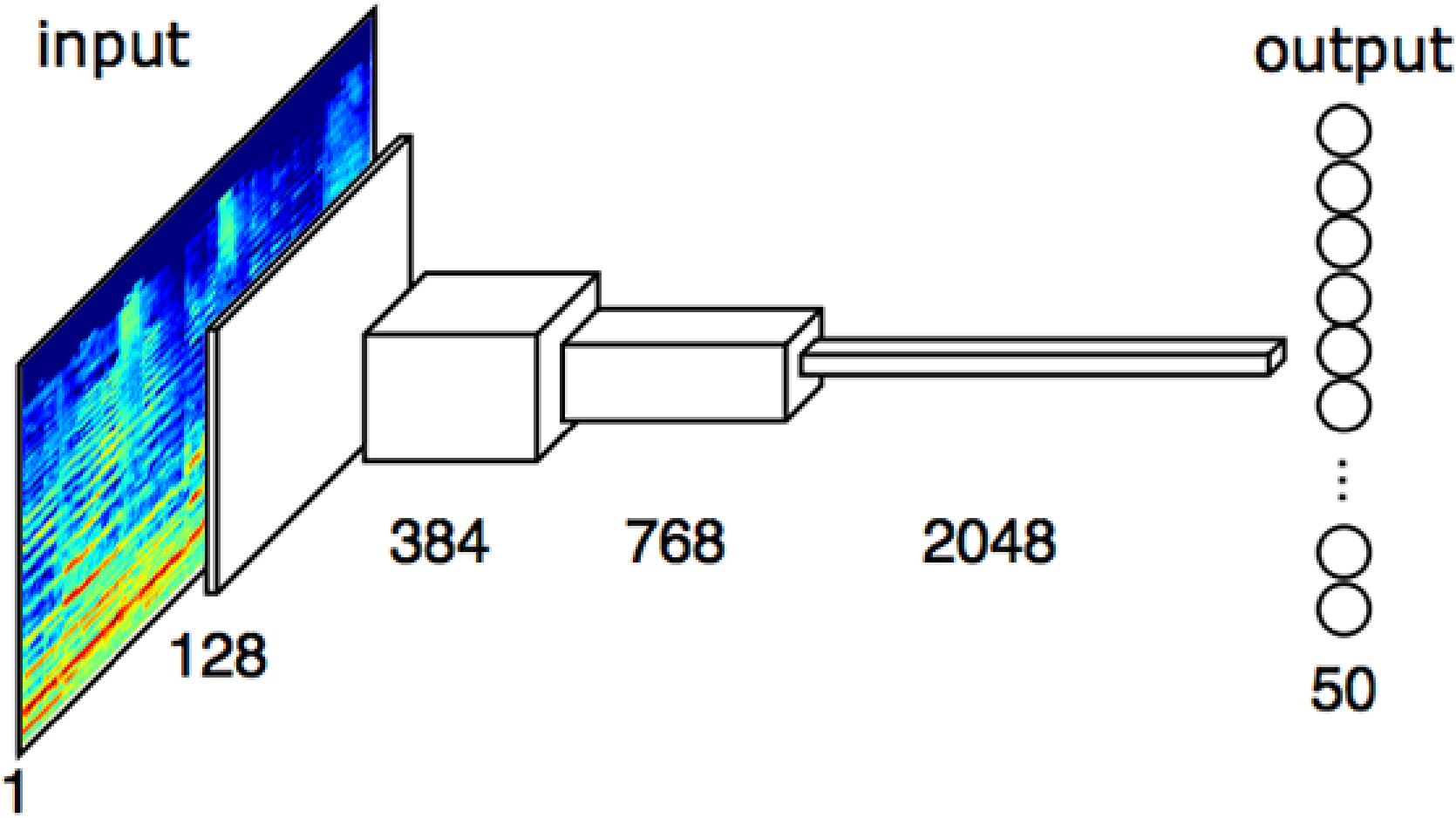


An adapted shape of 2D filters for 1st layer

2) 3x3 or 5x5 filters, used as in image processing

VGG-net on spectrum

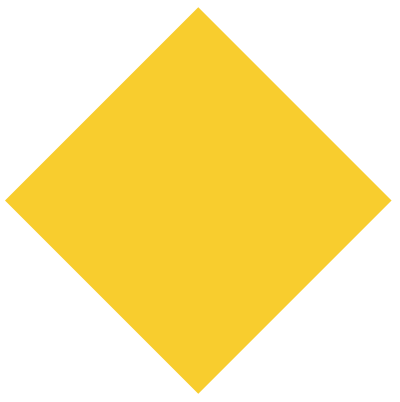
- for music automatic tagging (multi-label task)
- using STFT, MFCC and Mel-gram representation as input



FCN-5	FCN-6	FCN-7
Mel-spectrogram (input: $96 \times 1366 \times 1$)		
Conv $3 \times 3 \times 128$		
MP (2, 4) (output: $48 \times 341 \times 128$)		
Conv $3 \times 3 \times 256$		
MP (2, 4) (output: $24 \times 85 \times 256$)		
Conv $3 \times 3 \times 512$		
MP (2, 4) (output: $12 \times 21 \times 512$)		
Conv $3 \times 3 \times 1024$		
MP (3, 5) (output: $4 \times 4 \times 1024$)		
Conv $3 \times 3 \times 2048$		
MP (4, 4) (output: $1 \times 1 \times 2048$)		
.	Conv $1 \times 1 \times 1024$	Conv $1 \times 1 \times 1024$
	.	Conv $1 \times 1 \times 1024$
Output 50×1 (sigmoid)		

K. Choi, G. Fazekas and M. Sandler. Automatic tagging using deep convolutional neural networks. ISMIR 2016.
K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. 2015.

APPROACHES FOR FEATURE LEARNING

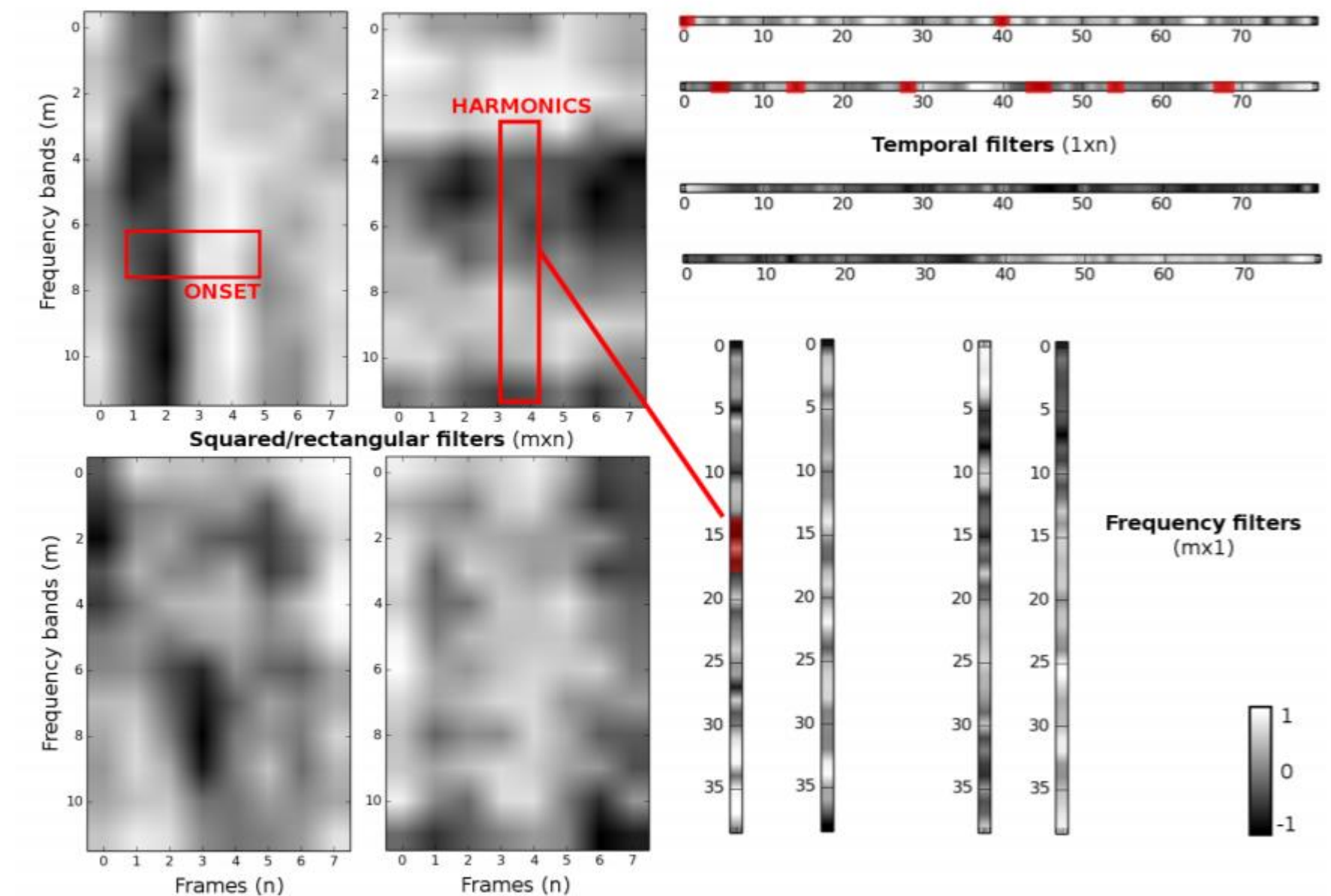
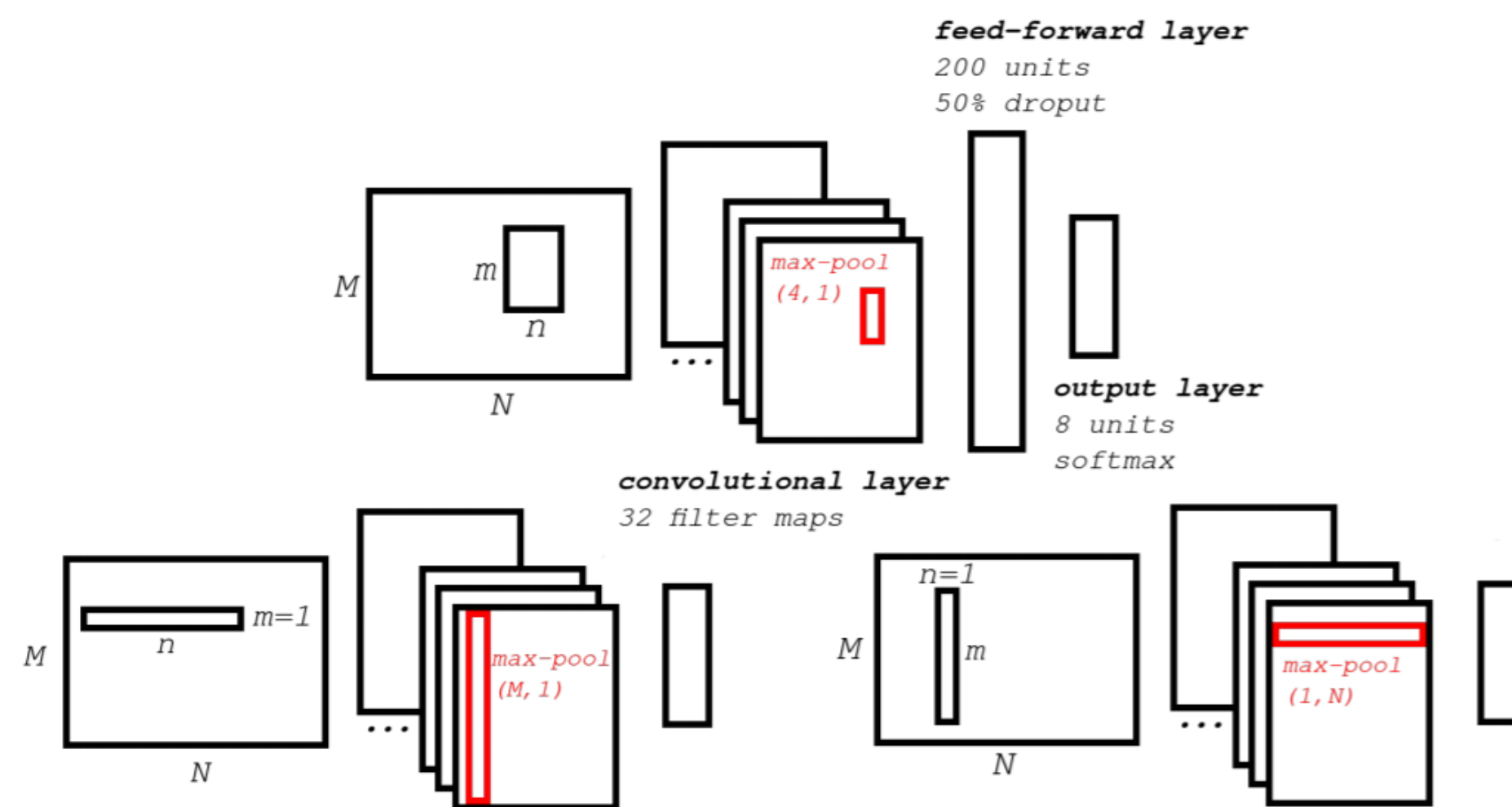


An adapted shape of 2D filters for 1st layer

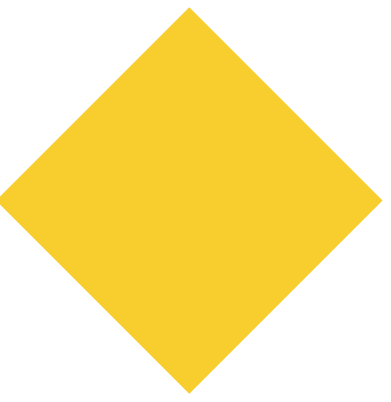
3) Adapted filtered forms to highlight specific properties

Musically-motivated CNN

→ the filters shape is suited to represent timbre (vertically)
and to represent ryhtm (horizontally)



APPROACHES FOR FEATURE LEARNING



1D representation

Representation

- Raw waveform (end-to-end learning)

Idea

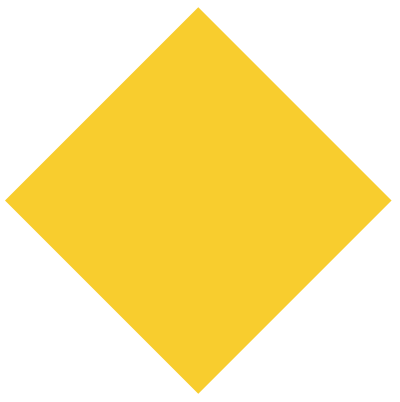
- Learn the filter to be applied directly on the waveform to get the most appropriate representation for a given task

Problem

- How to model time invariance ?

Research still going on to explore this approach.

APPROACHES FOR FEATURE LEARNING

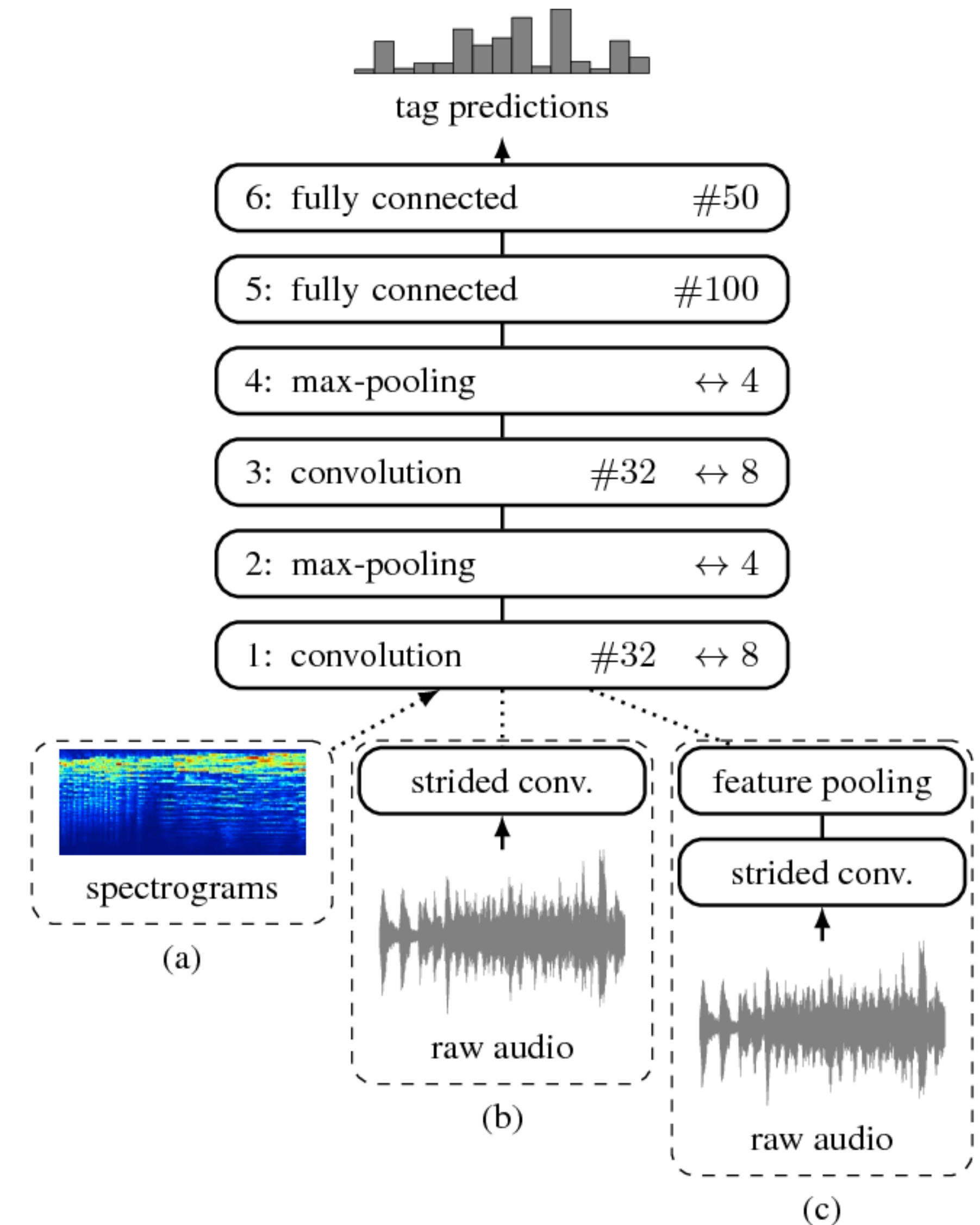


An adapted shape of 1D filters for 1st layer

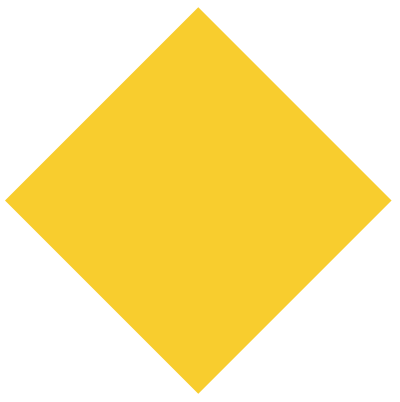
1) Filter with unique size

Using waveform in input of a convolutional network (Conv1D, TDC)
→ filter size and stride: parameter of STFT

problem: is temporal invariance respected ?

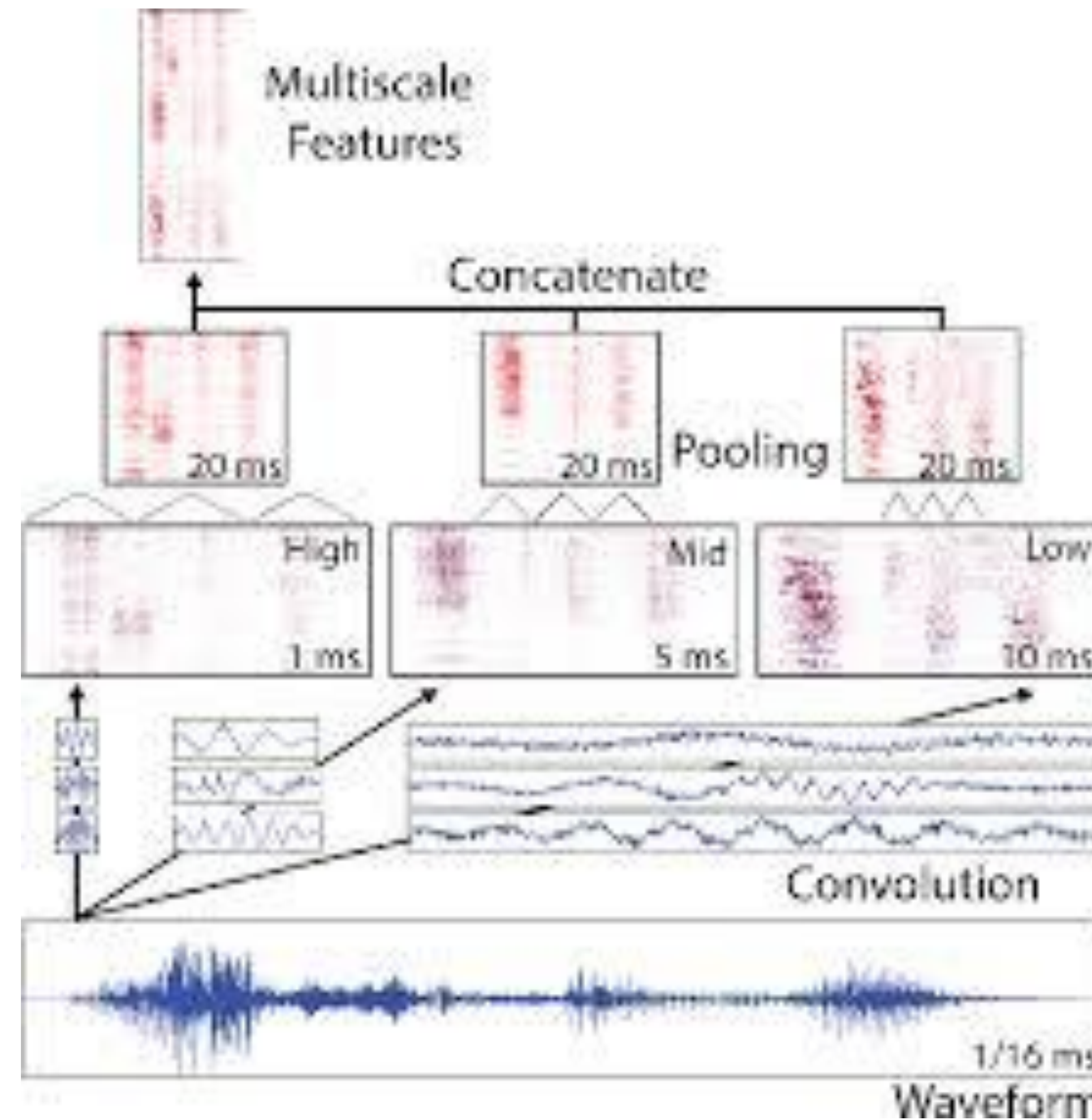


APPROACHES FOR FEATURE LEARNING

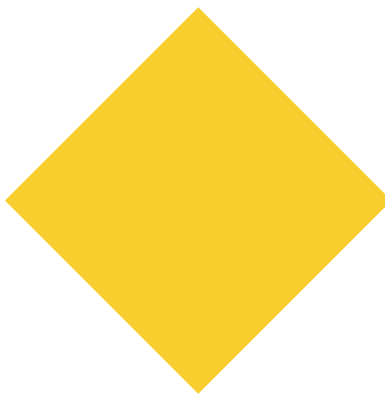


An adapted shape of 1D filters for 1st layer

- 2) Filters of different size
Multi-scale approach

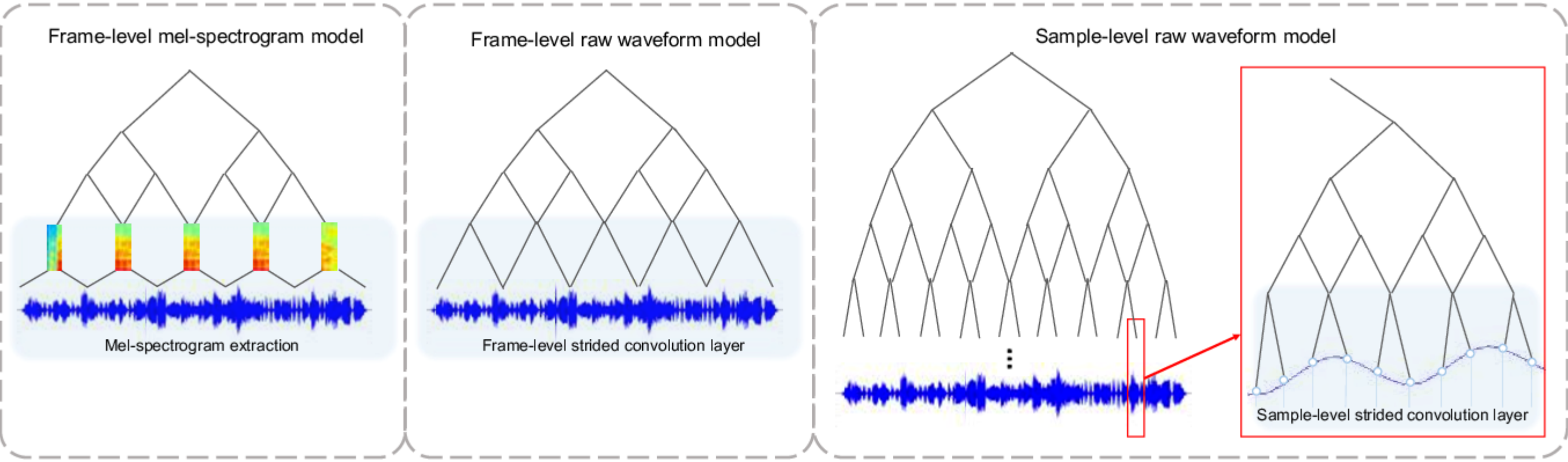


APPROACHES FOR FEATURE LEARNING



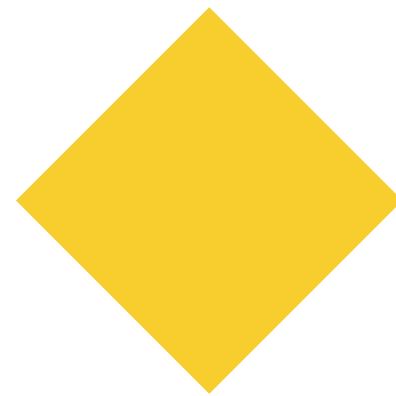
An adapted shape of 1D filters for 1st layer

- 3)
Sample CNN: VGG-net on waveform
→ make it easier to ensure time invariance


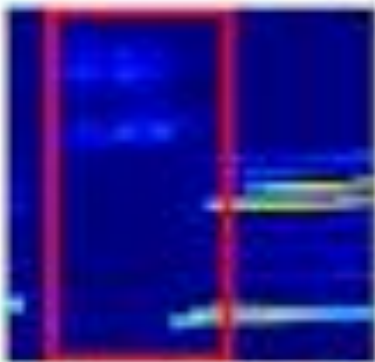
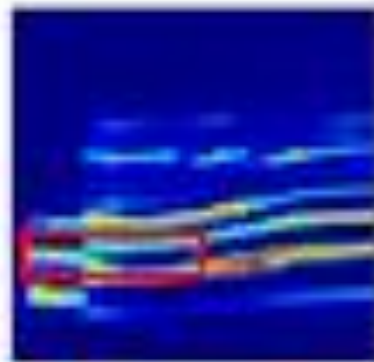

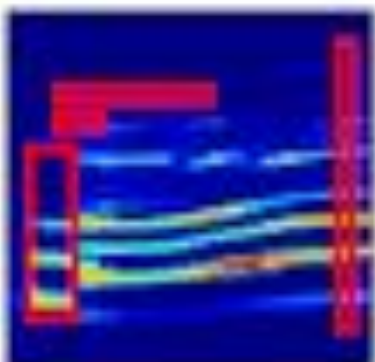
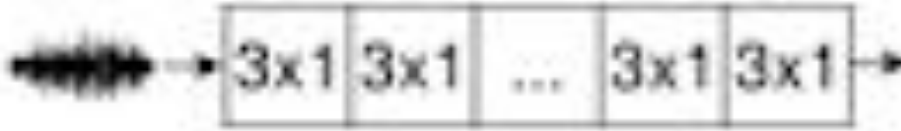
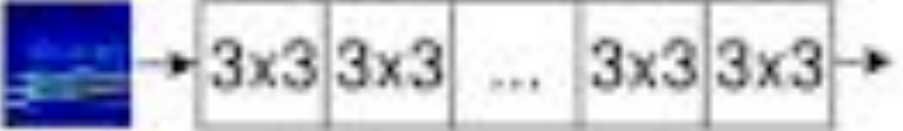


3 ⁹ model, 19683 frames 59049 samples (2678 ms) as input			
layer	stride	output	# of params
conv 3-128	3	19683 × 128	512
conv 3-128	1	19683 × 128	49280
maxpool 3	3	6561 × 128	
conv 3-128	1	6561 × 128	49280
maxpool 3	3	2187 × 128	
conv 3-256	1	2187 × 256	98560
maxpool 3	3	729 × 256	
conv 3-256	1	729 × 256	196864
maxpool 3	3	243 × 256	
conv 3-256	1	243 × 256	196864
maxpool 3	3	81 × 256	
conv 3-256	1	81 × 256	196864
maxpool 3	3	27 × 256	
conv 3-256	1	27 × 256	196864
maxpool 3	3	9 × 256	
conv 3-256	1	9 × 256	196864
maxpool 3	3	3 × 256	
conv 3-512	1	3 × 512	393728
maxpool 3	3	1 × 512	
conv 1-512	1	1 × 512	262656
dropout 0.5	—	1 × 512	
sigmoid	—	50	25650
Total params			1.9 × 10 ⁶

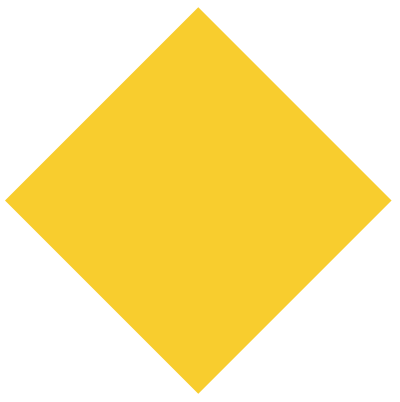
APPROACHES FOR FEATURE LEARNING



An adapted shape of 1D filters for 1st layer

DESIGN BASED ON DOMAIN KNOWLEDGE?	FILTERS CONFIG?	INPUT SIGNAL? <i>waveform</i> <i>end-to-end learning in the strictest sense</i>	<i>pre-processed waveform</i> <i>which is generally formatted in 2D i.e.: time-frequency representation</i>
yes	single filter shape in 1st CNN layer	FRAME-LEVEL  filter length: 512 stride: 256 (Dieleman et al., 2014)	VERTICAL OR HORIZONTAL  filter shape: 7x90 (Lee et al., 2009) OR  filter shape: 7x3 (Schlüter et al., 2014)
yes	many filter shapes in 1st CNN layer	FRAME-LEVEL  filter lengths: 512, 256, 128 stride: 64 (Zhu et al., 2016)	VERTICAL AND/OR HORIZONTAL  vertical filter shapes: 3x40, 1x75. horizontal filter shapes: 1x3, 1x10. (Pons et al., 2017)
no	minimal filter expression	SAMPLE-LEVEL  stack of 3x1 filters (Lee et al., 2017)	SMALL RECTANGULAR FILTERS  stack of 3x3 filters (Choi et al., 2016)

APPROACHES FOR FEATURE LEARNING



Signal models

1) Source-filter + harmonic

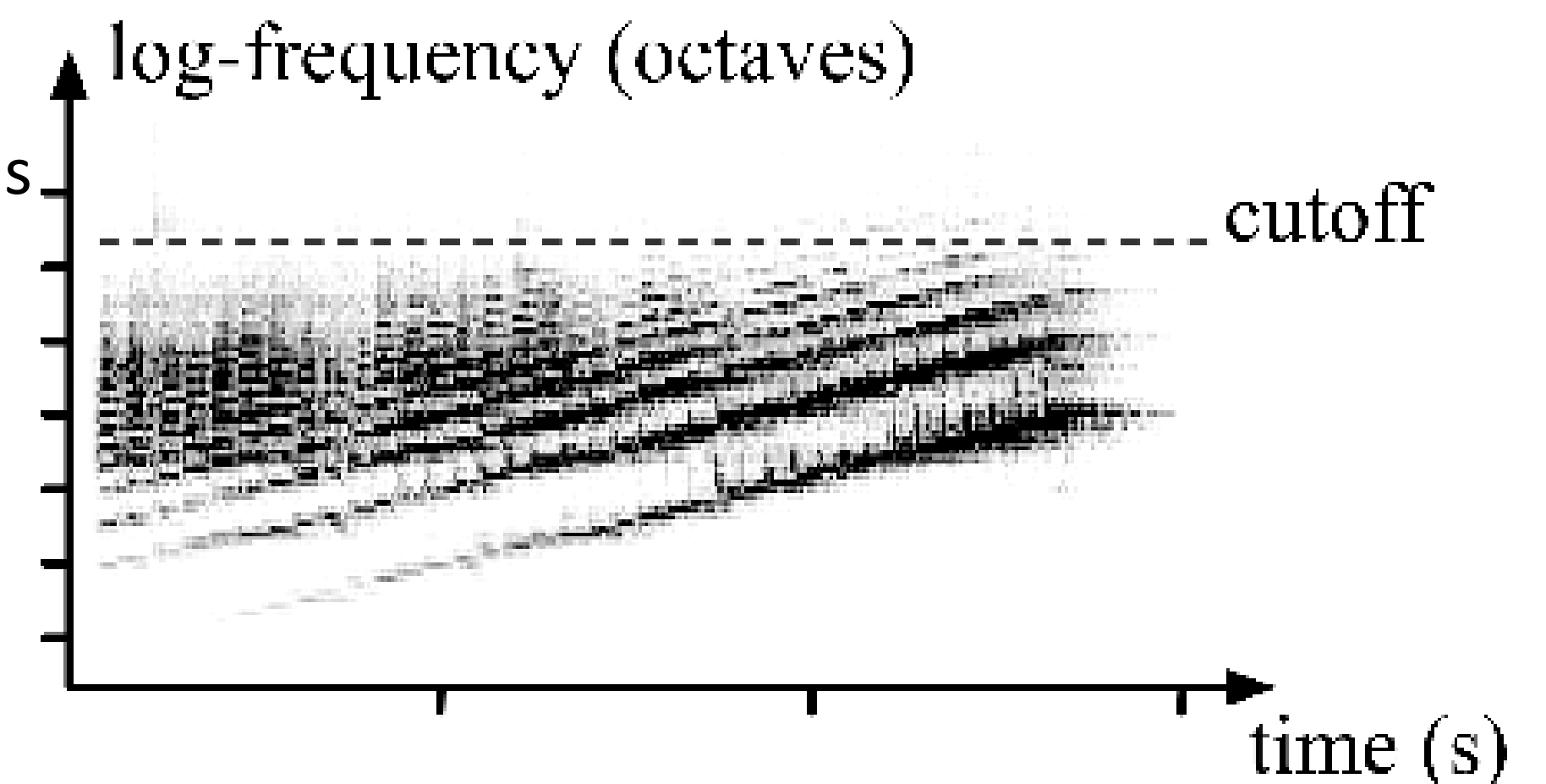
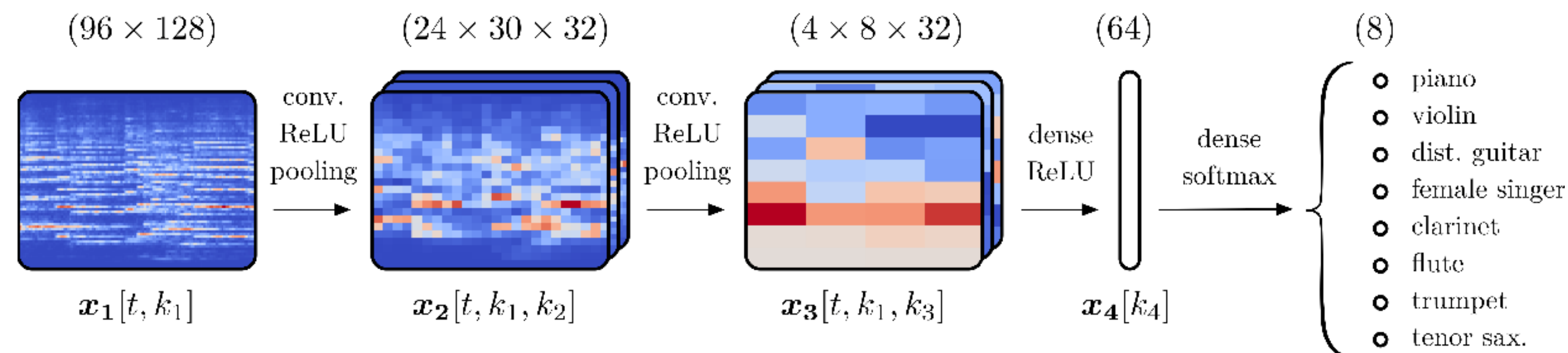
Application : instruments recognition with ConvNet

CQT: co-variant to transpositions

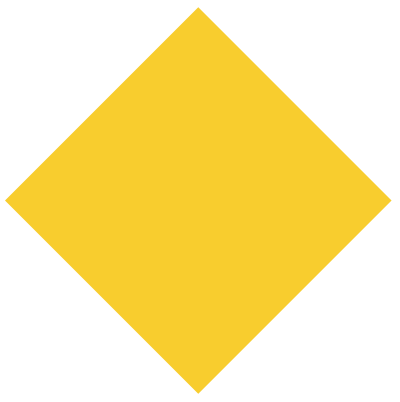
- chromatic scales are parallel diagonals
- low energy for frequencies over cut-off frequency

In high frequencies ($f > \text{cut-off}$):

- harmonics closed to each others, regularly distributed over frequency axis
- transposed sounds have similar spectra
- high correlation between CQT of different pitches
- the energy in a definite bandwidth is pitch-independent
→ 1D filter (convolution over time only)



APPROACHES FOR FEATURE LEARNING

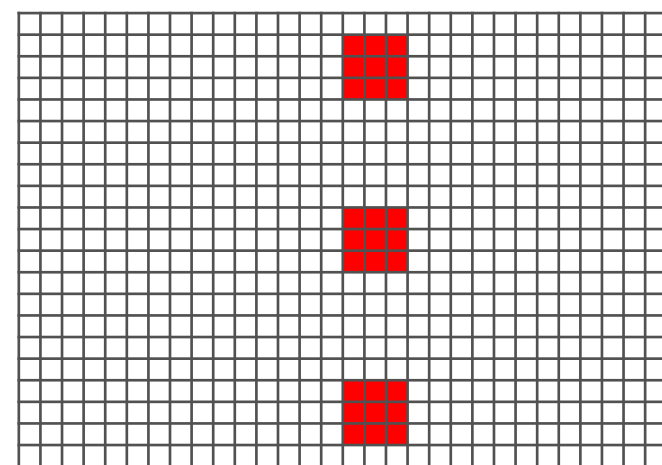
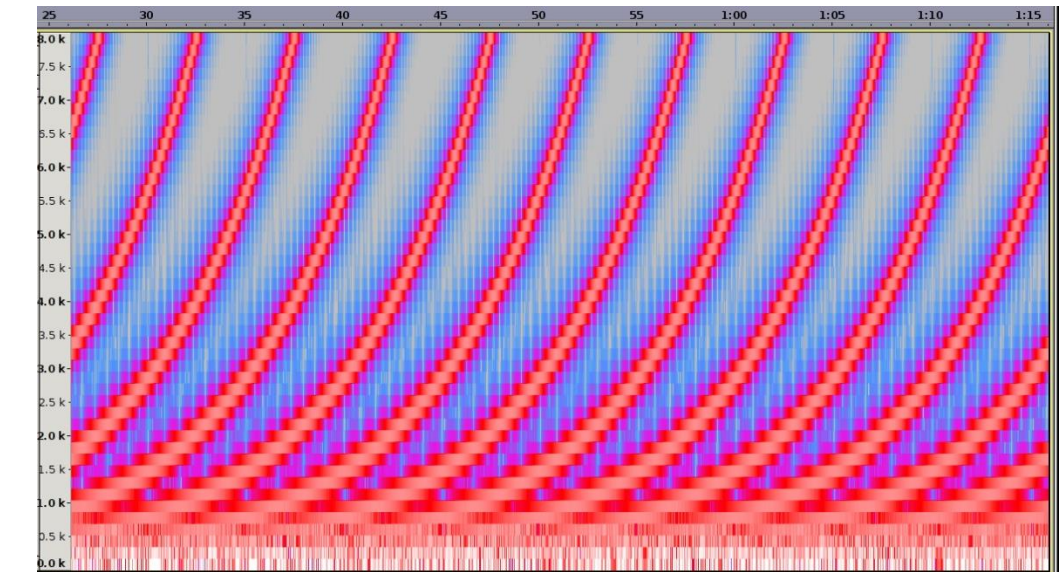


Signal models

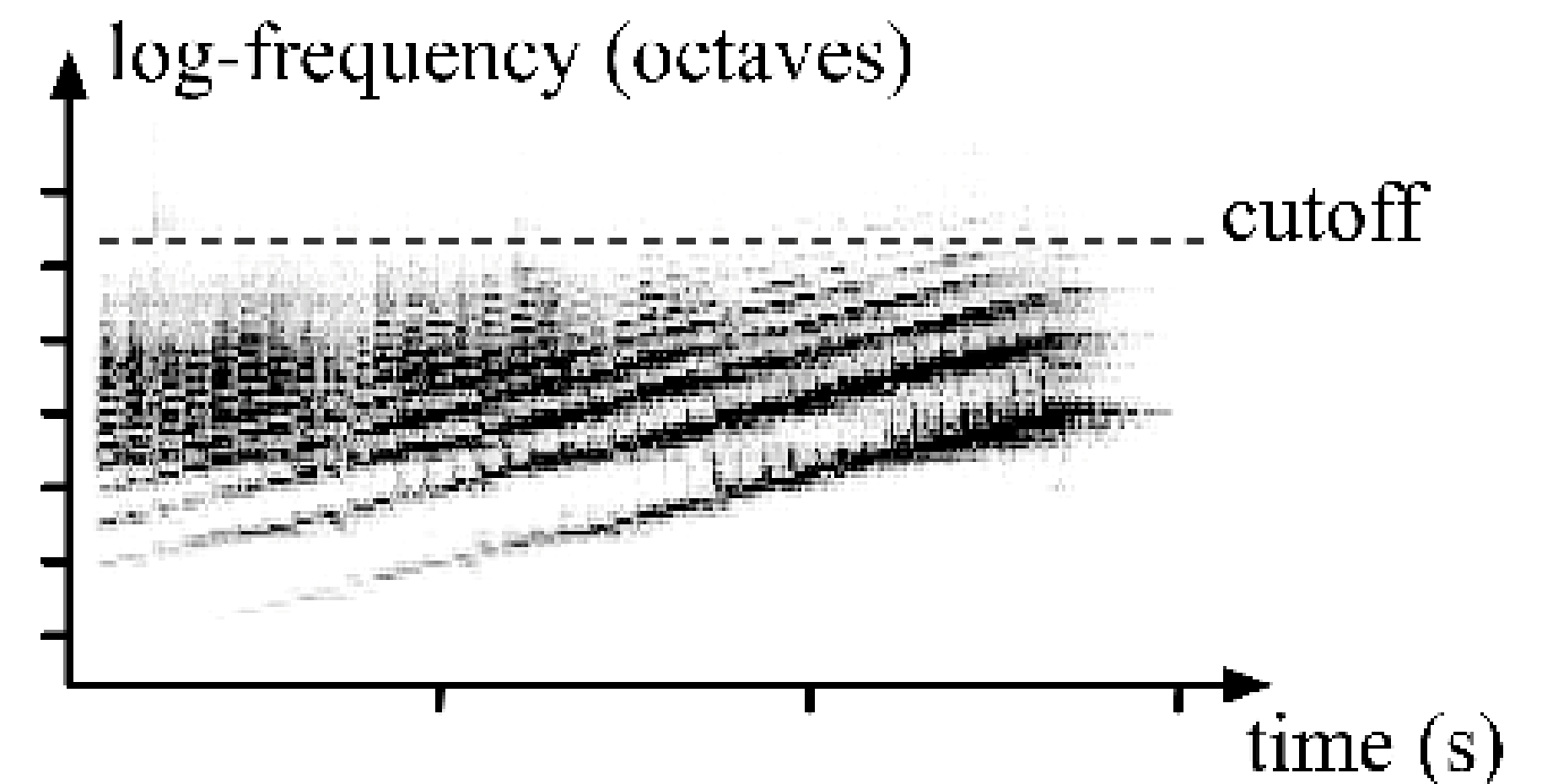
1) Source-filter + harmonic

In low frequencies (f < cut-off):

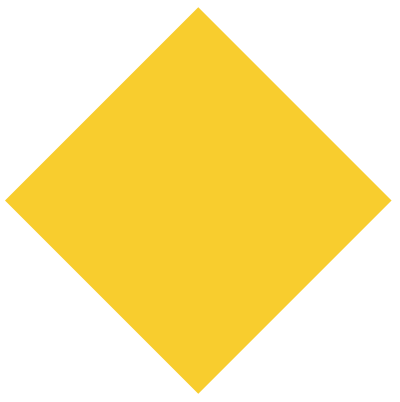
- harmonic comb is sparse and co-variant with pitch
- harmonic structure is well described, measuring correlation between harmonics themselves
- log-frequency axis is rolled on Shepard's spiral
→ time-frequency filters have a 1-octave differences over frequencies



$$y_2[t, k_1, k_2] = b_2[k_2] + \sum_{\tau, \kappa_1, j_1} W_2[\tau, \kappa_1, j_1, k_2] \times x_1[t - \tau, k_1 - \kappa_1 - Qj_1]$$



APPROACHES FOR FEATURE LEARNING



Signal models

2) Source-filter

Application: main melody estimation

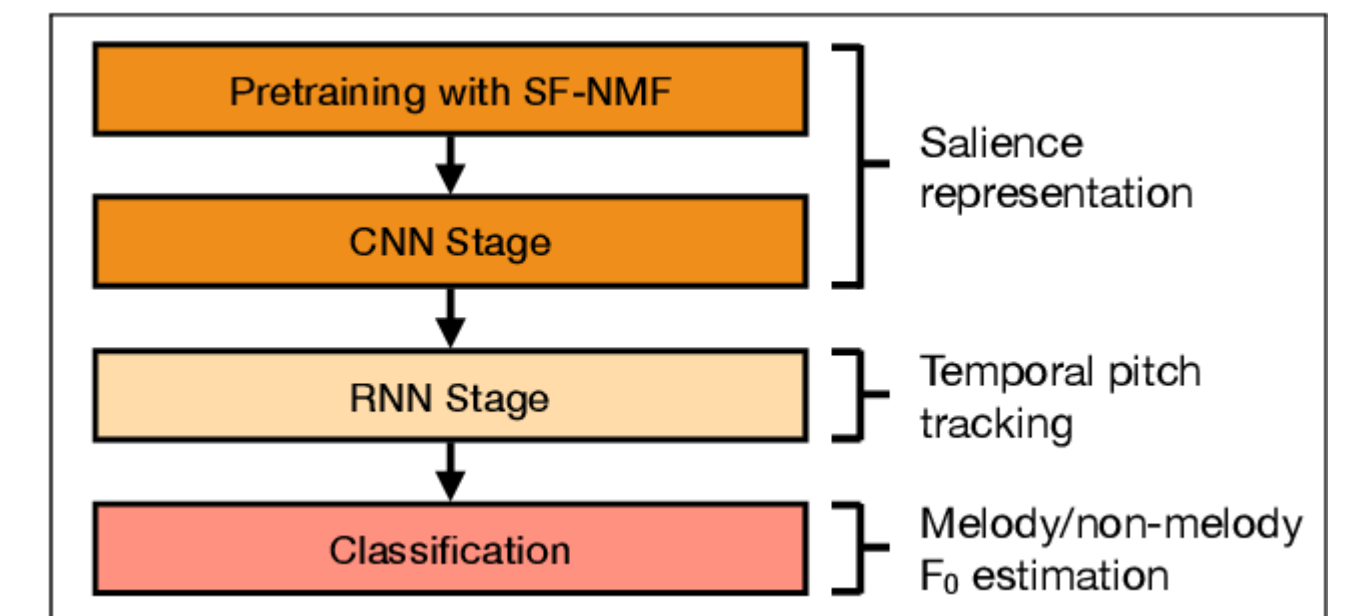
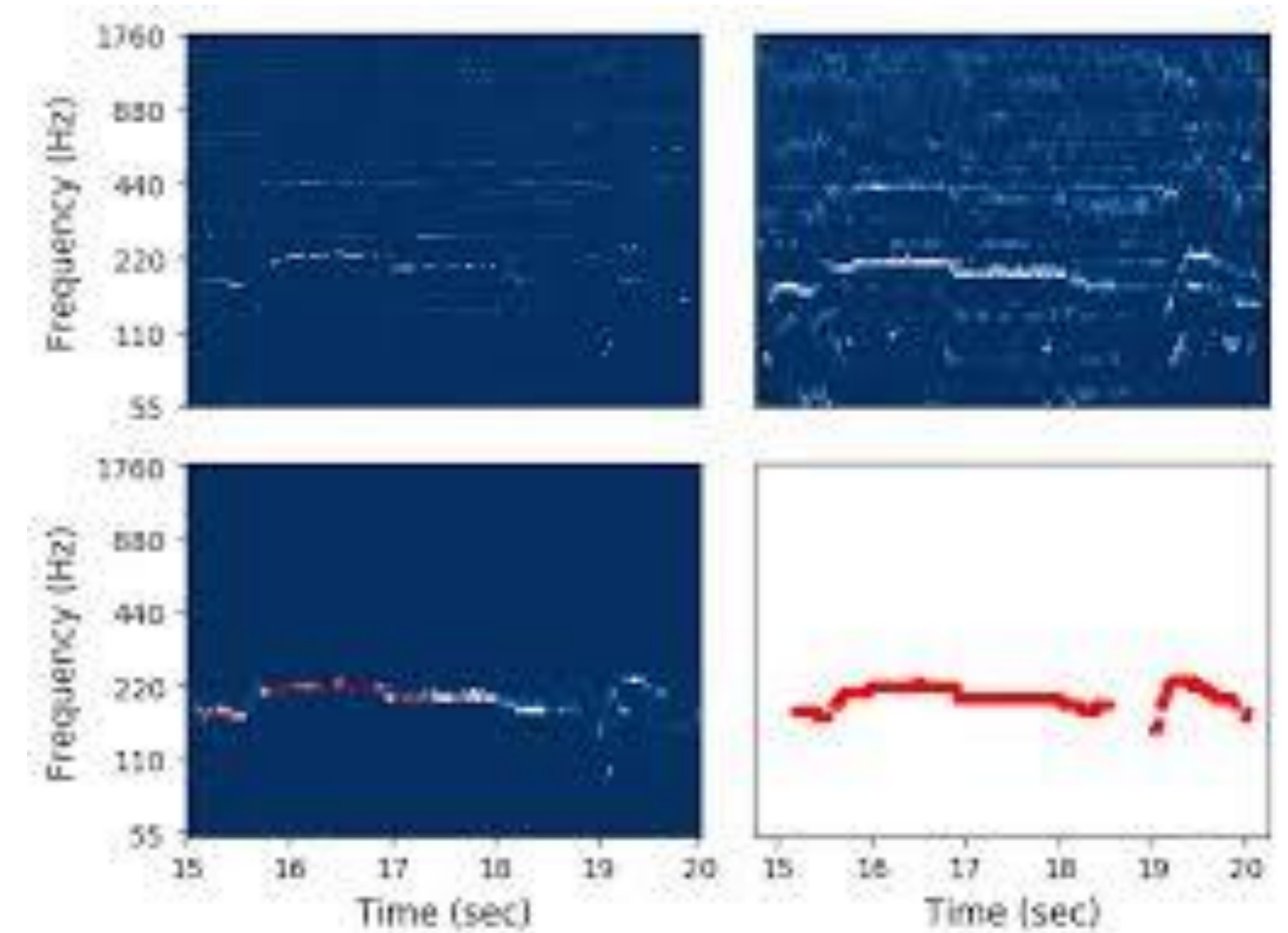
Idea:

- Desomposition of signal according to NMF model [Durrieu, 2010]

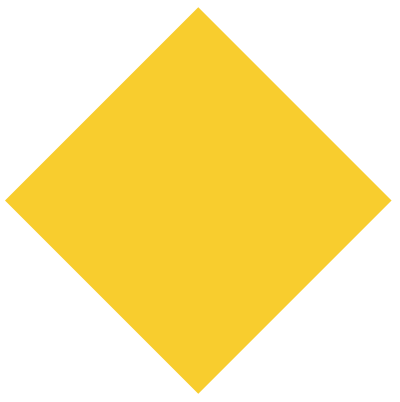
$$\begin{aligned}
 V &\approx \hat{V} = \overset{\text{source}}{V^{F_0}} \odot \overset{\text{filter}}{V^\Phi} + \overset{\text{background}}{V^B} \\
 &= \underset{\substack{\uparrow \\ \text{NMF}}}{W^{F_0}} \underset{\substack{\uparrow \\ \text{source}}}{H^{F_0}} \odot \underset{\substack{\uparrow \\ \text{NMF}}}{W^\Phi} \underset{\substack{\uparrow \\ \text{filter}}}{H^\Phi} + \underset{\substack{\uparrow \\ \text{NMF}}}{W^B} \underset{\substack{\uparrow \\ \text{background}}}{H^B}
 \end{aligned}$$

F0 basis F0 activations

- Estimated activations F_0 are used as inputs of a CNN or RNN



APPROACHES FOR FEATURE LEARNING

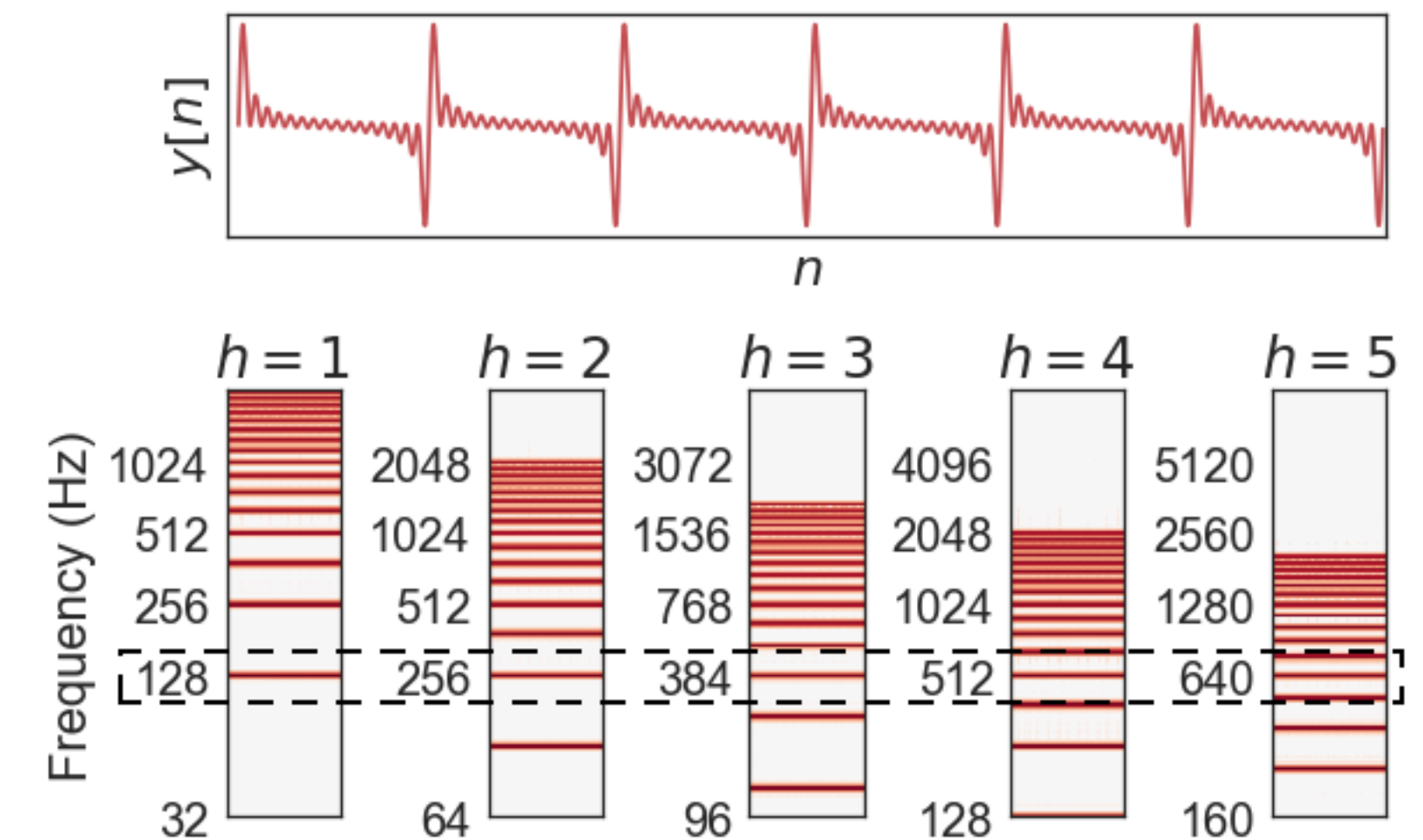
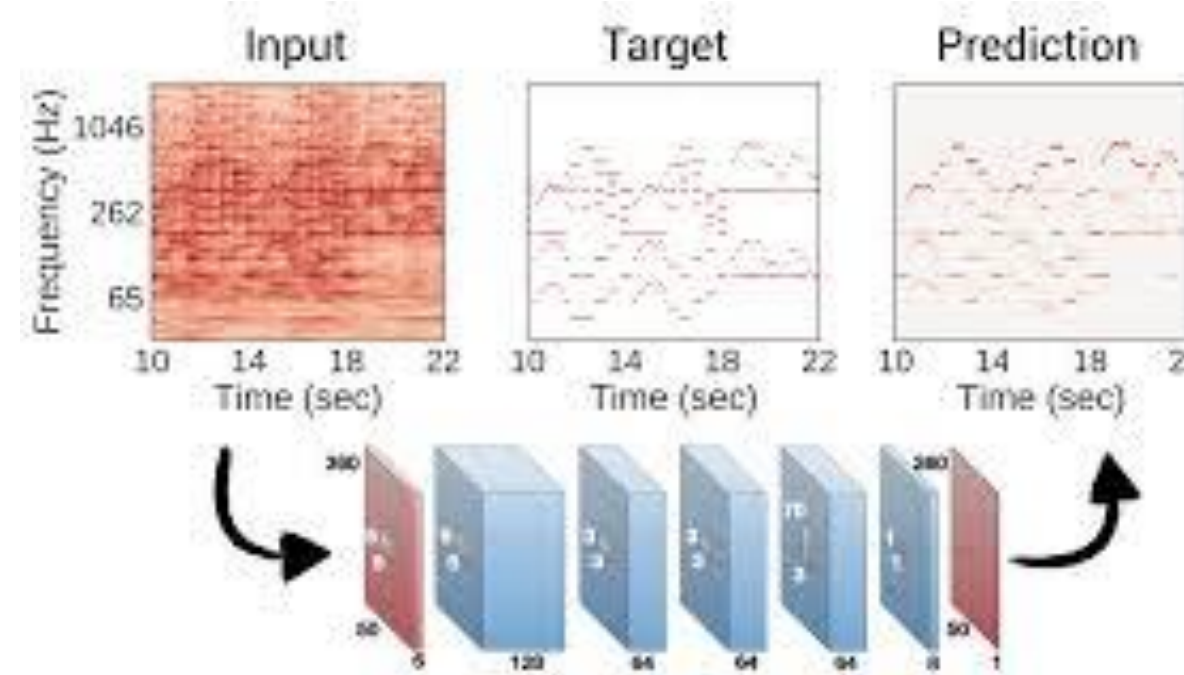


Signal models

3) Harmonic-Constant-Q-Transform (HCQT):

- Representation of audio signal using Constant-Q-Transform (CQT)
- Several CQT are computed for several minimal frequencies hf_{\min}
 - harmonics of hf_0 at same position in different CQTs
 - CQTs are stacked in input layer (RGB) depth

Application: multi-pitch estimation

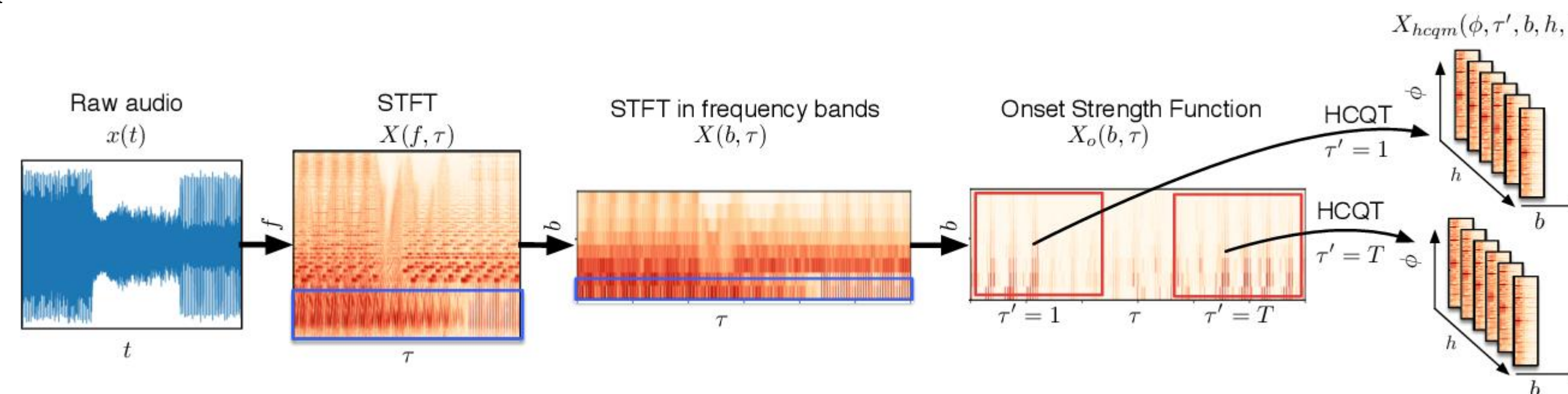


APPROACHES FOR FEATURE LEARNING

Signal models

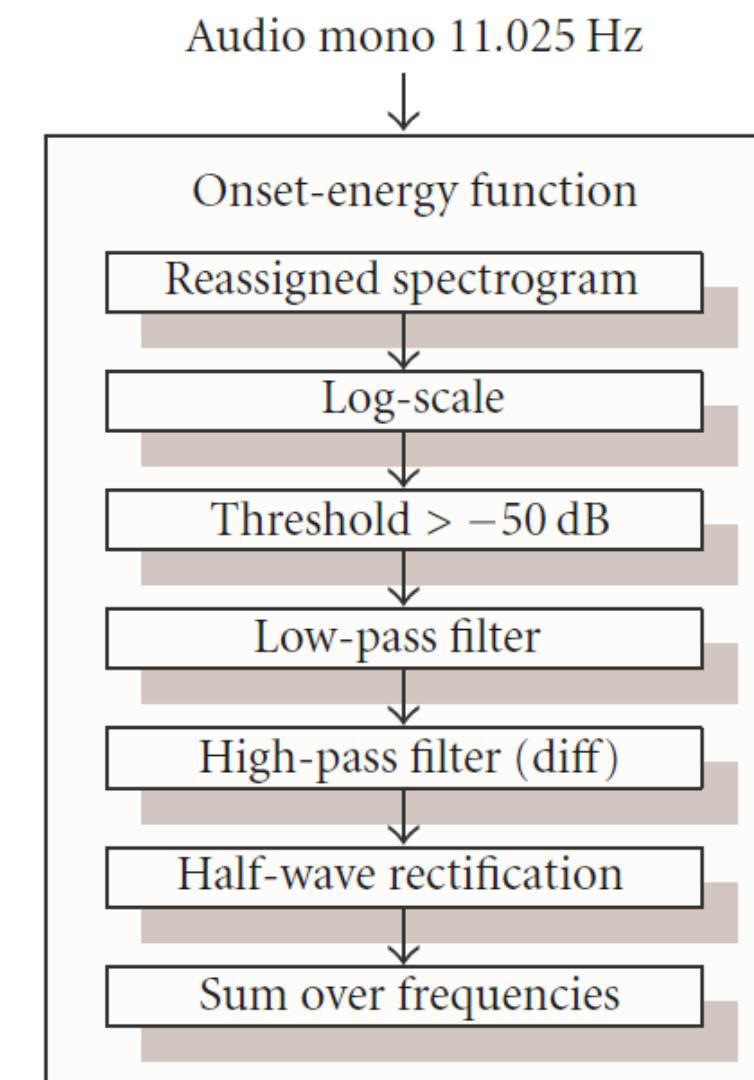
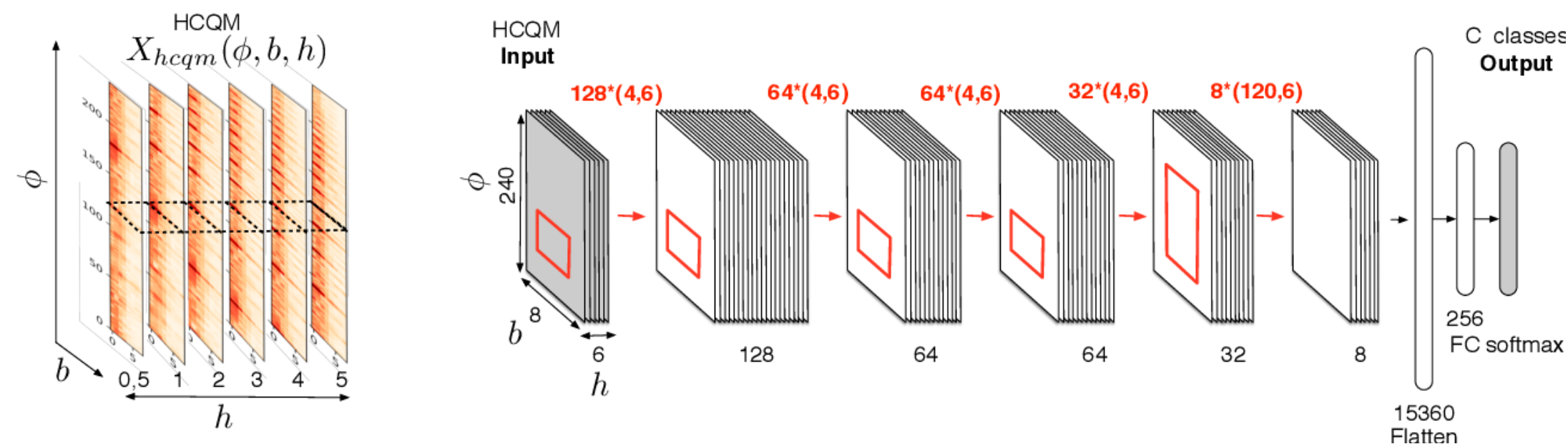
3) Harmonic-Constant-Q-Modulation (HCQM):

- In each percetive frequency bandwidth b , an onset function is computed $X_0(b, \tau)$
- The periodical content (tempo, metric, rythm) of each onset function is represented using an HCQT $X_{hcqm}(\phi, \tau', b, h)$ with ϕ modulation, h harmonic

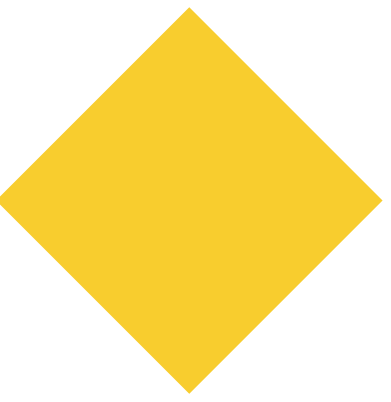


Application:

- Tempo estimation
- rythm classification



APPROACHES FOR FEATURE LEARNING



Generative models

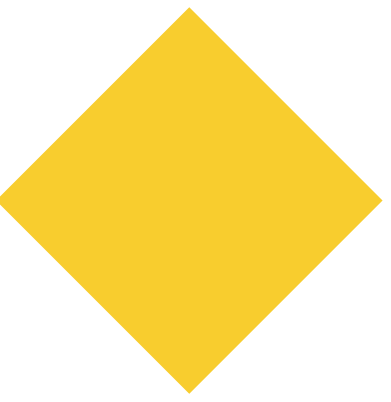
A generative model unable to generate an audio signal $x(m)$ using a z -representation

Usually:

- if z is a complex STFT
→ use DFT inverse
- If z is a spectrogram (magnitude of STFT)
→ use DFT inverse and try to reconstruct phase with Griffin and Lim algorithm
(many artifacts)

Otherwise, do not use Fourier Transform anymore for that purpose

APPROACHES FOR FEATURE LEARNING



Generative models

1) Neural-Autoregressive Models: WaveNet

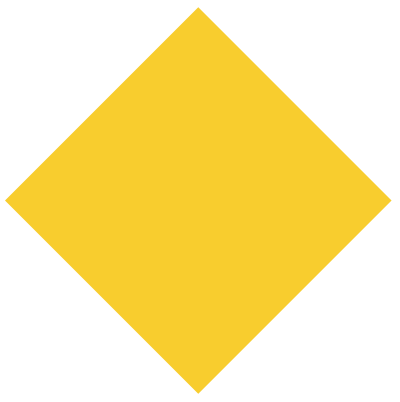
Generative model operating directly on audio samples

- « raw audio → challenging models
- Based on PixelCNN
- High resolution and long-term dependancies

Autoregressive model

next sample is almost reconstructed from linear convolution of past samples

APPROACHES FOR FEATURE LEARNING



Generative models

1) Neural-Autoregressive Models: WaveNet

Causal convolution

requires many layers of large filters to increase receptive field

Dilated convolution (wholes)

increase the receptive field by orders of magnitude

Stacked dilated convolution

dilation doubled

Input/output signal representation

Softmax layer

Conditional wavenet

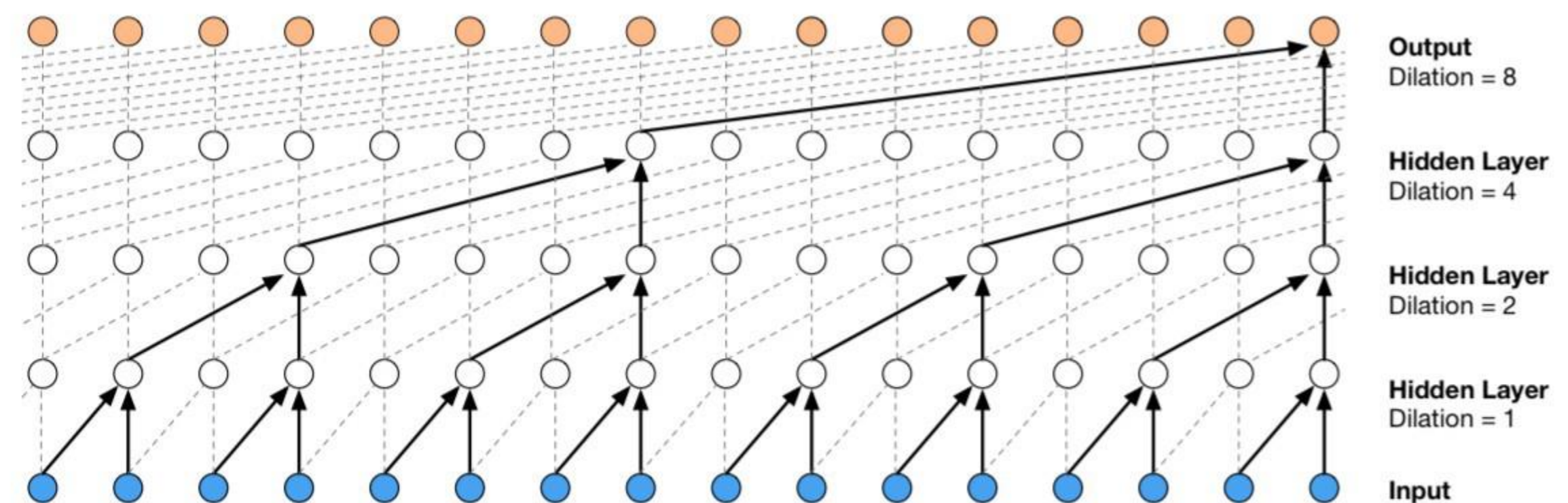
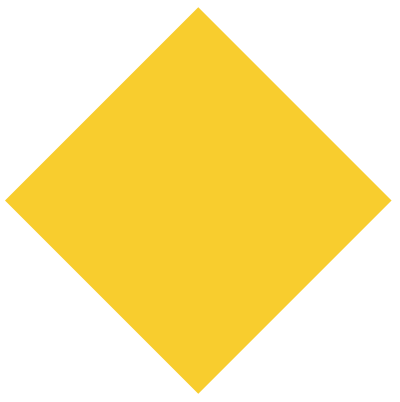


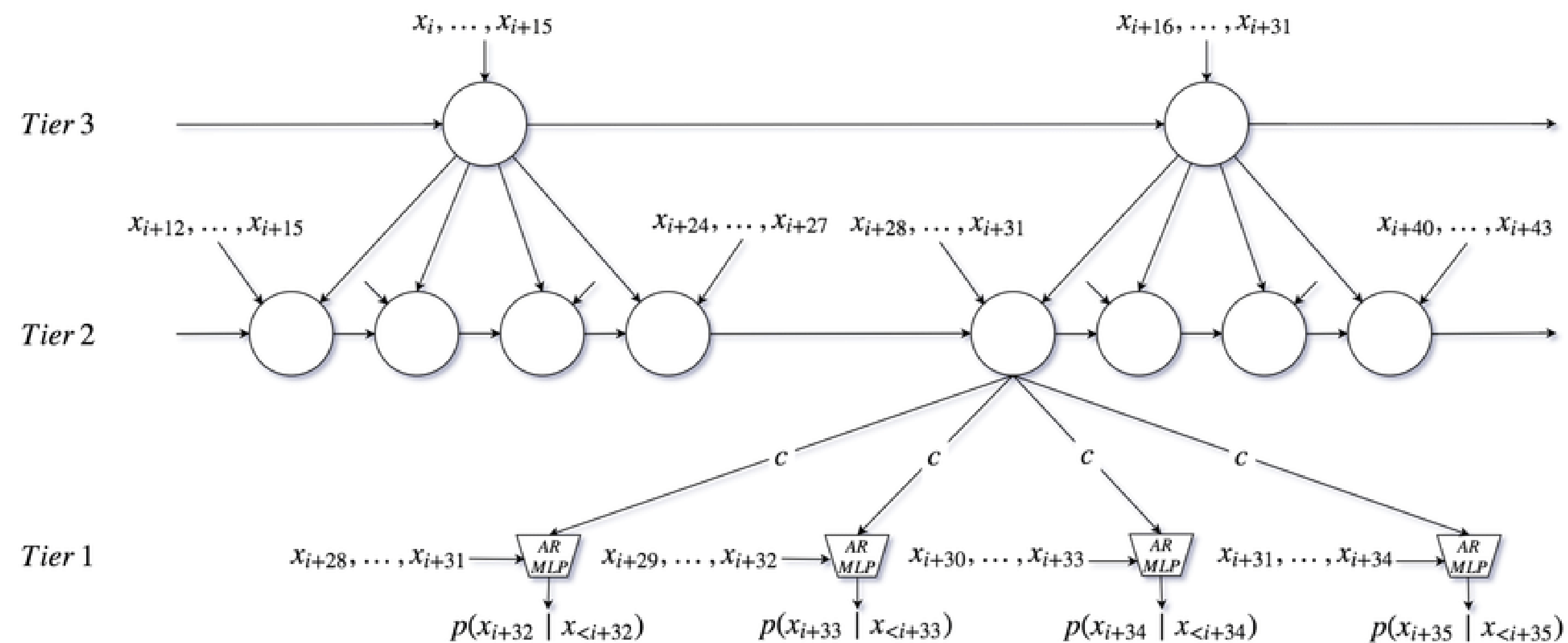
Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

APPROACHES FOR FEATURE LEARNING



Generative models

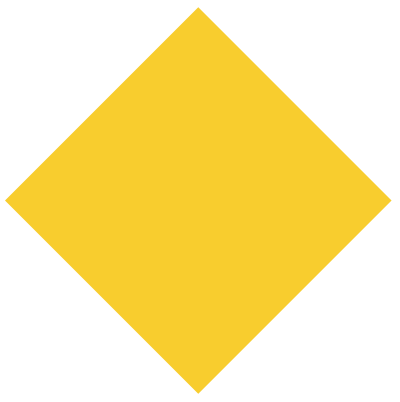
2) Neural-Autoregressive Models: SampleRNN



DEEP LEARNING FOR AUDIO AND SPEECH PROCESSING.

New learning paradigms

NEW LEARNING PARADIGMS



Classification

- Binary classification

$$x^{(i)} \rightarrow \hat{y}^{(i)} = f_{\theta}(x^i)$$

$$\theta^* = \min_{\theta} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

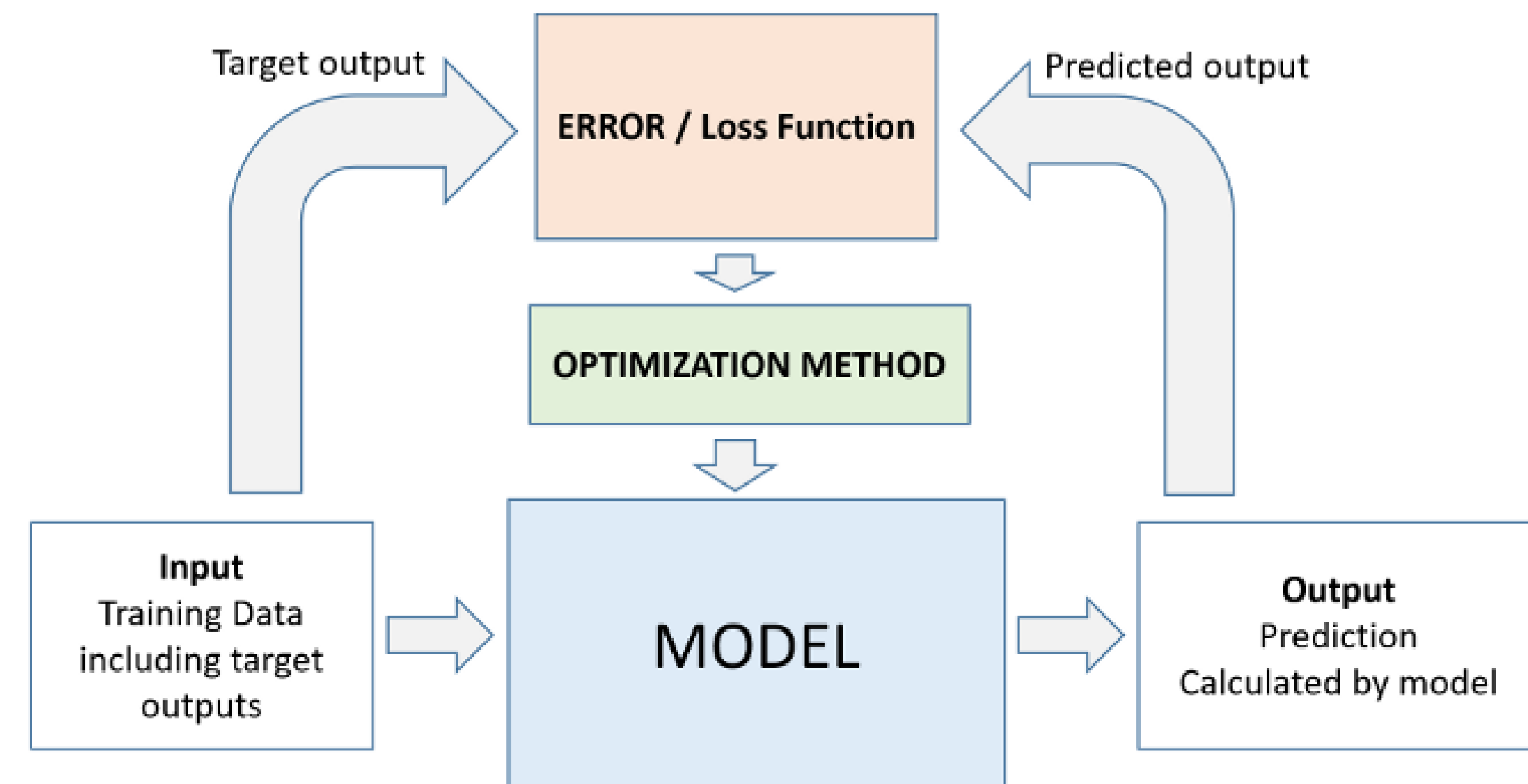
$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log \hat{y}^{(i)}) + (1 - y^{(i)} \log(1 - \hat{y}^{(i)}))$$

- Multi-class classification (single-label)

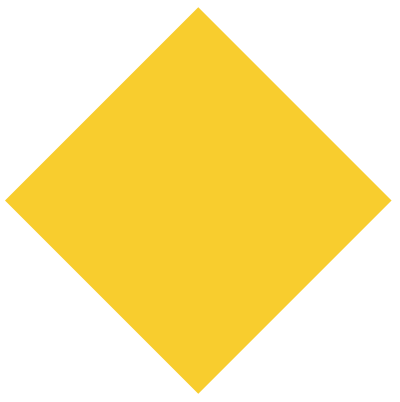
$$x^{(i)} \rightarrow \hat{y}^{(i)} = f_{\theta}(x^i)$$

$$\theta^* = \min_{\theta} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = - \sum_{k=1}^K y_k^{(i)} \log \hat{y}_k^{(i)}$$



NEW LEARNING PARADIGMS



Encoder/Decoder (Auto-encoder/VAE)

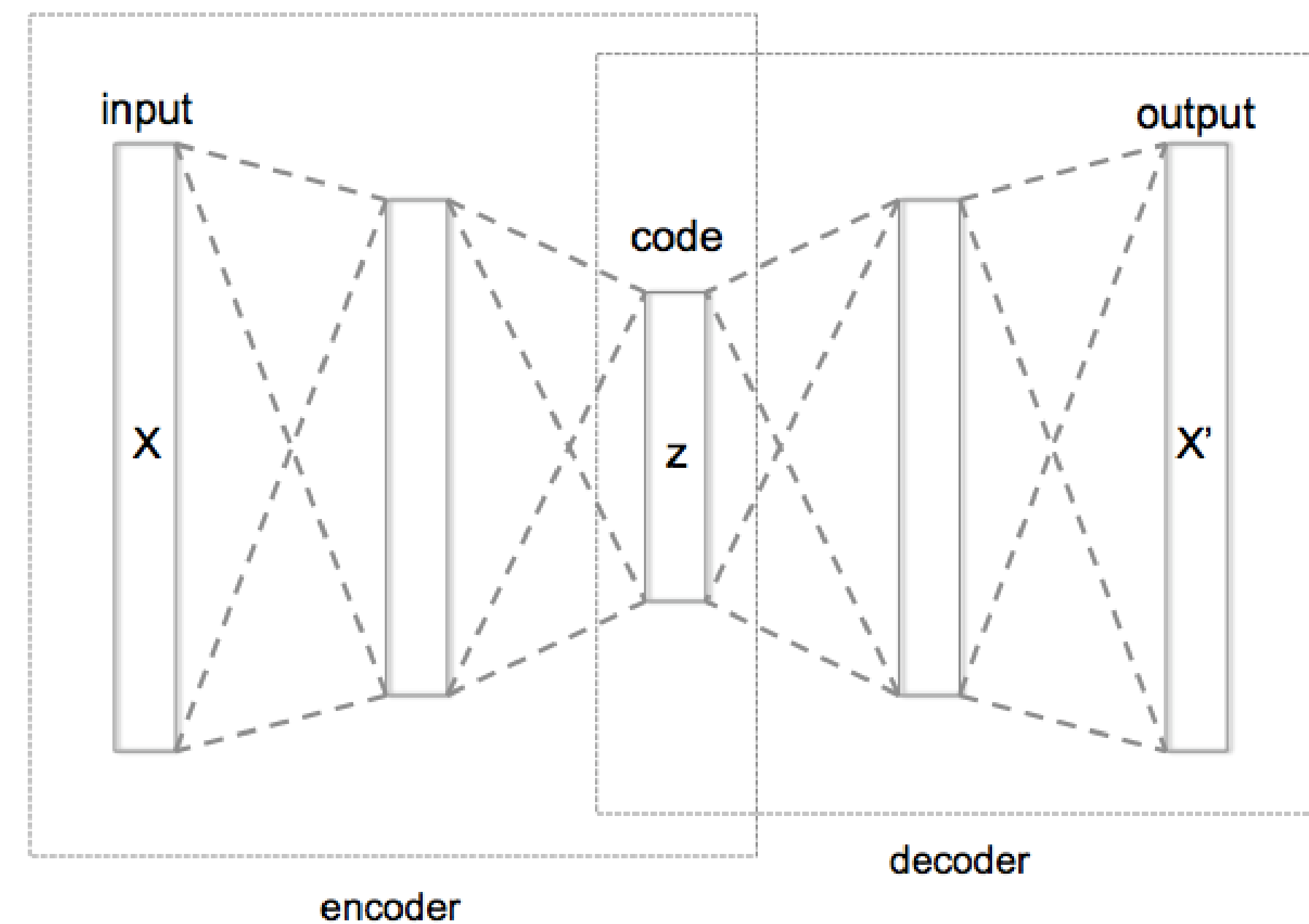
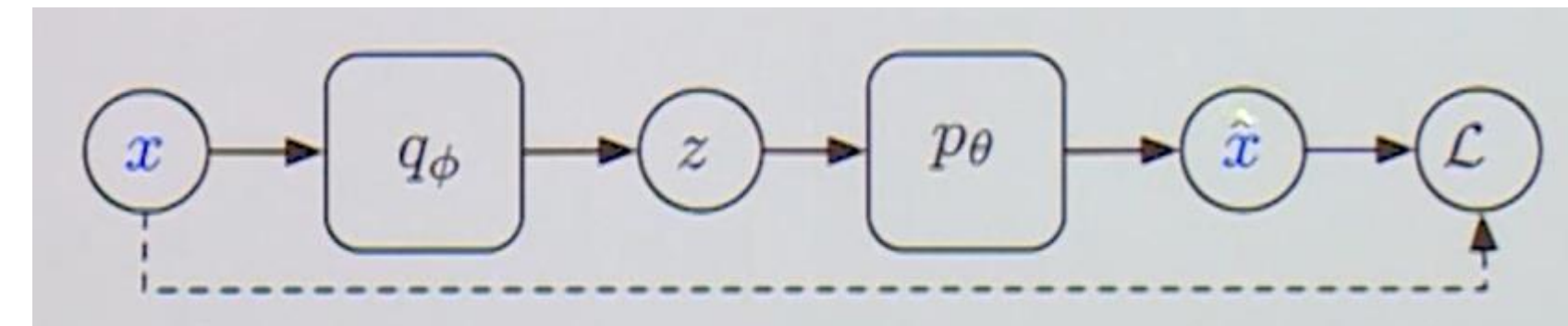
1) Auto-encoder

$$z = q_{\phi}(x)$$

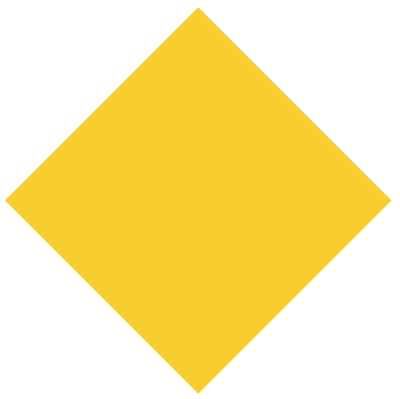
$$\hat{x} = p_{\theta}(z)$$

$$\mathcal{L} = \|x - \hat{x}\|_2^2$$

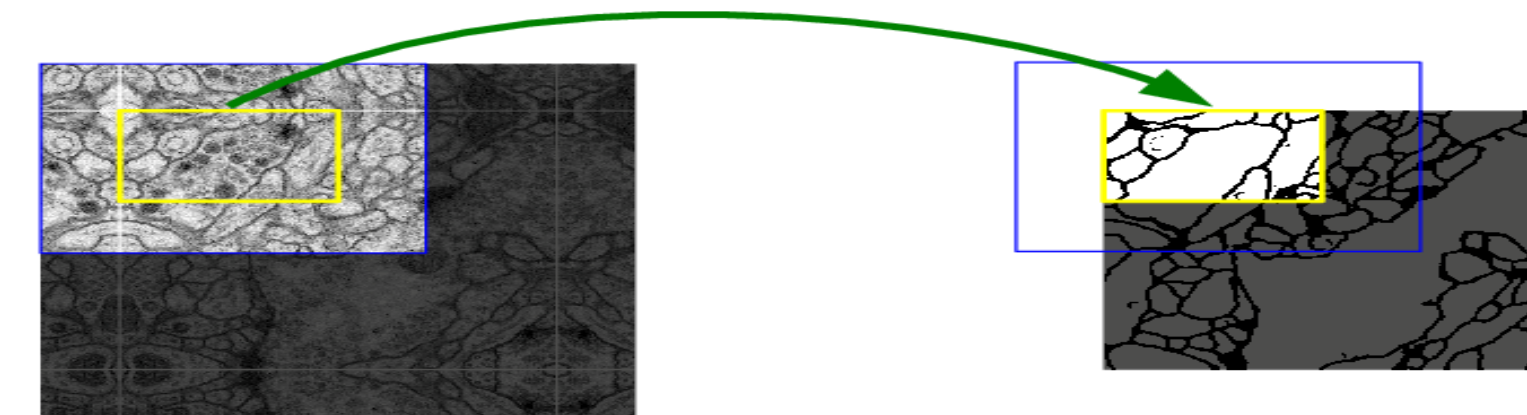
z refers to code, latent variable or latent representation
(projection of x in a variety/manifold)



NEW LEARNING PARADIGMS



Encoder/Decoder (Auto-encoder/VAE)



2) Denoising auto-encoder

U-Net

- Contracting path to capture context and a symmetric expanding path that enables precise localization

Separation of vocal and instrumental parts

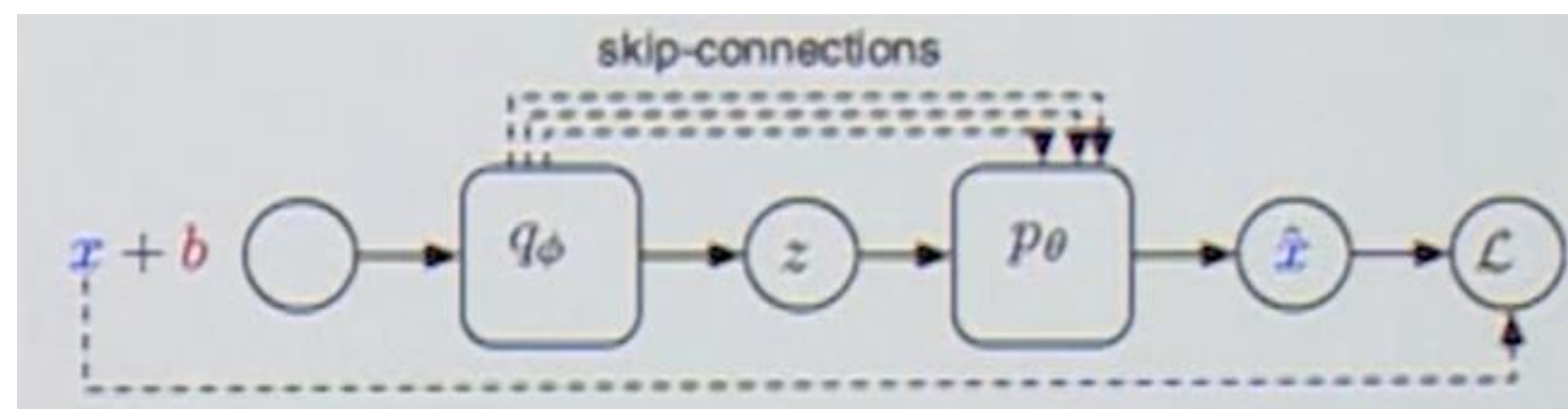
- A time-frequency mask is learnt so that:

$$X = X_v + X_i$$

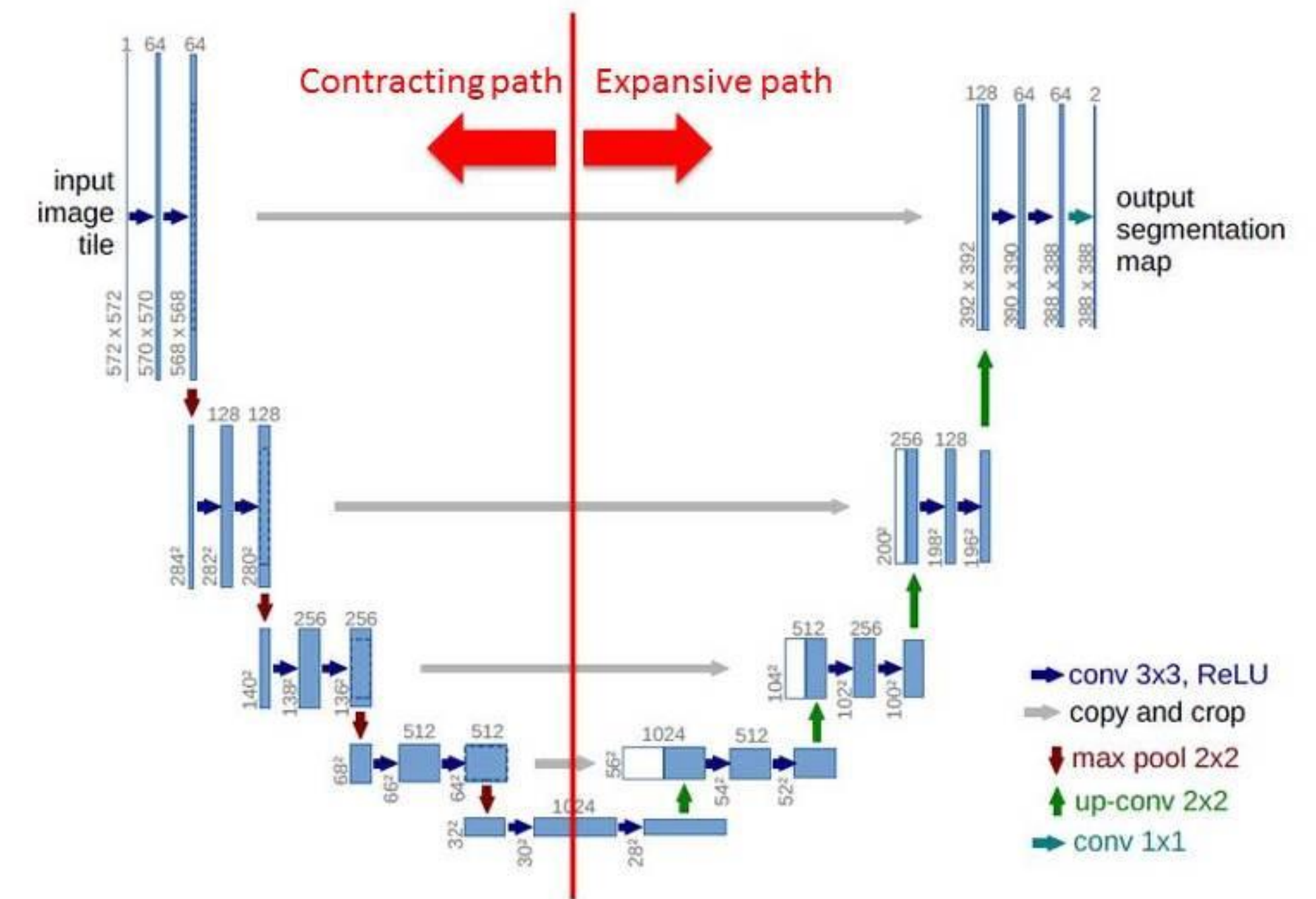
$$\hat{X}_v = X \otimes M$$

$$\mathcal{L} = \|X_v - \hat{X}_v\|_1$$

demo



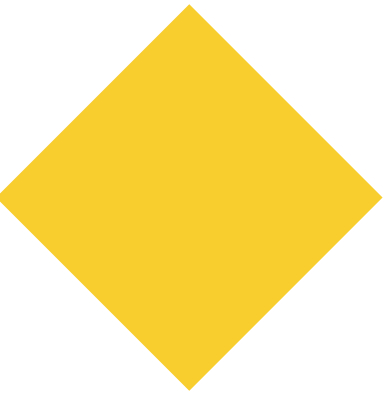
Network Architecture



A. Jansson, E. Humphrey, N. Montecchio, R. Bittner, A. Kumar, T. Weyde. Singing voice separation with deep U-Net convolutional networks. 2017.

O. Ronneberger, P. Fischer, T. Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. 2015.

NEW LEARNING PARADIGMS



Encoder/Decoder (Auto-encoder/VAE)

3) Complex input / complex network

- Deep Complex U-Net for source separation (complex mask)
Complex convolution

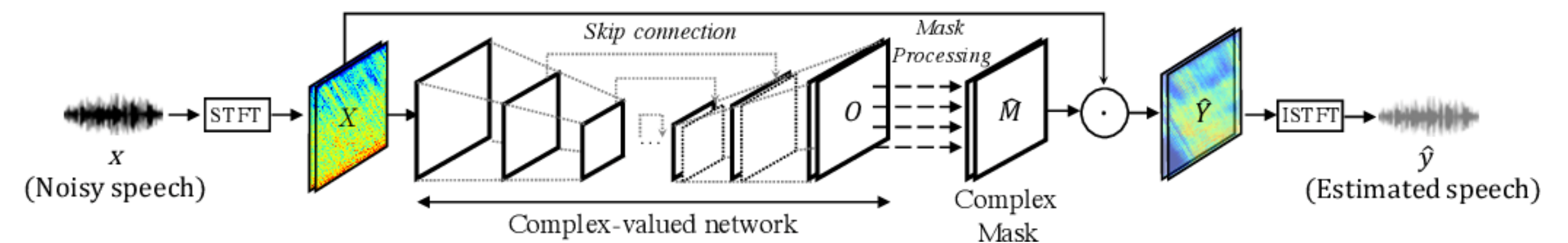
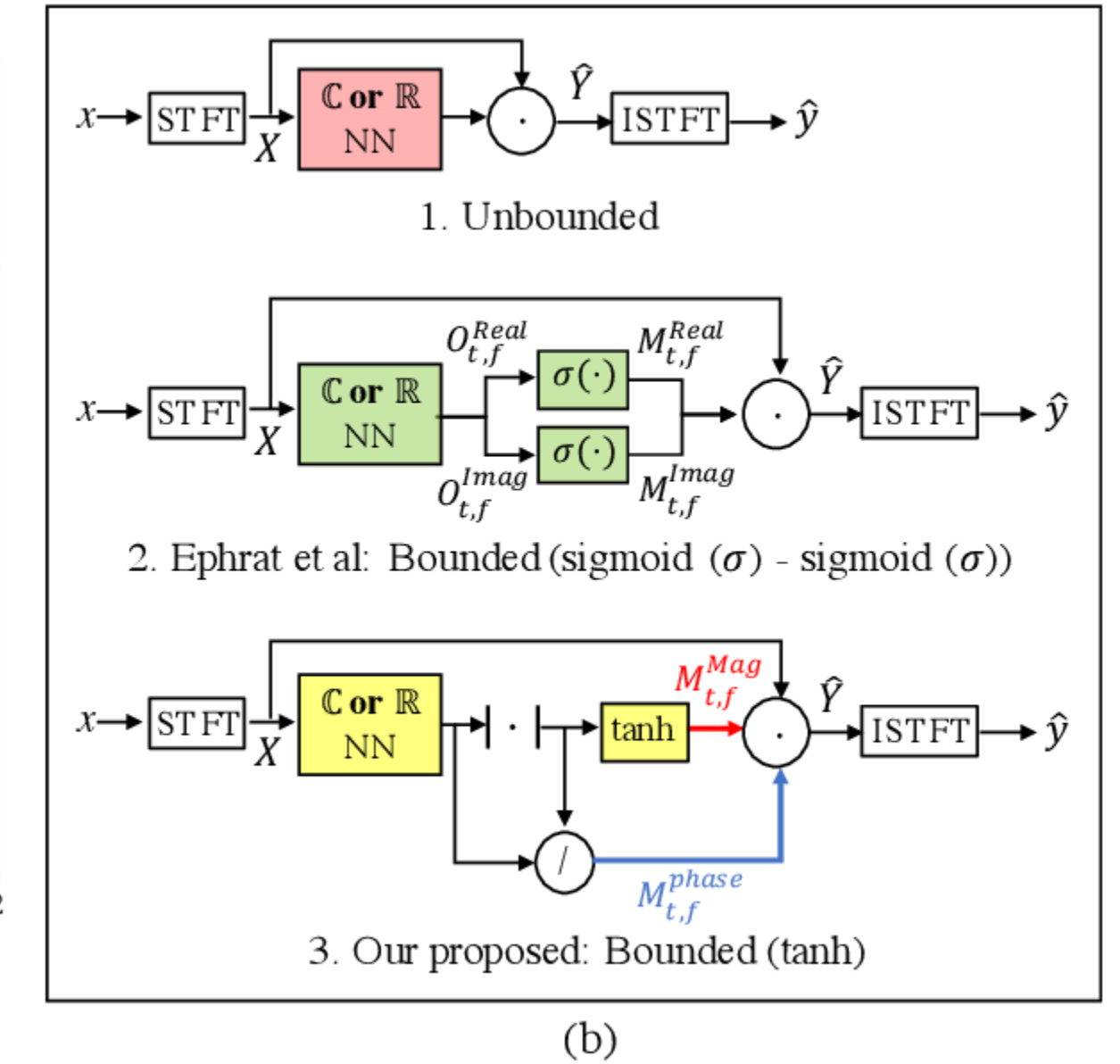
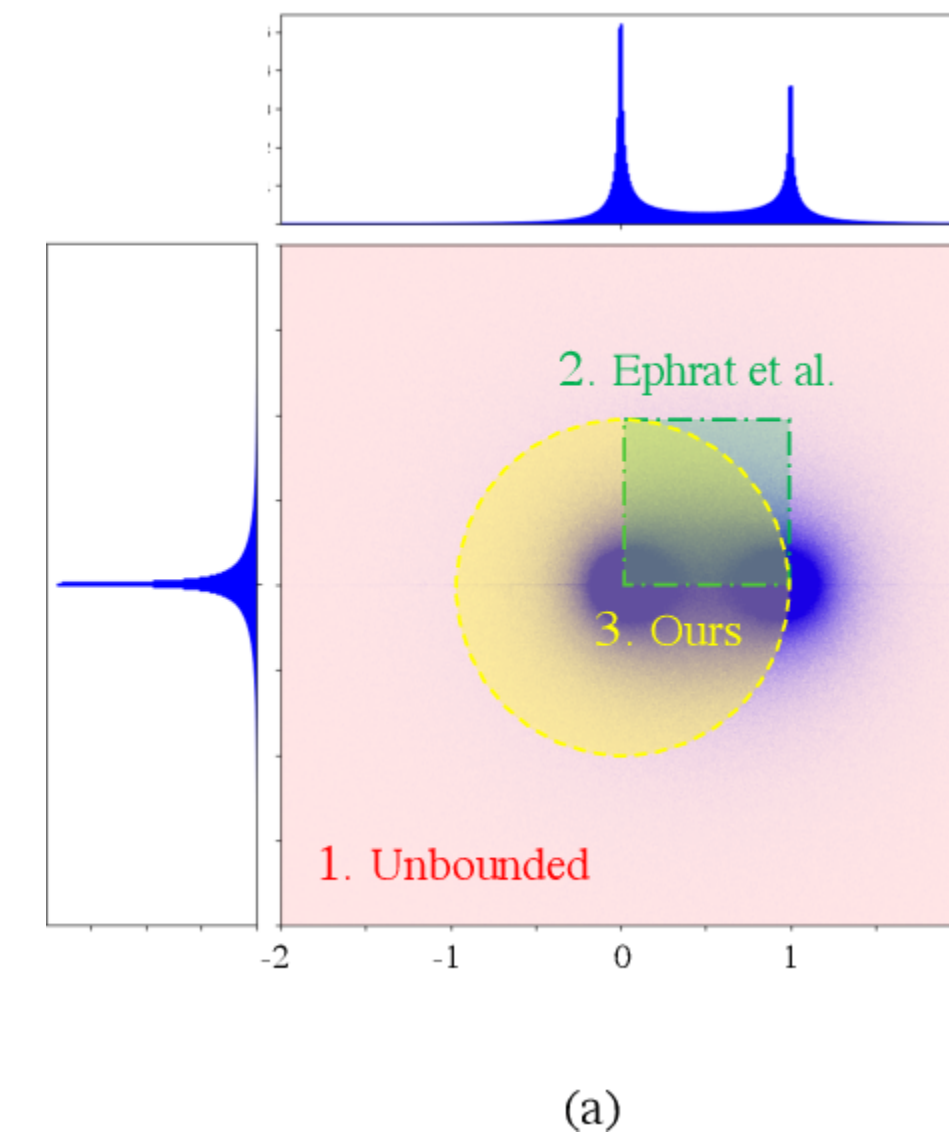
$$W = A + iB$$

$$h = x + iy$$

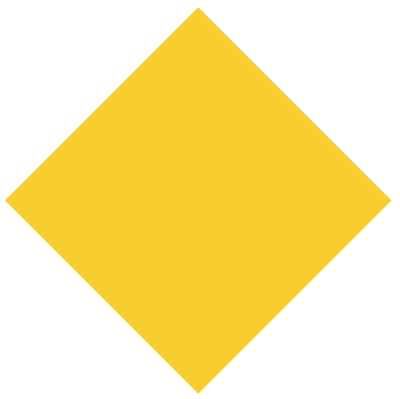
$$W * h = (A * x - B * y) + i \cdot (B * x + a * y)$$

Complex masking

$$\begin{aligned}\hat{Y}_{t,f} &= \hat{M}_{t,f} \cdot X_{t,f} \\ &= |\hat{M}_{t,f}| \cdot |X_{t,f}| \cdot e^{i\phi_{\hat{M}_{t,f}} + \phi_{X_{t,f}}}\end{aligned}$$



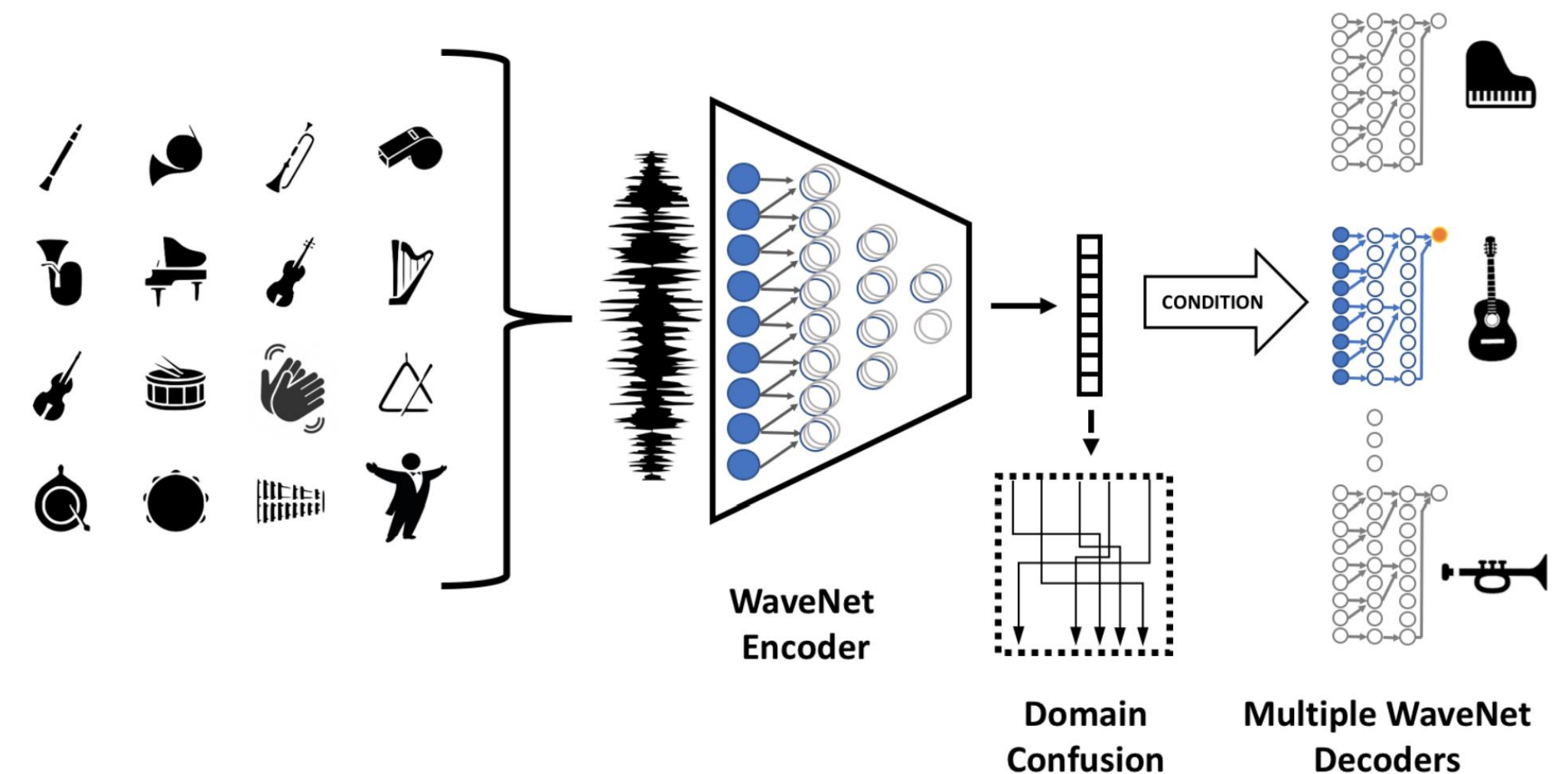
NEW LEARNING PARADIGMS



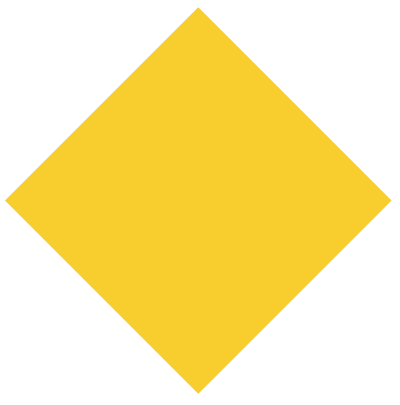
Encoder/Decoder (Auto-encoder/VAE)

4) Translating music across musical instruments, genres and styles

- $O(s, r)$:
random augmentation of input s with seed r
- Encoder E :
shared wavenet
- Disentangled latent space
Domain classification C
- Decoder D^j (domain j):
multiple wavenet, conditioned on the latent representation produced by E
- Adversarial loss
minimize reconstruction loss
maximize domain classification
→ prevent the latent space to learn domain characteristic



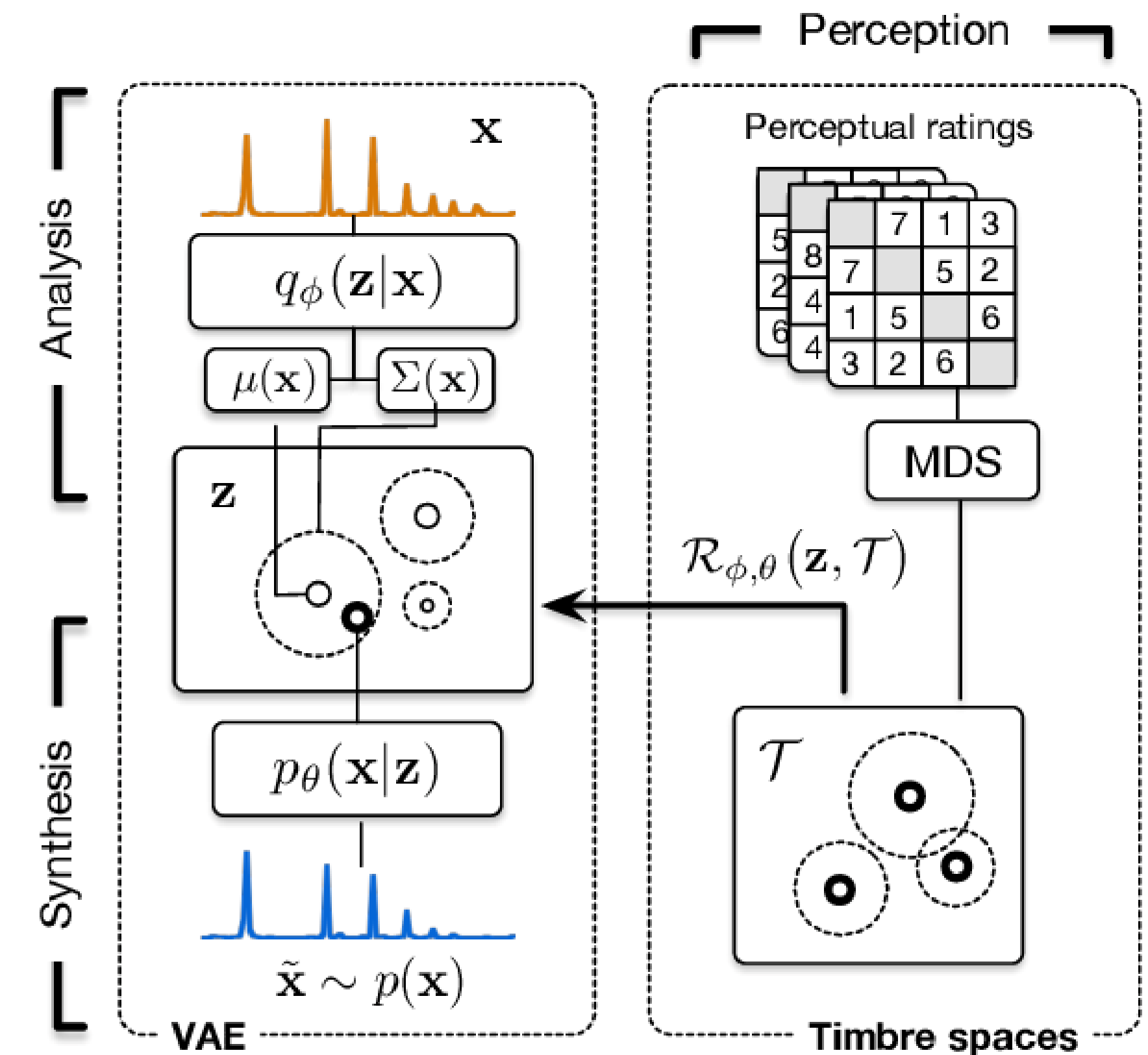
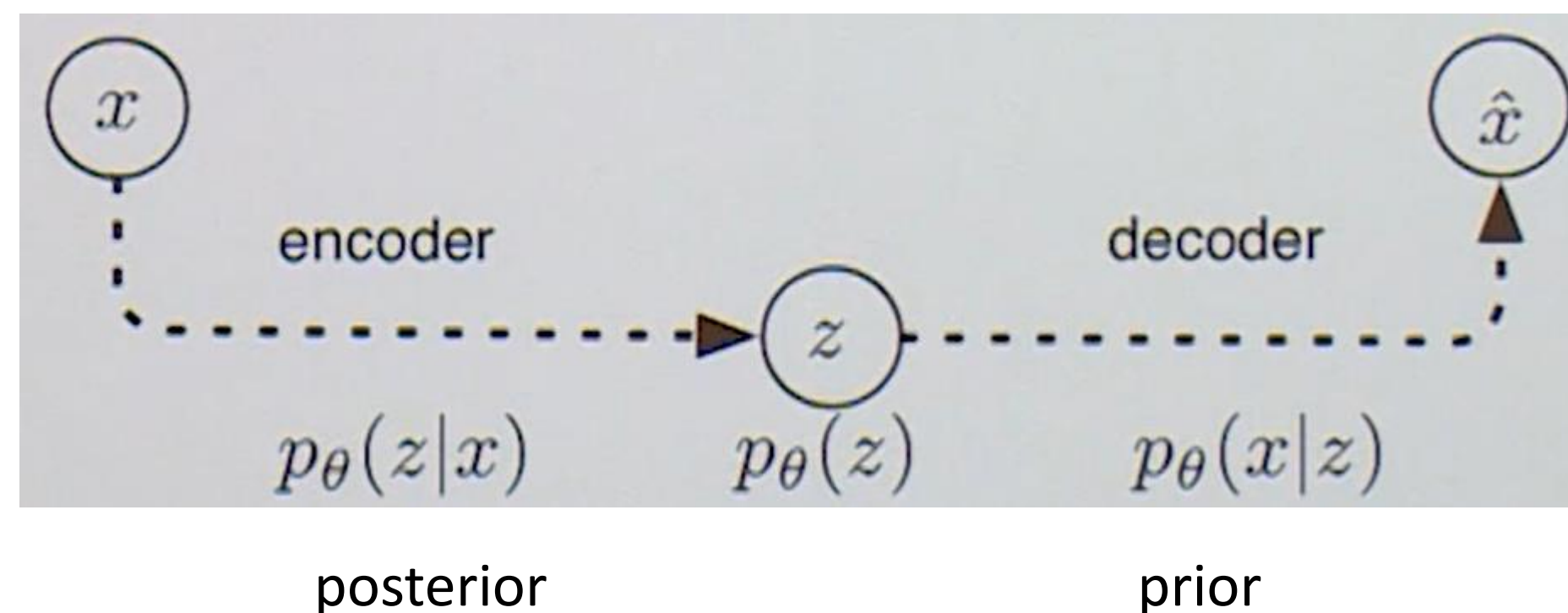
NEW LEARNING PARADIGMS



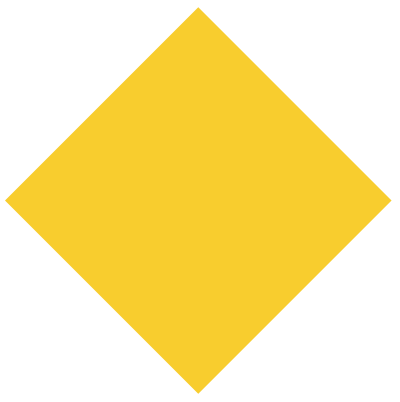
Encoder/Decoder (Auto-encoder/VAE)

5) Variational Auto-Encoder

- Latent variables are drawn from a prior
 $z_i \sim p(z)$
- data x have a likelihood that is conditioned on latent variables z :
 $x_i \sim p(x|z)$
- likelihood and prior: $p(x, z) = p(x|z)p(z) = p(xz)p(z)$



NEW LEARNING PARADIGMS



Metric learning

1) Triplet Loss

We train the network for a triplet of data anchor, positive, negative

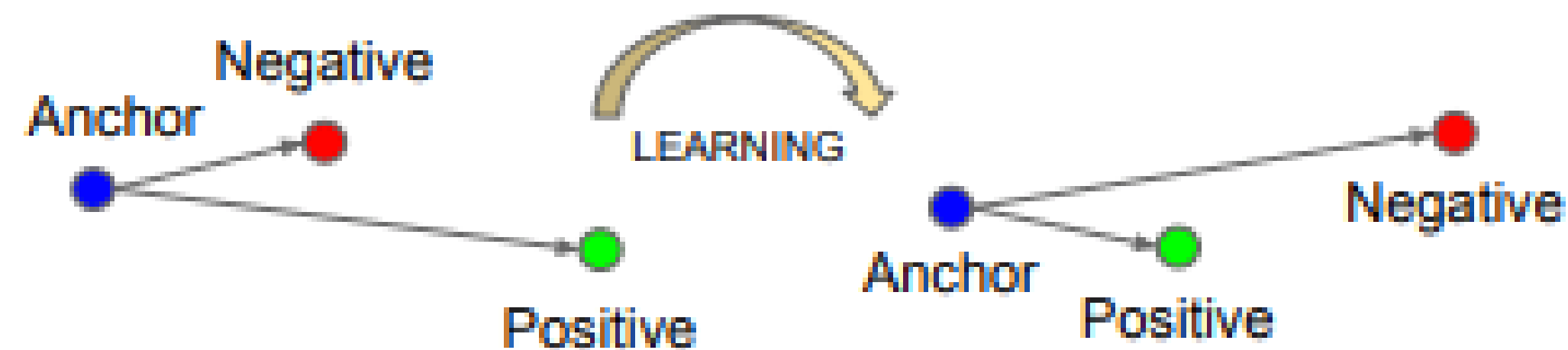


Figure 3. The **Triplet Loss** minimizes the distance between an *anchor* and a *positive*, both of which have the same identity, and maximizes the distance between the *anchor* and a *negative* of a different identity.

DEEP LEARNING FOR AUDIO AND SPEECH PROCESSING.

Thank you for your attention.

Referencences:

- Geoffroy Peeters, Telecom Paris Tech