



Scenarios:

Know

5(+) un known

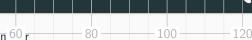
5(+)

estima 40 ; $\tau_3^{(r)} \rightarrow \mathsf{TOAs}$ estimation

Source signal is



Arrival time of acoustic events







nes recordings
$$\{ ilde{x}_i\}_i o \{ au_i^{(r)}, lpha_i^{(r)}\}_{i,r}$$



Our case: signal source and passive system of (I microphones)





Scenarios:

Know

5(+) un known

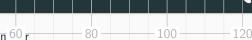
5(+)

estima 40 ; $\tau_3^{(r)} \rightarrow \mathsf{TOAs}$ estimation

Source signal is



Arrival time of acoustic events







nes recordings
$$\{ ilde{x}_i\}_i o \{ au_i^{(r)}, lpha_i^{(r)}\}_{i,r}$$



Our case: signal source and passive system of (I microphones)





Scenarios:

Know

5(+) un known

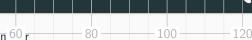
5(+)

estima 40 ; $\tau_3^{(r)} \rightarrow \mathsf{TOAs}$ estimation

Source signal is



Arrival time of acoustic events







nes recordings
$$\{ ilde{x}_i\}_i o \{ au_i^{(r)}, lpha_i^{(r)}\}_{i,r}$$



Our case: signal source and passive system of (I microphones)







Acoustic Echo Retrieval (AER)

Estir mm; early (str 40;) reflection 60 r

microphones recordings
$$\{\widetilde{x}_i\}_i \longrightarrow \{\tau_i^{(r)},\alpha_i^{(r)}\}_{i,r}$$

estima 40 ;
$$\tau_i^{(r)} \rightarrow \mathsf{TOAs}$$
 estimation

Source signal is

h(+)

Scenarios:

5(+)

5(+)

1/12

un known

Active intrusive or specific setup



Passive

passive and more common setups

Arrival time of acoustic events

blind inverse problem (harder)

(Applications: passive listening, smart speakers, etc.)



Know

5(+)

5(+)

un known

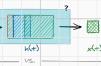


Estir mm, early (str 40) reflection 60 r microphones recordings

 $\{\tilde{x}_i\}_i \longrightarrow \{\tau_i^{(r)}, \alpha_i^{(r)}\}_{i,r}$ estima 40; $\tau_3^{(r)} \rightarrow \mathsf{TOAs}$ estimation

Scenarios: Source signal is

h(+)





x(+)

x(+)

RIR -agriostic



Estimation is

Arrival time of acoustic events





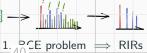






Our case: signal source and passive system of (I microphones)





60 rany frameworks, solver, >20

Pros years of literature)

Cons Peak picking

2/12 • On-grid estimation

Full RIR is estimated

1.
$${}^{\circ}_{40}$$
CE problem \implies RIRs
2. Peak picking \implies Echoes

exploratory (no solver)

• easy ill-conditioned

no direct translation of BCE

1. Estimation in
$$\Theta = \{\tau_i^{(r)}\}_{i,r}$$







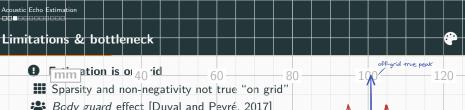












Body guard effect [Duval and Peyré, 2017] \rightarrow low recall \Longrightarrow low accuracy

How 60 ut higher F_* ? → Increase Precision

the higher the sampling frequency the

→ echo labeling (NP-hard problem)

-80 memory usage

by: at best
$$\mathcal{O}(F_s^2)$$
 1. Learning-based

• Computational complexity: at best $\mathcal{O}(F_s^2)$ 2. Analytical per iteration

How to solve this?

on-grid estimated peaks

RIR-agnostic + off-grid

Learning-based off-grid AER

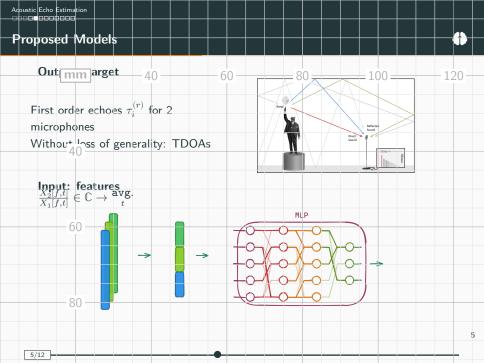
43

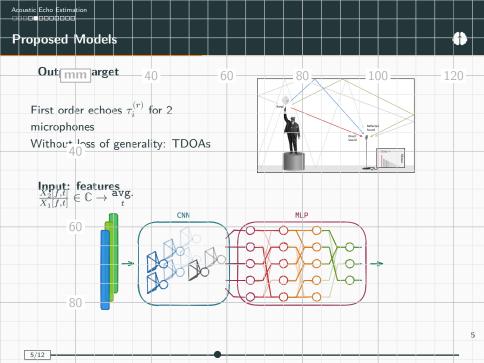
- **Rec:** $\underset{\longrightarrow}{\text{mm}} \mathbb{R} \Leftrightarrow \{\widetilde{x}_{i,40} \xrightarrow{?} \{\tau_{i}^{(r)}, \phi_{60}^{(r)}\}$
- Observations:
- This direct mapping is difficult, the inverse "is not"
 - \rightarrow acoustic simulators: mic/src/room \rightarrow $\tau_i^{(r)}, \alpha_i^{(r)}, \tilde{h}_i, \tilde{x}_i$ = 440 istic simulator are "simple", versatile and fast
 - F 40 istic simulator are "simple", versatile and last → many data

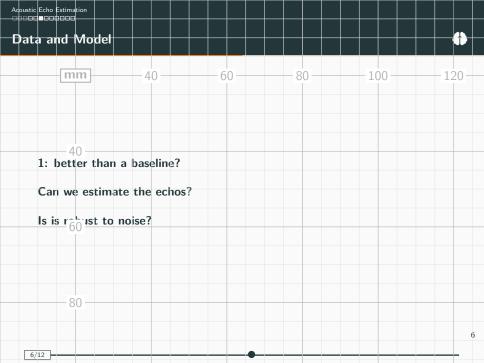
This approach is successfully in Sound Source Localization

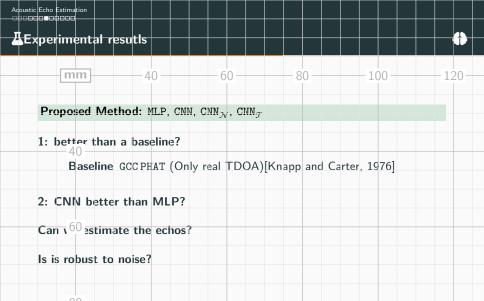
- → position is related to echoes

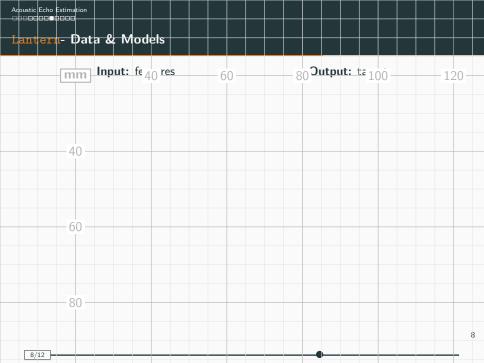
 [Navyen et al., 2018, Perotin et al., 2019] ▲ Not only DNN
- Idea: (Deep) Learning-based AER
- 1. Extend virtually learning-based \$SL to AER
- 2. Estimate first echo estimation (simple but important)
- 3. C⁸⁰, 2 microphones

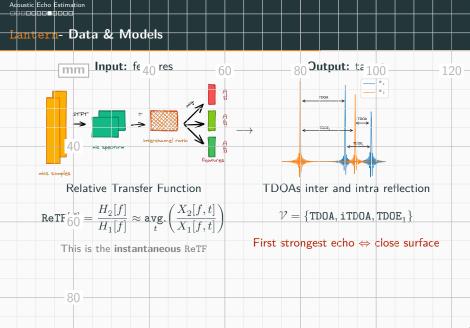






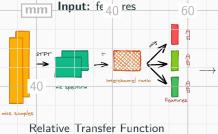






8/12





$$\begin{aligned} \text{ReTF}_{60}^{\text{r}} &= \frac{H_2[f]}{H_1[f]} \approx \text{avg.} \left(\frac{X_2[f,t]}{X_1[f,t]} \right) \\ \text{This is the instantaneous ReTF} \end{aligned}$$

$$\mathcal{V} = \{ \texttt{TDOA}, \texttt{iTDOA}, \texttt{TDOE}_1 \}$$
First strongest echo \Leftrightarrow close surface

TDOAs inter and intra reflection

80 Dutput: ta 100

Model



mics samples Relative Transfer Function

$$\text{ReTF}_{60}^{\text{l.s.}} = \frac{H_2[f]}{H_1[f]} \approx \underset{t}{\text{avg.}} \left(\frac{X_2[f,t]}{X_1[f,t]} \right)$$

This is the instantaneous ReTF

TDOAs inter and intra reflection

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

80 Dutput: ta 100

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

Model

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



mics samples Relative Transfer Function

$$\text{ReTF}_{60}^{\text{l.s.}} = \frac{H_2[f]}{H_1[f]} \approx \underset{t}{\text{avg.}} \left(\frac{X_2[f,t]}{X_1[f,t]} \right)$$

This is the instantaneous ReTF

TDOAs inter and intra reflection

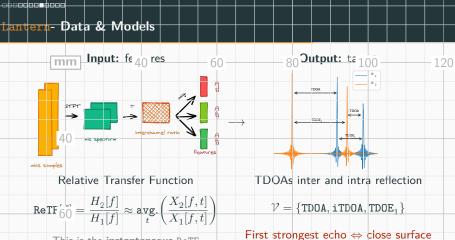
 $\mathcal{V} = \{ \text{TDOA}, \text{iTDOA}, \text{TDOE}_1 \}$

80 Dutput: ta 100

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

Model

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



for data fusion

NADAL ID:- Land 10041

Generative models ← similar

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

Adoustic Echo Estimation

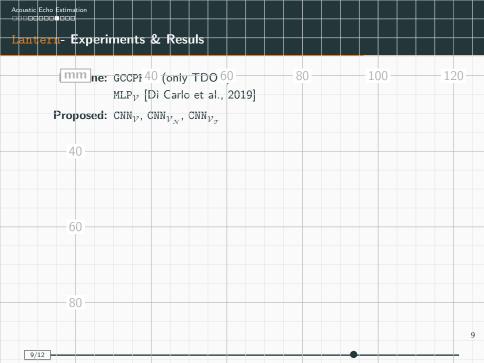
Loss Function:

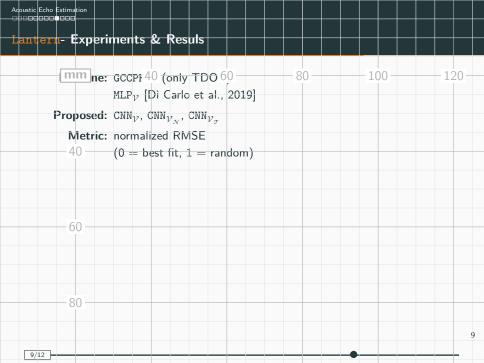
1. RMSE (Multi-label regression) on \mathcal{V}

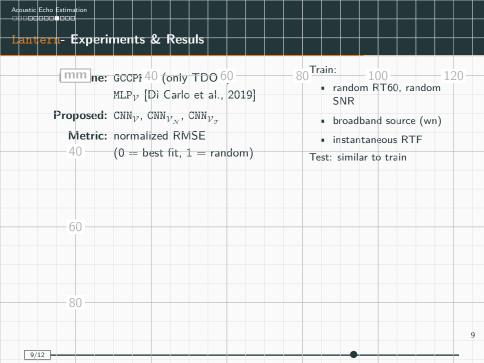
2. Gaussian log-likelihood $\rightarrow \{\mu, \sigma^2\}$

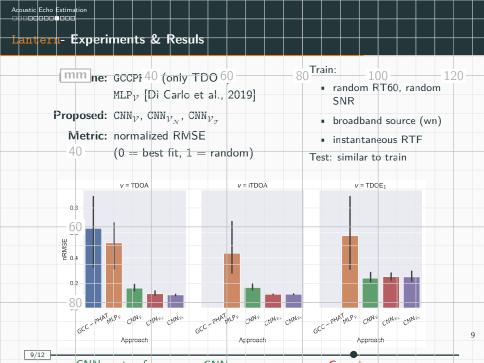
3. Student log-likelihood $\rightarrow \{\mu, \lambda, \nu\}$

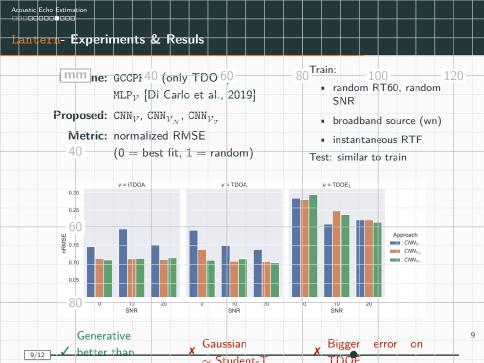
This is the instantaneous ReTF













mm

Bishop, C. M. (1994).

Mixture density networks.

tworks. Habets, I

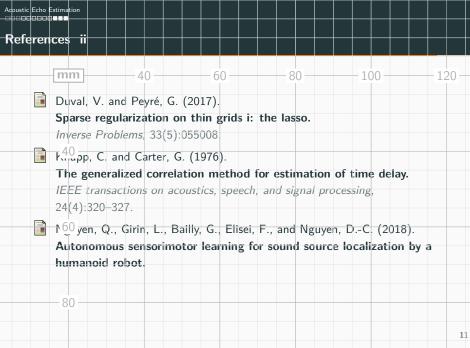
Chakrabarty, S. and Habets, E. A. (2017). \mathbf{F}_{40} adband doa estimation using convolutional neural networks trained with noise signals.

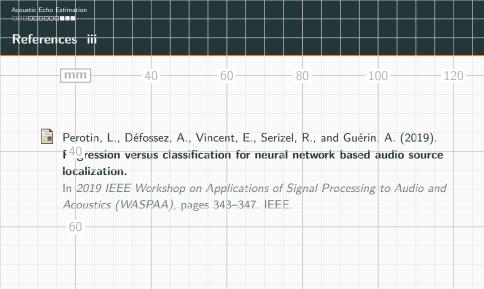
In 2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), pages 135–140. IEEE.

10/12

Nirage: 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.





12