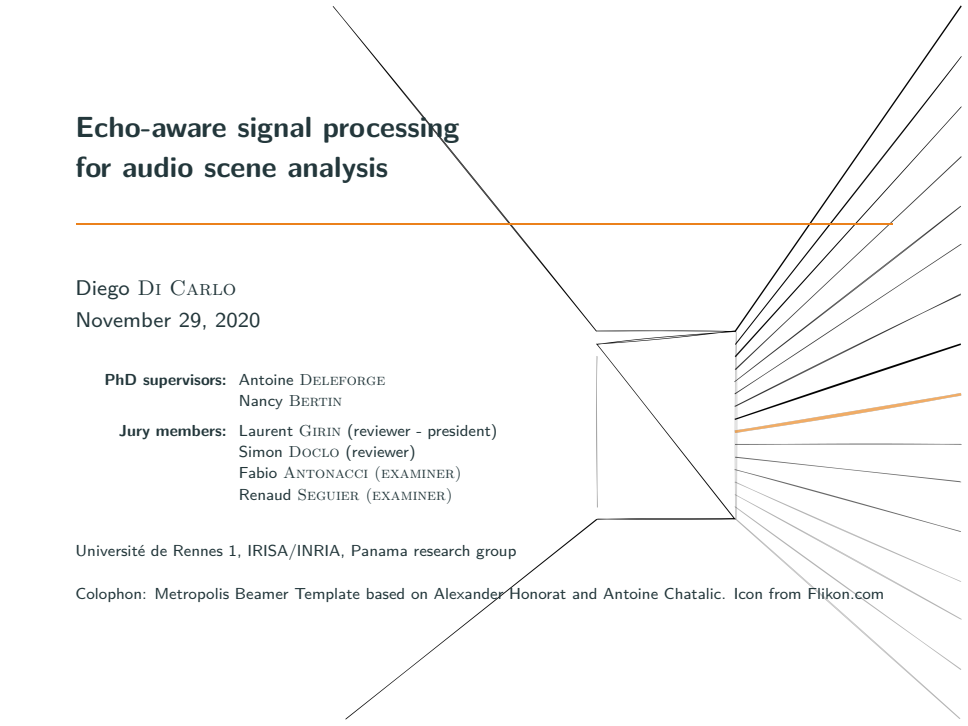# Echo-aware signal processing for audio scene analysis

Diego Di Carlo

November 29, 2020

**PhD supervisors:** Antoine Deleforge
Nancy Bertin

**Jury members:** Laurent Girin (reviewer - president)
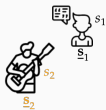Simon Doclo (reviewer)
Fabio Antonacci (examiner)
Renaud Seguier (examiner)

# Introduction

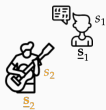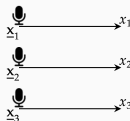# Echo-aware signal processing for audio scene analysis

Current Scenario



**Sound**

- produced by sources

# Echo-aware signal processing for audio scene analysis

Current Scenario

**Sound**
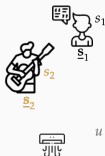
- produced by sources

- recorded by (array of) microphones

$s_1$

$\underline{s}_1$

$s_2$

$\underline{s}_2$

$\underline{x}_1$      $x_1$

$\underline{x}_2$      $x_2$

$\underline{x}_3$      $x_3$

1

# Echo-aware signal processing for audio scene analysis
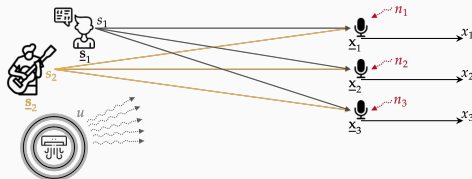
Current Scenario



**Sound**

- produced by sources
- recorded by (array of) microphones
- corrupted by noise

1

# Echo-aware signal processing for audio scene analysis

Current Scenario



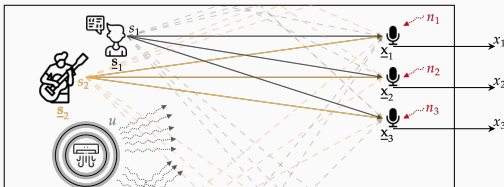**Sound**

- produced by sources

- recorded by (array of) microphones

- corrupted by noise

- propagates in the space

1

# Echo-aware signal processing for audio scene analysis

Current Scenario



**Sound**

- produced by sources

- recorded by (array of) microphones

- corrupted by noise

- propagates in the space

- interacts with the room
  ↪ reverberation

1

# Echo-aware signal processing for audio scene analysis

Semantic information



on nature and content

# Echo-aware signal processing for audio scene analysis

Semantic information

Spatial information



on nature and content

on position and geometry
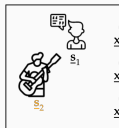
# Echo-aware signal processing for audio scene analysis

Semantic information



on nature and content

Spatial information



on position and geometry

Temporal information



on events activity

# Echo-aware signal processing for audio scene analysis

Semantic information



on nature and content

Spatial information



on position and geometry

Temporal information



on events activity

**Audio Scene Analysis**

Extraction and organization of all the information in the sound

2

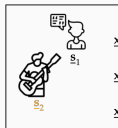# Echo-aware signal processing for audio scene analysis

### Semantic information



on nature and content

### Spatial information



on position and geometry

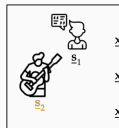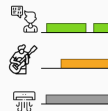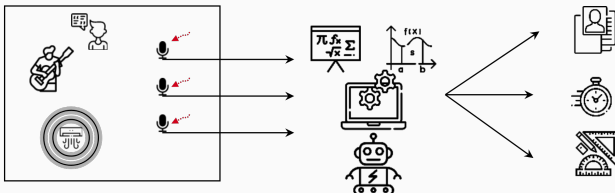### Temporal information



on events activity

**Audio Scene Analysis**

Extraction and organization of all the information in the sound
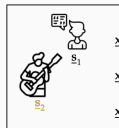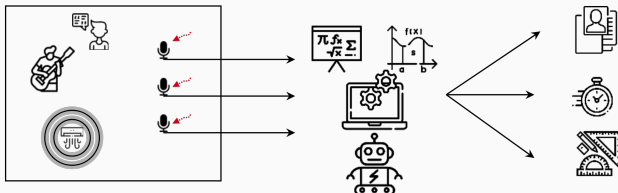


**Can computer do it?**

2

# Echo-aware signal processing for audio scene analysis

# Echo-aware signal processing for audio scene analysis



## Signal Processing

Mathematical models, frameworks and tools to tackle and solve such problems

# Echo-aware signal processing for audio scene analysis



> **Signal Processing**
>
> Mathematical models, frameworks and tools to tackle and solve such problems

Some (inverse) problems

- Speaker Identification
- Sound Source Separation (SSS)
- Speech Enhancement (SE)
- Automatic Speech Recognition (ASR)
- Sound Source Localization (SSL)

- Voice Activity Detection
- Diarization
- $RT_{60}$ estimation
- Acoustic Channel Estimation
- Wall Absorption Estimation

3

# Echo-aware signal processing for audio scene analysis



### Signal Processing

Mathematical models, frameworks and tools to tackle and solve such problems

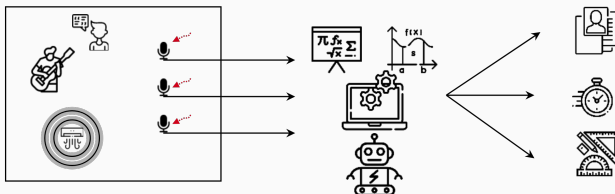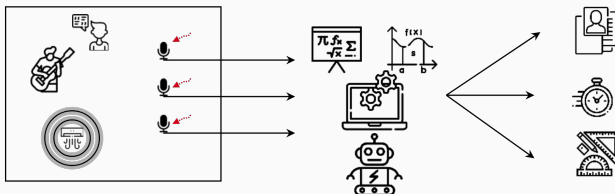Some (inverse) problems

- Speaker Identification

- Sound Source Separation (SSS)

- Speech Enhancement (SE)

- Automatic Speech Recognition (ASR)

- Voice Activity Detection

- Diarization

- $RT_{60}$ estimation

- Acoustic Channel Estimation

- Wall Absorption Estimation

3

# Echo-aware signal processing for audio scene analysis

**Sound interacts with indoor environment:**

it is reflected,
　specularly and diffusely

+ it is absorbed,
$\Bigg\}$ = all reverberation

+ it is transmitted,

+ and other.

---

**Acoustic Echoes**

- Elements of reverberation
- Specular reflection standing out for time and strength
- Repetition of a sound but after
  - same content
  - delay $\Leftrightarrow$ distance

**Everyday examples:**

## Outline and contributions

**Thesis title**
Audio Scene Analysis
↓
context and problems

## Outline and contributions

**Thesis title**

Audio Scene Analysis

↓

context and problems

Signal Processing

↓

models and frameworks

5

## Outline and contributions

**Thesis title**

Audio Scene Analysis

↓

context and problems

Signal Processing

↓

models and frameworks

Echo-aware

↓

better processing

## Outline and contributions

**Thesis title**

Audio Scene Analysis          Signal Processing          Echo-aware

↓                    ↓                    ↓

context and problems      models and frameworks      better processing

**Thesis content:**

How to estimate them?                    How to use them?

- Analytical method                    - Source Localization

- Learning-based method                - Source Separation

  no parameter tuning                  - Speech Enhancement

  no full sound modeling               - Room Geometry Estimation

Where to find them?
Echo-aware database for
estimation and application

5

# Problem Statement

## Signal model

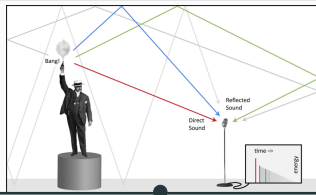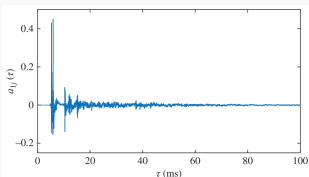Sound propagation process $\Leftrightarrow$ Source $\rightarrow$ Filter $\rightarrow$ Receiver model

source signal

microphone signal $\longleftarrow$ $\tilde{x}_i(t) = (\tilde{h}_i * \tilde{s})(t) + +\tilde{n}(t)$ $\longrightarrow$ noise term

continuous-time convolution

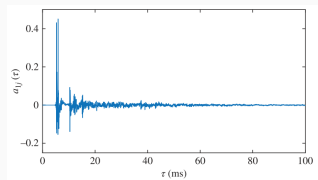⚠ continuous time

### Room Impulse Response (RIR)

- linear filtering effect of the sound
- acoustic response of a room to a (prefect) impulsive sound
- depends on spatial properties (room geometry, mic/src position)

6

## Echoes in the RIR

RIR model

$$\tilde{h}_i(t) = \tilde{h}_i^{\mathsf{d}}(t) + \tilde{h}_i^{\mathsf{e}}(t) + \tilde{h}_i^{\mathsf{lrev}}(t) + \varepsilon_i(t)$$



Echoes can be modeled as sum of Dirac's delta

$$\tilde{h}_i^{\mathsf{echoes}} = \tilde{h}_i^{\mathsf{d}}(t) + \tilde{h}_i^{\mathsf{e}}(t) \approx \sum_{r=0}^{R} \alpha_i^{(r)} \delta(t - \tau_i^{(r)})$$

**Goal:** estimated the $\tau_{i_{i,r}}$

**Challenges:**

- $\alpha$ distortion (even if we know it $\implies$ labeling)
- $\alpha \to \alpha(t)$ (sum of diracs $\to$ sum of filters)
- $h_l$ reverberation is included in the noise term
- depends on the scene geometry (room, source and mic position)

# References i