

mm

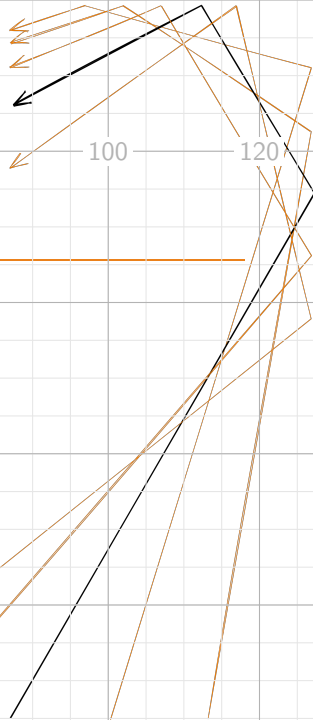
Echo-aware signal processing for audio scene analysis

Diego L. Carlo
December 2, 2020

PhD supervisors: Antoine Deleforge
Nancy Bertin

members: Renaud Seghier (president, examiner)
Simon Doclo (reviewer)
Laurent Girin (reviewer)
Fabio Antonacci (examiner)

Université de Rennes 1, IRISA/INRIA, Panama research group



mm

40

60

80

100

120

40

Echo-aware Application

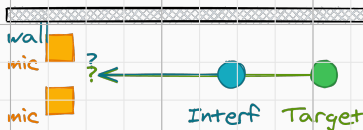
60

80

Echo-aware Application



40



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farruggian et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes ← image

- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation

How?

Echoes ∈ sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farruggian et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes ← image

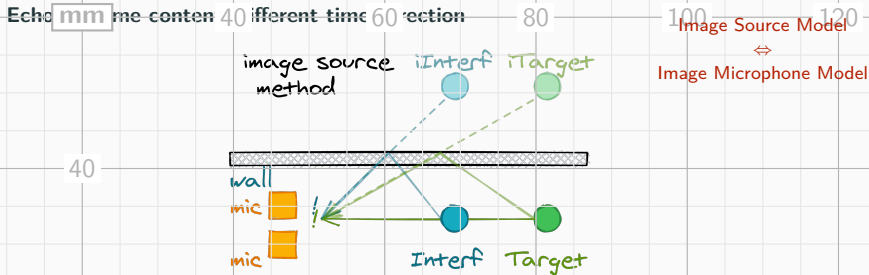
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation

How?

Echoes ∈ sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farrington et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes ← image

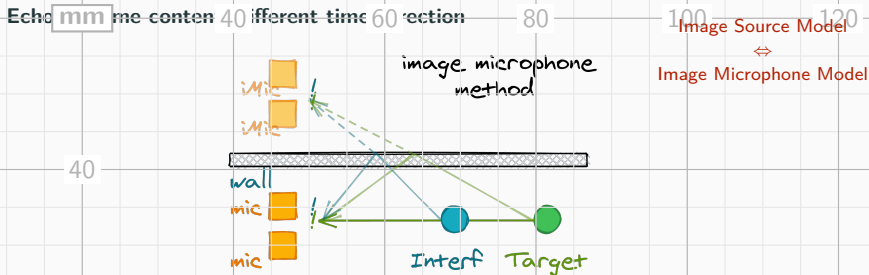
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation

How?

Echoes ∈ sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farrington et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes \leftarrow image

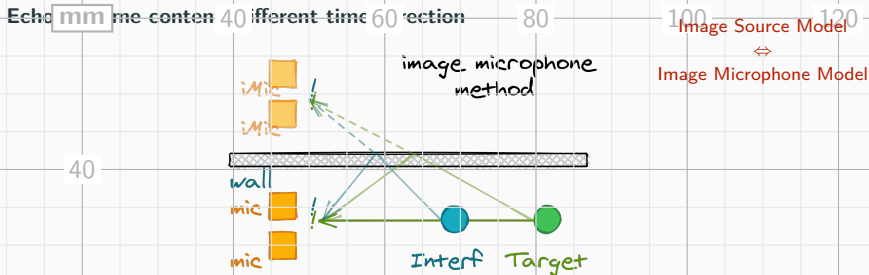
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation

How?

Echoes \in sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farrington et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes \leftarrow image

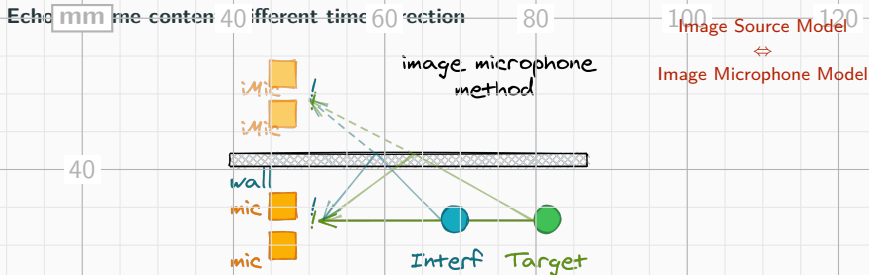
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation

How?

Echoes \in sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farrar et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes \leftarrow image

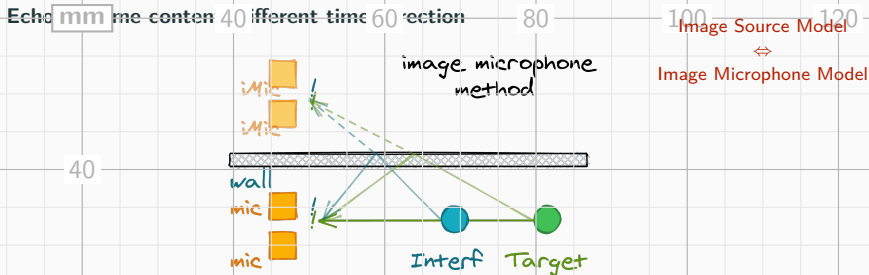
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation [Z. Crocco et al., 2017]

How?

Echoes \in sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 opy

- Sound Source Separation [Leglaive et al., 2016]
- Speech Enhancement [Farrington et al., 1993, Dokmanić et al., 2015, Kowalczyk, 2019]

Where?

Echoes \leftarrow image

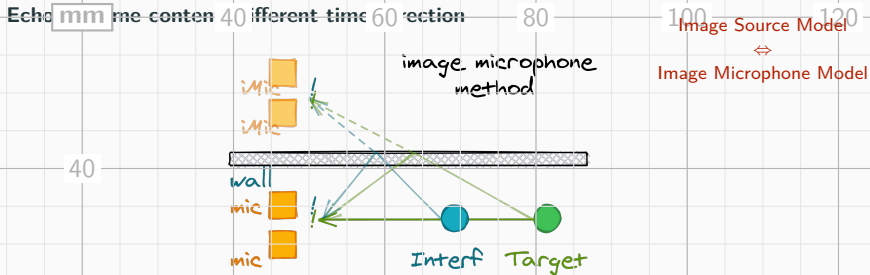
- Sound Source Localization [Ribeiro et al., 2010, Jensen et al., 2019]
- Microphone Calibration [Dokmanić et al., 2015, Salvati et al., 2016]
- Room Geometry Estimation [Z. Crocco et al., 2017]

How?

Echoes \in sound propagation

- Blind Channel Estimation [Lin et al., 2007, Crocco et al., 2017]
- Acoustic Measurements [Eaton et al., 2015, Kuttruff, 2016]

Echo-aware Application



What?

Echoes 60 copy

- Sound Source Separation

[Leglaive et al., 2016]

- Speech Enhancement

[Farruggian et al., 1993,

Dokmanić et al., 2015,

Kowalczyk, 2019]

Where?

Echoes \leftarrow image

- Sound Source Localization

[Ribeiro et al., 2010,

Jensen et al., 2019]

- Microphone Calibration

[Dokmanić et al., 2015,

Salvati et al., 2016]

- Room Geometry

Estimation

[Z. Crocco et al., 2017]

How?

Echoes \in sound propagation

- Blind Channel Estimation

[Lin et al., 2007,

Crocco et al., 2017]

- Acoustic Measurements

[Eaton et al., 2015,

Kuttruff, 2016]

Sound Source Localization (SSL)

(common knowledge) 

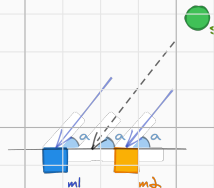
We consider h distance estimation.

SSL with 2 microphones

- Only one angle of arrival (AOA) 
- can be approximated from TDOA using e.g.

GC-RFAT¹

(known limitation, but good in practice)



² [DiBiase et al., 2001]

¹ [Knapp and Carter, 1976]

Sound Source Localization (SSL)

(common knowledge) 

We consider ≤ 40 distance estimation on.

80

100

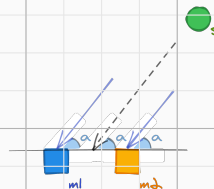
120

SSL with 2 microphones

- Only one angle of arrival (AOA) \uparrow
- can be approximated from TDOA using e.g.

GC-RFAT¹

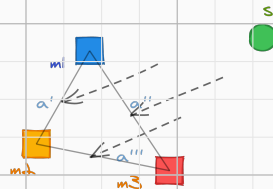
(known limitation, but good in practice)



SSL with more microphones

- Direction of Arrival (DoA): azimuth (\leftrightarrow) and elevation (\updownarrow)
- AOA for each pair can be “fused” together (e.g. angular spectra in SRP-PHAT²)

(known limitation, but good in practice)



² [DiBiase et al., 2001]

¹ [Knapp and Carter, 1976]

Sound Source Localization with Echoes

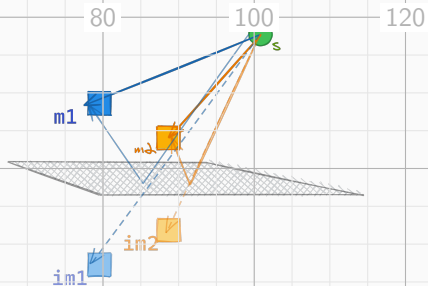


mm

The Picnic Scenario:

- One source
- Two microphones
- passive scenario

40 generalizable to any array geometry



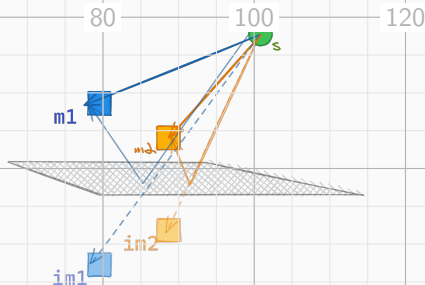
Sound Source Localization with Echoes



mm

The Picnic Scenario:

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



Sound Source Localization with Echoes



mm

40

60

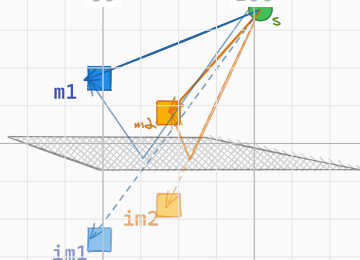
80

100

120

The Picnic Scenario:

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

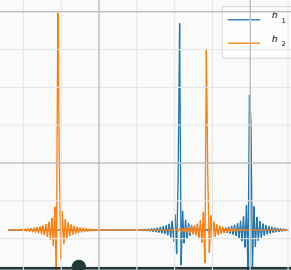


Each is augmented with echoes

Mirage Array

(Microphone Array Augmentation with Echoes)

How to process the *image* microphones?



Sound Source Localization with Echoes



mm

40

60

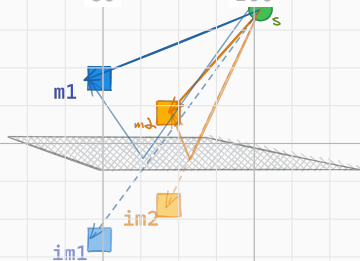
80

100

120

The Picnic Scenario:

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

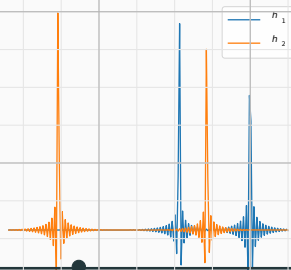


Each is augmented with echoes

Mirage Array

(Microphone Array Augmentation with Echoes)

How to process the *image* microphones?



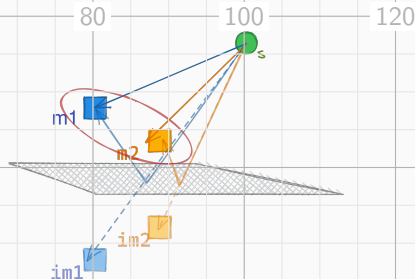
Sound Source Localization with Echoes



mm

The Picnic Scenario:

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

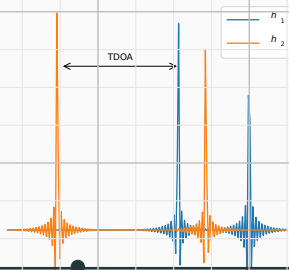


Each is augmented with echoes

Mirage Array

(Microphone Array Augmentation with Echoes)

How to process the *image* microphones?



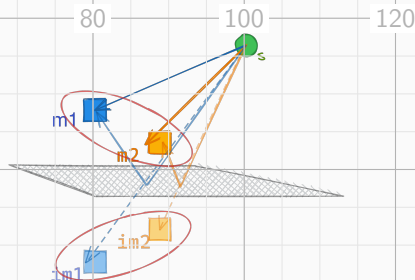
Sound Source Localization with Echoes



mm

The Picnic Scenario:

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

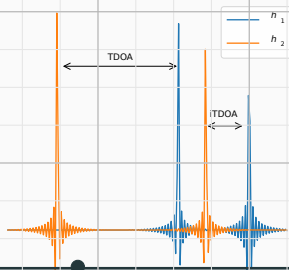


Each is augmented with echoes

Mirage Array

(Microphone Array Augmentation with Echoes)

How to process the *image* microphones?



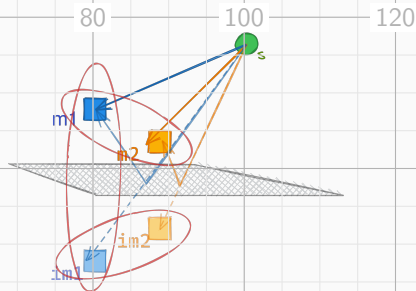
Sound Source Localization with Echoes



mm

The Picnic Scenario:

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

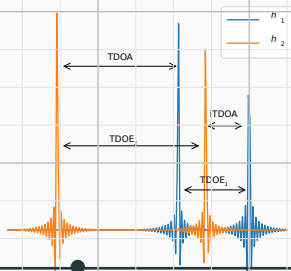


Each is augmented with echoes

Mirage Array

(Microphone Array Augmentation with Echoes)

How to process the *image* microphones?



Sound Source Localization with Echoes



mm

40

60

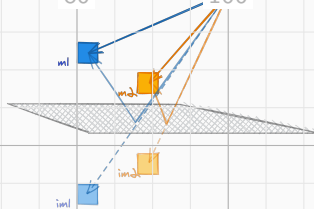
80

100

120

Idea: DoA estimate on the MIRAGE array.

Recall: these TDOAs are the same of the DNN-based method



40

Proposed Approach:

1. use proposed learning-based approach (MLP) for TDOAs estimation for each pair
2. fuse together the estimation ...
 - 2.1 use the error on a validation set as measure of uncertainty.
 - 2.2 DDDDDDDOOOODOO DDOO BETTER HERE for each pair of the Mirage array (similar to SRP-PHAT¹),
 - 2.3 knowing the position of the microphones,

Baseline GCC PHAT on true microphones²

² [DiBiase et al., 2001]

¹ [Knapp and Carter, 1976]

Experimental results



Proposed: MLP with **Mirage**

mm

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↩ also error in the manuscript)

AOI	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
GCC PHAT	wn	81	97

Observation

✓ $\alpha < 80^\circ$ arable to baseline when white noise source in noiseless case

Experimental results



Proposed: MLP with **Mirage**

mm

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↩ also error in the manuscript)

AOI 40	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
Mirage	wn+n	26	54
GCC PHAT	wn	81	97
GCC PHAT	wn+n	65	83

Observation

✓ $\alpha < 80^\circ$ arable to baseline when white noise source in noiseless case

Experimental results



Proposed: MLP with **Mirage**

mm

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↩ also error in the manuscript)

AO: 40	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
Mirage	wn+n	26	54
GCC PHAT	wn	81	97
GCC PHAT	wn+n	65	83
<hr/>			
Mir: 60	sp	63	82
GCC PHAT	sp	82	97

Observation

✓ $\alpha < 80$ arable to baseline when white noise source in noiseless case

Experimental results



Proposed: MLP with **Mirage**

mm

40

60

80

100

120

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↩ also error in the manuscript)

AOI 40	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
Mirage	wn+n	26	54
GCC PHAT	wn	81	97
GCC PHAT	wn+n	65	83
<hr/>			
Mir 60	sp	63	82
Mirage	sp+n	16	35
GCC PHAT	sp	82	97
GCC PHAT	sp+n	19	32

Observation

✓ $\alpha < 80$ arable to baseline when white noise source in noiseless case

✗ not generalize to noisy and speech data

Experimental results



Proposed: MLP with **Mirage**

mm

40

60

80

100

120

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↔ also error in the manuscript)

AOI 40	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
Mirage	wn+n	26	54
GCC PHAT	wn	81	97
GCC PHAT	wn+n	65	83

Mir 60	sp	63	82
Mirage	sp+n	16	35
GCC PHAT	sp	82	97
GCC PHAT	sp+n	19	32

DoA ↕	Input	ACCURACY			
		$< 10^\circ$		$< 20^\circ$	
		$\theta \leftrightarrow$	$\phi \updownarrow$	$\theta \leftrightarrow$	$\phi \updownarrow$
Mirage	wn	59	71	79	88
Mirage	wn+n	18	26	35	66
Mirage	sp	45	59	71	83
Mirage	sp+n	17	12	38	43

Observation

- ✓ α_{80} arable to baseline when white noise source in noiseless case
- ✗ not generalize to noisy and speech data
- ✓ Takled “impossible” localization

Experimental results



Proposed: MLP with **Mirage**

mm

40

60

80

100

120

Baseline: GCC PHAT¹

Data: 200 synthetic stereophonic recordings for close-surface scenario

Metric: accuracy in % ($<10^\circ$, $<20^\circ$) (↔ also error in the manuscript)

AO: 40	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	wn	77	97
Mirage	wn+n	26	54
GCC PHAT	wn	81	97
GCC PHAT	wn+n	65	83

Mir: 60	Input	ACCURACY	
		$\alpha < 10^\circ$	$\alpha < 20^\circ$
Mirage	sp	63	82
Mirage	sp+n	16	35
GCC PHAT	sp	82	97
GCC PHAT	sp+n	19	32

DoA: ↕	Input	ACCURACY			
		$< 10^\circ$		$< 20^\circ$	
		$\theta \leftrightarrow$	$\phi \updownarrow$	$\theta \leftrightarrow$	$\phi \updownarrow$
Mirage	wn	59	71	79	88
Mirage	wn+n	18	26	35	66
Mirage	sp	45	59	71	83
Mirage	sp+n	17	12	38	43

Observation

✓ $\alpha < 80$ arable to baseline when white noise source in noiseless case

✗ not generalize to noisy and speech data

✓ Takled “impossible” localization

⚠ Performance depending on echo estimation methods (work in progress)

References



mm, M., Truc, A., and Del, A. (2017). **Uncalibrated 3d room geometry estimation from sound impulse responses.** *Journal of the Franklin Institute*, 354(18):8678–8709.



DiBiase, J. H., Silverman, H. F., and Brandstein, M. S. (2001). **Robust localization in reverberant rooms.** *Microphone Arrays*, pages 157–180. Springer.



Dokmanić, I., Scheibler, R., and Vetterli, M. (2015). **Raking the cocktail party.** *IEEE journal of selected topics in signal processing*, 9(5):825–836.



En, J., Gaubitch, N. D., Moore, A. H., and Naylor, P. A. (2015). **The ace challenge—corpus description and performance evaluation.** In *2015 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 1–5. IEEE.



Evers, C. and Naylor, P. A. (2018). **Acoustic slam.** *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(9):1484–1498.

References ii



Elkann, J. L., Soudran, A. C., and Jan, E.-E. (1993).

Spatially selective sound capture for speech and audio processing.

Speech Communication, 13(1-2):207–222.



Jensen, J. R., Saqib, U., and Gannot, S. (2019).

An em method for multichannel toa and doa estimation of acoustic echoes.

In *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 120–124. IEEE.



Knapp, C. and Carter, G. (1976).

The generalized correlation method for estimation of time delay.

IEEE transactions on acoustics, speech, and signal processing, 24(4):320–327.



Kowalczyk, K. (2019).

Raking early reflection signals for late reverberation and noise reduction.

The Journal of the Acoustical Society of America, 145(3):EL257–EL263.



Kreković, M., Dokmanić, I., and Vetterli, M. (2016).

Ech-slam: Simultaneous localization and mapping with acoustic echoes.

In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 11–15. IEEE.



mm

H. (2014)

Room acoustics.

CRC Press.



Leglaive, S., Badeau, R., and Richard, G. (2016).

Multichannel audio source separation with probabilistic reverberation priors.*ACM Transactions on Audio, Speech, and Language Processing*, 24(12):2453–2465.

Lin, Y., Chen, J., Kim, Y., and Lee, D. D. (2007).

Blind sparse-nonnegative (bsn) channel identification for acoustic time-difference-of-arrival estimation.*IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pages 106–109. IEEE.

Ribeiro, F., Ba, D., Zhang, C., and Florêncio, D. (2010).

Turning enemies into friends: Using reflections to improve sound source localization.*IEEE International Conference on Multimedia and Expo*, pages 731–736. IEEE.

mm

40

60

80

100

120



40

Salvati, D., Drioli, C., and Foresti, G. L. (2016).

Sound source and microphone localization from acoustic impulse responses.

IEEE Signal Processing Letters, 23(10):1459–1463.

60

80