Clarity Workshop

Machine Listening in Dynamic Environments

Dr Christine Evers

University of Southampton, School of Electronics & Computer Science 17 September 2021

Clarity Workshop, Sep' 2021

Data Challenge



Sound Source LOCAlisation & TrAcking (LOCATA)

	Static Loudspeakers		Moving Human Talkers	
Array	Single	Multiple	Single	Multiple
Fixed	Task 1	Task 2	Task 3	Task 4
Moving	-	-	Task 5	Task 6

Development dataset

- Multichannel recordings + close-talking mouth signals
- Ground-truth positions & orientations for all sources and microphones
- Voice activity labels for each source

Evaluation dataset

- Multichannel recordings for all tasks and arrays
- Ground-truth data for array positions and orientations

C. Evers et al., "The LOCATA Challenge: Acoustic Source Localization and Tracking," in IEEE/ACM Trans. Audio, Speech, and Lang. Proc., Apr. 2020 C. Evers et al., "Data Corpus for the IEEE-AASP Challenge on Acoustic Source Localization and Tracking (LOCATA)," Data set. Zenodo. http://doi.org/10.5281/zenodo.3630471

Clarity Workshop, Sep' 2021

LOCATA Data Challenge

Sound Sources:

- ► Tasks 1 & 2: Static loudspeakers
- ► Tasks 2 6: Moving human talkers

Microphone Arrays:

- Pseudospherical robot head (12 mics)
- Spherical Eigenmike (32 mics)
- Siemens Signia hearing aids in head-torso simulator (4 mics)
- Planar DICIT array (15 mics)

Close-talking mouth signals:

► DPA d:screet SC4060



C. Evers et al., "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 C. Evers et al., "Data Corpus for the IEEE-AASP Challenge on Acoustic Source Localization and Tracking (LOCATA)," Data set. Zenodo. http://doi.org/10.5281/zenodo.3630471

Clarity Workshop, Sep' 2021

Christine Evers: "Machine Listening in Dynamic Environments"

Southampton

LOCATA Data Challenge



OptiTrack system (10 synchronized IR cameras)

Position & orientation of all sources & arrays across time and space

Synchronization:

Audio data + tracking data: 120 Hz



C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 C. Evers *et al.*, "Data Corpus for the IEEE-AASP Challenge on Acoustic Source Localization and Tracking (LOCATA)," Data set. Zenodo. http://doi.org/10.5281/zenodo.3630471

Clarity Workshop, Sep' 2021





LOCATA - Moving Source

Southampton



C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: *"Machine Listening in Dynamic Environments"*





Avg human walking speed: 1.4 m/s

Clarity Workshop, Sep' 2021





Clarity Workshop, Sep' 2021





Southampton



Clarity Workshop, Sep' 2021

LOCATA: Impact of source inactivity Southampton



C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: "*Machine Listening in Dynamic Environments*" 11 / 26



C. Evers, E. Habets, S. Gannot, P. Naylor, "DoA Reliability for Distributed Acoustic Tracking", in IEEE Signal Proc. Letters, 2018.

Clarity Workshop, Sep' 2021







C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: *"Machine Listening in Dynamic Environments"*

LOCATA: Impact of source ambiguity Southampton



Task 4, Recording 4, Submission ID 4 - Multiple moving sources

C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: *"Machine Listening in Dynamic Environments"* 16 / 26



A. Hogg *et al.*, "Overlapping speaker segmentation using multiple hypothesis tracking of fundamental frequency", IEEE Trans. Audio, Speech, Lang. Proc., 2021 A. Hogg, C. Evers, P. Naylor, "Multichannel Overlapping Speaker Segmentation Using Multiple Hypothesis Tracking Of Acoustic And Spatial Features," *ICASSP* 2021

Clarity Workshop, Sep' 2021

LOCATA: Impact of source cardinality Southampton



Task 6, Recording 2, Submission ID 4 - Multiple moving sources

C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.*, Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: *"Machine Listening in Dynamic Environments"*



C. Evers & P. Naylor, "Acoustic SLAM," in IEEE/ACM Trans. Audio, Speech, and Lang. Proc., Sep. 2018. C. Evers & P. Naylor, "Optimized Self-Localization for SLAM in Dynamic Scenes using Probability Hypothesis Density Filters," in IEEE Trans. Signal Proc., Feb. 2018.

Clarity Workshop, Sep' 2021



Stabilising perception in ego-motion southempton



X. Alameda-Pineda et al., "RAVEL: An Annotated Corpus for Training Robots with Audiovisual Abilities," Journal of Multimodal User Interfaces, 2013.

Clarity Workshop, Sep' 2021



C. Evers *et al.*, "The LOCATA Challenge: Acoustic Source Localization and Tracking," in *IEEE/ACM Trans. Audio, Speech, and Lang. Proc.,* Apr. 2020 Clarity Workshop, Sep' 2021 Christine Evers: *"Machine Listening in Dynamic Environments"* 2



Y. Ban, X. Alameda-Pineda, C. Evers, R. Horaud, "Tracking Multiple Audio Sources with the von Mises Distribution and Variational EM", in IEEE Signal Proc. Letters, 2019Clarity Workshop, Sep' 2021Christine Evers: "Machine Listening in Dynamic Environments"23 / 26



 $m{S}_t$: Set of Cartesian source states at frame t; $m{s}_t^{(n)}$: Cartesian position of source n; $m{P}$: Set of partitions of $m{S}_t$

C. Evers & P. Naylor, "Acoustic SLAM," in IEEE/ACM Trans. Audio, Speech, and Lang. Proc., Sep. 2018. C. Evers & P. Naylor, "Optimized Self-Localization for SLAM in Dynamic Scenes using Probability Hypothesis Density Filters," in IEEE Trans. Signal Proc., Feb. 2018.

Clarity Workshop, Sep' 2021

Summary



Challenges impacting on machine listening in dynamic scenarios:

• Source motion:

- Spatio-temporal variations, prohibiting batch-processing of recordings
- Source inactivity, leading to estimation bias and false estimates

Competing sources:

- Overlapping speech, resulting in ambiguity in the source identities
- Uncertainty in the number of sources, resulting in a combinatorial problem

• Ego-motion:

- Unknown source-sensor distance, required for reference frame transformations
- Unknown self-position, prohibitive for coherent integration of long-term memory

