



Institut Mines-Télécom

"*MACHINE LISTENING: AI FOR SOUNDS AND MUSIC*"

COLLOQUE IMT - L'INTELLIGENCE
ARTIFICIELLE AU COEUR DES
MUTATIONS INDUSTRIELLES.
APRIL 4TH, 2019

Gaël RICHARD

Professor, Head of the Image, Data, Signal department

A well established domain for speech ...

Speech Recognition



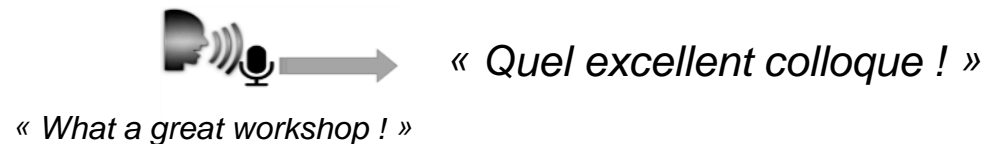
Speaker Recognition



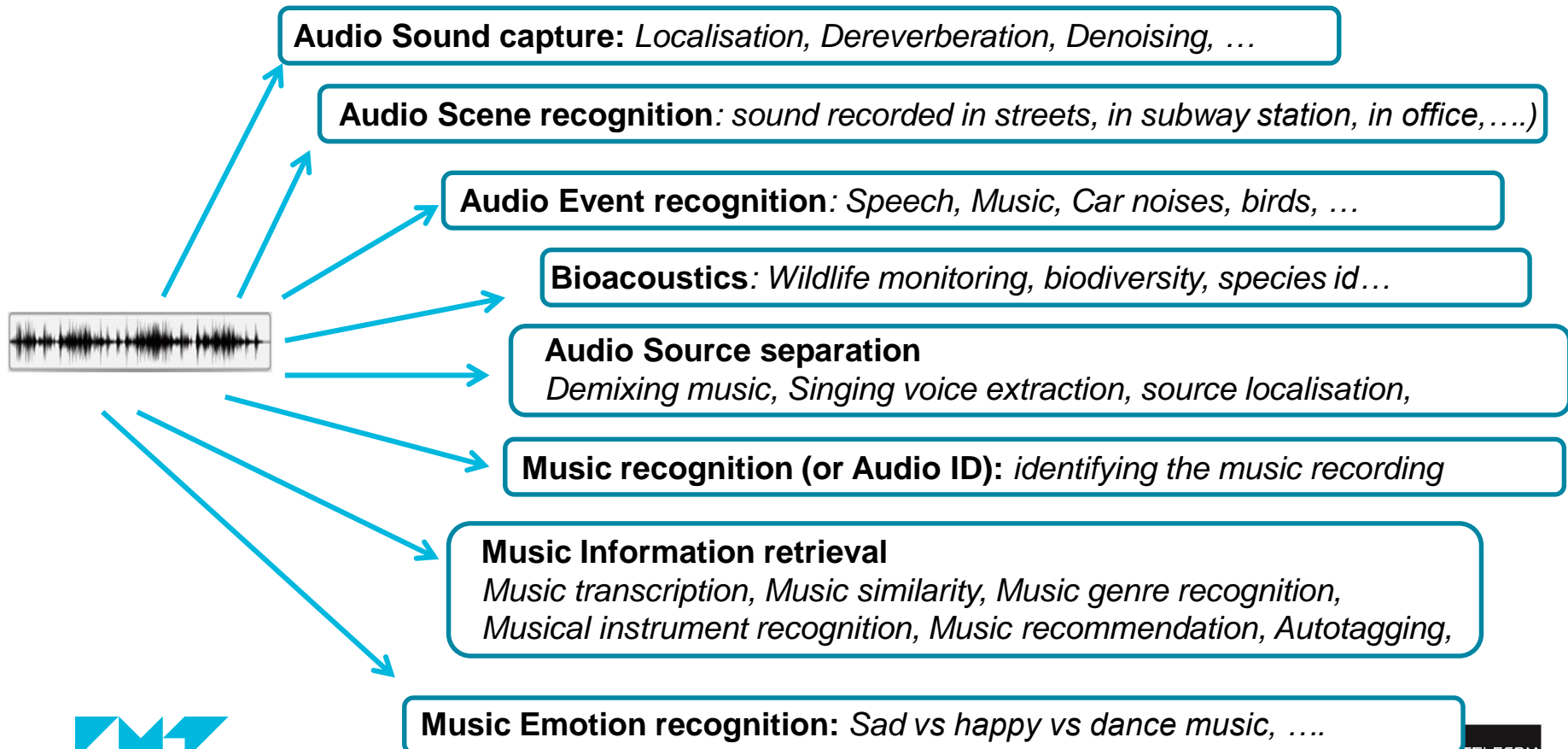
Speech Emotion Recognition



Speech Translation



Goal: to extract « **information** » from the audio signal



Music streaming,
music recommendation



Vocal separation, music separation



Music education



Bioacoustics



Music Identification, Audio Fringerprint



Karaoke, speech to rap conversion

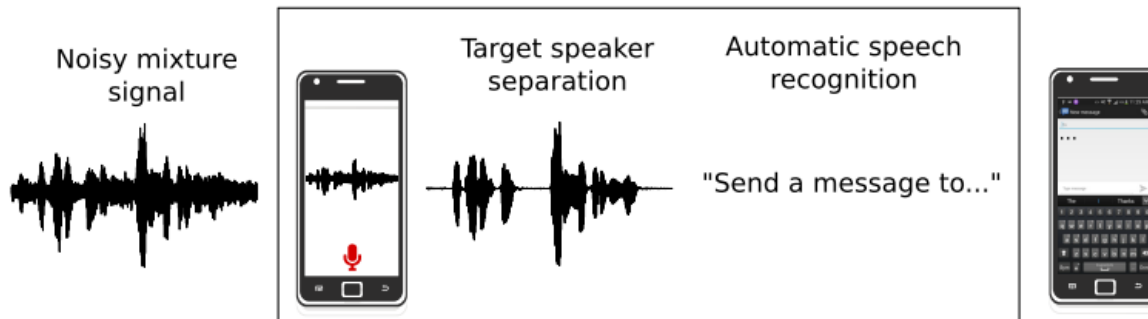


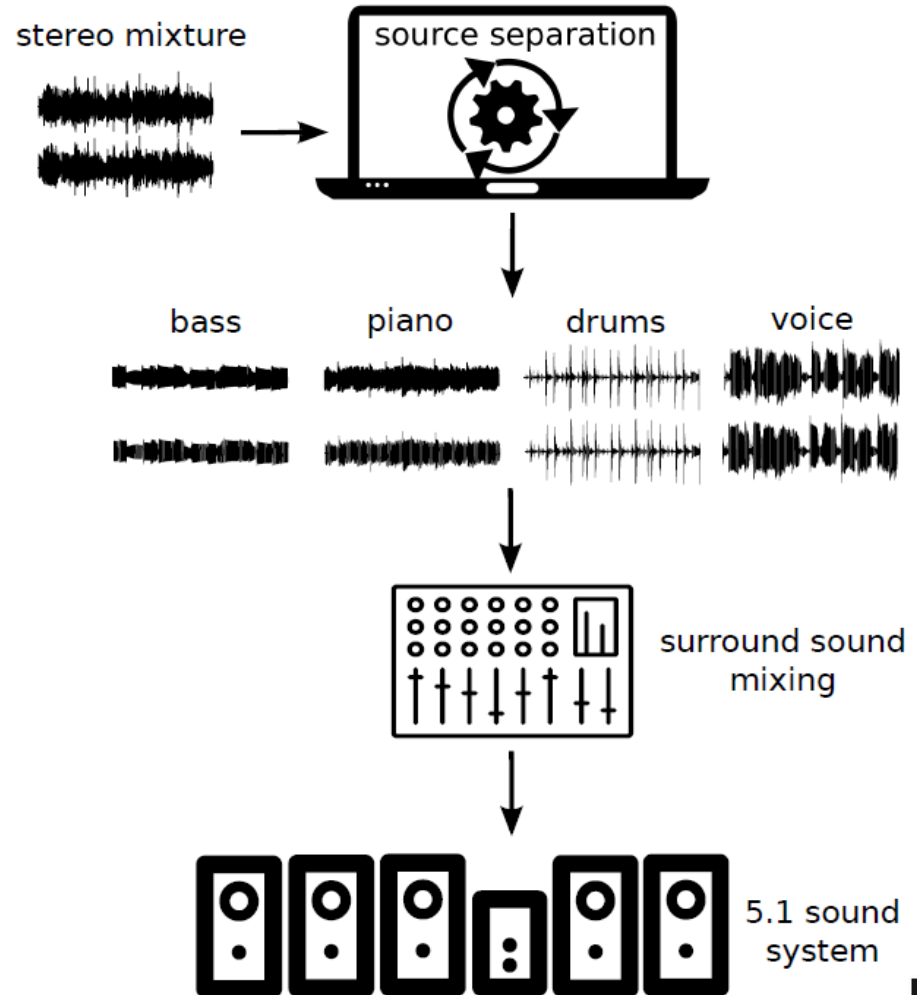
Sound recognition,
smarthomes, smart hearables



Stratégie de marque
musicale, Supervision
musicale (pub.; films)



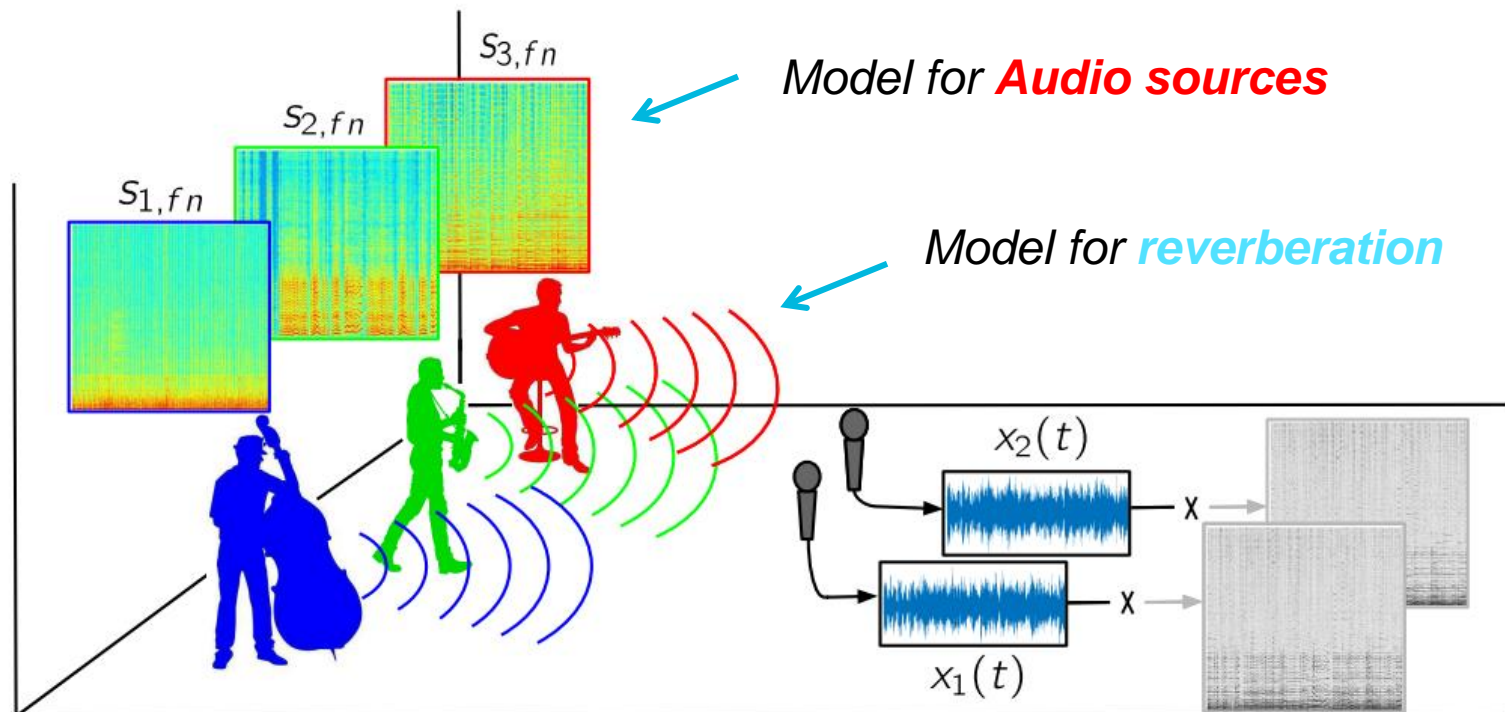




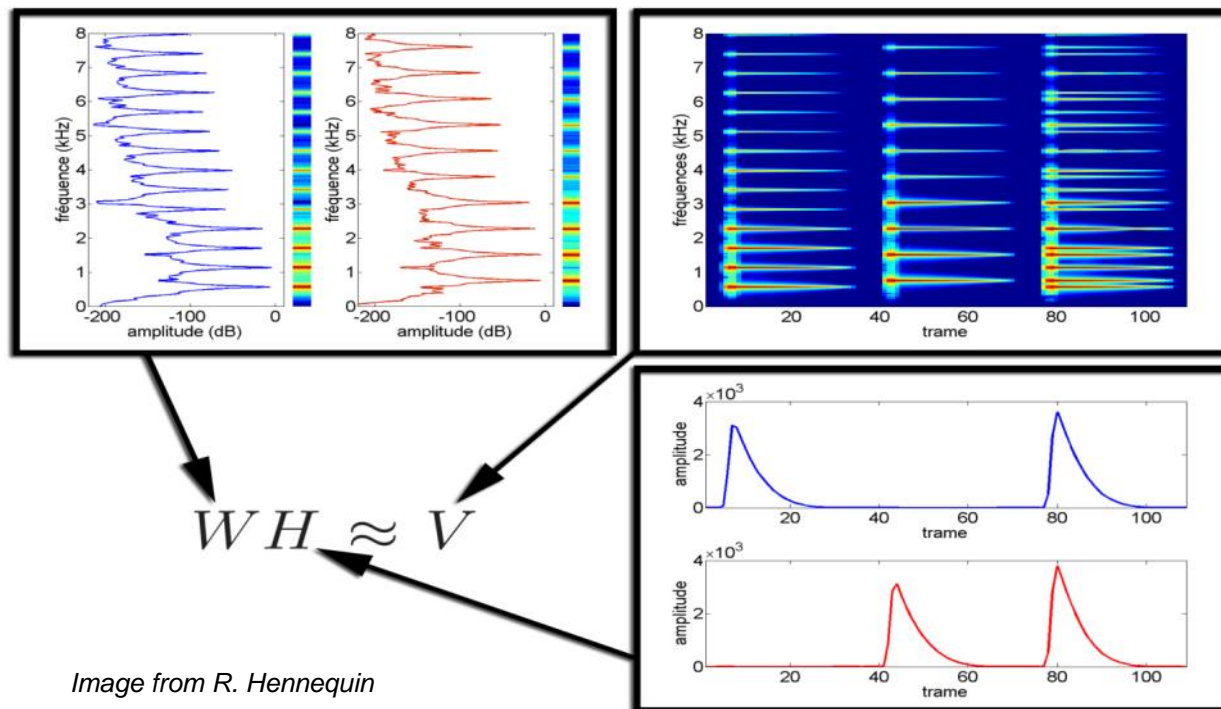
Use case of the ANR project

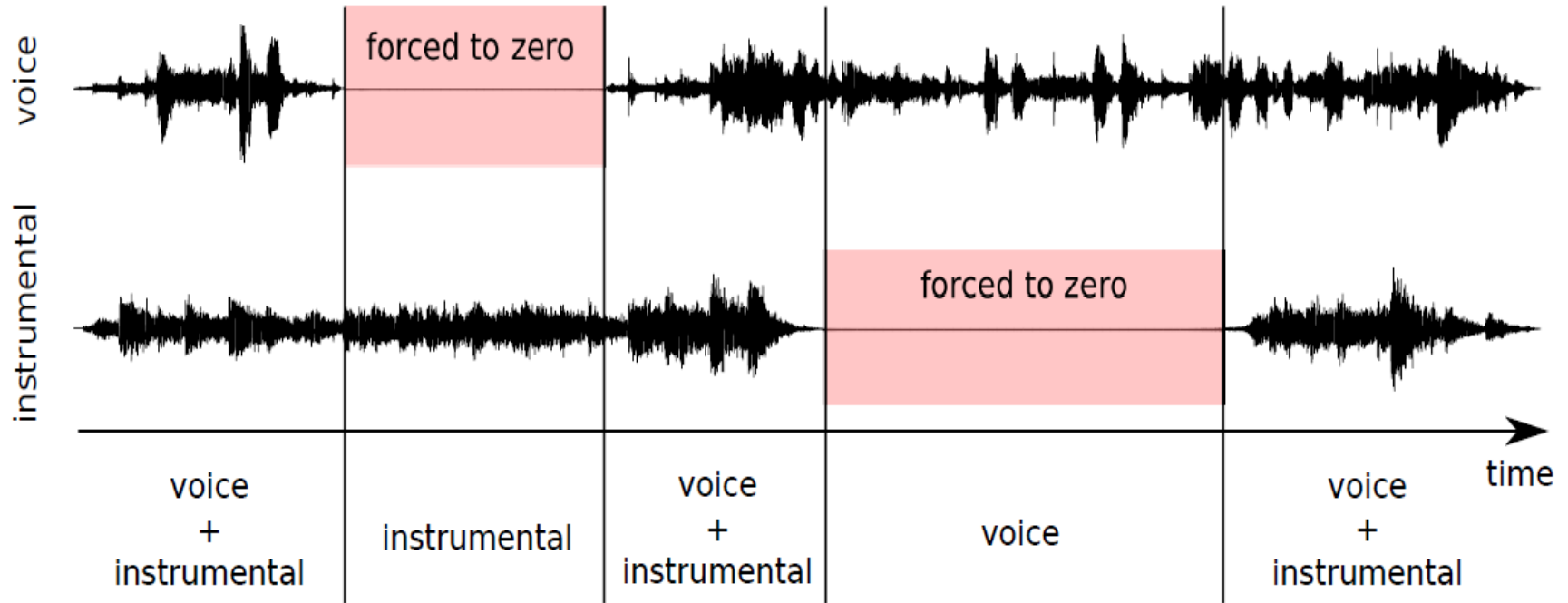
3DISON
EDISON 3D

- ▶ **Time-domain** mixture representation: $x_i(t) = \sum_{j=1}^J [a_{ij} \star s_j](t)$
- ▶ **Time-frequency** source representation: $s_j(t) = \mathcal{T}^{-1}(\{s_{j,fn}\}_{f,n})$



Non Negative Matrix factorization





S. Leglaive, R. Badeau, G. Richard, "Multichannel Audio Source Separation with Probabilistic Reverberation Priors", IEEE/ACM Transactions on Audio, Speech, and Language Processing, Vol. 24, no. 12, December 2016

Simon Leglaive, Roland Badeau, Gaël Richard, Separating Time-Frequency Sources from Time-Domain

Convulsive Mixtures Using Non-negative Matrix Factorization. WASPAA , Oct. 2017 New Paltz, US.

Audio ID = find high-level metadata from a music recording



Challenges:

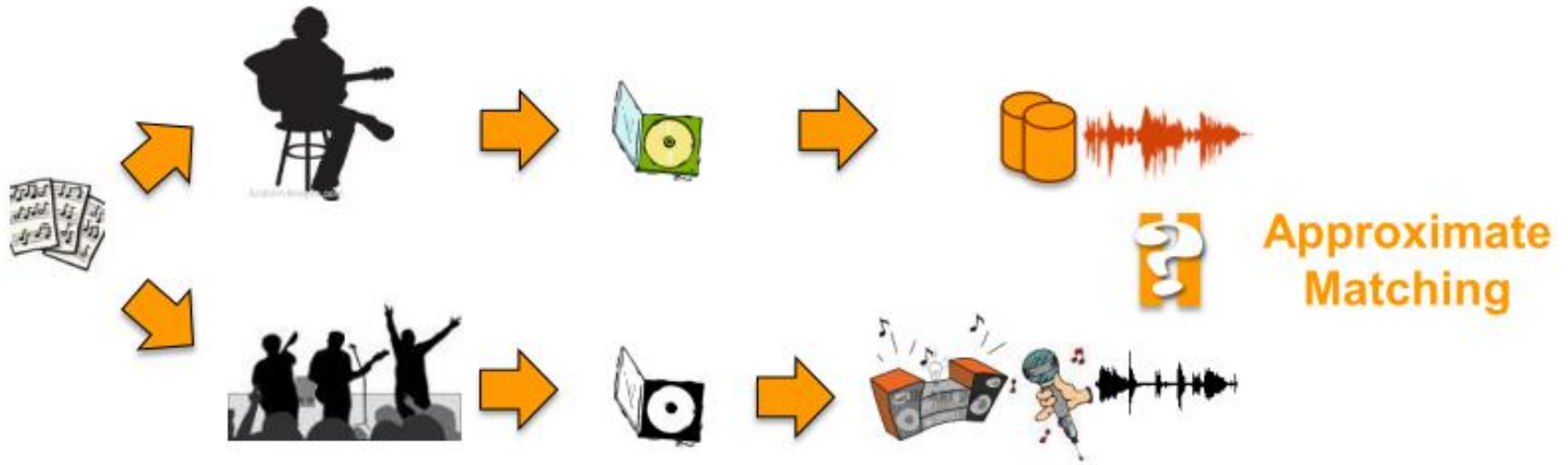
Efficiency in adverse conditions (distorsion, noises,..)

Scale to "Big data" (bases > millions of titles)

Rapidity / Real time

Product example :



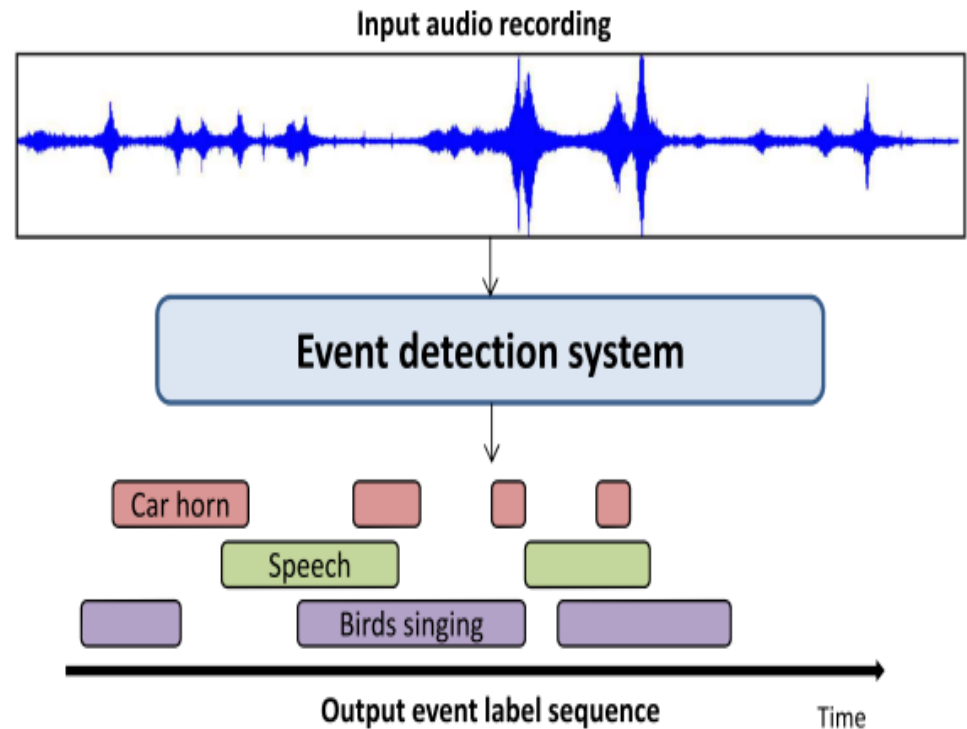
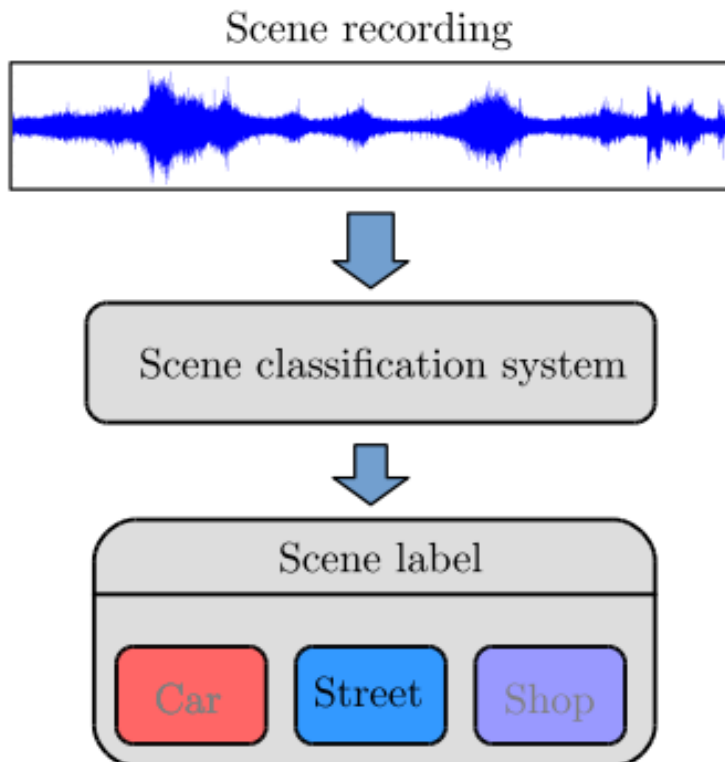


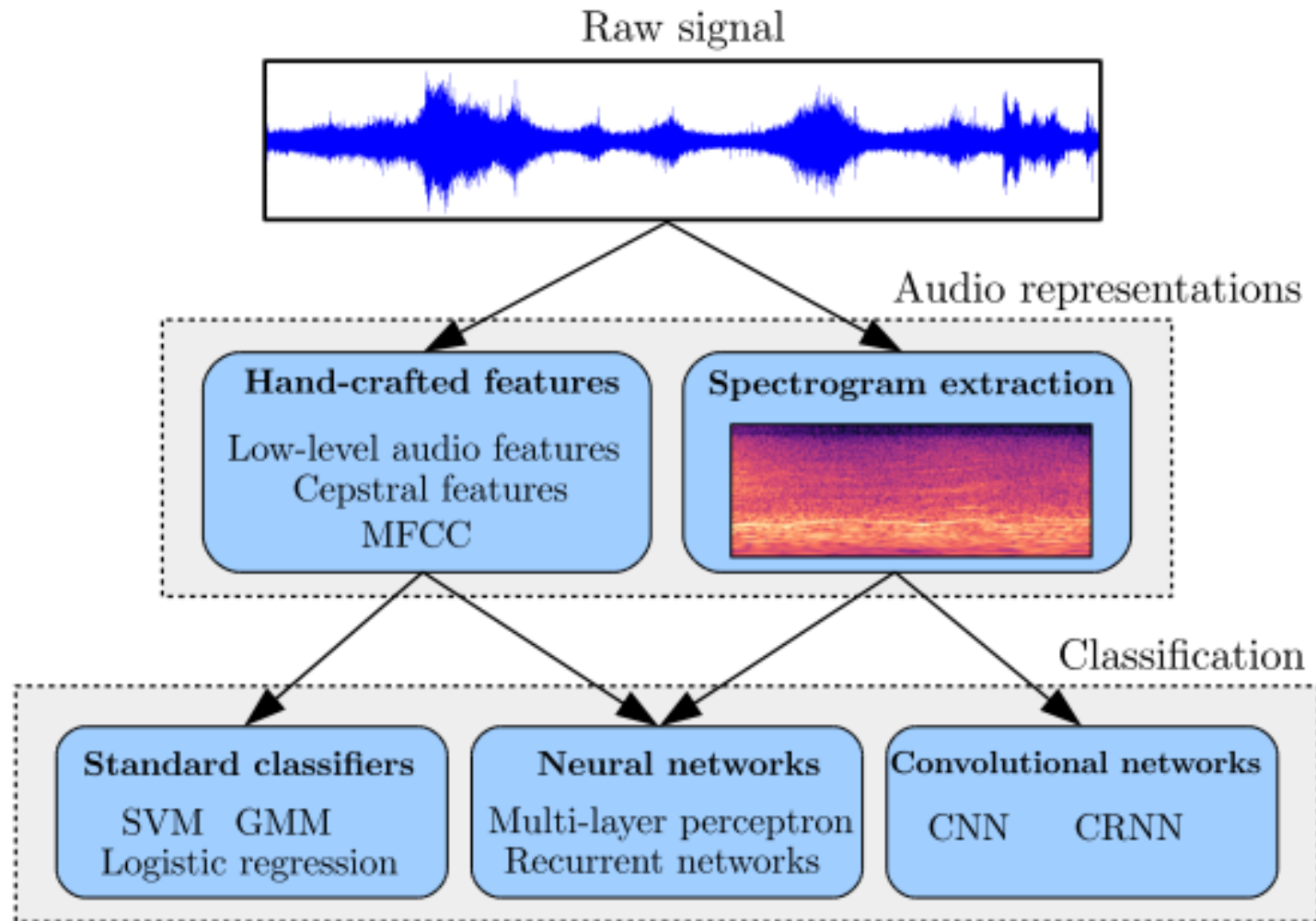
Audio recordings recognition

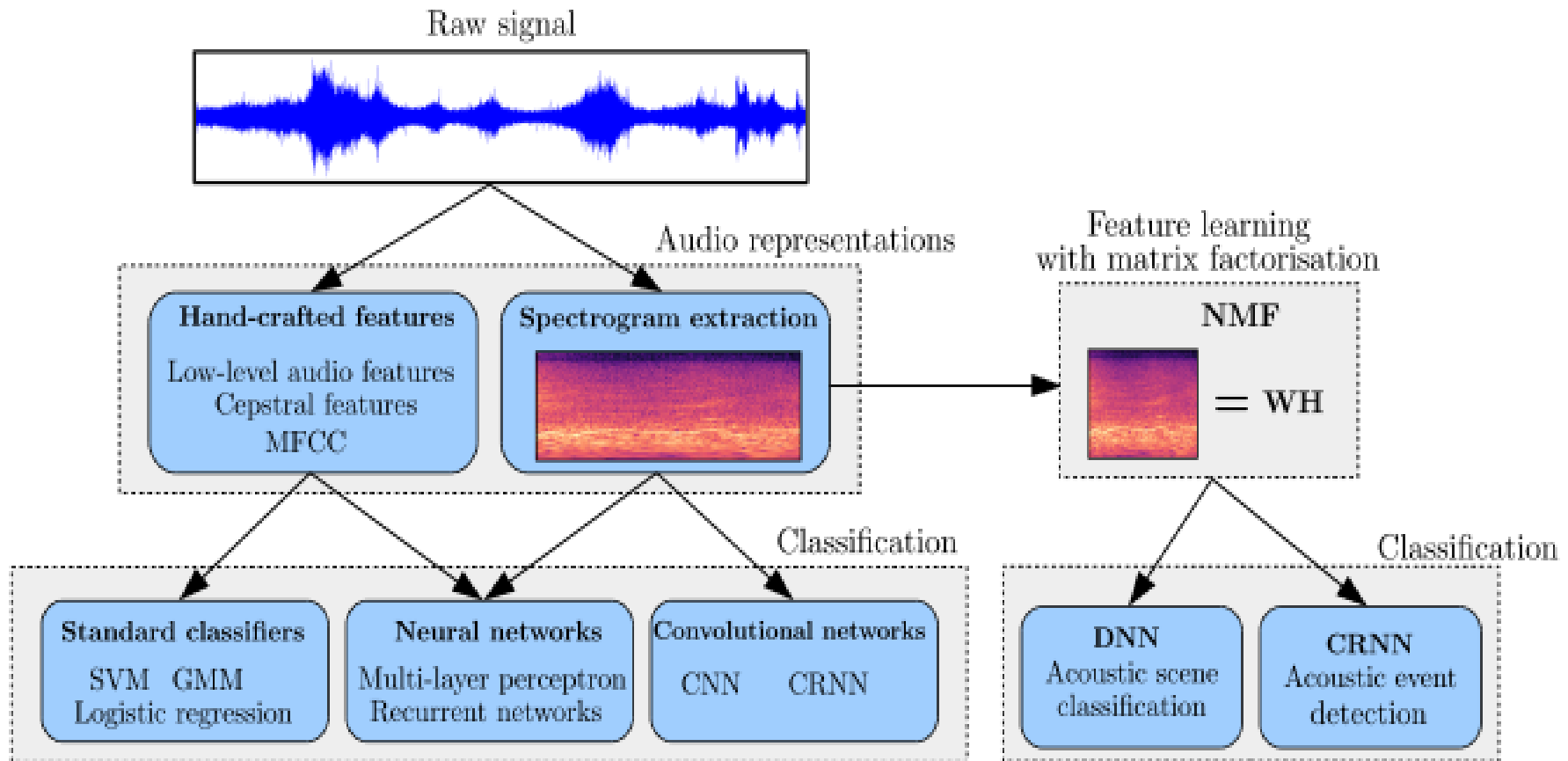
- Identical
- Approximate (live vs studio)
- Real time demonstrator
- For music recommendation, second screen applications, ...



SOUND EVENTS AND ACOUSTIC SCENE RECOGNITION



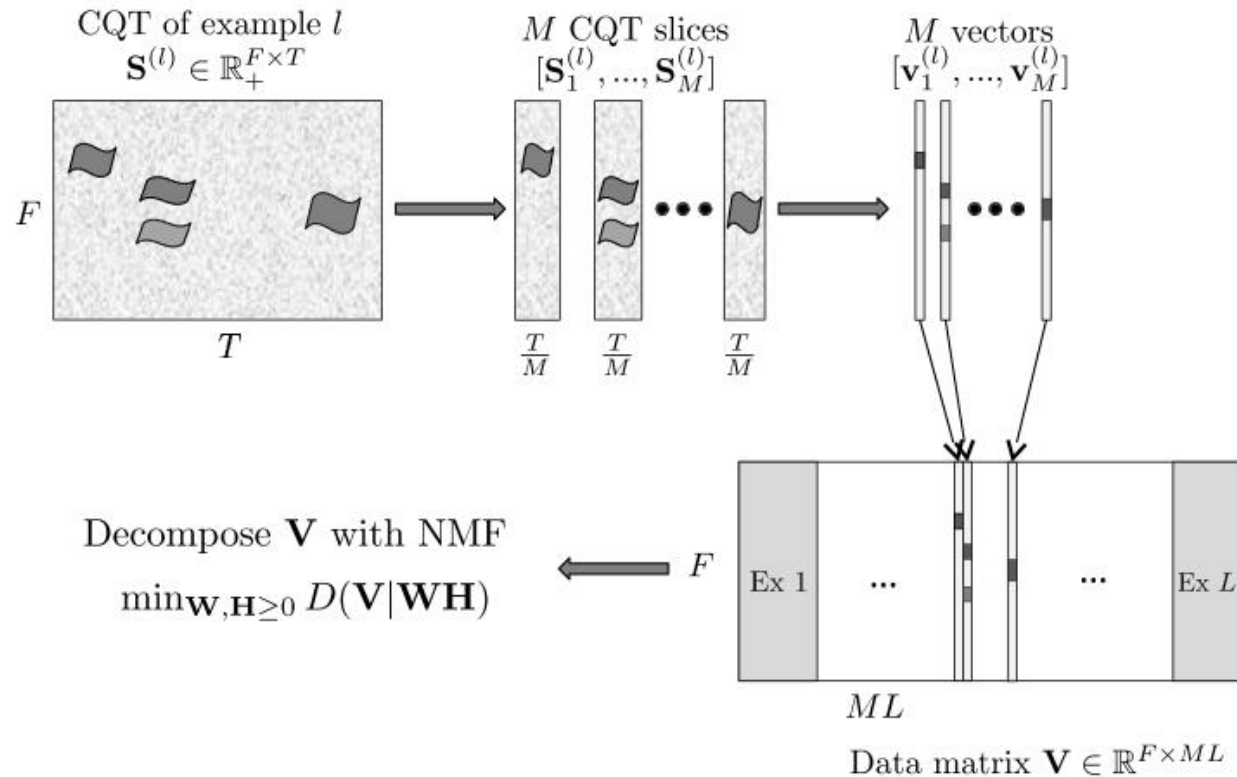




V. Bisot, R. Serizel, S. Essid, G. Richard, "Feature Learning with Matrix Factorization Applied to Acoustic Scene Classification", *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, (2017), Special Issue on *Sound Scene and Event Analysis*.

EXAMPLE FOR SCENE CLASSIFICATION

From time-frequency representations to dictionary learning

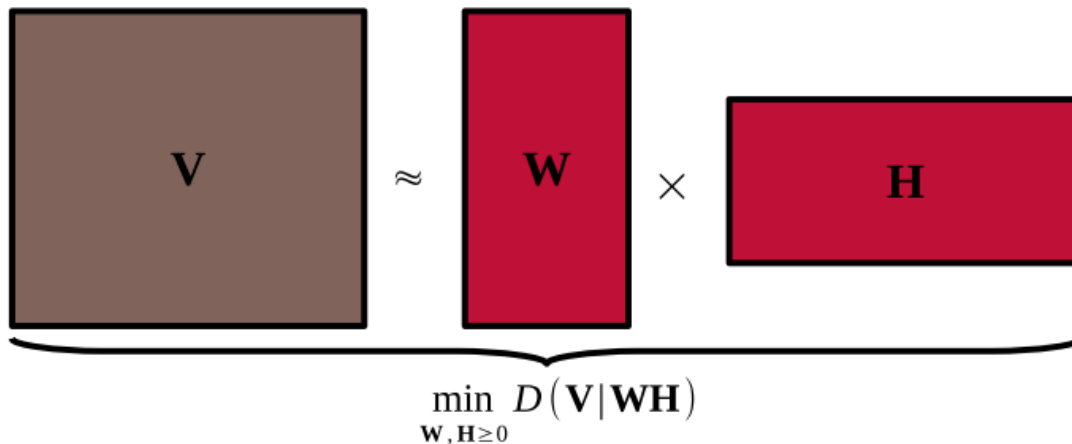


UNSUPERVISED NMF FOR ACOUSTIC SCENE RECOGNITION

Nonnegative matrix factorization

$$\min_{\mathbf{W}, \mathbf{H} \geq 0} D(\mathbf{V} | \mathbf{W}\mathbf{H}) \text{ with } \mathbf{W} \in \mathbb{R}_+^{F \times K} \text{ and } \mathbf{H} \in \mathbb{R}_+^{K \times N}$$

Dictionary learning with NMF

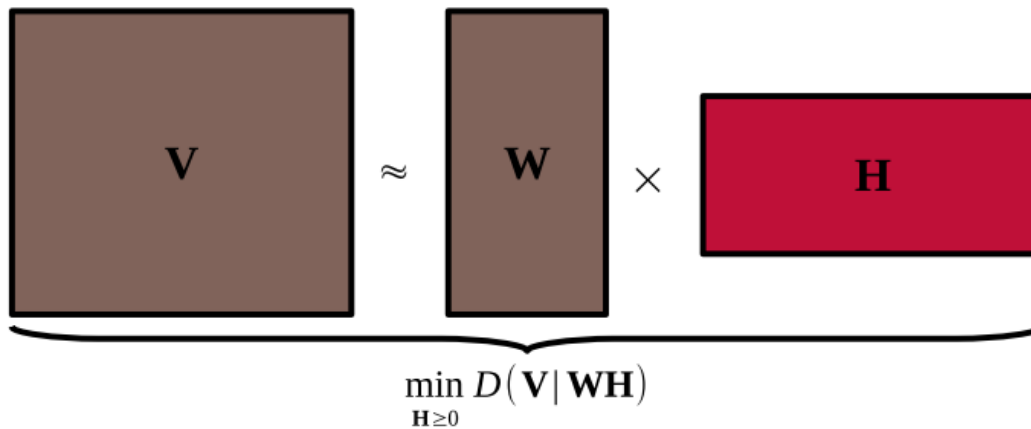


UNSUPERVISED NMF FOR ACOUSTIC SCENE RECOGNITION

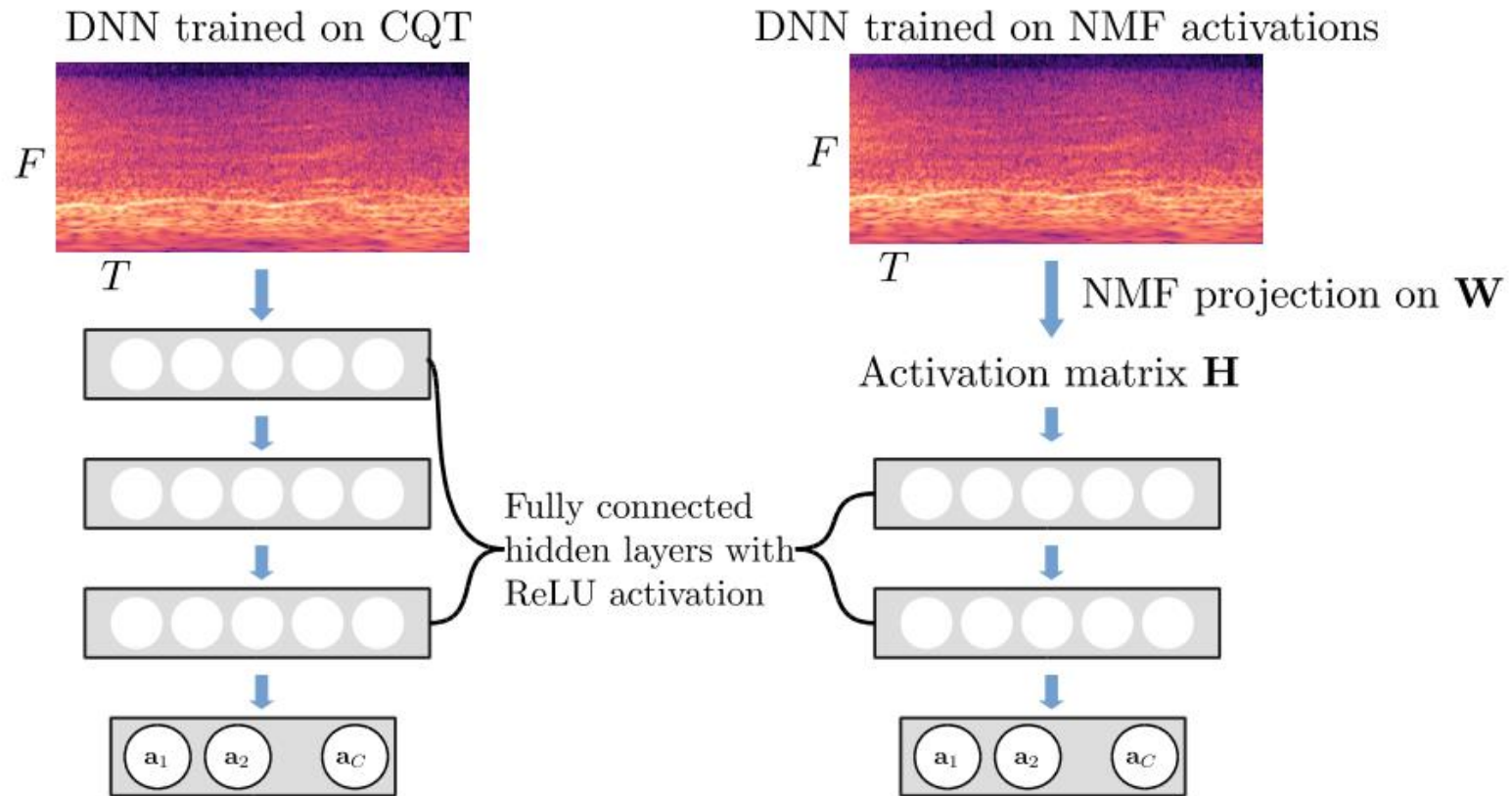
Nonnegative matrix factorization

$$\min_{\mathbf{W}, \mathbf{H} \geq 0} D(\mathbf{V} | \mathbf{W}\mathbf{H}) \text{ with } \mathbf{W} \in \mathbb{R}_+^{F \times K} \text{ and } \mathbf{H} \in \mathbb{R}_+^{K \times N}$$

Feature extraction \rightarrow project on learned dictionary



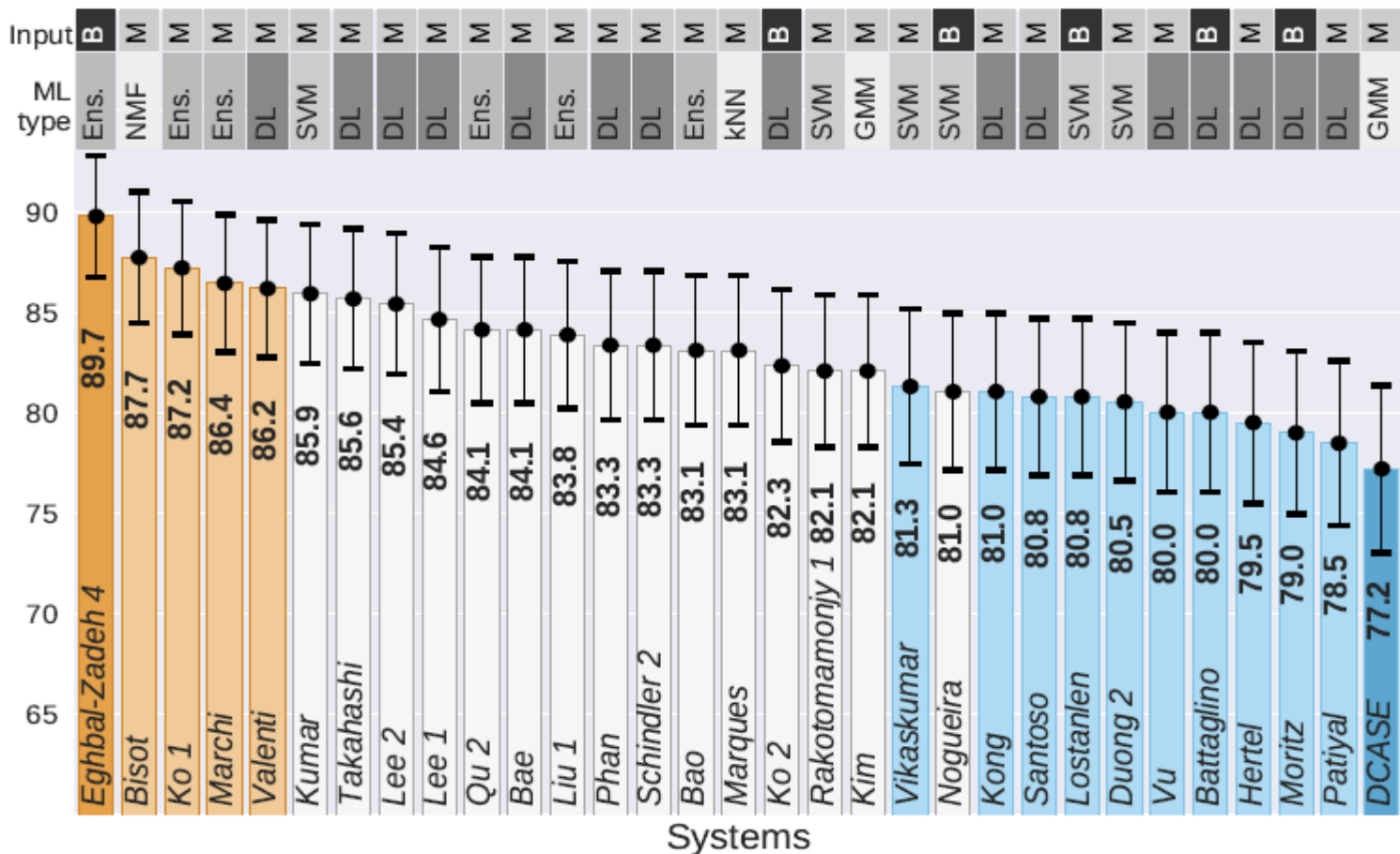
EXAMPLE WITH DNN: ACOUSTIC SCENE RECOGNITION



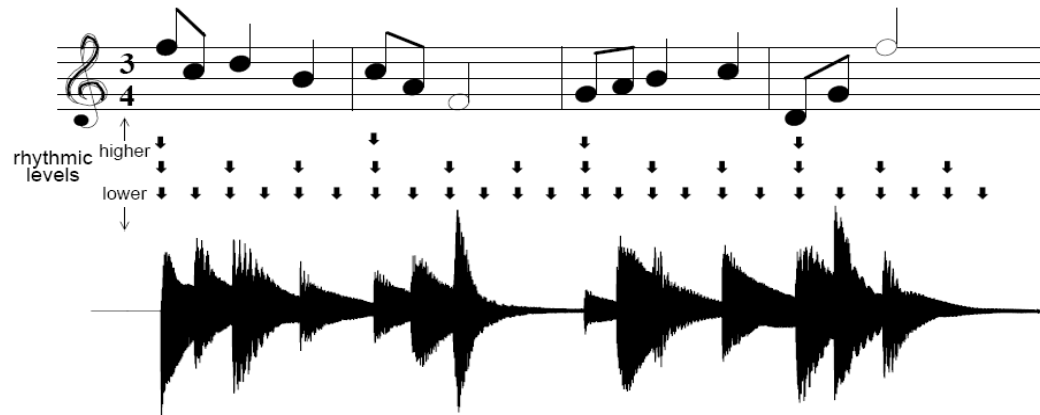
V. Bisot & al., "Feature Learning with Matrix Factorization Applied to Acoustic Scene Classification", *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, (2017),

V. Bisot & al., Leveraging deep neural networks with nonnegative representations for improved environmental sound classification *IEEE International Workshop on Machine Learning for Signal Processing MLSP*, Sep 2017, Tokyo

TYPICAL PERFORMANCES OF ACOUSTIC SCENE RECOGNITION (CHALLENGE DCASE 2016)



■ A Mesaros & al. Detection and Classification of Acoustic Scenes and Events: Outcome of the DCASE 2016 challenge IEEE/ACM Transactions on Audio, Speech, and Language Processing 26 (2), 379-393

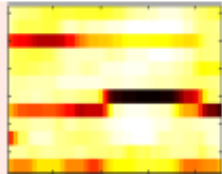
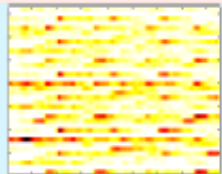
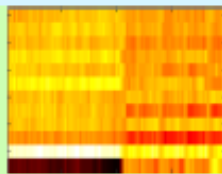
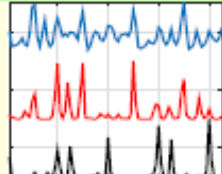
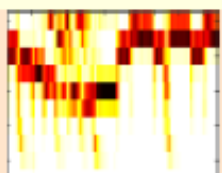


Major topics:

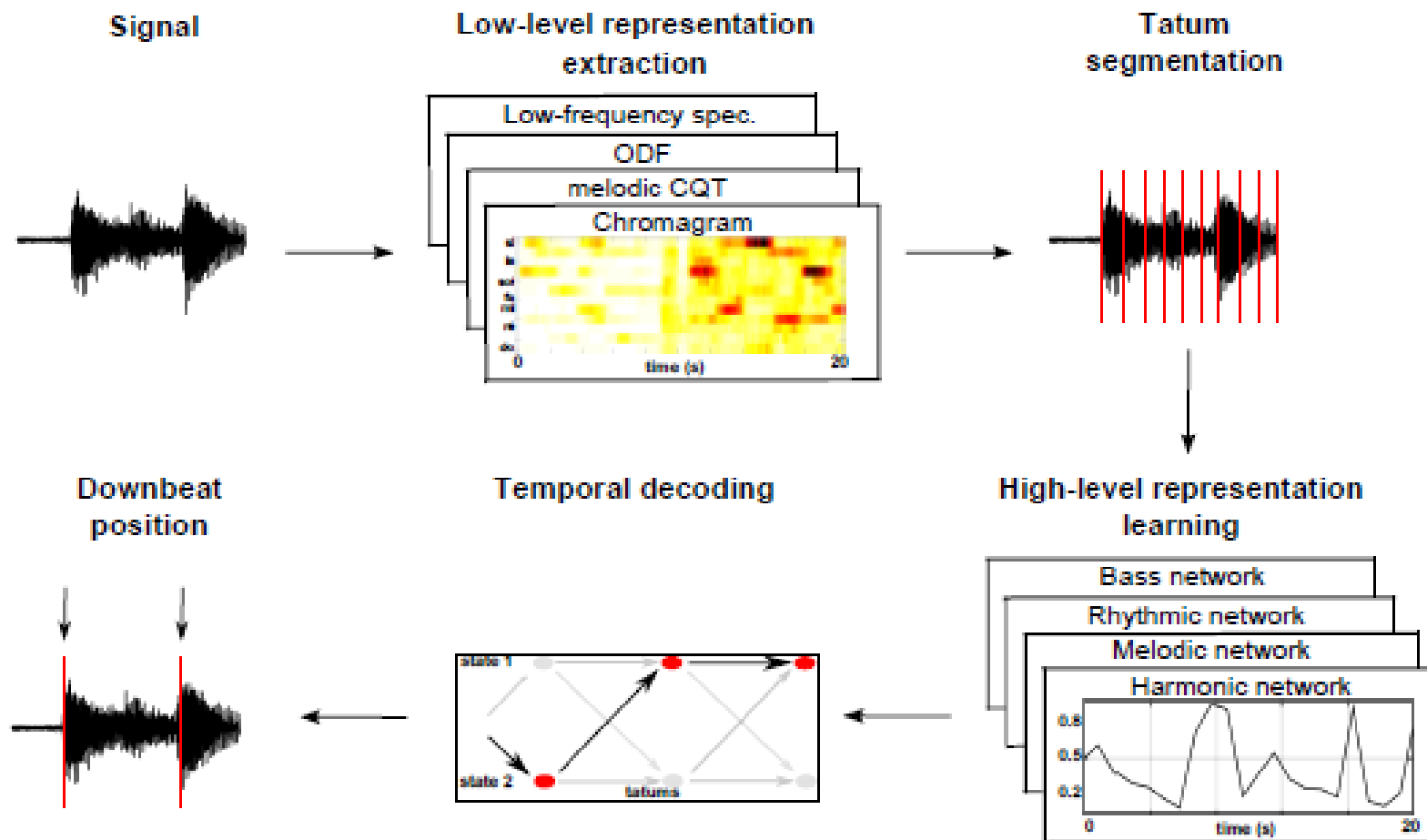
- . Music transcription (Multiple F0 estimation, Beat/Downbeat detection, instrument classification, ...),
- . Music recommendation
- . Source separation,
- . Multimodal music processing

MIR: AN EXAMPLE WITH DOWNBEAT ESTIMATION

(DURAND & AL. 2017)

Cue	Examples	Input
Harmony	Chord change, Cadence	
Melody	Melodic pattern, pivot notes	
Timbre	Section change, new instrument	
Rhythm	Bar-length rhythm patterns	
Bass content	Bass, Double bass and kick drum highlight downbeats	

MIR: AN EXAMPLE WITH DOWNBEAT ESTIMATION (DURAND & AL. 2017)



Examples at the output of each network

https://simondurand.github.io/dnn_audio.html

Other audio example

JBB (Tatum)



JBB (Downbeat)



Exemple (Tatum)



Exemple (Downbeat)



Some Current trends:

**Context-driven
Style transfer**

With  technicolor

**EEG-driven
music processing**

With  technicolor

**Conditional audio
generation**

With  Sony CSL

**Context-aware music
Recommendation**

with  DEEZER

**Text-Informed Lead
Vocal Extraction**

with  **Audionamix**
SEPARATE2CREATE®