

# ECHO-AWARE signal processing for audio scene analysis

---

Diego DI CARLO

November 24, 2020

supervisors: Antione DELEFORGE, Nancy BERTIN

collaborators: Clément ELVIRA, Robin SCHEIBLER, Ivan DOKMANIĆ, Sharon GANNOT, Pini A

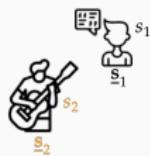
INRIA IRISA

## Introduction

---

# Scenario

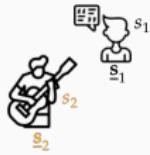
## Sound



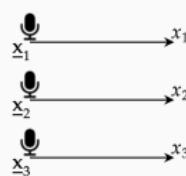
- produced by **sources**

Attention: artificial sound vs (natural) microphone recordings

# Scenario



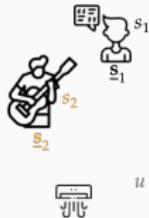
## Sound



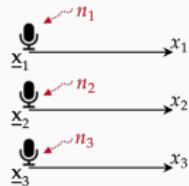
- produced by **sources**
- recorded by **microphones**

Attention: artificial sound vs (natural) microphone recordings

# Scenario



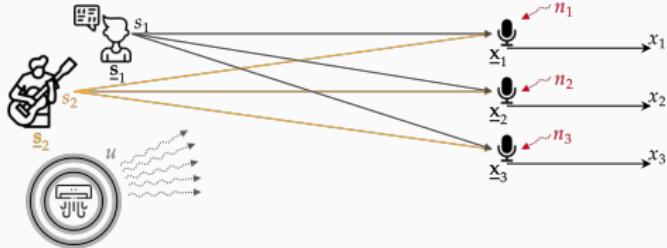
## Sound



- produced by **sources**
- recorded by **microphones**
- corrupted by **noise**

Attention: artificial sound vs (natural) microphone recordings

# Scenario

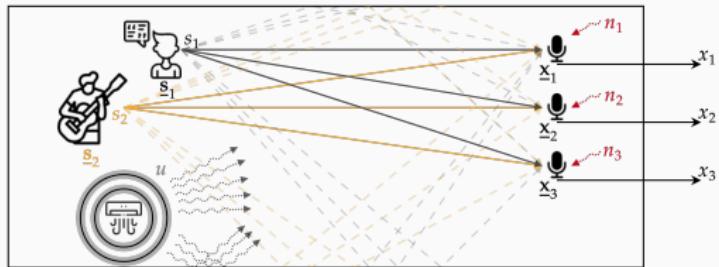


## Sound

- produced by **sources**
- recorded by **microphones**
- corrupted by **noise**
- propagates in the **space**

Attention: artificial sound vs (natural) microphone recordings

# Scenario



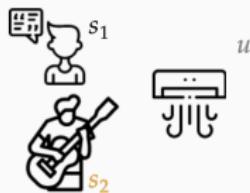
## Sound

- produced by **sources**
- recorded by **microphones**
- corrupted by **noise**
- propagates in the **room**  
     $\hookrightarrow$  **reverberation**

Attention: artificial sound vs (natural) microphone recordings

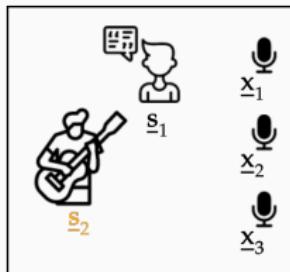
# Echo-aware signal processing for audio scene analysis

## Semantic information



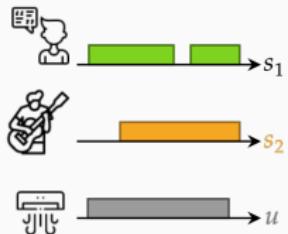
about source nature and semantic content

## Spatial information



about source position and room geometry

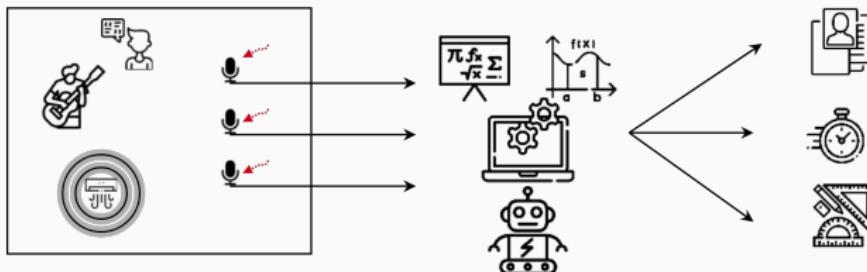
## Temporal information

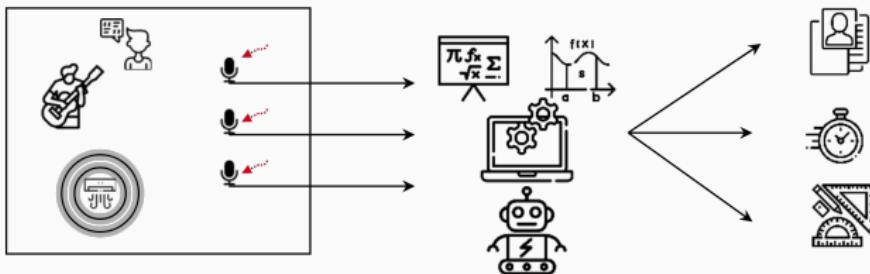


about events activity

## Audio Scene Analysis

Extraction and organization of all the information in the sound





## Signal Processing

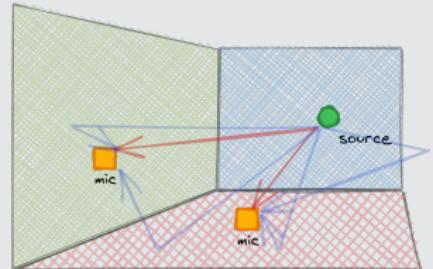
Mathematical models, frameworks and tools to tackle and solve such problems

Some (inverse) problems

- Speaker Identification
  - Sound Source Separation (SSS)
  - Speech Enhancement (SE)
  - Automatic Speech Recognition (ASR)
  - Sound Source Localization (SSL)
  - Room Geometry Estimation (RooGE)
- Who?
- What?
- Where?
- Voice Activity Detection
  - Diarization
  - $RT_{60}$  estimation
  - Acoustic Channel Estimation
  - Wall Absorption Estimation
  - *and many many other*
- When?
- How?

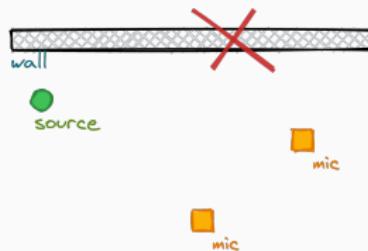
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



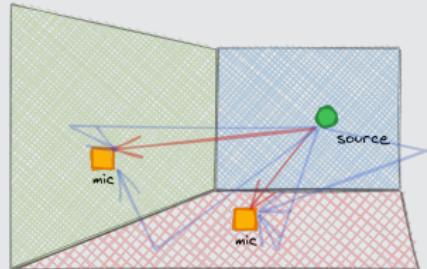
## Audio signal processing methods

- ignore it



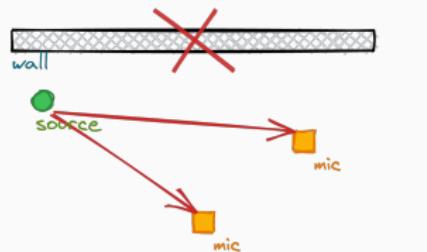
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



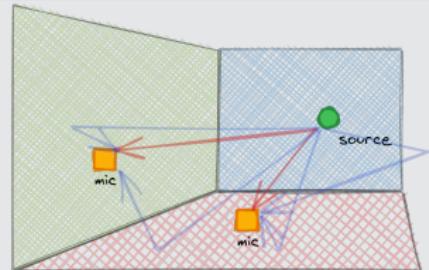
## Audio signal processing methods

- ignore it
- assume it free-field



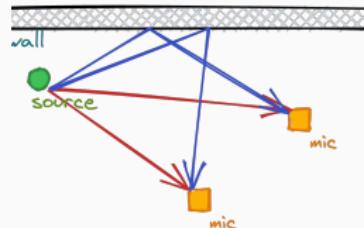
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



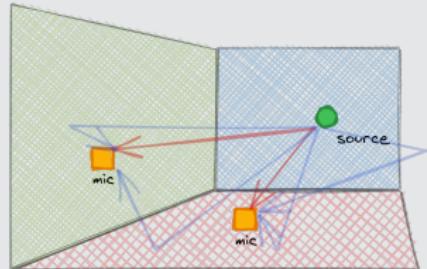
## Audio signal processing methods

- ignore it
- assume it free-field
- model it entirely



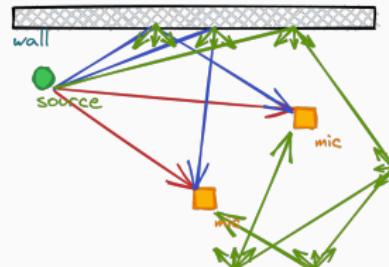
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



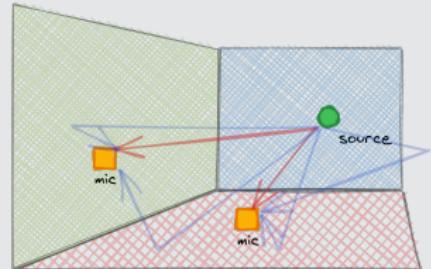
## Audio signal processing methods

- ignore it
- assume it free-field
- model it entirely
- model as few reflection



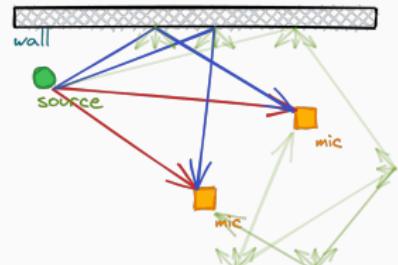
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



## Audio signal processing methods

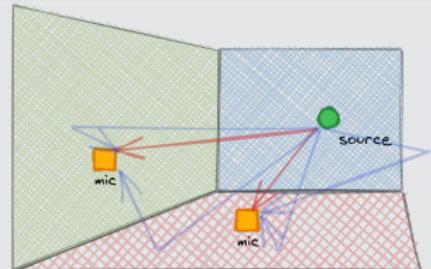
- ignore it
- assume it free-field
- model it entirely
- model as few reflection
- model it as early and late parts



# Echo-aware signal processing for audio scene analysis

## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor

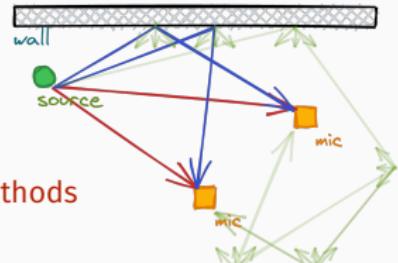


## Audio signal processing methods

- ignore it
- assume it free-field
- model it entirely
- model as few reflection
- model it as early and late parts

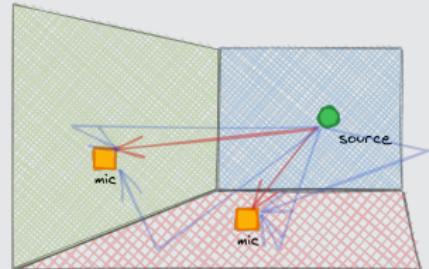
}

**Echo-aware methods**



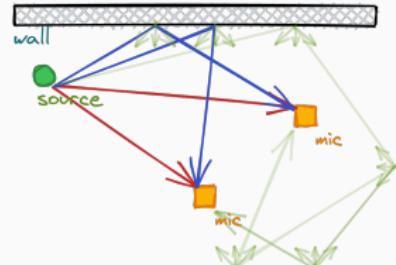
## Acoustic Echoes

- Elements of the sound propagation
- Standing out for time and strength
- Repetition of the source sound but later
- Both outdoor and indoor



## Audio signal processing methods

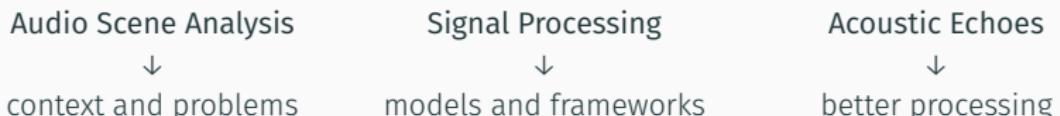
- ignore it
- assume it free-field
- model it entirely
- model as few reflection
- model it as early and late parts



## Model sound field

- as free field, then reverberation is noise
- entirely, is a too challenging

# Thesis goal and contribution



## Challenges and Objective

1. How to estimate acoustic echoes?
2. How to extend methods for echo-aware audio scene analysis

### 1. Estimation

- Knowledge-based echo estimation  
↪ **Blaster**
- Learning-based echo estimation  
↪ **Lantern**

### 2. Application

- Echo-aware Source Separation  
↪ **Separake**
- Echo-aware Source Localization  
↪ **Mirage**
- Echo-aware Speech Enhancement
- Echo-aware Room Geometry Estimation

### 3. Data:

Echo-aware database → **dEchorate**

## Modeling

---

# Acoustic Impulse Response

Sound propagates

Sound source

→

?

→

microphone

Sound propagates

$s$

→

$h$

→

$x$

# Acoustic Impulse Response

Sound propagates  
Sound source → environment → microphone

Sound propagates  
 $s \rightarrow h \rightarrow x$

# Acoustic Impulse Response

Sound propagates  
Sound source → room → microphone

Sound propagates  
 $s \rightarrow h \rightarrow x$

# Echoes and Room Impulse Response

RIRs can be modeled with the Image Methods

- specular reflection only
- “playing billiard in a concert hall”
- for shoebox room it is the solution for physics
- in frequency domain it writes as

RIRs accounts for  
the **geometry** of the room

- Room shape and size
- Mic and Source position
- presence of objects

the acoustic properties of the audio scene

- surface materials
- objects materials

examples



## Acoustic Echo Estimation

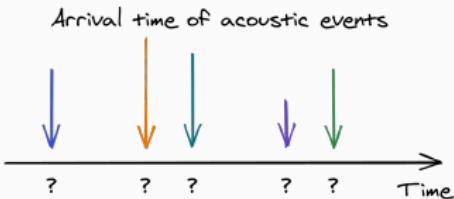
---

# Acoustic Echo Retrieval

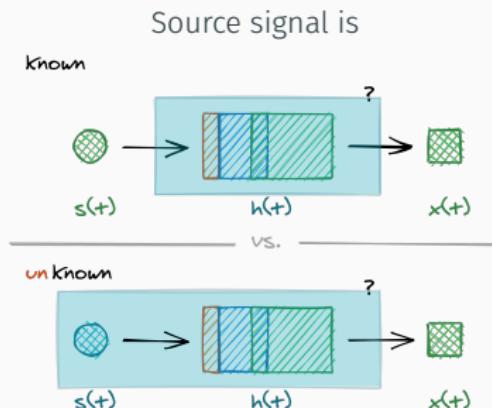
## The acoustic echoes retrieval (AER) problem

Estimating early (strong) acoustic reflections:

- their time of arrivals → TOAs Estimation  
↪ sufficient sometimes
- their amplitude  
↪ closed-form form TOA



## Approaches

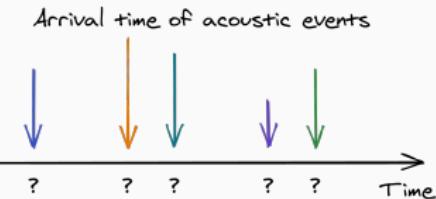


# Acoustic Echo Retrieval

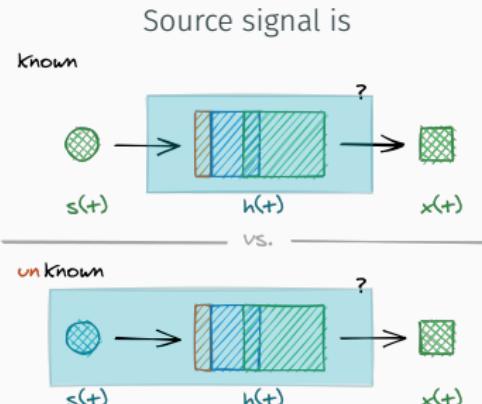
## The acoustic echoes retrieval (AER) problem

Estimating early (strong) acoustic reflections:

- their time of arrivals → TOAs Estimation  
↪ sufficient sometimes
- their amplitude  
↪ closed-form form TOA



Approaches

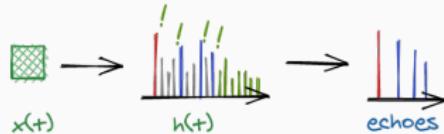


Scenario: signal source, only TOAs and passive system

# Passive Acoustic Echo Estimation

## Passive Acoustic Echo Estimation:

### RIR-based approaches



1. SIMO BCE problem  $\Rightarrow$  RIRs
2. Peak picking and *disambiguation*  $\Rightarrow$  Echoes

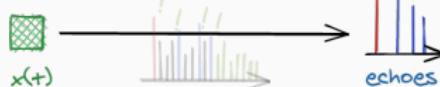
#### Pros

- SIMO BCE is well studied (elegant framework)
- It works well in some scenarios and in practice  
↪ if not limitation

#### Cons

- Full RIR
- dependent of manually tuned peak picking
- Pathological issue (sampling and body-guard)
- Complexity

### RIR-agnostic approaches



1. Estimation directly in the echoes parameters space  $\{\tau, \alpha\}$  and direction of arrivals can be used instead

#### Performed with

- Cross-correlation on-grid, eg. EM, Acoustic Cameras
- Cross-relation with super-resolution off-grid, [?, ?]

#### Pro

- No need for full RIRs
- Sub-sampling accuracy
- Low complexity
- Sparsity and Non-negativity are respected

#### Cons

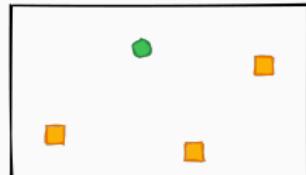
- Exploratory

# AER as discrete SIMO BCE

Key ingredient – *Cross relation identity*

$$x_i = h_i * s$$

$$h_2 * x_1 = h_2 * h_1 * s = h_1 * h_2 * s = h_1 * x_2$$



Ideas:

1. Sampled version of  $x_1, x_2$  are available ( $\mathbf{x}_1, \mathbf{x}_2$ )
2. echoes' TOAs  $\propto$  sampling frequency
3. Find echoes  $\rightarrow$  find sparse vectors  $\mathbf{h}_1, \mathbf{h}_2$  of length  $L$
4. Modeled as Lasso-like problem

$$\widehat{\mathbf{h}}_1, \widehat{\mathbf{h}}_2 \in \arg \min_{\mathbf{h}_1, \mathbf{h}_2 \in \mathbf{R}^n} \|\mathbf{x}_1 * \mathbf{h}_2 - \mathbf{x}_2 * \mathbf{h}_1\|_2^2 + \lambda \mathcal{P}(\mathbf{h}_1, \mathbf{h}_2) \quad \text{s.t.} \quad \mathcal{C}(\mathbf{h}_1, \mathbf{h}_2)$$

$\mathcal{P}(\mathbf{h}_1, \mathbf{h}_2) \rightarrow$  sparse promoting regularizer

$\mathcal{C}(\mathbf{h}_1, \mathbf{h}_2) \rightarrow$  non-negativity anchor constraints

$\mathbf{x}_i * \mathbf{h}_j$  computed as  $\mathcal{T}(\mathbf{x}_i)\mathbf{h}_j \in \mathcal{O}(L^2)$

3 [Tong et al., 1994]      3 [Lin et al., 2007, Lin et al., 2008]      3 [Aissa-El-Bey and Abed-Meraim, 2008]

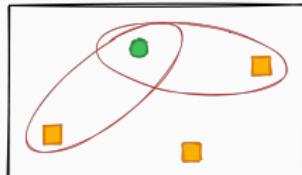
3 [Kowalczyk et al., 2013]      3 [Crocco and Del Bue, 2015, Crocco and Del Bue, 2016]

# AER as discrete SIMO BCE

Key ingredient – *Cross relation identity*

$$x_i = h_i * s$$

$$h_2 * x_1 = h_2 * h_1 * s = h_1 * h_2 * s = h_1 * x_2$$



Ideas:

1. Sampled version of  $x_1, x_2$  are available ( $\mathbf{x}_1, \mathbf{x}_2$ )
2. echoes' TOAs  $\propto$  sampling frequency
3. Find echoes  $\rightarrow$  find sparse vectors  $\mathbf{h}_1, \mathbf{h}_2$  of length  $L$
4. Modeled as Lasso-like problem

$$\widehat{\mathbf{h}}_1, \widehat{\mathbf{h}}_2 \in \arg \min_{\mathbf{h}_1, \mathbf{h}_2 \in \mathbf{R}^n} \|\mathbf{x}_1 * \mathbf{h}_2 - \mathbf{x}_2 * \mathbf{h}_1\|_2^2 + \lambda \mathcal{P}(\mathbf{h}_1, \mathbf{h}_2) \quad \text{s.t.} \quad \mathcal{C}(\mathbf{h}_1, \mathbf{h}_2)$$

$\mathcal{P}(\mathbf{h}_1, \mathbf{h}_2) \rightarrow$  sparse promoting regularizer

$\mathcal{C}(\mathbf{h}_1, \mathbf{h}_2) \rightarrow$  non-negativity anchor constraints

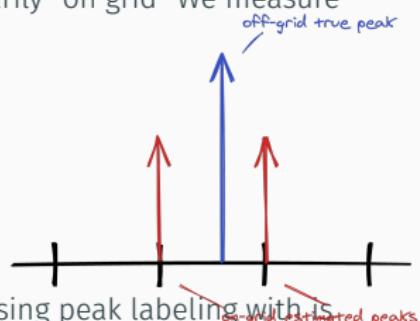
$\mathbf{x}_i * \mathbf{h}_j$  computed as  $\mathcal{T}(\mathbf{x}_i)\mathbf{h}_j \in \mathcal{O}(L^2)$

3 [Tong et al., 1994]      3 [Lin et al., 2007, Lin et al., 2008]      3 [Aissa-El-Bey and Abed-Meraim, 2008]

3 [Kowalczyk et al., 2013]      3 [Crocco and Del Bue, 2015, Crocco and Del Bue, 2016]

## 1. Estimation is on-grid

- Sparsity and non-negativity Echoes are not necessarily “on grid” We measure filters, not diracs
- *Body guard effect* [Duval and Peyré, 2017]
  - low recall  $\Rightarrow$  low accuracy
  - slow convergence



## ... and Pick Picking

- Manually tuned peaking or peak disambiguation (using peak labeling with is NP-hard and need other prior knowledge)

Increase the sampling frequency,  $F_s$

- Increase Precision

## Computational bottleneck

- Bigger vectors and matrices
  - memory usage
- Computational complexity: at best  $\mathcal{O}(F_s^2)$  per iteration
- the higher the sampling frequency, the more ill-conditioned

## State Of The Art

1. discrete (sparse) SIMO BCE  
based on time-domain XREL
2. Peak-picking

## State Of The Art

1. discrete (sparse) SIMO BCE  
based on time-domain XREL
2. Peak-picking

⇒ however

- Estimation in the RIR space  
memory issue
- Echoes are “off-grid”  
accuracy issues and mismatch
- Peak picking and labeling  
tuned and NP-hard

## State Of The Art

1. discrete (sparse) SIMO BCE based on time-domain XREL
2. Peak-picking

⇒ however

- Estimation in the RIR space memory issue
- Echoes are “off-grid” accuracy issues and mismatch
- Peak picking and labeling tuned and NP-hard

⇒ we propose  
**Blaster**

1. Knowledge-based approach
2. BCE + Continuous Dictionary based on XREL
3. Iterative-like approach
4. Inputs:
  - stereo mic recordings
  - # echoes
5. Output:  $\tau_i^{(r)}, \alpha_{i,r}^{(r)}$

[Di Carlo et al., 2020] Collaboration

## Lantern

1. Learning-based regression
2. Deep Learning used for SSL
3. Inputs: stereo audio feature
4. Output in the TDOA space ( $\neq$  Echo space)

[Di Carlo et al., 2019]

# Blaster- Knowledge-based Off-grid AER

Observation 1: the cross relation remains true in the frequency domain

$$\mathcal{F}x_1 \cdot \mathcal{F}h_2(n/F_s) = \mathcal{F}x_2 \cdot \mathcal{F}h_1(n/F_s) \quad n = 0 \dots N - 1$$

Observation 2:  $\mathcal{F}\delta_{\text{echo}}$  is known in closed-form

Observation 3:  $\mathcal{F}x_i$  can be (well) approximated by DFT

$$\mathbf{X}_i = \text{DFT}(\mathbf{x}_i) \simeq \mathcal{F}\mathbf{x}_i(nF_s) \quad n = 0 \dots N - 1$$

Idea: Recover echoes by matching a finite number of frequencies

$$\arg \min_{h_1, h_2 \in \underset{\text{measure}}{\text{space}}} \frac{1}{2} \|\mathbf{X}_1 \cdot \mathcal{F}h_2(f) - \mathbf{X}_2 \cdot \mathcal{F}h_1(f)\|_2^2 + \lambda \|h_1 + h_2\|_{\text{TV}} \quad \text{s.t.} \begin{cases} h_1(\{0\}) = 1 \\ h_l \geq 0 \end{cases}$$

Looks like a Lasso problem, but  $\mathcal{F}h_2(f)$  is a continuous function.

Instance of a BLasso problem [Bredies and Carioni, 2020]

Solved with Sliding Frank-Wolfe algorithm [Denoyelle et al., 2019]

✓ no Toeplitz matrix

✓ Solutions is  
a train of Dirac

✓ anchor prevents  
trivial solution

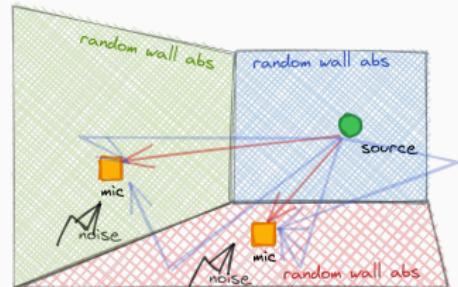
In the manuscript:

# Blaster- Experimental Results

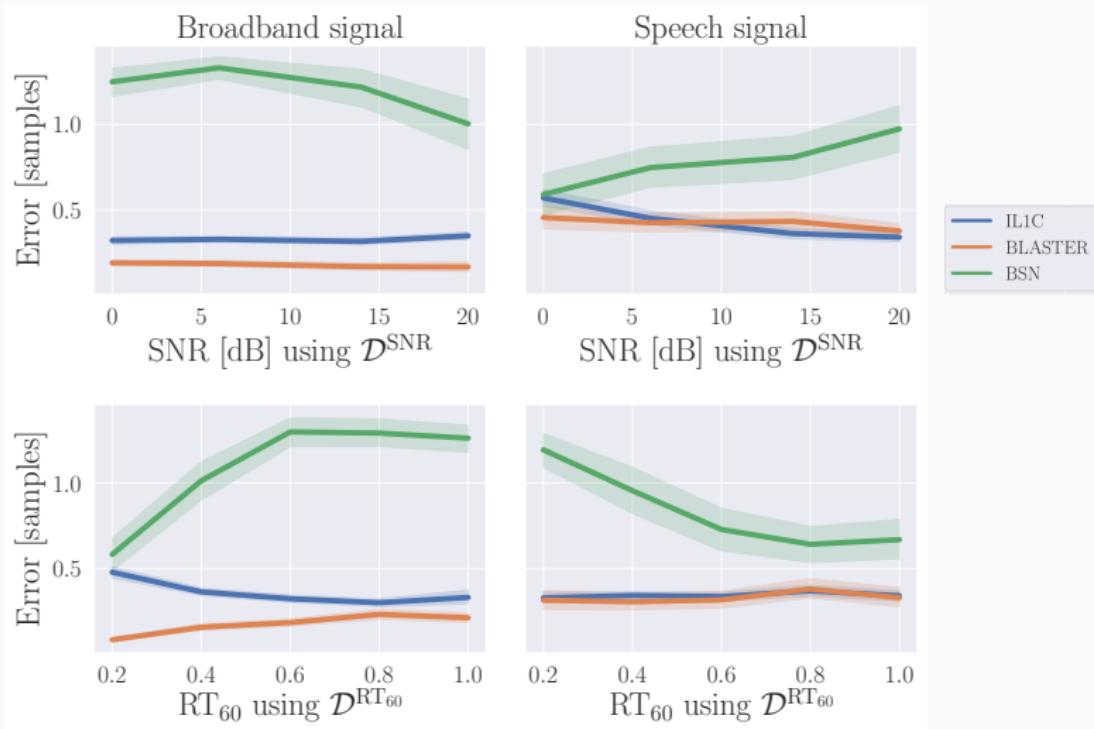
## Methods

- BSN [Lin et al., 2007]
- IL1C: iteratively-weighted  $\ell_1$  constraint SIMO BCE [Crocco and Del Bue, 2015]
- **Blaster**: Proposed off-grid approach

Baseline method are xvalidated on other dataset



# Error per Dataset/Signal while recovering 7 echoes



✓ Lower RMSE

✓ Robustness  
to SNR and  $\text{RT}_{60}$

✗ Source signal  
dependent

## Precision per threshold in typical scenario

$\tau_{\text{thr}}$ [samples]	Precision [%]									
	R = 2 echoes					R = 7 echoes				
	0.5	1	2	3	10	0.5	1	2	3	10
BSN	8	9	27	46	62	5	8	38	54	73
IL1C	51	55	55	56	58	42	53	55	56	58
BLASTER	68	73	74	75	75	46	53	56	57	61

Table 1:  $RT_{60} = 200$  ms and SNR = 20 dB.

✓ Invariant  
to threshold

✗ Sensitive  
to # echoes

# Performance per # of echoes

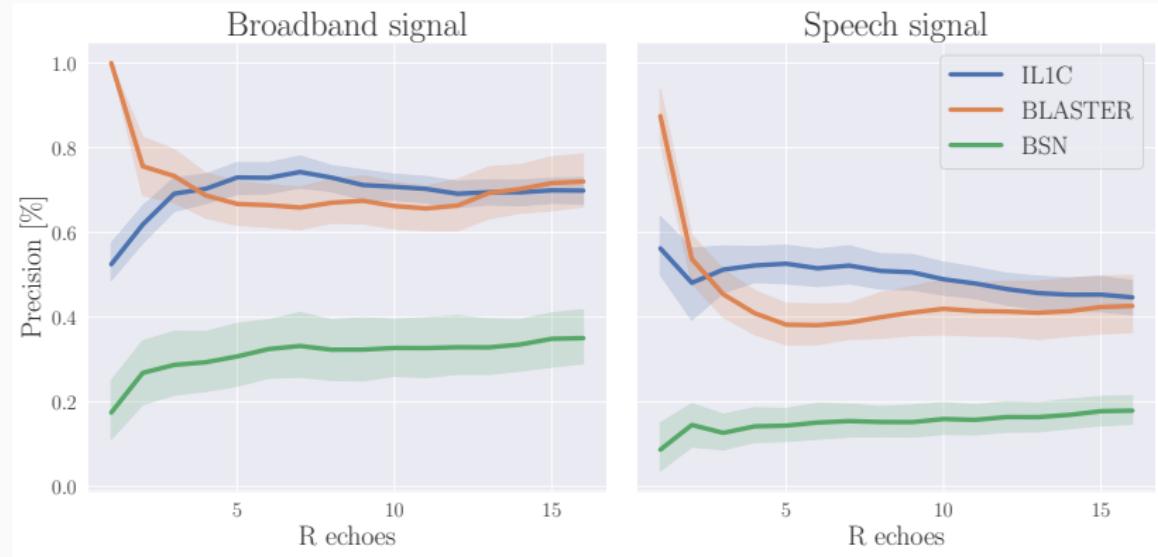


Figure 1:  $RT_{60} = 400$  ms and SNR = 20 dB.

**x** Sensitive  
to # echoes

**x** Sensitive  
source signal

**✓** Good  
for 2 echoes

# Performance per # of echoes

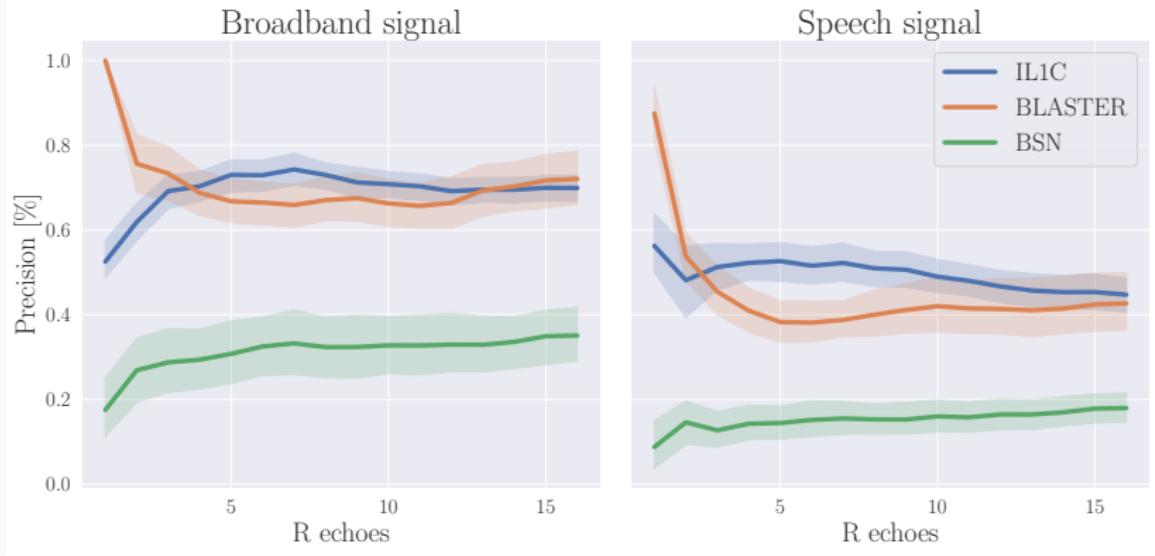


Figure 1:  $RT_{60} = 400$  ms and SNR = 20 dB.

**x** Sensitive  
to # echoes

**x** Sensitive  
source signal

Good  
for 2 echoes  
[Scheibler et al., 2018,  
Di Carlo et al., 2019]

**Observation 1:** Mapping from observation to echo is extremely difficult  
Later echoes are not considered, may help

**Observation 2:** We have acoustic simulators  
Acoustic simulators based on ISM  
source position, room ← reverberation elements ←  
annotation for free

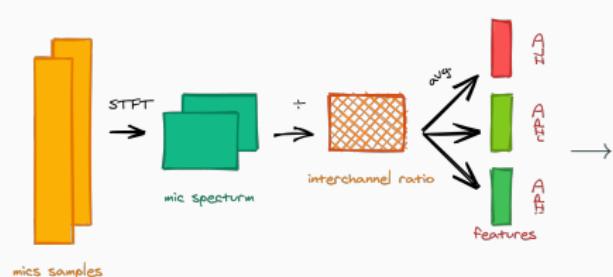
**Observation 3:** (Deep) Learning-based methods successful for localization  
Echoes are strongly related to the source position

## Idea: Use Deep Learning for AER

- Extend previous work on source localization for Echo Estimation
- Estimate the first echo TOA  
    ↪ simple case, but with important application in SSL

Which mapping?

Input: features



Relative Transfer Function

$$\text{ReTF}[f] = \frac{H_2[f]}{H_1[f]} \approx \text{avg.}_t \left( \frac{X_2[f, t]}{X_1[f, t]} \right)$$

This is the instantaneous ReTF

Which Model?

- Architecture: CNN [Chakrabarty and Habets, 2017, Nguyen et al., 2018]

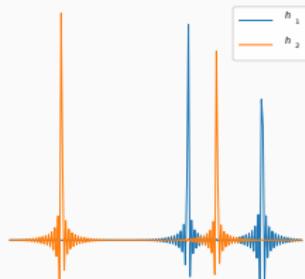
- Loss Function:

1. RMSE (Multi-label regression) on  $\mathcal{V}$
2. Gaussian log-likelihood  $\rightarrow \{\mu, \sigma^2\}$
3. Student log-likelihood  $\rightarrow \{\mu, \lambda, \nu\}$

} Generative models  $\leftarrow$  for data fusion  
similar to MDN [Bishop, 1994]

- Virtually Supervised Learning (= data from acoustic simulator)

Output: target



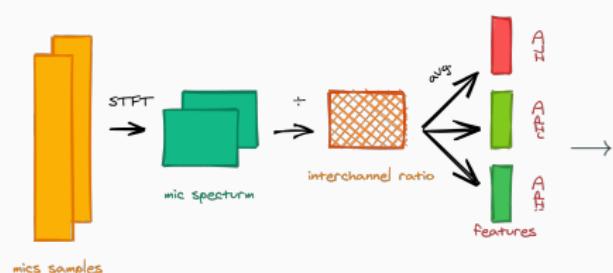
TDOAs inter and intra reflection

$$\mathcal{V} = \{\text{TDOA}, \text{iTDOA}, \text{TDOE}_1\}$$

First strongest echo  $\Leftrightarrow$  close surface

Which mapping?

Input: features



Relative Transfer Function

$$\text{ReTF}[f] = \frac{H_2[f]}{H_1[f]} \approx \text{avg.}_t \left( \frac{X_2[f, t]}{X_1[f, t]} \right)$$

This is the instantaneous ReTF

Which Model?

- Architecture: CNN [Chakrabarty and Habets, 2017, Nguyen et al., 2018]

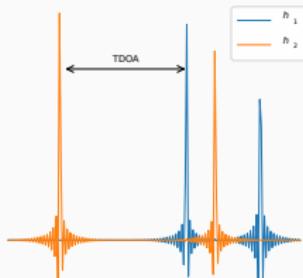
Loss Function:

1. RMSE (Multi-label regression) on  $\mathcal{V}$
2. Gaussian log-likelihood  $\rightarrow \{\mu, \sigma^2\}$
3. Student log-likelihood  $\rightarrow \{\mu, \lambda, \nu\}$

Generative models  $\leftarrow$  for data fusion  
similar to MDN [Bishop, 1994]

- Virtually Supervised Learning (= data from acoustic simulator)

Output: target



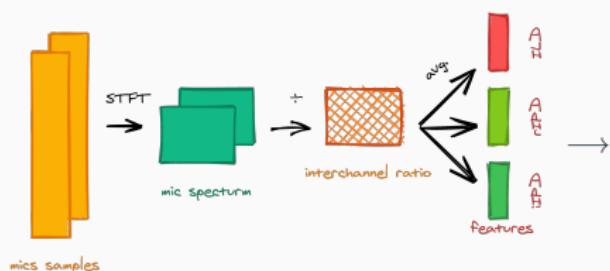
TDOAs inter and intra reflection

$$\mathcal{V} = \{\text{TDOA}, \text{iTDOA}, \text{TDOE}_1\}$$

First strongest echo  $\Leftrightarrow$  close surface

Which mapping?

Input: features



Relative Transfer Function

$$\text{ReTF}[f] = \frac{H_2[f]}{H_1[f]} \approx \text{avg.}_t \left( \frac{X_2[f, t]}{X_1[f, t]} \right)$$

This is the instantaneous ReTF

Which Model?

- Architecture: CNN [Chakrabarty and Habets, 2017, Nguyen et al., 2018]

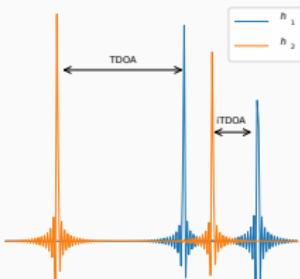
Loss Function:

1. RMSE (Multi-label regression) on  $\mathcal{V}$
2. Gaussian log-likelihood  $\rightarrow \{\mu, \sigma^2\}$
3. Student log-likelihood  $\rightarrow \{\mu, \lambda, \nu\}$

} Generative models  $\leftarrow$  for data fusion  
similar to MDN [Bishop, 1994]

- Virtually Supervised Learning (= data from acoustic simulator)

Output: target



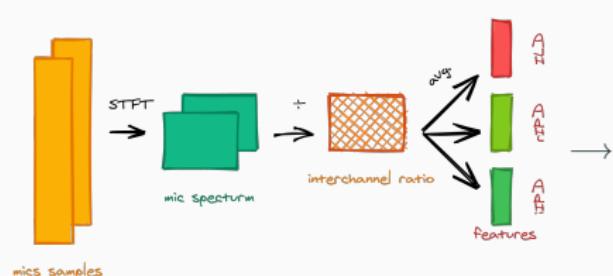
TDOAs inter and intra reflection

$$\mathcal{V} = \{\text{TDOA}, \text{iTDOA}, \text{TDOE}_1\}$$

First strongest echo  $\Leftrightarrow$  close surface

Which mapping?

Input: features

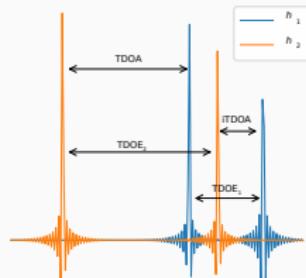


Relative Transfer Function

$$\text{ReTF}[f] = \frac{H_2[f]}{H_1[f]} \approx \text{avg.}_t \left( \frac{X_2[f, t]}{X_1[f, t]} \right)$$

This is the instantaneous ReTF

Output: target



TDOAs inter and intra reflection

$$\mathcal{V} = \{\text{TDOA}, \text{iTDOA}, \text{TDOE}_1\}$$

First strongest echo  $\Leftrightarrow$  close surface

Which Model?

- Architecture: CNN [Chakrabarty and Habets, 2017, Nguyen et al., 2018]

- Loss Function:

1. RMSE (Multi-label regression) on  $\mathcal{V}$
2. Gaussian log-likelihood  $\rightarrow \{\mu, \sigma^2\}$
3. Student log-likelihood  $\rightarrow \{\mu, \lambda, \nu\}$

Generative models  $\leftarrow$  for data fusion  
similar to MDN [Bishop, 1994]

- Virtually Supervised Learning (= data from acoustic simulator)

# Lantern- Experiments & Results

Baseline: GCCPHAT (only TDOA),  
 $\text{MLP}_\mathcal{V}$  [Di Carlo et al., 2019]

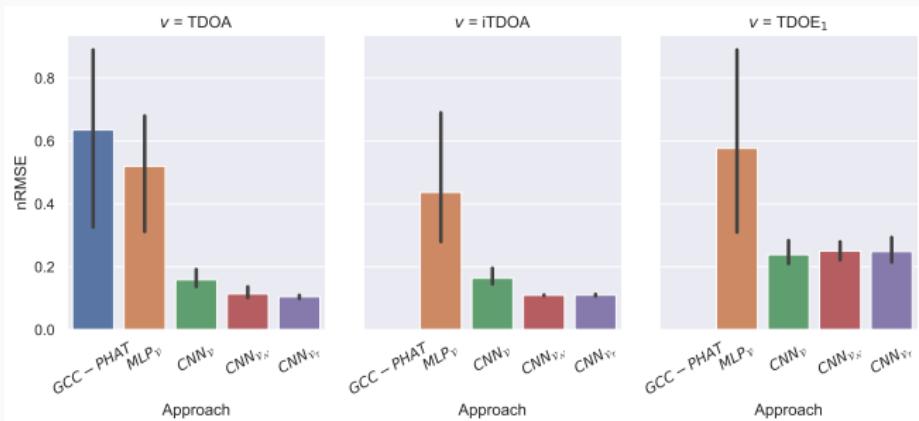
Proposed:  $\text{CNN}_\mathcal{V}$ ,  $\text{CNN}_{\mathcal{V}_N}$ ,  $\text{CNN}_{\mathcal{V}_T}$

Metric: normalized RMSE  
(0 = best fit, 1 = random)

Train:

- RT60, SNR
- white noise
- instantaneous RTF

Test: similar to train



✓ CNNs outperform  
GCC-PHAT and  
 $\text{MLP}$

✓ CNNs more robust  
to noise

✗ Gaussian  
 $\sim \text{Student-T}$

# Lantern- Experiments & Results

Baseline: GCCPHAT (only TDOA),  
 $\text{MLP}_{\mathcal{V}}$  [Di Carlo et al., 2019]

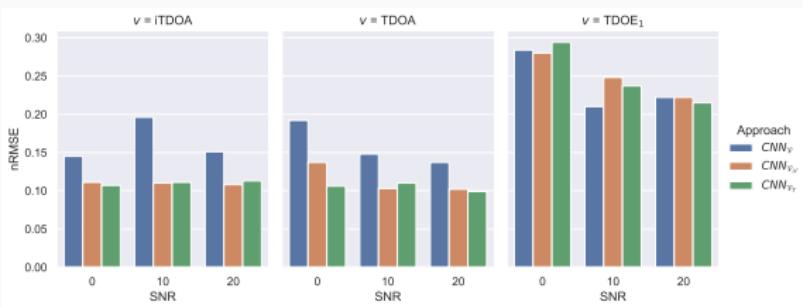
Proposed:  $\text{CNN}_{\mathcal{V}}$ ,  $\text{CNN}_{\mathcal{V}_N}$ ,  $\text{CNN}_{\mathcal{V}_T}$

Metric: normalized RMSE  
(0 = best fit, 1 = random)

Train:

- RT60, SNR
- white noise
- instantaneous RTF

Test: similar to train



✓ Generative  
better than  
Normal

✗ Gaussian  
Student-T

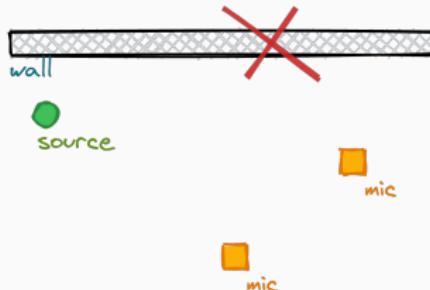
~

## Echo-aware Application

---

# Audio signal processing and sound propagation

Sound propagation is



- completely ignored
- assumed free-field case (*anechoic*)
- model it full (*reverberant*)
- *learned* it full (*reverberant*)
- model few early echoes (*multipath*)

$$x_i(t) = (h_i * s)(t)$$

$$h_i(t) = h_i^d(t) + h_i^e(t) + h_i^r(t)$$

Recall

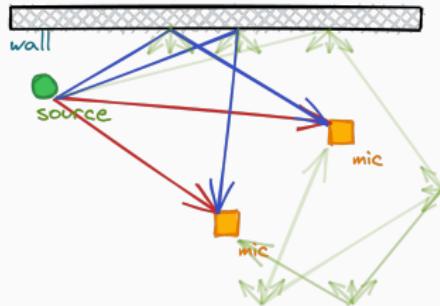
- anechoic case: easy mapping, but incoherence or wrong processing
- reverberant case: difficult mapping and estimation, but coherent processing

What can we do with echoes?

# Audio signal processing and sound propagation

Sound propagation is

- completely ignored
- assumed free-field case (*anechoic*)
- model it full (*reverberant*)
- *learned* it full (*reverberant*)
- model few early echoes (*multipath*)



$$x_i(t) = (h_i * s)(t)$$

$$h_i(t) = h_i^d(t) + h_i^e(t) + h_i^r(t)$$

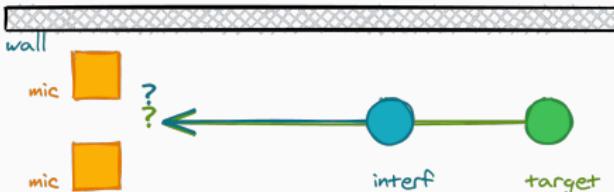
Recall

- anechoic case: easy mapping, but incoherence or wrong processing
- reverberant case: difficult mapping and estimation, but coherent processing

What can we do with echoes?

# Echo-aware Application

Echoes = same content, different time/direction



Echoes helps indoor processing:

## What?

Echoes = repetitions

- Sound Source Separation
- Speech Enhancement  
(Dereverberation,  
Denoising, Room  
Equalization)

## Where?

Echoes ∈ indoor propagation

- Sound Source Localization
- Microphone Calibration
- Room Geometry  
Reconstruction

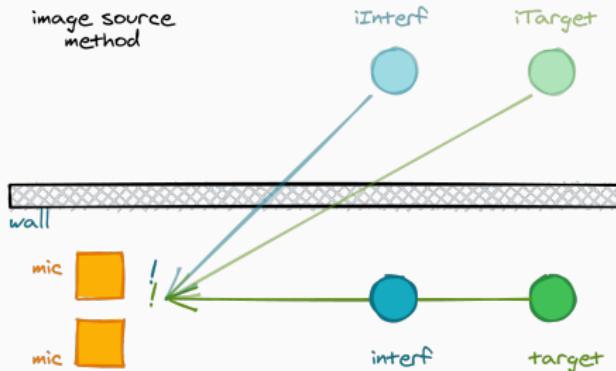
## How?

Echoes ∈ sound propagation

- Blind Channel Estimation
- Acoustic Measurements

# Echo-aware Application

Echoes = same content, different time/direction



Echoes helps indoor processing:

## What?

Echoes = repetitions

- Sound Source Separation
- Speech Enhancement  
(Dereverberation,  
Denoising, Room  
Equalization)

## Where?

Echoes ∈ indoor propagation

- Sound Source Localization
- Microphone Calibration
- Room Geometry  
Reconstruction

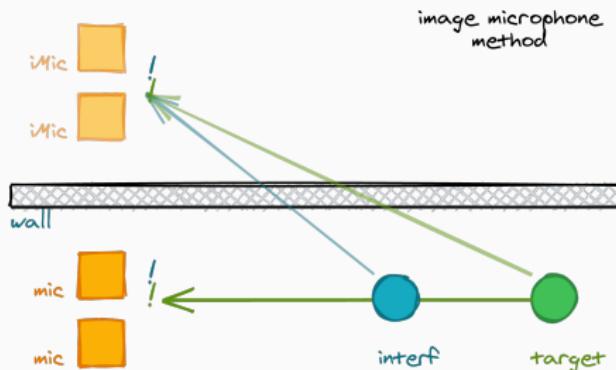
## How?

Echoes ∈ sound propagation

- Blind Channel Estimation
- Acoustic Measurements

# Echo-aware Application

Echoes = same content, different time/direction



Echoes helps indoor processing:

## What?

Echoes = repetitions

- Speech Enhancement  
(Dereverberation,  
Denoising, Room  
Equalization)

## Where?

Echoes ∈ indoor propagation

- Sound Source Localization
- Microphone Calibration
- Room Geometry Reconstruction

## How?

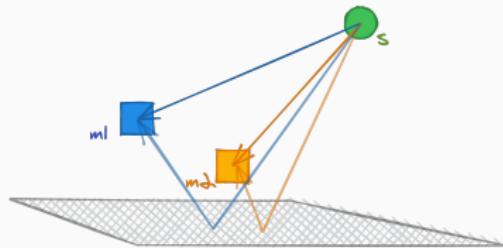
Echoes ∈ sound propagation

- Blind Channel Estimation
- Acoustic Measurements

# Mirage- Sound Source Localization with Echoes

The **Picnic** Scenario:

- One source
- Two microphones
  - passive scenario
  - generalizable later
- Close to a very reflective surface
  - First echo = Strongest echo
  - $\alpha_{\text{picnic}}$  const.  $\forall f$
  - table-top device

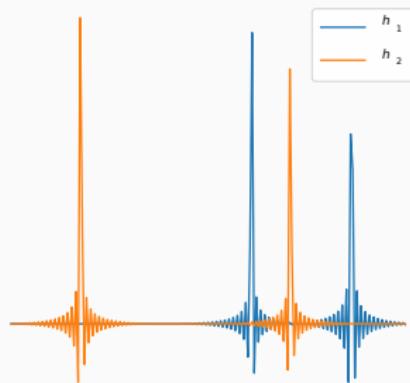


How to access the *image* microphones?  
⇒ each pair is augmented with echoes

## Mirage Array

idea: use SSL algorithm on this augmented array

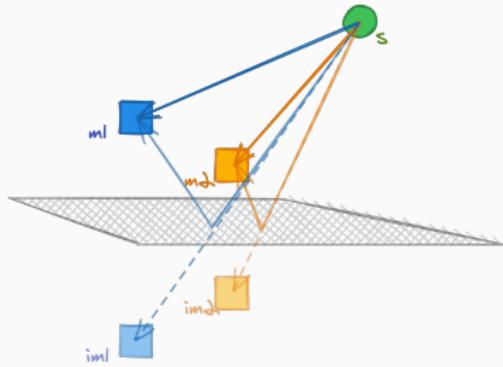
recall: echoes are known



# Mirage- Sound Source Localization with Echoes

The **Picnic** Scenario:

- One source
- Two microphones
  - passive scenario
  - generalizable later
- Close to a very reflective surface
  - First echo = Strongest echo
  - $\alpha_{\text{picnic}}$  const.  $\forall f$
  - table-top device

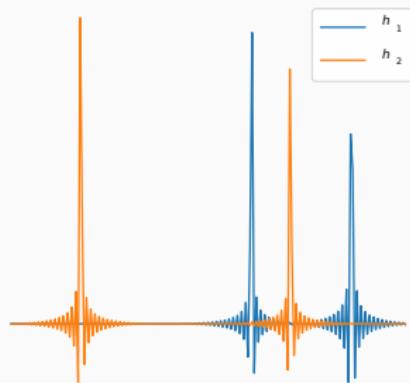


How to access the *image* microphones?  
⇒ each pair is augmented with echoes

## Mirage Array

idea: use SSL algorithm on this augmented array

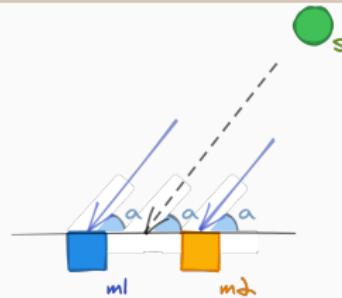
recall: echoes are known



# Mirage- Sound Source Localization with Echoes

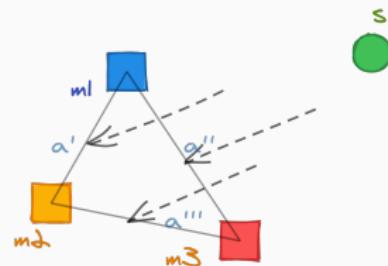
## SSL with 2 microphones

1. only Angle of Arrival (AOA) w.r.t. the frame of the pair
  - e.g. GCC-PHAT for TDOA estimation (known limitation, but good in practice)
  - TDOA to AoA known frame distance



## SSL with more microphones

1. For each pair  $m$ :  
 $\text{AOA}_m \leftarrow \text{TDOA-based 2-mic-SSL}$
2. "Aggregate" together all the observation  
(Angular spectra, Probability distributions)  
eg. SRP-PHAT



**Baseline:**  
GCC-PHAT on true microphones

**Proposed Approach:**  
Using **Lantern** (DNN-based TDOA estimation)  
problem: real value not estimation → Generative Model given the TDOA axis

# Mirage- Results

## Data

Virtually generated dataset as for **Lantern**

### AOA estimation normalized nRMSE

### Angular Error — mean and accuracy

### Azimuth Elevation Estimation

### Angular Error — mean and accuracy

	Input	nRMSE			ACCURACY	
		TDOA	iTDOA	TDOE	$\theta < 10$	$\theta < 20$
MIRAGE	wn	0.18	0.28	0.25	4.10 (77)	5.97 (97)
	wn+n	0.68	0.69	0.89	5.00 (26)	9.89 (54)
	sp	0.31	0.34	0.56	4.83 (63)	7.26 (82)
	sp+n	0.99	0.98	1.48	4.60 (16)	9.88 (35)
	GCC-PHAT	0.21	-	-	4.22 (81)	6.19 (97)
	wn	0.68	-	-	4.03 (65)	5.34 (83)
	wn+n	0.32	-	-	4.08 (82)	5.34 (97)
	sp	1.38	-	-	4.70 (19)	8.38 (32)
DoA		ACCURACY			ACCURACY	
		< 10		< 20		
		θ	ϕ	θ	ϕ	
MIRAGE	wn	4.5 (59)	3.9 (71)	6.8 (79)	5.9 (88)	
	wn+n	4.4 (18)	5.5 (26)	9.4 (35)	11.1 (66)	
	sp	4.6 (45)	4.8 (59)	8.1 (71)	7.2 (83)	
	sp+n	5.2 (17)	5.9 (12)	10.7 (38)	12.3 (43)	

✓ Solved “impossible”  
localization

✗ Performance depending on  
echo estimation

## Echo-aware Dataset

---

## Data in audio signal processing

1. are necessary for validating (and learning) models
2. collecting real data is not always possible  
annotation and recording require expertise, equipment and time
3. dataset of real data cannot be easily shared  
they do not generalize to different use-cases and scenarios (array, recording scenario)
4. simulated data are used instead: quantity, versatility, annotation easiness and “quality”

## Echo-aware Data in audio signal processing

For SE: strong echoes, but not annotated

[Szöke et al., 2019, Bertin et al., 2019, Remaggi et al., 2016]

For RooGE: good geo. annotation, but no variety of acoustic scenarios

[Dokmanić et al., 2013, Crocco et al., 2017, Remaggi et al., 2019]

A good echo-aware dataset should allow SE, RooGE and AER

HOW?

signal annotation  $\leftrightarrow$  geometric annotation

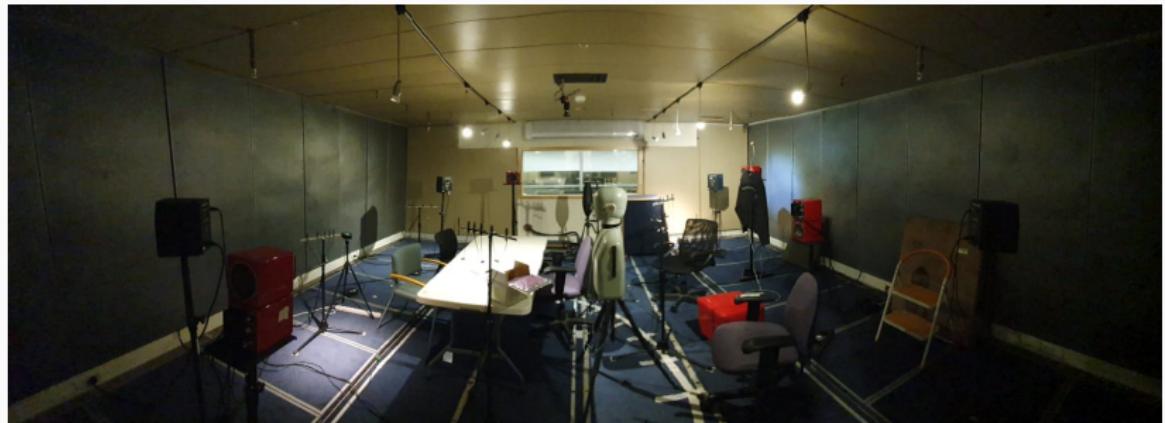
# dEchorate realization

**dEchorate:** echo-aware dataset

Recorded Acoustic lab of Bar'Ilan (Shoebox)

Annotated during confinement COVID-2020

Collaboration with prof. Sharon Gannot and ing. Pinchas Tandeitnik



# dEchorate realization

## dEchorate: echo-aware dataset

Recorded Acoustic lab of Bar'Ilan (Shoebox)

Annotated during confinement COVID-2020

Collaboration with prof. Sharon Gannot and ing. Pinchas Tandeitnik

### Key features:

- many acoustic environments (revolving panels)
- 6 nULA with 5 mics and 4 sound sources
- geometry annotated & echo annotated
- measured RIRs  $\xrightarrow{\text{matching}}$  simulated RIRs



1. RIR estimation with chirps signal [Farina, 2007, Szöke et al., 2019]
2. IPS with beacon → mic and src positioning ( $\pm 2$  cm)
3. GUI for echo annotation  
Skyline, Matched Filter, Assisted Peak Picking
4. Refined position with Least Square optimization
5. iterate including ceiling (perfectly flat)

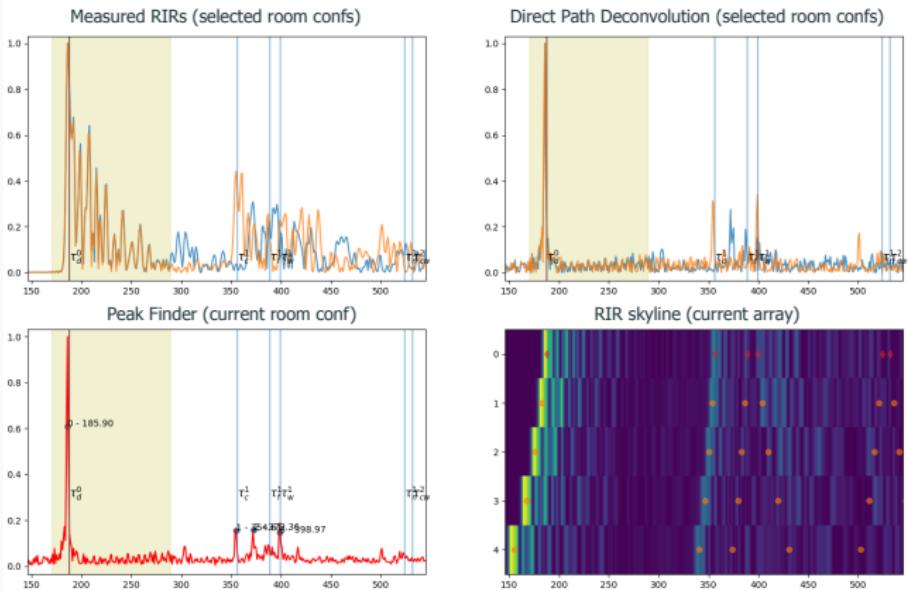
1. RIR estimation with chirps signal [Farina, 2007, Szöke et al., 2019]
2. IPS with beacon → mic and src positioning ( $\pm 2$  cm)
3. GUI for echo annotation  
Skyline, Matched Filter, Assisted Peak Picking
4. Refined position with Least Square optimization
5. iterate including ceiling (perfectly flat)

	Metrics	bIPS	dMDS	dcMDS
Geom.	Max.	0	6.1	1.07
	Avg. $\pm$ Std.	0	$1.8 \pm 1.4$	$0.39 \pm 0.2$
Signal	Max.	5.86	1.20	1.86
	Avg. $\pm$ Std.	$1.85 \pm 1.5$	$0.16 \pm 0.2$	$0.41 \pm 0.3$
Mismatch	GoM (1.0 ms)	97.9%	93.4%	98.1%
	GoM (0.1 ms)	26.6%	44.8%	53.1%
	GoM (0.05 ms)	12.5%	14.4%	30.2%

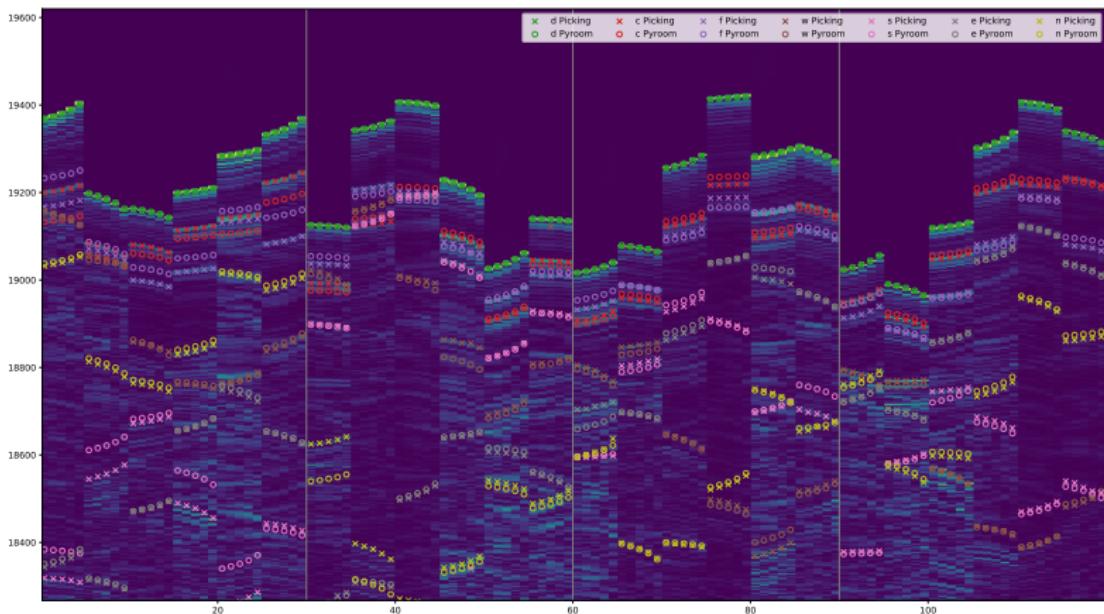
GoM = Goodness of Match ( $\neq$  error, because no groundtruth)

# dEchorate— Annotation

1. RIR estimation with chirps signal [Farina, 2007, Szöke et al., 2019]
2. IPS with beacon → mic and src positioning ( $\pm 2$  cm)
3. GUI for echo annotation  
Skyline, Matched Filter, Assisted Peak Picking
4. Refined position with Least Square optimization
5. iterate including ceiling (perfectly flat)



# dEchorate— Annotation



Estimating the room geometry: shape, volume or reflector position from signal or form TOAs and labels

## RooGE as Image Source Inversion

If TOAs annotation (label and value) are available:

For each wall/label:

1. TOA → image source position via 3D multilateration
2. image source position → reflector estimation via geometric reasoning

other methods differ for priors and setup [Filos et al., 2011, Antonacci et al., 2012, Crocco et al., 2017]

# Room Geometry Estimation (RooGE) with dEchorate

Estimating the room geometry: shape, volume or reflector position from signal or form TOAs and labels

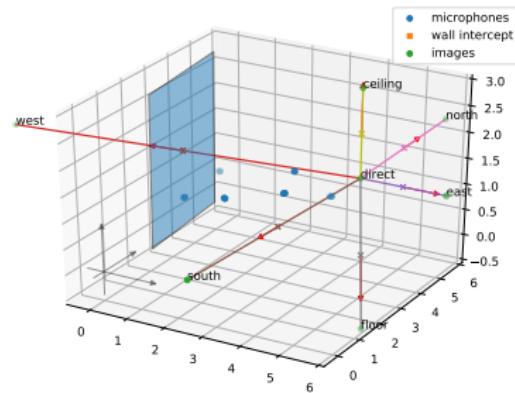
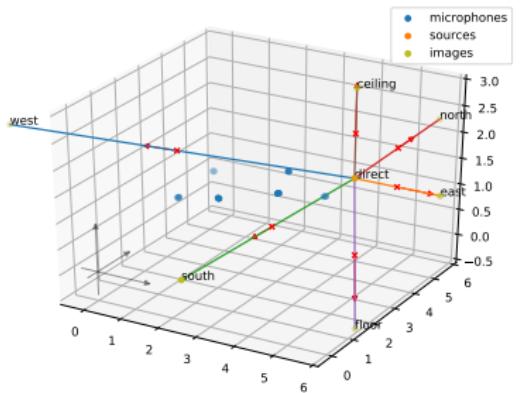
## RooGE as Image Source Inversion

If TOAs annotation (label and value) are available:

For each wall/label:

1. TOA → image source position via 3D multilateration
2. image source position → reflector estimation via geometric reasoning

other methods differ for priors and setup [Filos et al., 2011, Antonacci et al., 2012, Crocco et al., 2017]



# Room Geometry Estimation (RooGE) with dEchorate

Estimating the room geometry: shape, volume or reflector position from signal or form TOAs and labels

## RooGE as Image Source Inversion

If TOAs annotation (label and value) are available:

For each wall/label:

1. TOA → image source position via 3D multilateration
2. image source position → reflector estimation via geometric reasoning

other methods differ for priors and setup [Filos et al., 2011, Antonacci et al., 2012, Crocco et al., 2017]

source id wall	1		2		3		4	
	DE	AE	DE	AE	DE	AE	DE	AE
west	0.74	8.99	4.59	8.32	5.89	5.75	<b>0.05</b>	<b>2.40</b>
east	<b>0.81</b>	<b>0.08</b>	0.9	0.50	<i>69.51</i>	<i>55.70°</i>	0.31	0.21
south	3.94	<i>16.08°</i>	<b>0.18</b>	1.77	<i>14.37</i>	<i>18.55°</i>	0.82	<b>1.65</b>
north	1.34	0.76	1.40	8.94	<b>0.63</b>	<b>0.17</b>	2.08	1.38
floor	<b>5.19</b>	<b>1.76</b>	7.27	2.66	7.11	2.02	5.22	1.90
ceiling	1.16	0.28	0.67	0.76	<b>0.24</b>	1.16	0.48	<b>0.26</b>

Distance Error (DE) [cm] and Angular Error (AE)

## Speech Enhancement

Improve the quality of a *target* sound source with respect:

- interferences, i.e. from other sources  $\rightsquigarrow$  sound source separation
- background noise  $\rightsquigarrow$  denoising
- reverberation  $\rightsquigarrow$  dereverberation, room equalization

## Spatial filtering via Beamformers

- Is a speech enhancement techniques for multichannel
- vs. Wiener Filtering, the target is distortionless
- in anechoic case, it correspond to delay-and-sum beamformer
- physical interpretation with steering vector based on DOA
- both in time and frequency domain

### Beamformer: closed-form solution

$$\mathbf{w} = \Sigma_n^{-1} \mathbf{C} (\mathbf{C}^H \Sigma_n^{-1} \mathbf{C}) \mathbf{g}$$

## Speech Enhancement

Improve the quality of a *target* sound source with respect:

- interferences, i.e. form other sources  $\rightsquigarrow$  sound source separation
- background noise  $\rightsquigarrow$  denoising
- reverberation  $\rightsquigarrow$  dereverberation, room equalization

## Spatial filtering via Beamformers

- Is a speech enhancement techniques for multichannel
- vs. Wiener Filtering, the target is distortionless
- in anechoic case, it correspond to delay-and-sum beamformer
- physical interpretation with steering vector based on DOA
- both in time and frequency domain

### Beamformer: closed-form solution

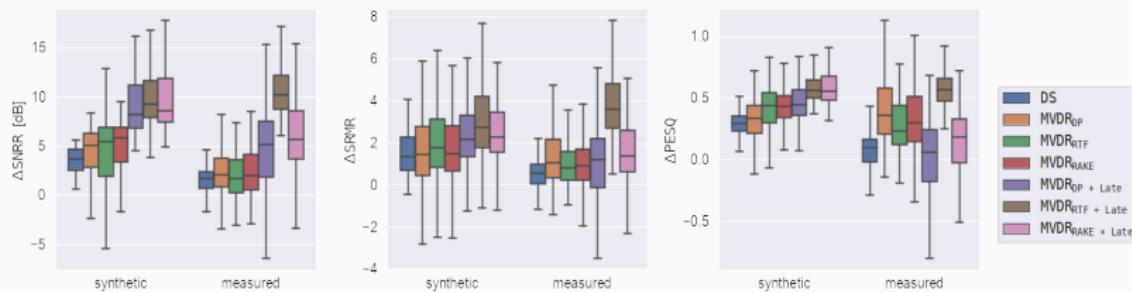
$$\mathbf{w} = \Sigma_n^{-1} \mathbf{C} (\mathbf{C}^H \Sigma_n^{-1} \mathbf{C}) \mathbf{g}$$

# Echo-aware Speech Enhancement

The PSD of various components  
asd

## Different Criteria and Solution

- DS
- MVDR - DP
- MVDR ReTF



## Conclusion

---

# Summary of contributions

## Audio Scene Analysis

↓  
context and problems

## Signal Processing

↓  
models and frameworks

## Acoustic Echoes

↓  
better processing

### Challenges and Objective

1. How to estimate acoustic echoes?
2. How to extend methods for echo-aware audio scene analysis

#### 1. Estimation

- Knowledge-based echo estimation  
↪ **Blaster**
- Learning-based echo estimation  
↪ **Lantern**

#### 2. Application

- Echo-aware Source Separation  
↪ **Separake**
- Echo-aware Source Localization  
↪ **Mirage**
- Echo-aware Speech Enhancement
- Echo-aware Room Geometry Estimation

#### 3. Data:

Echo-aware database → **dEchorate**

Directions for future work:

- Estimation
  - RTF and Multichannel **Blaster**
  - Physic-based Learning
- Application
  - Other field of echoes (Sismology, Surveliance)
- dEchorate
  - Synthetic vs. Real RIRs

# List of publications

---

- Estimation
  - **Lantern** [Di Carlo et al., 2019]
  - **dEchorate** [Di Carlo et al., 2020]
- Application
  - **Mirage** [Di Carlo et al., 2019]
  - **Separake** [Scheibler et al., 2018]
- Data
  - **dEchorate** (Unpublished)
- Other
  - Signal Processing CUP 2019 [Deleforge et al., 2019]
  - Locata Challenge 2019 [Lebarbenchon et al., 2018]
  - Collaboration with Honda [Di Carlo and Deleforge, ]

Code

dEchorate:

Risotto:

Brioche:

pyMBSSLocate:

Separake:

-  Aissa-El-Bey, A. and Abed-Meraim, K. (2008).  
**Blind simo channel identification using a sparsity criterion.**  
In *2008 IEEE 9th Workshop on Signal Processing Advances in Wireless Communications*, pages 271–275. IEEE.
-  Antonacci, F., Filos, J., Thomas, M. R., Habets, E. A., Sarti, A., Naylor, P. A., and Tubaro, S. (2012).  
**Inference of room geometry from acoustic impulse responses.**  
*IEEE Transactions on Audio, Speech, and Language Processing*, 20(10):2683–2695.
-  Bertin, N., Camberlein, E., Lebarbenchon, R., Vincent, E., Sivasankaran, S., Illina, I., and Bimbot, F. (2019).  
**Voicehome-2, an extended corpus for multichannel speech processing in real homes.**  
*Speech Communication*, 106:68–78.
-  Bishop, C. M. (1994).  
**Mixture density networks.**

-  Bredies, K. and Carioni, M. (2020).  
Sparsity of solutions for variational inverse problems with finite-dimensional data.  
*Calculus of Variations and Partial Differential Equations*, 59(1):14.
-  Chakrabarty, S. and Habets, E. A. (2017).  
Broadband doa estimation using convolutional neural networks trained with noise signals.  
In *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 136–140. IEEE.
-  Crocco, M. and Del Bue, A. (2015).  
**Room impulse response estimation by iterative weighted l 1-norm.**  
In *2015 23rd European Signal Processing Conference (EUSIPCO)*, pages 1895–1899. IEEE.
-  Crocco, M. and Del Bue, A. (2016).  
**Estimation of tdoa for room reflections by iterative weighted l 1 constraint.**  
In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3201–3205. IEEE.

-  Crocco, M., Trucco, A., and Del Bue, A. (2017).  
**Uncalibrated 3d room geometry estimation from sound impulse responses.**  
*Journal of the Franklin Institute*, 354(18):8678–8709.
-  Deleforge, A., Di Carlo, D., Strauss, M., Serizel, R., and Marcenaro, L. (2019).  
**Audio-based search and rescue with a drone: Highlights from the ieee signal processing cup 2019 student competition [sp competitions].**  
*IEEE Signal Processing Magazine*, 36(5):138–144.
-  Denoyelle, Q., Duval, V., Peyré, G., and Soubies, E. (2019).  
**The sliding frank-wolfe algorithm and its application to super-resolution microscopy.**  
*Inverse Problems*, 36(1):014001.
-  Di Carlo, D. and Deleforge, A.  
**Hri-jf collaboration - final phase ii deliverable.**  
Technical report, Inria Nancy - Grand Est.

-  Di Carlo, D., Deleforge, A., and Bertin, N. (2019).  
**Mirage: 2d source localization using microphone pair augmentation with echoes.**  
In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 775–779. IEEE.
-  Di Carlo, D., Elvira, C., Deleforge, A., Bertin, N., and Gribonval, R. (2020).  
**Blaster: An off-grid method for blind and regularized acoustic echoes retrieval.**  
In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 156–160. IEEE.
-  Dokmanić, I., Parhizkar, R., Walther, A., Lu, Y. M., and Vetterli, M. (2013).  
**Acoustic echoes reveal room shape.**  
*Proceedings of the National Academy of Sciences*, 110(30):12186–12191.
-  Duval, V. and Peyré, G. (2017).  
**Sparse regularization on thin grids i: the lasso.**  
*Inverse Problems*, 33(5):055008.
-  Farina, A. (2007).  
**Advancements in impulse response measurements by sine sweeps.**  
In *Audio Engineering Society Convention 122*. Audio Engineering Society.

-  Filos, J., Canclini, A., Thomas, M. R., Antonacci, F., Sarti, A., and Naylor, P. A. (2011). Robust inference of room geometry from acoustic measurements using the hough transform. In *2011 19th European Signal Processing Conference*, pages 161–165. IEEE.
-  Kowalczyk, K., Habets, E. A., Kellermann, W., and Naylor, P. A. (2013). Blind system identification using sparse learning for tdoa estimation of room reflections. *IEEE Signal Processing Letters*, 20(7):653–656.
-  Lebarbenchon, R., Camberlein, E., Di Carlo, D., Gaultier, C., Deleforge, A., and Bertin, N. (2018). Evaluation of an open-source implementation of the srp-phat algorithm within the 2018 locata challenge. *Proc. of LOCATA Challenge Workshop-a satellite event of IWAENC*.
-  Lin, Y., Chen, J., Kim, Y., and Lee, D. D. (2007). Blind sparse-nonnegative (bsn) channel identification for acoustic time-difference-of-arrival estimation. In *2007 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pages 106–109. IEEE.

-  Lin, Y., Chen, J., Kim, Y., and Lee, D. D. (2008).  
**Blind channel identification for speech dereverberation using l1-norm sparse learning.**  
In *Advances in Neural Information Processing Systems*, pages 921–928.
-  Nguyen, Q., Girin, L., Bailly, G., Elisei, F., and Nguyen, D.-C. (2018).  
**Autonomous sensorimotor learning for sound source localization by a humanoid robot.**
-  Remaggi, L., Jackson, P. J., Coleman, P., and Wang, W. (2016).  
**Acoustic reflector localization: novel image source reversion and direct localization methods.**  
*IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(2):296–309.
-  Remaggi, L., Jackson, P. J., and Wang, W. (2019).  
**Modeling the comb filter effect and interaural coherence for binaural source separation.**  
*IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(12):2263–2277.

-  Scheibler, R., Di Carlo, D., Deleforge, A., and Dokmanić, I. (2018). **Separake: Source separation with a little help from echoes.** In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6897–6901. IEEE.
-  Szöke, I., Skácel, M., Mošner, L., Paliesek, J., and Černocký, J. H. (2019). **Building and evaluation of a real room impulse response dataset.** *IEEE Journal of Selected Topics in Signal Processing*, 13(4):863–876.
-  Tong, L., Xu, G., and Kailath, T. (1994). **Blind identification and equalization based on second-order statistics: A time domain approach.** *IEEE Transactions on Information Theory*, 40(2):340–349.